5. Orthogonal Matrices and Subspaces

Orthogonal vectors x and y

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

• Orthogonal basis for a subspace

orthogonal: $v_i^T v_j = 0 \xrightarrow{v_i / \|v_i\|}$ orthonormal: $v_i^T v_i = 1$

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

- 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}



Applied Mathematics for Deep

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

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

Pb is the orthogonal projection of **b** onto the column space of **P**



• Orthogonal matrices are square with orthonormal columns: $Q^T = Q^{-1}$ $Q^T = Q^{-1}$ $Q^T = Q^{-1} = Q^{-1}$

$$\mathbf{Q} \text{ is square} \rightarrow \left\{ \mathbf{Q} \mathbf{Q}^{T} = \mathbf{I} \right\} \rightarrow \mathbf{Q}^{T} = \mathbf{Q}$$

$$\mathbf{Q}_{\text{rotate}} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}, \quad \mathbf{Q}_{\text{reflect}} = \begin{bmatrix} \cos \theta & \sin \theta \\ \sin \theta & -\cos \theta \end{bmatrix}$$

$$\underbrace{\operatorname{rotation through an angle } \theta \quad \operatorname{reflection across the} \frac{\theta}{2} \text{ line}}_{\text{reflection across the} \frac{\theta}{2} \text{ line}}$$

 $\mathbf{Q}_1, \mathbf{Q}_2$: orthogonal $\rightarrow \mathbf{Q}_1 \mathbf{Q}_2$: orthogonal



Highlights of Linear Algebra - 14

All reflection matrices have eigenvalues -1 and 1

Examples

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

6. Eigenvalues and Eigenvectors

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

$$\mathbf{x}: \text{ eigenvector of } \mathbf{A} \\ \lambda: \text{ eigenvector of } \mathbf{A} \\ \lambda: \text{ eigenvector of } \mathbf{A} \\ \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} \\ n \times n \text{ matrices } \rightarrow n \text{ independent eigenvectors } \mathbf{x}_{1} \text{ to } \mathbf{x}_{n} \text{ with } n \text{ different eigenvalues } \lambda_{1} \text{ to } \lambda_{n} \\ \mathbf{v} = c_{1}\mathbf{x}_{1} + \dots + c_{n}\mathbf{x}_{n} \rightarrow \mathbf{A}\mathbf{v} = c_{1}\lambda_{1}\mathbf{x}_{1} + \dots + c_{n}\lambda_{n}\mathbf{x}_{n} \rightarrow \mathbf{A}^{k}\mathbf{v} = c_{1}\lambda_{1}^{k}\mathbf{x}_{1} + \dots + c_{n}\lambda_{n}^{k}\mathbf{x}_{n} \\ \text{How useful?} \begin{cases} \text{solution of differential equations} \\ \text{similar matrices } \rightarrow \text{ same eigenvalues} \\ \text{diagonalize a matrix} \end{cases} \\ \text{Four properties: matrix } \mathbf{A}(\text{real}), \mathbf{S}(\text{symmetric}), \quad \mathbf{Q}(\text{orthogonal}) \\ \text{like real numbers: } \lambda, \quad \mathbf{Q}(\text{orthogonal}) \\ \text{like complex numbers: } e^{i\theta} \end{cases} \\ \text{(Trace of } \mathbf{S}) \sum_{i=1}^{n} \lambda_{i} = \text{trace of matrix} \\ \text{(Determinant)} \quad \prod \lambda_{i} = \text{determinant of matrix} \end{cases} \\ \text{(Real eigenvalues of } \mathbf{S}) \mathbf{S}: \text{ real eigenvalues, orthogonal eigenvectors}} \\ \text{(Orthogonal eigenvectors)} \text{ if } \lambda_{1} \neq \lambda_{3}, \text{ then } \mathbf{x}_{1} \cdot \mathbf{x}_{2} = 0, \text{ eigenvectors of } \mathbf{A} \text{ 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 \\ 1 \end{bmatrix} = 3 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \text{ and } \mathbf{S} \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} = 1 \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}, \mathbf{Q} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \rightarrow \left\{ \mathbf{Q} \begin{bmatrix} 1 \\ -i \\ -i \end{bmatrix} = i \begin{bmatrix} 1 \\ -i \\ 1 \end{bmatrix} \text{ and } \mathbf{Q} \begin{bmatrix} 1 \\ i \\ i \end{bmatrix} = -i \begin{bmatrix} 1 \\ i \\ i \end{bmatrix} \end{cases}$$

Applied Mathematics for Deep Learning

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

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

Make **B** gradually into a triangular matrix \rightarrow Gradually show up on the main diagonal $\mathbf{B}\mathbf{y} = \lambda \mathbf{y} \rightarrow \mathbf{M}\mathbf{A}\mathbf{M}^{-1}\mathbf{y} = \lambda \mathbf{y} \rightarrow \mathbf{A}(\mathbf{M}^{-1}\mathbf{y}) = \lambda(\mathbf{M}^{-1}\mathbf{y})$

(3) diagonalize a matrix

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

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] $\begin{cases}
\left(\text{Geometric Multiplicity} = \text{GM}\right): \text{ count the independent eigenvectors, } \dim N(\mathbf{A} - \lambda \mathbf{I}) \\
\left(\text{Algebraic Multiplicity} = \text{AM}\right): \text{ count the repetitions of eigenvalues, } \det(\mathbf{A} - \lambda \mathbf{I}) = 0 \\
\mathbf{A} = \begin{bmatrix} 5 & 1 \\ 0 & 5 \end{bmatrix}, \begin{bmatrix} 6 & -1 \\ 1 & 4 \end{bmatrix}, \begin{bmatrix} 7 & 2 \\ -2 & 3 \end{bmatrix}
\end{cases}$

8. Singular Value Decomposition (SVD)

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

(

key point: two sets of singular vectors

$$\begin{bmatrix} n & \text{right singular vectors } (\mathbf{v}_1, \dots, \mathbf{v}_n) & \text{orthogonal in } \mathbf{R}^n \\ m & \text{left singular vectors } (\mathbf{u}_1, \dots, \mathbf{u}_m) & \text{orthogonal in } \mathbf{R}^m \end{bmatrix}$$

connection between $n \mathbf{v}$'s and $m \mathbf{u}$'s

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

$$\mathbf{A} \begin{bmatrix} \mathbf{v}_{1} \cdots \mathbf{v}_{r} \cdots \mathbf{v}_{n} \\ \mathbf{v}_{1} \cdots \mathbf{v}_{r} \cdots \mathbf{v}_{n} \end{bmatrix}_{r=1}^{r=1} = \begin{bmatrix} \mathbf{u}_{1} \cdots \mathbf{u}_{r} & \mathbf{u}_{n} \\ \mathbf{u}_{1} \cdots \mathbf{u}_{r} & \mathbf{u}_{n} \end{bmatrix}_{r=1}^{\sigma_{1}} \begin{bmatrix} \sigma_{1} & \cdots & \sigma_{r} \\ \mathbf{u}_{n} & \mathbf{u}_{r} & \mathbf{u}_{n} \end{bmatrix}_{r=1}^{\sigma_{1}} \begin{bmatrix} \mathbf{v}_{1} & \mathbf{v}_{r} \mathbf{v}_{r} \\ \mathbf{v}_{r} & \mathbf{v}_{r} \end{bmatrix}_{r=1}^{r=1} = \begin{bmatrix} \mathbf{u}_{1} & \mathbf{v}_{r} \mathbf{v}_{r} \\ \mathbf{v}_{r} & \mathbf{v}_{r} \end{bmatrix}_{r=1}^{\sigma_{1}} \begin{bmatrix} \sigma_{1} & \cdots & \sigma_{r} \\ \mathbf{v}_{r} & \mathbf{v}_{r} \end{bmatrix}_{r=1}^{\sigma_{1}} \begin{bmatrix} \sigma_{1} & \cdots & \sigma_{r} \end{bmatrix}_{r=1}^{r=1} = \begin{bmatrix} \mathbf{u}_{1} & \mathbf{v}_{r} \mathbf{u}_{r} \end{bmatrix}_{r=1}^{\sigma_{1}} \begin{bmatrix} \sigma_{1} & \cdots & \sigma_{r} \end{bmatrix}_{r=1}^{\sigma_{1}} \begin{bmatrix} \sigma_{1}$$

Applied Mathematics for Deep Learning

Proof of SVD

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

$$\begin{bmatrix} 3 & 0\\ 4 & 5\\ \end{bmatrix} \underbrace{\frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1\\ 1 & 1\\ \end{bmatrix}}_{\mathbf{V}^{T} = \mathbf{V}^{-1}} = \underbrace{\frac{1}{\sqrt{10}} \begin{bmatrix} 1 & -3\\ 3 & 1\\ \end{bmatrix}}_{\mathbf{U}^{T} = \mathbf{U}^{-1}} \underbrace{\frac{\sqrt{5}}{\sqrt{5}}}_{rank(\mathbf{A}) = 2 \to \sigma_{1}, \sigma_{2}}$$

$$\sigma_{1} \mathbf{u}_{1} \mathbf{v}_{1}^{T} + \sigma_{2} \mathbf{u}_{2} \mathbf{v}_{2}^{T} = \frac{3\sqrt{5}}{\sqrt{10}\sqrt{2}} \begin{bmatrix} 1\\ 3\\ \end{bmatrix} \begin{bmatrix} 1 & 1\\ \end{bmatrix} + \frac{\sqrt{5}}{\sqrt{10}\sqrt{2}} \begin{bmatrix} -3\\ 1\\ \end{bmatrix} \begin{bmatrix} -1 & 1\\ \end{bmatrix} = \frac{3}{2} \begin{bmatrix} 1 & 1\\ 3 & 3\\ \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 3 & -3\\ -1 & 1\\ \end{bmatrix} = \begin{bmatrix} 3 & 0\\ 4 & 5\\ \end{bmatrix} = \mathbf{A}$$

$$\begin{cases} \mathbf{V} \text{ contains orthonormal eigenvectors of } \mathbf{A}^{T} \mathbf{A} \\ \mathbf{U} \text{ contains orthonormal eigenvectors of } \mathbf{A} \mathbf{A}^{T} \\ \sigma_{1}^{2} \text{ to } \sigma_{r}^{2} \text{ are the nonzero eigenvalues of both } \mathbf{A}^{T} \mathbf{A} \text{ and } \mathbf{A} \mathbf{A}^{T} \end{cases}$$

$$SVD \text{ requires that } \mathbf{Av}_{k} = \sigma_{k} \mathbf{u}_{k} : \mathbf{v} \text{'s } \left(\mathbf{A}^{T} \mathbf{Av}_{k} = \sigma_{k}^{2} \mathbf{v}_{k}\right) \rightarrow \mathbf{u} \text{'s } (\mathbf{u}_{k} = \frac{\mathbf{Av}_{k}}{\sigma_{k}} : \text{ sign, multiple eigenvalues)}$$

$$\text{check } 1: \mathbf{u} \text{'s are eigenvectors of } \mathbf{A} \mathbf{A}^{T} \rightarrow \mathbf{A} \mathbf{A}^{T} \mathbf{u}_{k} =$$

$$\text{check } 2: \mathbf{u} \text{'s are also orthonormal } \rightarrow \mathbf{u}_{j}^{T} \mathbf{u}_{k} =$$

$$\text{choose } (n-r) \mathbf{v} \text{'s in } N(\mathbf{A}) \text{ and } (m-r) \mathbf{u} \text{'s in } N(\mathbf{A}^{T})$$

- Columns of V are orthogonal eigenvectors of A^TA
- Av=σu gives orthonormal eigenvectors u of AA^T
- σ^2 = eigenvalue of $A^T A$ = 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 Λ 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 $\|\mathbf{A} - \mathbf{A}_{k}\| \le \|\mathbf{A} - \mathbf{B}\|$

Example

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

- 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 $|\lambda_1| \le \sigma_1$

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



4 parameters: two angles, two numbers

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), $\|\mathbf{v}_1\| = 1 \rightarrow \|\mathbf{A}\mathbf{v}_1\| = \sigma_1$ $\Rightarrow \text{ Find the maximum value } \lambda \text{ of } \frac{\|\mathbf{A}\mathbf{x}\|^2}{\|\mathbf{x}\|^2} = \frac{(\mathbf{A}\mathbf{x})^T \mathbf{A}\mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\mathbf{x}^T \mathbf{S}\mathbf{x}}{\mathbf{x}^T \mathbf{x}}$ $\frac{\partial}{\partial x_i} \left(\frac{\mathbf{x}^T \mathbf{S} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \right) = \left(\mathbf{x}^T \mathbf{x} \right) 2 \left(\mathbf{S} \mathbf{x} \right)_i - \left(\mathbf{x}^T \mathbf{S} \mathbf{x} \right) 2 \left(\mathbf{x} \right)_i = 0 \text{ for } i = 1, \dots, n$ $\rightarrow (\mathbf{S}\mathbf{x})_i = \left(\frac{\mathbf{x}^T \mathbf{S}\mathbf{x}}{\mathbf{x}^T \mathbf{x}}\right) (\mathbf{x})_i \rightarrow \mathbf{S}\mathbf{x} = \lambda \mathbf{x}$ Maximize $\frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{y}\|}$ under the conditon $\mathbf{v}_1^T \mathbf{x} = 0 \rightarrow$ The maximum is σ_2 at $\mathbf{x} = \mathbf{v}_2$

Polar decomposition

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

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)
- Eckart-Young low rank approximation theorem
 - The norm of $A-A_k$ is below the norm of all other $A-B_k$

$$- A_k = \sigma_1 u_1 v_1^{\mathsf{T}} + \ldots + \sigma_k u_k v_k^{\mathsf{T}}$$

Eckart-Young: If **B** has rank k, then $\|\mathbf{A} - \mathbf{B}\| \ge \|\mathbf{A} - \mathbf{A}_k\|$ $\mathbf{A}_k = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \dots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T$: the closest rank k matrix to **A** $\begin{cases}
\text{Spectral norm: } \|\mathbf{A}\|_2 = \max_{\mathbf{x}\neq 0} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \sigma_1 (\ell^2 \text{ norm}) \\
\text{Frobenius norm: } \|\mathbf{A}\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2} \\
\text{Nuclear norm: } \|\mathbf{A}\|_N = \sigma_1 + \dots + \sigma_r \text{ (the trace norm)} \\
\|\mathbf{I}\|_2 = 1 = \|\mathbf{Q}\|_2, \|\mathbf{I}\|_F = \sqrt{n} = \|\mathbf{Q}\|_F, \|\mathbf{I}\|_N = n = \|\mathbf{Q}\|_N
\end{cases}$

Applied Mathematics for Deep Learning

Eckart-Young Theorem

Best approximation by A_k

Eckart-Young in L^2 :

If rank
$$(\mathbf{B}) \leq k$$
, then $\|\mathbf{A} - \mathbf{B}\| = \max_{\mathbf{x} \neq 0} \frac{\|(\mathbf{A} - \mathbf{B})\mathbf{x}\|}{\|\mathbf{x}\|} \geq \sigma_{k+1}$

Eckart-Young in the Frobenius norm:

If **B** is closest to **A**, then $\mathbf{U}^T \mathbf{B} \mathbf{V}$ is closest to $\mathbf{U}^T \mathbf{A} \mathbf{V}$

$$\mathbf{B} = \mathbf{U} \begin{bmatrix} \mathbf{D} & \mathbf{0} \\ {}^{k \times k} & \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}^{T}, \mathbf{A} = \begin{bmatrix} \mathbf{L} + \mathbf{E} + \mathbf{R} & \mathbf{F} \\ \mathbf{G} & \mathbf{H} \end{bmatrix}$$

The matrix **D** must be the same as $\mathbf{E} = diag(\sigma_1, ..., \sigma_k)$

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 $n \times n$ (large matrix)

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

- Statistics behind PCA
 - Variances: diagonal entries of the matrix AA^{T}
 - Covariances: off- diagonal entries of the matrix AA^{T}
 - Sample covariance matrix: S=AA^T/(n-1)
- Geometry behind PCA
 - Sum of squared distances from the data points to the line is a minimum
- Linear algebra behind PCA
 - Singular values σ_i and singular vectors u_i of A
 - Total variance:

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

11. Norms of Vectors and Matrices

- The norm of a nonzero vector v is a positive number ||v||
- That number measures the "length" of the vector

every norm for vectors or functions or matrice must share these two properties of the absolute value |c| of a number

All norms
$$\begin{cases} \text{multiply } \mathbf{v} \text{ by } c \text{ (rescaling)} \rightarrow \|c\mathbf{v}\| = |c|\|\mathbf{v}\| \\ \text{add } \mathbf{v} \text{ to } \mathbf{w} \text{ (Triangle inequality)} \rightarrow \|\mathbf{v} + \mathbf{w}\| \leq \|\mathbf{v}\| + \|\mathbf{w}\| \\ \\ \ell^2 \text{ norm } = \text{Euclidean norm: } \|\mathbf{v}\|_2 = \sqrt{|v_1|^2 + \dots + |v_n|^2} \\ \\ \ell^1 \text{ norm } = 1 \text{ norm: } \|\mathbf{v}\|_1 = |v_1| + \dots + |v_n| \\ \\ \ell^\infty \text{ norm } = \text{ max norm: } \|\mathbf{v}\|_\infty = \text{maximum of } |v_1|, \dots, |v_n| \\ \\ \|\mathbf{v}\|_p = \left(|v_1|^p + \dots + |v_n|^p\right)^{1/p} \end{cases}$$

Important vector norms and a failure



• Minimum of $||v||_p$ on the line $a_1v_1 + a_2v_2 = 1$



Inner products and S=norm

Inner product = length squared : $\mathbf{v} \cdot \mathbf{v} = \mathbf{v}^T \mathbf{v} = \|\mathbf{v}\|^2$ Angle θ between vector v and w: $\mathbf{v}^T \mathbf{w} = \|\mathbf{v}\| \|\mathbf{w}\| \cos \theta$ \rightarrow $\begin{cases} \text{Cauchy-Schwarz:} \|\mathbf{v}^T \mathbf{w}\| \le \|\mathbf{v}\| \|\mathbf{w}\| \\ \text{Triangle Inequality:} \|\mathbf{v} + \mathbf{w}\| \le \|\mathbf{v}\| + \|\mathbf{w}\| \end{cases}$ Choose any symmetric positive definite matrix \mathbf{S} $\|\mathbf{v}\|_{\mathbf{S}}^2 = \mathbf{v}^T \mathbf{S} \mathbf{v}$ gives a norm for \mathbf{v} in \mathfrak{R}^n (called the S - norm) $(\mathbf{v}, \mathbf{w})_{\mathbf{S}} = \mathbf{v}^T \mathbf{S} \mathbf{w}$ gives the S - inner product for \mathbf{v}, \mathbf{w} in \mathfrak{R}^n

Norm of Matrices

- Frobenius Norm
- Matrix Norm ||A|| from vector norm ||v||
- Nuclear Norm

$$\begin{cases} \left\|\mathbf{A}\right\|_{F} = \sqrt{\left|a_{11}\right|^{2} + \dots + \left|a_{1n}\right|^{2} + \dots + \left|a_{mn}\right|^{2}} \\ \left\|\mathbf{A}\right\|_{F} = \left\|\mathbf{U}\mathbf{\Sigma}\mathbf{V}^{T}\right\|_{F} = \left\|\mathbf{\Sigma}\mathbf{V}^{T}\right\|_{F} = \left\|\mathbf{\Sigma}\right\|_{F} = \sqrt{\sigma_{1}^{2} + \dots + \sigma_{r}^{2}} \\ \left\|\mathbf{A}\right\|_{F}^{2} = \text{trace of } \mathbf{A}^{T}\mathbf{A} = \text{sum of eigenvalues} = \sigma_{1}^{2} + \dots + \sigma_{r}^{2} \\ \left\|\mathbf{A}\right\| = \max_{\mathbf{v}\neq 0} \frac{\left\|\mathbf{A}\mathbf{v}\right\|}{\left\|\mathbf{v}\right\|} = \text{largest growth factor} \\ \ell^{2} \text{ norm : } \left\|\mathbf{A}\right\|_{2} = \text{largest singular value } \sigma_{1} \text{ of } \mathbf{A} \\ \ell^{1} \text{ norm: } \left\|\mathbf{A}\right\|_{1} = \text{largest } \ell^{1} \text{ norm of the columns of } \mathbf{A} \\ \ell^{\infty} \text{ norm : } \left\|\mathbf{A}\right\|_{\infty} = \text{largest } \ell^{1} \text{ norm of the rows of } \mathbf{A} \\ \left\|\mathbf{A}\right\|_{N} = \sigma_{1} + \dots + \sigma_{r} = \text{trace norm} \end{cases}$$