I.5 Orthogonal Matrices and Subspaces

• 1. Orthogonal vectors x and y

$$\begin{array}{l} \left(\text{test} \right) \\ \mathbf{x}^{T} \mathbf{y} = 0 \\ \overline{\mathbf{x}}^{T} \mathbf{y} = 0 \end{array} \right\} \rightarrow \begin{cases} \text{Pytagoras Law of right triangles: } \|\mathbf{x} - \mathbf{y}\|^{2} = \|\mathbf{x}\|^{2} + \|\mathbf{y}\|^{2} \\ \text{Law of cosines: } \|\mathbf{x} - \mathbf{y}\|^{2} = \|\mathbf{x}\|^{2} + \|\mathbf{y}\|^{2} - 2\|\mathbf{x}\|\|\mathbf{y}\|\cos\theta \end{cases}$$

- 2. Orthogonal basis for a subspace
 - Standard basis is orthogonal (even orthonormal) in \mathbb{R}^n (*i*, *j*, *k* in \mathbb{R}^3)
 - Hadamard matrices H_n containing orthogonal bases of Rⁿ
 - Are those orthogonal matrices? (square, orthonormal vectors)

Applied Mathematics for Deep Learning

- Every subspace of **R**ⁿ has an orthogonal basis: Gram-Schmidt idea
 - Two independent vectors a and b in the plane: $\mathbf{c} = \mathbf{b} \frac{\mathbf{a}^T \mathbf{b}}{\mathbf{a}^T \mathbf{c}} \mathbf{a} \rightarrow \mathbf{a}^T \mathbf{c} = \mathbf{0}$

$$\mathbf{q}_{1} = \frac{\mathbf{a}_{1}}{\|\mathbf{a}_{1}\|}$$

$$\mathbf{A}_{2} = \mathbf{a}_{2} - (\mathbf{a}_{2}^{T}\mathbf{q}_{1})\mathbf{q}_{1} \rightarrow \mathbf{q}_{2} = \frac{\mathbf{A}_{2}}{\|\mathbf{A}_{2}\|} \rightarrow \mathbf{q}_{1}^{T}\mathbf{A}_{2} = 0?$$

$$\mathbf{A}_{3} = \mathbf{a}_{3} - (\mathbf{a}_{3}^{T}\mathbf{q}_{1})\mathbf{q}_{1} - (\mathbf{a}_{3}^{T}\mathbf{q}_{2})\mathbf{q}_{2} \rightarrow \mathbf{q}_{3} = \frac{\mathbf{A}_{3}}{\|\mathbf{A}_{3}\|}$$

$$\rightarrow \begin{bmatrix} \mathbf{a}_{1} & \mathbf{a}_{2} & \mathbf{a}_{3} \end{bmatrix} = \begin{bmatrix} \mathbf{q}_{1} & \mathbf{q}_{2} & \mathbf{q}_{3} \end{bmatrix} \begin{bmatrix} \|\mathbf{a}_{1}\| & \mathbf{a}_{2}^{T}\mathbf{q}_{1} & \mathbf{a}_{3}^{T}\mathbf{q}_{1} \\ \|\mathbf{A}_{2}\| & \mathbf{a}_{3}^{T}\mathbf{q}_{2} \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{q}_{1} & \mathbf{q}_{2} & \mathbf{q}_{3} \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{22} & r_{23} \\ r_{33} \end{bmatrix} \Leftrightarrow \mathbf{A} = \mathbf{Q}\mathbf{R}$$

- 3. Orthogonal subspace R (row space) and N (null space)
 - Ax=0: The row space of A is orthogonal to the nullspace of A
 - $A^{T}y=0$: The column space of A is orthogonal to the nullspace of A^{T}





4. Tall thin matrices Q with orthonormal columns: Q^TQ=I

 $\begin{cases} \text{if } \mathbf{Q} \text{ multiplies any vector } \mathbf{x}, \text{ the length of the vector does not change: } \|\mathbf{Q}\mathbf{x}\| = \|\mathbf{x}\| \\ \text{if } m > n \text{ then } m \text{ rows cannot be orthogonal in } \mathbf{R}^n : \mathbf{Q}\mathbf{Q}^T \neq \mathbf{I} \\ \mathbf{Q}_1 = \frac{1}{3} \begin{bmatrix} 2 \\ 2 \\ -1 \end{bmatrix}, \mathbf{Q}_2 = \frac{1}{3} \begin{bmatrix} 2 & 2 \\ 2 & -1 \\ -1 & 2 \end{bmatrix}, \mathbf{Q}_3 = \frac{1}{3} \begin{bmatrix} 2 & 2 & -1 \\ 2 & -1 & 2 \\ -1 & 2 & 2 \end{bmatrix} \rightarrow \mathbf{Q}_i \mathbf{Q}_i^T = \mathbf{I}? \\ \mathbf{P} = \mathbf{Q}\mathbf{Q}^T : \text{ projection matrix } \rightarrow \mathbf{P}^2 = \mathbf{P} = \mathbf{P}^T \xrightarrow{\text{"least squares"}} \end{cases}$

Pb is the orthogonal projection of **b** onto the column space of **P** : P_1b , P_2b , P_3b



Highlights of Linear Algebra - 5

$$\mathbf{Q}_{1} = \frac{1}{3} \begin{bmatrix} 2\\2\\-1 \end{bmatrix}, \mathbf{Q}_{2} = \frac{1}{3} \begin{bmatrix} 2 & 2\\2 & -1\\-1 & 2 \end{bmatrix}, \mathbf{Q}_{3} = \frac{1}{3} \begin{bmatrix} 2 & 2 & -1\\2 & -1 & 2\\-1 & 2 & 2 \end{bmatrix} \rightarrow \mathbf{Q}_{i} \mathbf{Q}_{i}^{T} = \mathbf{I}?$$

$$\mathbf{P}_{1} = \mathbf{Q}_{1} \mathbf{Q}_{1}^{T} = \frac{1}{9} \begin{bmatrix} 2\\2\\-1 \end{bmatrix} \begin{bmatrix} 2 & 2 & -1 \end{bmatrix} = \frac{1}{9} \begin{bmatrix} 4 & 4 & -2\\4 & 4 & -2\\-2 & -2 & 1 \end{bmatrix} \rightarrow \mathbf{P}_{1} \mathbf{b} = \frac{1}{9} \begin{bmatrix} 18\\18\\-9 \end{bmatrix} = \begin{bmatrix} 2\\2\\-1 \end{bmatrix}$$

$$\mathbf{P}_{2} = \mathbf{Q}_{2} \mathbf{Q}_{2}^{T} = \frac{1}{9} \begin{bmatrix} 2 & 2\\2 & -1\\-1 & 2 \end{bmatrix} \begin{bmatrix} 2 & 2 & -1\\2 & -1 & 2 \end{bmatrix} = \frac{1}{9} \begin{bmatrix} 8 & 2 & 2\\2 & 5 & -4\\2 & -4 & 5 \end{bmatrix} \rightarrow \mathbf{P}_{2} \mathbf{b} = \frac{1}{9} \begin{bmatrix} 36\\9\\9\\9 \end{bmatrix} = \begin{bmatrix} 4\\1\\1 \end{bmatrix}$$

$$\mathbf{P}_{3} = \mathbf{Q}_{3} \mathbf{Q}_{3}^{T} = \frac{1}{9} \begin{bmatrix} 2 & 2 & -1\\2 & -1 & 2\\-1 & 2 & 2 \end{bmatrix} \begin{bmatrix} 2 & 2 & -1\\2 & -1 & 2\\-1 & 2 & 2 \end{bmatrix} = \frac{1}{9} \begin{bmatrix} 9 & 0 & 0\\0 & 9 & 0\\0 & 0 & 9 \end{bmatrix} = \mathbf{I} \rightarrow \mathbf{P}_{3} \mathbf{b} = \mathbf{b}$$

Applied Mathematics for Deep Learning



Highlights of Linear Algebra - 7

All reflection matrices have eigenvalues -1 and 1

Householder Reflections

$$\mathbf{H} = \mathbf{I} - 2\frac{\mathbf{v}\mathbf{v}^{T}}{\|\mathbf{v}\|^{2}} = \mathbf{I} - 2\mathbf{u}\mathbf{u}^{T}$$
key point: if $\mathbf{v} = \mathbf{a} - \mathbf{r}$ and $\|\mathbf{a}\| = \|\mathbf{r}\|$, then $\mathbf{H}\mathbf{a} = \mathbf{a} - 2\frac{(\mathbf{a} - \mathbf{r})(\mathbf{a} - \mathbf{r})^{T}}{(\mathbf{a} - \mathbf{r})^{T}}\mathbf{a} = \mathbf{r}$

$$\mathbf{H}_{k} [\text{column } k] = \begin{bmatrix} \mathbf{I} & \\ \mathbf{I} - 2\mathbf{u}\mathbf{u}^{T} \end{bmatrix} \begin{bmatrix} \mathbf{a}_{\text{upper}} \\ \mathbf{a}_{\text{lower}} \end{bmatrix} = r_{k}$$

$$\mathbf{H}_{n-1} \cdots \mathbf{H}_{2}\mathbf{H}_{1}\mathbf{A} = \begin{bmatrix} \mathbf{r}_{1} & \mathbf{r}_{2} & \cdots & \mathbf{r}_{n} \end{bmatrix} \rightarrow \mathbf{Q}^{T}\mathbf{A} = \mathbf{R}$$
keep a record of the \mathbf{H}_{j} by storing only the vectors $\mathbf{v}_{j} = \mathbf{a}_{j} - \mathbf{r}_{j}$, not the matrix

 $\mathbf{A}\mathbf{x} = \mathbf{b} \xrightarrow{\text{multiply by all the H's}} \mathbf{R}\mathbf{x} = \mathbf{Q}^T \mathbf{b}$

$$\mathbf{A} = \begin{bmatrix} 4 & x \\ 3 & x \end{bmatrix} \rightarrow \mathbf{a} = \begin{bmatrix} 4 \\ 3 \end{bmatrix}, \mathbf{r} = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \rightarrow \mathbf{v} = \mathbf{a} - \mathbf{r} = \begin{bmatrix} -1 \\ 3 \end{bmatrix} \rightarrow \mathbf{u} = \frac{1}{\sqrt{10}} \begin{bmatrix} -1 \\ 3 \end{bmatrix}$$
$$\mathbf{H} = \mathbf{I} - 2\mathbf{u}\mathbf{u}^{T} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 2\frac{1}{10} \begin{bmatrix} 1 & -3 \\ -3 & 9 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 4 & 3 \\ 3 & -4 \end{bmatrix} = \mathbf{Q}^{T} \rightarrow \mathbf{H}\mathbf{A} = \begin{bmatrix} 5 & x \\ 0 & x \end{bmatrix} = \mathbf{R}$$

Applied Mathematics for Deep Learning

Examples

- Rotations
- Reflections
- Hadamard matrices
- Haar wavelets
- Discrete Fourier Transform (DFT)
- Complex inner product



I.6 Eigenvalues and Eigenvectors

eigenvectors of A don't change direction when you multiply them by A

 $\begin{array}{l} \mathbf{x}: \text{ eigenvector of } \mathbf{A} \\ \lambda: \text{ eigenvector of } \mathbf{A} \end{array} \rightarrow \mathbf{A}\mathbf{x} = \lambda \mathbf{x} \rightarrow \mathbf{A}(\mathbf{A}\mathbf{x}) = \mathbf{A}(\lambda \mathbf{x}) = \lambda^2 \mathbf{x} \rightarrow \mathbf{A}^k \mathbf{x} = \lambda^k \mathbf{x}, \ \mathbf{A}^{-1} \mathbf{x} = \frac{1}{\lambda} \mathbf{x} \text{ if } \lambda \neq 0 \end{array}$

 $n \times n$ matrices \rightarrow n independent eigenvectors \mathbf{x}_1 to \mathbf{x}_n with n different eigenvalues λ_1 to λ_n

$$\mathbf{v} = c_1 \mathbf{x}_1 + \dots + c_n \mathbf{x}_n \to \mathbf{A} \mathbf{v} = c_1 \lambda_1 \mathbf{x}_1 + \dots + c_n \lambda_n \mathbf{x}_n \to \mathbf{A}^k \mathbf{v} = c_1 \lambda_1^k \mathbf{x}_1 + \dots + c_n \lambda_n^k \mathbf{x}_n$$

How useful?
$$\begin{cases} (1) \text{ solution of differential equations} \\ (2) \text{ similar matrices} \to \text{ same eigenvalues} \\ (3) \text{ diagonalize a matrix} \end{cases}$$

$$|\lambda_1| > 1: c_1 \lambda_1^n \mathbf{x}_1 \text{ will grow as } n \text{ increases} \end{cases}$$

$$\Big\} \to \text{follow each eigenvector separately!}$$

 $|\lambda_2| < 1: c_2 \lambda_2^n \mathbf{x}_2$ will steadily disappear as *n* increases

 $eig(\mathbf{A} + \mathbf{B}) \neq eig(\mathbf{A}) + eig(\mathbf{B})$ $eig(\mathbf{A}\mathbf{B}) \neq eig(\mathbf{A})eig(\mathbf{B})$ $\lambda_{1} = \lambda_{2}$ might or might not have two independent eigenvectors Four properties: matrix $\mathbf{A}(\text{real})$, $\underbrace{\mathbf{S}(\text{symmetric})}_{\text{like real numbers: }\lambda}$, $\underbrace{\mathbf{Q}(\text{orthogonal})}_{\text{like complex numbers: }e^{i\theta}}$ every $|\lambda|=1$ powers of \mathbf{Q} don't grow or decay (Trace of \mathbf{S}) $\sum_{i=1}^{n} \lambda_i$ = trace of matrix (Determinant) $\prod \lambda_i$ = determinant of matrix (Real eigenvalues of \mathbf{S}) \mathbf{S} : real eigenvalues, orthogonal eigenvectors (Orthogonal eigenvectors) if $\lambda_1 \neq \lambda_2$, then $\mathbf{x}_1 \cdot \mathbf{x}_2 = 0$, eigenvectors of \mathbf{A} are orthogonal iff $\mathbf{A}^T \mathbf{A} = \mathbf{A}\mathbf{A}^T$

$$\mathbf{S} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \rightarrow \left\{ \mathbf{S} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 3 \begin{bmatrix} 1 \\ 1 \end{bmatrix} \text{ and } \mathbf{S} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 1 \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right\}$$
$$\mathbf{Q} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \rightarrow \left\{ \mathbf{Q} \begin{bmatrix} 1 \\ -i \end{bmatrix} = i \begin{bmatrix} 1 \\ -i \end{bmatrix} \text{ and } \mathbf{Q} \begin{bmatrix} 1 \\ i \end{bmatrix} = -i \begin{bmatrix} 1 \\ i \end{bmatrix} \right\}$$
$$\mathbf{A} = \begin{bmatrix} 8 & 3 \\ 2 & 7 \end{bmatrix} \rightarrow \left\{ \mathbf{A} \begin{bmatrix} 3 \\ 2 \end{bmatrix} = 10 \begin{bmatrix} 3 \\ 2 \end{bmatrix} \text{ and } \mathbf{A} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 5 \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right\}$$
$$\mathbf{3} \rightarrow \mathbf{30? \text{ complex } \lambda}$$

(1) A controls a system of linear differential equations: $\frac{d\mathbf{u}}{dt} = \mathbf{A}\mathbf{u}$ with $\mathbf{u}(0)$ $\mathbf{u}(0) = c_1 \mathbf{x}_1 + \dots + c_n \mathbf{x}_n$ $\mathbf{u}(t) = c_1 e^{\lambda_1 t} \mathbf{x}_1 + \dots + c_n e^{\lambda_n t} \mathbf{x}_n \xrightarrow{\lambda = a + ib} \begin{cases} e^{at} = \begin{cases} \operatorname{Re} \lambda > 0 : \text{ grow} \\ \operatorname{Re} \lambda < 0 : \text{ decay} \end{cases}$ $e^{ibt} = \cos bt + i \sin bt : \text{ oscillate} \end{cases}$

shift in $\mathbf{A} \rightarrow \text{shift}$ in λ : $(\mathbf{A} + \mathbf{sI})\mathbf{x} = \lambda \mathbf{x} + \mathbf{sx} = (\lambda + \mathbf{s})\mathbf{x}$

(2) **B** similar to $\mathbf{A} \to \mathbf{B} = \underset{\text{invertible}}{\mathbf{M}} \mathbf{A}\mathbf{M}^{-1} \to eig(\mathbf{B}) = eig(\mathbf{A})$: compute eigenvalues of large matrices

Make **B** gradually into a triangular matrix \rightarrow Gradually show up on the main diagonal

$$By = \lambda y \rightarrow MAM^{-1}y = \lambda y \rightarrow A(M^{-1}y) = \lambda(M^{-1}y)$$
$$Ax = \lambda x \rightarrow MAM^{-1}(Mx) = M\lambda x = \lambda(Mx)$$

(3) diagonalize a matrix

$$\mathbf{A}\begin{bmatrix}\mathbf{x}_{1} & \cdots & \mathbf{x}_{n}\end{bmatrix} = \begin{bmatrix}\mathbf{A}\mathbf{x}_{1} & \cdots & \mathbf{A}\mathbf{x}_{n}\end{bmatrix} = \begin{bmatrix}\lambda_{1}\mathbf{x}_{1} & \cdots & \lambda_{n}\mathbf{x}_{n}\end{bmatrix} = \begin{bmatrix}\mathbf{x}_{1} & \cdots & \mathbf{x}_{n}\end{bmatrix}\begin{bmatrix}\lambda_{1} & & \\ & \ddots & \\ & & \lambda_{n}\end{bmatrix} \rightarrow \begin{bmatrix}\mathbf{A}\mathbf{X} = \mathbf{X}\mathbf{A}\\\mathbf{A} = \mathbf{X}\mathbf{A}\mathbf{X}^{-1}\\\mathbf{A}^{2} = \mathbf{X}\mathbf{A}^{2}\mathbf{X}^{-1}\\\mathbf{A}^{2} = \mathbf{X}\mathbf{A}^{2}\mathbf{X}^{-1}$$

$$\mathbf{A}^{k}\mathbf{v} = \mathbf{X}\mathbf{A}^{k}\mathbf{X}^{-1}\mathbf{v}: \mathbf{v} = \mathbf{X}\mathbf{c} \rightarrow \underbrace{\mathbf{c}} = \underbrace{\mathbf{X}^{-1}\mathbf{v}}_{c_{i}} \rightarrow \underbrace{\mathbf{A}^{k}\mathbf{X}^{-1}\mathbf{v}}_{c_{i}\lambda_{i}^{k}} \rightarrow \underbrace{\mathbf{X}\mathbf{A}\mathbf{X}^{-1}\mathbf{v}}_{\sum c_{i}\lambda_{i}^{k}x_{i}}$$
Example:
$$\mathbf{A} = \begin{bmatrix}\mathbf{8} & \mathbf{3}\\ 2 & 7\end{bmatrix} \xrightarrow{\text{divide by 10}} \underbrace{\mathbf{A}}_{\text{markov matrix with positive columns adding to 1} \rightarrow \begin{bmatrix}\mathbf{A}^{k}\mathbf{v} = c_{1}(1)^{k}\mathbf{x}_{1} + c_{1}\left(\frac{1}{2}\right)^{k}\mathbf{x}_{2}$$
as k increases $\mathbf{A}^{k}\mathbf{v}$ approaches to $c_{1}\mathbf{x}_{1}$

the action of the whole matrix A is broken into simple actions (just muliply by λ)

[nondiagonalizable matrices: when GM < AM, **A** is not diagonalizable] $\begin{cases}
(Geometric Multiplicity = GM): count the independent eigenvectors, dim <math>N(\mathbf{A} - \lambda \mathbf{I}) \\
(Algebraic Multiplicity = AM): count the repetitions of eigenvalues, det(\mathbf{A} - \lambda \mathbf{I})=0
\end{cases}$

$$\mathbf{A} = \begin{bmatrix} 5 & 1 \\ 0 & 5 \end{bmatrix}, \begin{bmatrix} 6 & -1 \\ 1 & 4 \end{bmatrix}, \begin{bmatrix} 7 & 2 \\ -2 & 3 \end{bmatrix} \rightarrow \begin{cases} \det(\mathbf{A} - \lambda \mathbf{I}) = (\lambda - 5)^2 = 0 \rightarrow AM = 2\\ rank(\mathbf{A} - 5\mathbf{I}) = 1 \rightarrow GM = 1 \end{cases}$$

Applied Mathematics for Deep Learning

I.8 Singular Value Decomposition (SVD)

best matrices (real symmetric matrices S): real eigenvalues and orthogonal eigenvectors other matrices (A is not square, $m \times n$, matrix of data): complex eigenvalues and not orthogonal eigenvectors

key point: two sets of singular vectors {

$$\begin{bmatrix} n & n \\ m & n \\ n & n \\ n$$

connection between n v's and m u's

$$\underbrace{\operatorname{Av}_{1} = \sigma_{1}\mathbf{u}_{1}, \dots, \operatorname{Av}_{r} = \sigma_{r}\mathbf{u}_{r}}_{r=\operatorname{rank}(\mathbf{A})}, \underbrace{\operatorname{Av}_{r+1} = \mathbf{0}, \dots, \operatorname{Av}_{n} = \mathbf{0}}_{(n-r) \text{ v's in } N(\mathbf{A})} \leftarrow \operatorname{rank}(\mathbf{A}) = 2 \cdot \operatorname{age, height}_{n = 1000} \right\}$$

$$\mathbf{A} \begin{bmatrix} \mathbf{v}_{1} & \cdots & \mathbf{v}_{r} & \cdots & \mathbf{v}_{n} \\ \mathbf{v}_{1} & \cdots & \mathbf{v}_{r} & \cdots & \mathbf{v}_{n} \end{bmatrix}_{r=1}^{r=1} = \begin{bmatrix} \mathbf{u}_{1} & \cdots & \mathbf{u}_{r} & \cdots & \mathbf{u}_{m} \end{bmatrix} \begin{bmatrix} \sigma_{1} & \cdots & \sigma_{r} & \sigma_{r} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \rightarrow \begin{cases} \mathbf{AV} = \mathbf{U\Sigma} \\ \mathbf{A} = \mathbf{U\Sigma}\mathbf{V}^{T} \\ \mathbf{AV} = \mathbf{U\Sigma} \\ \mathbf{AV} = \mathbf{U} \\ \mathbf{AV} \\ \mathbf{AV} = \mathbf{U} \\ \mathbf{AV} = \mathbf{U} \\ \mathbf{AV} \\ \mathbf{AV} = \mathbf{U} \\ \mathbf{AV} \\ \mathbf{AV} = \mathbf{U} \\ \mathbf{AV} = \mathbf{U} \\ \mathbf{AV} \\ \mathbf{AV} = \mathbf{U} \\ \mathbf{AV} \\ \mathbf{AV} = \mathbf{U} \\ \mathbf{AV} \\ \mathbf{$$

- Columns of V are orthogonal eigenvectors of A^TA
- Av=σu gives orthonormal eigenvectors u of AA^T
- σ^2 = eigenvalue of A^TA = eigenvalue of AA^T \neq 0
- Why is the SVD so important?
 - It separates the matrix into rank one pieces like the other factorizations A=LU, A=QR, S=Q∧Q^T
 - Those pieces come in order of importance
 - First piece $\sigma_1 u_1 v_1^T$ is the closest rank one matrix to A
 - Sum of the first k pieces is best possible for rank k

 $\mathbf{A}_{k} = \sigma_{1} \mathbf{u}_{1} \mathbf{v}_{1}^{T} + \dots + \sigma_{k} \mathbf{u}_{k} \mathbf{v}_{k}^{T} \text{ is the best rank } k \text{ approximation to } \mathbf{A}:$ If **B** has rank k then $\|\mathbf{A} - \mathbf{A}_{k}\| \le \|\mathbf{A} - \mathbf{B}\|$

Proof of SVD

$$\mathbf{AX} = \mathbf{XA} \Leftrightarrow \mathbf{AV} = \mathbf{U\Sigma}$$

$$\begin{bmatrix} 3 & 0\\ 4 & 5 \end{bmatrix} \xrightarrow{1} \sqrt{2} \begin{bmatrix} 1 & -1\\ 1 & 1 \end{bmatrix} = \frac{1}{\sqrt{10}} \begin{bmatrix} 1 & -3\\ 3 & 1 \end{bmatrix} \xrightarrow{\sqrt{5}} \sqrt{5} \\ \xrightarrow{rank(\mathbf{A})=2 \to \sigma_1, \sigma_2} \\ \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T = \frac{3\sqrt{5}}{\sqrt{10}\sqrt{2}} \begin{bmatrix} 1\\ 3 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} + \frac{\sqrt{5}}{\sqrt{10}\sqrt{2}} \begin{bmatrix} -3\\ 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \end{bmatrix} = \frac{3}{2} \begin{bmatrix} 1 & 1\\ 3 & 3 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 3 & -3\\ -1 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0\\ 4 & 5 \end{bmatrix} = \mathbf{A}$$

$$\begin{cases} \mathbf{V} \text{ contains orthonormal eigenvectors of } \mathbf{A}^T \mathbf{A} \\ \mathbf{U} \text{ contains orthonormal eigenvectors of } \mathbf{A}\mathbf{A}^T \\ \sigma_1^2 \text{ to } \sigma_r^2 \text{ are the nonzero eigenvalues of both } \mathbf{A}^T \mathbf{A} \text{ and } \mathbf{A}\mathbf{A}^T \end{cases}$$

$$SVD \text{ requires that } \mathbf{Av}_k = \sigma_k \mathbf{u}_k : \mathbf{v} \cdot \mathbf{v} \cdot \mathbf{A}^T \mathbf{Av}_k = \sigma_k^2 \mathbf{v}_k \end{pmatrix} \rightarrow \mathbf{u} \cdot \mathbf{v} \cdot \mathbf{u} \cdot \mathbf{u} = \frac{\mathbf{Av}_k}{\sigma_k} \text{ for } k = 1, \dots, r \\ \text{ sign, multiple eigenvalues} \end{bmatrix}$$

(check 1) **u**'s are eigenvectors of
$$\mathbf{A}\mathbf{A}^T \to \mathbf{A}\mathbf{A}^T\mathbf{u}_k = \sigma_k^2\mathbf{u}_k$$

(check 2) **u**'s are also orthonormal $\to \mathbf{u}_j^T\mathbf{u}_k = \frac{\sigma_k}{\sigma_j}\mathbf{v}_j^T\mathbf{v}_k = \begin{cases} 1 \text{ if } j = k \\ 0 \text{ if } j \neq k \end{cases}$
choose $(n-r)$ **v**'s in $N(\mathbf{A})$ and $(m-r)$ **u**'s in $N(\mathbf{A}^T)$
Applied Mathematics for Deep Learning

Example

Find the matrices
$$\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}$$
 for $\mathbf{A} = \begin{bmatrix} 3 & 0 \\ 4 & 5 \end{bmatrix} \rightarrow \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 25 & 20 \\ 20 & 25 \end{bmatrix}, \mathbf{A}\mathbf{A}^T = \begin{bmatrix} 9 & 12 \\ 12 & 41 \end{bmatrix}$
 $\rightarrow \begin{bmatrix} 25 & 20 \\ 20 & 25 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 45 \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 25 & 20 \\ 20 & 25 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = 5 \begin{bmatrix} -1 \\ 1 \end{bmatrix}$
 $\mathbf{v}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \rightarrow \mathbf{A}\mathbf{v}_1 = \begin{bmatrix} 3 & 0 \\ 4 & 5 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \frac{3}{\sqrt{2}} \begin{bmatrix} 1 \\ 3 \end{bmatrix} = \sigma_1 \mathbf{u}_1 = \sqrt{45} \frac{1}{\sqrt{10}} \begin{bmatrix} 1 \\ 3 \end{bmatrix}$
 $\mathbf{v}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \rightarrow \mathbf{A}\mathbf{v}_2 = \begin{bmatrix} 3 & 0 \\ 4 & 5 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} -3 \\ 1 \end{bmatrix} = \sigma_1 \mathbf{u}_1 = \sqrt{5} \frac{1}{\sqrt{10}} \begin{bmatrix} -3 \\ 1 \end{bmatrix}$
 $\mathbf{U} = \frac{1}{\sqrt{10}} \begin{bmatrix} 1 & -3 \\ 3 & 1 \end{bmatrix}, \mathbf{\Sigma} = \begin{bmatrix} \sqrt{45} \\ \sqrt{5} \end{bmatrix}, \mathbf{V} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$
 $\sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T = \frac{\sqrt{45}}{\sqrt{20}} \begin{bmatrix} 1 & 3 \\ 1 & 3 \end{bmatrix} + \frac{\sqrt{5}}{\sqrt{20}} \begin{bmatrix} 3 & -3 \\ -3 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 4 & 5 \end{bmatrix} = \mathbf{A}$

Questions

- If $S=QAQ^T$ is symmetric positive definite, what is its SVD?
- If S=QΛQ^T has a negative eigenvalue(Sx=-αx), what is the singular value and what are the vectors v and u?
- If A=Q is an orthogonal matrix, why does every singular value equal 1?
- Why are all eigenvalues of a square matrix A less than or equal to σ_1 ?
- If $A=xy^T$ has rank 1, what are u_1 , v_1 , σ_1 ? Check that $|\lambda_1| \le \sigma_1$
- What is the Karhunen-Loève transform and its connection to SVD? stochastic (random) form of PCA

Answers

(1)
$$\mathbf{S} = \mathbf{Q}\mathbf{A}\mathbf{Q}^{T} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{T} \rightarrow \mathbf{U} = \mathbf{V} = \mathbf{Q}, \ \mathbf{\Sigma} = \mathbf{A}$$

(2) $\mathbf{S}\mathbf{x} = (-\alpha)\mathbf{x} \leftrightarrow \mathbf{S}\mathbf{v} = \sigma\mathbf{u} : \mathbf{u} (\text{or } \mathbf{v}) = -\mathbf{x}, \ \sigma = \alpha$
(3) $\mathbf{A}^{T}\mathbf{A} = \mathbf{Q}^{T}\mathbf{Q} = \mathbf{I} \rightarrow \text{all } \sigma = 1 \rightarrow \mathbf{A} = \mathbf{Q} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{T} \leftrightarrow \mathbf{\Sigma} = \mathbf{I}, \ \mathbf{U} = \mathbf{Q}, \ \mathbf{V} = \mathbf{I}$
(4) $\|\mathbf{A}\mathbf{x}\| = \|\mathbf{U}\mathbf{\Sigma}\mathbf{V}^{T}\mathbf{x}\| = \|\mathbf{\Sigma}\mathbf{V}^{T}\mathbf{x}\| \le \sigma_{1}\|\mathbf{V}^{T}\mathbf{x}\| = \sigma_{1}\|\mathbf{x}\|$
 $\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \rightarrow \|\mathbf{A}\mathbf{x}\| = |\lambda|\|\mathbf{x}\|$
(5) $\mathbf{A} = \mathbf{x}\mathbf{y}^{T} = \frac{\mathbf{x}}{\|\mathbf{x}\|} (\|\mathbf{x}\|\|\|\mathbf{y}\|) \frac{\mathbf{y}^{T}}{\|\mathbf{y}\|} = \mathbf{u}_{1}\sigma_{1}\mathbf{v}_{1}^{T}, \ \mathbf{A}\mathbf{x} = (\mathbf{x}\mathbf{y}^{T})\mathbf{x} = \mathbf{x}(\mathbf{y}^{T}\mathbf{x}) = \lambda\mathbf{x} \rightarrow \lambda_{1} = |\mathbf{y}^{T}\mathbf{x}| \le \sigma_{1} = \|\mathbf{x}\|\|\mathbf{y}\|$

(6) KL begins with a covariance matrix V of a zero-mean random process. In general V could be an infinite matrix or a covaraiance function. Then the KL expansion will be an infinite series. The eigenvectors of V, in order of $\sigma_1^2 \ge \sigma_2^2 \ge ... \ge 0$, are the basis functions \mathbf{u}_i for the KL transform.

The expansion of any vector **v** in an orthonormal basis $\mathbf{u}_1, \mathbf{u}_2, \dots$ is $\mathbf{v} = \sum (\mathbf{u}_i^T \mathbf{v}) \mathbf{u}_i$.

In this stochastic case, that transform decorrelates the random process: the \mathbf{u}_i are independent.

More than that, the ordering of the eigenvalues means that the first k terms, stopping at $(\mathbf{u}_k^T \mathbf{v})\mathbf{u}_k$, minimize the expected square error.

Geometry of SVD

- A = (rotation)(stretching)(rotation) $U\Sigma V^T$ for every A
- If A is m by n and B is n by m, then AB and BA have the same nonzero eigenvalues



Applied Mathematics for Deep Learning

First singular vector v₁

Maximize the ratio $\frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} \rightarrow$ The maximum is σ_1 at the vector $\mathbf{x} = \mathbf{v}_1$ maximizing **x** is \mathbf{v}_1 : $\mathbf{A}\mathbf{v}_1 = \sigma_1 \mathbf{u}_1$ (the longest axis of the ellipse), $\|\mathbf{v}_1\| = 1 \rightarrow \|\mathbf{A}\mathbf{v}_1\| = \sigma_1$ $\Rightarrow \text{ Find the maximum value } \lambda \text{ of } \frac{\|\mathbf{A}\mathbf{x}\|^2}{\|\mathbf{x}\|^2} = \frac{(\mathbf{A}\mathbf{x})^T \mathbf{A}\mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\mathbf{x}^T \mathbf{S}\mathbf{x}}{\mathbf{x}^T \mathbf{x}}$ $\frac{\partial}{\partial \mathbf{x}_{i}} \left(\frac{\mathbf{x}^{T} \mathbf{S} \mathbf{x}}{\mathbf{x}^{T} \mathbf{x}} \right) = \left(\mathbf{x}^{T} \mathbf{x} \right) 2 \left(\mathbf{S} \mathbf{x} \right)_{i} - \left(\mathbf{x}^{T} \mathbf{S} \mathbf{x} \right) 2 \left(\mathbf{x} \right)_{i} = 0 \text{ for } i = 1, \dots, n$ $\rightarrow (\mathbf{S}\mathbf{x})_i = \left(\frac{\mathbf{x}^T \mathbf{S}\mathbf{x}}{\mathbf{x}^T \mathbf{x}}\right) (\mathbf{x})_i \rightarrow \mathbf{S}\mathbf{x} = \lambda \mathbf{x}$ Maximize $\frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}$ under the conditon $\mathbf{v}_1^T \mathbf{x} = \mathbf{0} \rightarrow$ The maximum is σ_2 at $\mathbf{x} = \mathbf{v}_2$

Polar decomposition

$$\begin{aligned} x + iy &= re^{i\theta} \rightarrow \begin{cases} e^{i\theta} : \text{ orthogonal matrix } \mathbf{Q} \\ r \ge 0: \text{ positive semideinite matrix } \mathbf{S} \end{cases} \\ \mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \left(\mathbf{U}\mathbf{V}^T\right)\left(\mathbf{V}\mathbf{\Sigma}\mathbf{V}^T\right) = \mathbf{Q}\mathbf{S} \\ \text{if } \mathbf{A} \text{ is invertible, then } \mathbf{\Sigma} \text{ and } \mathbf{S} \text{ are also invertible} \\ \mathbf{S}^2 = \mathbf{V}\mathbf{\Sigma}^2\mathbf{V}^T = \mathbf{A}^T\mathbf{A} \rightarrow \begin{cases} \text{eigenvalues of } \mathbf{S} = \text{singular values of } \mathbf{A} \\ \text{eigenvectors of } \mathbf{S} = \text{singular vectors } \mathbf{v} \text{ of } \mathbf{A} \end{cases} \\ \mathbf{A} = \begin{bmatrix} 3 & 0 \\ 4 & 5 \end{bmatrix} \rightarrow \mathbf{U} = \frac{1}{\sqrt{10}} \begin{bmatrix} 1 & -3 \\ 3 & 1 \end{bmatrix}, \mathbf{\Sigma} = \begin{bmatrix} \sqrt{45} \\ \sqrt{5} \end{bmatrix}, \mathbf{V} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \\ \mathbf{Q} = \mathbf{U}\mathbf{V}^T = \frac{1}{\sqrt{10}} \begin{bmatrix} 1 & -3 \\ 3 & 1 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 & -1 \\ 1 & 2 \end{bmatrix} \\ \mathbf{S} = \mathbf{V}\mathbf{\Sigma}\mathbf{V}^T = \frac{\sqrt{5}}{2} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} = \sqrt{5} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \\ \mathbf{A} = \mathbf{Q} \quad \mathbf{S} \\ \text{rotation stretch} \end{aligned}$$

I.9 Principal Components and the Best Low Rank Matrix

- Major tool in understanding a matrix of data
 - Schmidt(1907)→ Eckart and Young(1936, ||A||_F)→Mirsky(1955, any norm ||A||)
- Eckart-Young low rank approximation theorem

- The norm of $A-A_k$ is below the norm of all other $A-B_k$ Eckart-Young: If **B** has rank *k*, then $||A-B|| \ge ||A-A_k||$

 $\mathbf{A}_{k} = \sigma_{1}\mathbf{u}_{1}\mathbf{v}_{1}^{T} + \dots + \sigma_{k}\mathbf{u}_{k}\mathbf{v}_{k}^{T}$: the closest rank k matrix to A

Spectral norm:
$$\|\mathbf{A}\|_2 = \max_{\mathbf{x}\neq 0} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \sigma_1(\ell^2 \text{ norm})$$

Frobenius norm: $\|\mathbf{A}\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2}$

Nuclear norm: $\|\mathbf{A}\|_{N} = \sigma_{1} + \dots + \sigma_{r}$ (the trace norm)

$$\|\mathbf{I}\|_{2} = 1 = \|\mathbf{Q}\|_{2}, \ \|\mathbf{I}\|_{F} = \sqrt{n} = \|\mathbf{Q}\|_{F}, \ \|\mathbf{I}\|_{N} = n = \|\mathbf{Q}\|_{N}$$

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{T} \xrightarrow{\text{no change in } \|\mathbf{A}\|} \Rightarrow \|\mathbf{Q}_{1}\mathbf{A}\mathbf{Q}_{2}^{T}\| = \|\mathbf{Q}_{1}\mathbf{U}\mathbf{\Sigma}\mathbf{V}^{T}\mathbf{Q}_{2}^{T}\| = \|(\mathbf{Q}_{1}\mathbf{U})\mathbf{\Sigma}(\mathbf{Q}_{2}\mathbf{V})^{T}\|_{T}$$

Applied Mathematics for Deep Learning

Highlights of Linear Algebra - 24

 $= \|\mathbf{A}\|$

Eckart-Young Theorem: Best approximation by A_k

Eckart-Young Theorem: Best approximation by A_k

Eckart-Young in the Frobenius norm: If **B** is closest to **A**, then $\mathbf{U}^T \mathbf{B} \mathbf{V}$ is closest to $\mathbf{U}^T \mathbf{A} \mathbf{V}$ rank $(\mathbf{B}) \le k$ is closest to $\mathbf{A} \rightarrow \mathbf{B} = \mathbf{A}_k$

$$\mathbf{B} = \mathbf{U} \begin{bmatrix} \mathbf{D} & \mathbf{0} \\ k \times k \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}^{T} \rightarrow \underbrace{\mathbf{U}, \mathbf{V} \text{ not necessarily}}_{\text{diagonalize } \mathbf{A}, \text{ rank}(\mathbf{C}) \leq k} \rightarrow \mathbf{A} = \mathbf{U} \begin{bmatrix} \mathbf{L} + \mathbf{E} + \mathbf{R} & \mathbf{F} \\ \mathbf{G} & \mathbf{H} \end{bmatrix} \mathbf{V}^{T}, \mathbf{C} = \mathbf{U} \begin{bmatrix} \mathbf{L} + \mathbf{D} + \mathbf{R} & \mathbf{F} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}^{T}$$
$$\|\mathbf{A} - \mathbf{B}\|_{F}^{2} = \|\mathbf{A} - \mathbf{C}\|_{F}^{2} + \|\mathbf{L}\|_{F}^{2} + \|\mathbf{R}\|_{F}^{2} + \|\mathbf{F}\|_{F}^{2} \xrightarrow{\text{as small as possible}}{\mathbf{L} = \mathbf{R} = \mathbf{F} = \mathbf{G} = \mathbf{0}} \rightarrow \mathbf{A} = \mathbf{U} \begin{bmatrix} \mathbf{E} & \mathbf{0} \\ \mathbf{0} & \mathbf{H} \end{bmatrix} \mathbf{V}^{T}$$

The matrix **D** must be the same as $\mathbf{E} = diag(\sigma_1, \dots, \sigma_k)$ The singular values of **H** must be the smallest (n-k) singular values of **A**

The smallest error $\|\mathbf{A} - \mathbf{B}\|_F$ must be $\|\mathbf{H}\|_F = \sqrt{\sigma_{k+1}^2 + \dots + \sigma_r^2}$

Principal Component Analysis

- Understand *n* sample points in *m*-dimensional space
- Data matrix A₀: *n* samples, *m* variables
 - Find the average (the sample mean) along each row of A_0
 - Subtract that mean from *m* entries in the row
 - Centered matrix $A=A_0$ -(mean)
 - How will linear algebra find that closest line through (0,0)? It is in the direction of the first singular vector u_1 of A?



A is $\mathbf{2} \times \mathbf{n}$ (large nullspace)

 AA^{T} is $\mathbf{2} \times \mathbf{2}$ (small matrix)

 $A^{\mathrm{T}}A$ is $\boldsymbol{n} \times \boldsymbol{n}$ (large matrix)

Two singular values $\sigma_1 > \sigma_2 > 0$

- Statistics behind PCA
 - Variances: diagonal entries of the matrix AA^T
 - sum of squares of distances from the mean
 - Covariances: off- diagonal entries of the matrix AA^{T}
 - Sample covariance matrix: S=AA^T/(n-1)
 - One DOF has already been used for mean=0

$$\mathbf{A} = \begin{bmatrix} 3 & -4 & 7 & 1 & -4 & -3 \\ 7 & -6 & 8 & -1 & -1 & 7 \end{bmatrix} \rightarrow \mathbf{S} = \frac{1}{n-1} \mathbf{A} \mathbf{A}^T = \frac{1}{6-1} \begin{bmatrix} 20 & 25 \\ 25 & 40 \end{bmatrix} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$$
$$\mathbf{U} = \begin{bmatrix} -0.5606 & -0.8281 \\ -0.8281 & 0.5606 \end{bmatrix}, \ \mathbf{\Sigma} = \begin{bmatrix} 56.9258 & 0 \\ 0 & 3.0742 \end{bmatrix}, \ \mathbf{V} = \begin{bmatrix} -0.5606 & -0.8281 \\ -0.8281 & 0.5606 \end{bmatrix}$$

- Geometry behind PCA
 - Minimize perpendicular distances: perpendicular least square, orthogonal regression
 - Sum of squared distances from the data points to the line is a minimum $n = \frac{n}{2} + \frac{n}{2}$

$$\sum_{\substack{j=1\\\text{fixed by data}}}^{n} \left\| \mathbf{a}_{j} \right\|^{2} = \sum_{\substack{j=1\\ \rightarrow \mathbf{u}_{1}^{T} \left(\mathbf{A} \mathbf{A}^{T} \right) \mathbf{u}_{1}}}^{n} \left\| \mathbf{a}_{j}^{T} \mathbf{u}_{2} \right\|^{2}$$

- Linear algebra behind PCA
 - Singular values σ_i and singular vectors u_i of A \leftarrow eigenvalues σ_i^2 and eigenvectors of S=AA^T/(n-1)
 - Total variance:

$$T = \frac{\|\mathbf{A}\|_{F}^{2}}{n-1} = \frac{\|\mathbf{a}_{1}\|^{2} + \dots + \|\mathbf{a}_{n}\|^{2}}{n-1} = \frac{\sigma_{1}^{2} + \dots + \sigma_{r}^{2}}{n-1}$$