



Matrices

DS-GA 1013 / MATH-GA 2824 Optimization-based Data Analysis

http://www.cims.nyu.edu/~cfgranda/pages/OBDA_fall17/index.html

Carlos Fernandez-Granda

Basic properties

Singular value decomposition

Denoising

Collaborative filtering

Principal component analysis

Probabilistic interpretation

Dimensionality reduction

Eigendecomposition

Column and row space

The column space $\text{col}(A)$ of a matrix A is the **span** of its columns

The row space $\text{row}(A)$ is the **span** of its rows.

Rank

For any matrix A

$$\dim(\text{col}(A)) = \dim(\text{row}(A))$$

This is the **rank** of A

Orthogonal column spaces

If the column spaces of $A, B \in \mathbb{R}^{m \times n}$ are orthogonal then

$$\langle A, B \rangle = 0$$

Proof:

$$\langle A, B \rangle := \operatorname{tr}(A^T B) = \sum_{i=1}^n \langle A_{:,i}, B_{:,i} \rangle = 0$$

Consequence:

$$\|A + B\|_F^2 = \|A\|_F^2 + \|B\|_F^2$$

Linear maps

Given two vector spaces \mathcal{V} and \mathcal{R} associated to the same scalar field, a linear map $f : \mathcal{V} \rightarrow \mathcal{R}$ is a map from vectors in \mathcal{V} to vectors in \mathcal{R} such that for any scalar α and any vectors $\vec{x}_1, \vec{x}_2 \in \mathcal{V}$

$$f(\vec{x}_1 + \vec{x}_2) = f(\vec{x}_1) + f(\vec{x}_2)$$

$$f(\alpha \vec{x}_1) = \alpha f(\vec{x}_1)$$

Matrix-vector product

The product of a matrix $A \in \mathbb{C}^{m \times n}$ and a vector $\vec{x} \in \mathbb{C}^n$ is a vector $A\vec{x} \in \mathbb{C}^m$, such that

$$(A\vec{x})[i] = \sum_{j=1}^n A_{ij} \vec{x}[j]$$

$A\vec{x}$ is a **linear combination of the columns** of A

$$A\vec{x} = \sum_{j=1}^n \vec{x}[j] A_{:,j}$$

Equivalence between matrices and linear maps

For finite m, n every linear map $f : \mathbb{C}^m \rightarrow \mathbb{C}^n$ can be uniquely represented by a matrix $F \in \mathbb{C}^{m \times n}$

Proof

The matrix is

$$F := [f(\vec{e}_1) \quad f(\vec{e}_2) \quad \cdots \quad f(\vec{e}_n)],$$

the columns are the result of applying f to the **standard basis**

For any vector $\vec{x} \in \mathbb{C}^n$

$$\begin{aligned} f(x) &= f\left(\sum_{i=1}^n \vec{x}[i] \vec{e}_i\right) \\ &= \sum_{i=1}^n \vec{x}[i] f(\vec{e}_i) \\ &= F\vec{x} \end{aligned}$$

Projecting and lifting

When a matrix $\mathbb{C}^{m \times n}$ is *fat*, i.e., $n > m$, we say that it **projects** vectors onto a lower dimensional space (this is **not** an orthogonal projection!)

When a matrix is *tall*, i.e., $m > n$, we say that it **lifts** vectors to a higher-dimensional space

Adjoint

Given two vector spaces \mathcal{V} and \mathcal{R} with inner products $\langle \cdot, \cdot \rangle_{\mathcal{V}}$ and $\langle \cdot, \cdot \rangle_{\mathcal{R}}$, the **adjoint** $f^* : \mathcal{R} \rightarrow \mathcal{V}$ of a linear map $f : \mathcal{V} \rightarrow \mathcal{R}$ satisfies

$$\langle f(\vec{x}), \vec{y} \rangle_{\mathcal{R}} = \langle \vec{x}, f^*(\vec{y}) \rangle_{\mathcal{V}}$$

for all $\vec{x} \in \mathcal{V}$ and $\vec{y} \in \mathcal{R}$

Conjugate transpose

The entries of the conjugate transpose $A^* \in \mathbb{C}^{n \times m}$ of $A \in \mathbb{C}^{m \times n}$ are

$$(A^*)_{ij} = \overline{A_{ji}}, \quad 1 \leq i \leq n, \quad 1 \leq j \leq m.$$

If the entries are all real, this is just the **transpose**

The adjoint $f^* : \mathbb{C}^n \rightarrow \mathbb{C}^m$ of a $f : \mathbb{C}^m \rightarrow \mathbb{C}^n$ represented by a matrix F corresponds to **F^***

Symmetric and Hermitian

A symmetric matrix is equal to its transpose

A Hermitian or self-adjoint matrix is equal to its conjugate transpose

Range

Let \mathcal{V} and \mathcal{R} be vector spaces associated to the same scalar field, the range of a map $f : \mathcal{V} \rightarrow \mathcal{R}$ is the set of vectors in \mathcal{R} reached by f

$$\text{range}(f) := \{\vec{y} \mid \vec{y} = f(\vec{x}) \text{ for some } \vec{x} \in \mathcal{V}\}$$

The range of a matrix is the range of its associated linear map

The range is the column space

For any matrix $A \in \mathbb{C}^{m \times n}$

$$\text{range}(A) = \text{col}(A)$$

Proof:

$\text{col}(A) \subseteq \text{range}(A)$ because $A_{:i} = A\vec{e}_i$ for $1 \leq i \leq n$

$\text{range}(A) \subseteq \text{col}(A)$ because $A\vec{x}$ is a linear combination of the columns of A for any $\vec{x} \in \mathbb{C}^n$

Null space

Let \mathcal{V} and \mathcal{R} be vector spaces, the null space of a map $f : \mathcal{V} \rightarrow \mathcal{R}$ is the set of vectors that are **mapped to zero**:

$$\text{null}(f) := \{ \vec{x} \mid f(\vec{x}) = \vec{0} \}$$

The null space of a matrix is the null space of its associated linear map

The null space of a linear map is a **subspace**

Null space and row space

For any matrix $A \in \mathbb{R}^{m \times n}$

$$\text{null}(A) = \text{row}(A)^\perp$$

Proof:

Any vector $\vec{x} \in \text{row}(A)$ can be written as $\vec{x} = A^T \vec{z}$, for some $\vec{z} \in \mathbb{R}^m$

If $y \in \text{null}(A)$ then

$$\langle \vec{y}, \vec{x} \rangle =$$

Null space and row space

For any matrix $A \in \mathbb{R}^{m \times n}$

$$\text{null}(A) = \text{row}(A)^\perp$$

Proof:

Any vector $\vec{x} \in \text{row}(A)$ can be written as $\vec{x} = A^T \vec{z}$, for some $\vec{z} \in \mathbb{R}^m$

If $y \in \text{null}(A)$ then

$$\langle \vec{y}, \vec{x} \rangle = \langle \vec{y}, A^T \vec{z} \rangle$$

Null space and row space

For any matrix $A \in \mathbb{R}^{m \times n}$

$$\text{null}(A) = \text{row}(A)^\perp$$

Proof:

Any vector $\vec{x} \in \text{row}(A)$ can be written as $\vec{x} = A^T \vec{z}$, for some $\vec{z} \in \mathbb{R}^m$

If $y \in \text{null}(A)$ then

$$\begin{aligned}\langle \vec{y}, \vec{x} \rangle &= \langle \vec{y}, A^T \vec{z} \rangle \\ &= \langle A\vec{y}, \vec{z} \rangle\end{aligned}$$

Null space and row space

For any matrix $A \in \mathbb{R}^{m \times n}$

$$\text{null}(A) = \text{row}(A)^\perp$$

Proof:

Any vector $\vec{x} \in \text{row}(A)$ can be written as $\vec{x} = A^T \vec{z}$, for some $\vec{z} \in \mathbb{R}^m$

If $y \in \text{null}(A)$ then

$$\begin{aligned}\langle \vec{y}, \vec{x} \rangle &= \langle \vec{y}, A^T \vec{z} \rangle \\ &= \langle A\vec{y}, \vec{z} \rangle \\ &= 0\end{aligned}$$

Null space and range

Let $A \in \mathbb{R}^{m \times n}$

$$\dim(\text{range}(A)) + \dim(\text{null}(A)) = n$$

Identity matrix

The identity matrix of dimensions $n \times n$ is

$$I = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ & & \cdots & \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

For any $\vec{x} \in \mathbb{C}^n$, $I\vec{x} = \vec{x}$

Matrix inverse

The inverse of $A \in \mathbb{C}^{n \times n}$ is a matrix $A^{-1} \in \mathbb{C}^{n \times n}$ such that

$$AA^{-1} = A^{-1}A = I$$

Orthogonal matrices

An orthogonal matrix is a square matrix such that

$$U^T U = U U^T = I$$

The columns $U_{:1}, U_{:2}, \dots, U_{:n}$ form an **orthonormal basis**

For any $\vec{x} \in \mathbb{R}^n$

$$\vec{x} = U U^T \vec{x} = \sum_{i=1}^n \langle U_{:i}, \vec{x} \rangle U_{:i}$$

Product of orthogonal matrices

If $U, V \in \mathbb{R}^{n \times n}$ are orthogonal matrices, then UV is also an orthogonal matrix

$$(UV)^T (UV) = V^T U^T UV = I$$

Orthogonal matrices

Orthogonal matrices change the **direction** of vectors, not their magnitude

$$\|U\vec{x}\|_2^2$$

Orthogonal matrices

Orthogonal matrices change the **direction** of vectors, not their magnitude

$$\|U\vec{x}\|_2^2 = \vec{x}^T U^T U \vec{x}$$

Orthogonal matrices

Orthogonal matrices change the **direction** of vectors, not their magnitude

$$\begin{aligned}\|U\vec{x}\|_2^2 &= \vec{x}^T U^T U \vec{x} \\ &= \vec{x}^T \vec{x}\end{aligned}$$

Orthogonal matrices

Orthogonal matrices change the **direction** of vectors, not their magnitude

$$\begin{aligned}\|U\vec{x}\|_2^2 &= \vec{x}^T U^T U \vec{x} \\ &= \vec{x}^T \vec{x} \\ &= \|\vec{x}\|_2^2\end{aligned}$$

Orthogonal-projection matrix

Given a subspace $\mathcal{S} \subseteq \mathbb{R}^n$ of dimension d , the matrix

$$P := UU^T$$

where the columns of $U_{:1}, U_{:2}, \dots, U_{:d}$ are an orthonormal basis of \mathcal{S} , maps any vector \vec{x} to $\mathcal{P}_{\mathcal{S}} \vec{x}$

$$\begin{aligned} P\vec{x} &= UU^T\vec{x} \\ &= \sum_{i=1}^d \langle U_{:i}, \vec{x} \rangle U_{:i} \\ &= \mathcal{P}_{\mathcal{S}} \vec{x} \end{aligned}$$

Basic properties

Singular value decomposition

Denoising

Collaborative filtering

Principal component analysis

Probabilistic interpretation

Dimensionality reduction

Eigendecomposition

Singular value decomposition

Every rank r real matrix $A \in R^{m \times n}$, has a singular-value decomposition (SVD) of the form

$$A = \begin{bmatrix} \vec{u}_1 & \vec{u}_2 & \cdots & \vec{u}_r \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ & & \ddots & \\ 0 & 0 & \cdots & \sigma_r \end{bmatrix} \begin{bmatrix} \vec{v}_1^T \\ \vec{v}_2^T \\ \vdots \\ \vec{v}_r^T \end{bmatrix}$$
$$= USV^T$$

Singular value decomposition

- ▶ The **singular values** $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ are positive real numbers
- ▶ The **left** singular vectors $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_r$ form an orthonormal set
- ▶ The **right** singular vectors $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_r$ also form an orthonormal set
- ▶ The SVD is **unique** if all the singular values are different
- ▶ If $\sigma_i = \sigma_{i+1} = \dots = \sigma_{i+k}$, then $\vec{u}_i, \dots, \vec{u}_{i+k}$ can be replaced by any orthonormal basis of their span (the same holds for $\vec{v}_i, \dots, \vec{v}_{i+k}$)
- ▶ The SVD of an $m \times n$ matrix with $m \geq n$ can be computed in $\mathcal{O}(mn^2)$

Column and row space

- ▶ The **left** singular vectors $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_r$ are a basis for the **column space**
- ▶ The **right** singular vectors $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_r$ are a basis for the **row space**

Proof:

$$\text{span}(\vec{u}_1, \dots, \vec{u}_r) \subseteq \text{col}(A) \text{ since } \vec{u}_i = A(\sigma_i^{-1} \vec{v}_i)$$

$$\text{col}(A) \subseteq \text{span}(\vec{u}_1, \dots, \vec{u}_r) \text{ because } A_{:i} = U(SV^T \vec{e}_i)$$

Singular value decomposition

$$A = \left[\underbrace{\vec{u}_1 \cdots \vec{u}_r}_{\text{Basis of range}(A)} \quad \vec{u}_{r+1} \cdots \vec{u}_n \right] \begin{bmatrix} \sigma_1 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \cdots & \sigma_r & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ & & \cdots & \cdots & & \\ 0 & \cdots & 0 & 0 & \cdots & 0 \end{bmatrix} \left[\underbrace{\vec{v}_1 \vec{v}_2 \cdots \vec{v}_r}_{\text{Basis of row}(A)} \quad \underbrace{\vec{v}_{r+1} \cdots \vec{v}_n}_{\text{Basis of null}(A)} \right]^T$$

Rank and numerical rank

The rank of a matrix is equal to the number of nonzero singular values

Given a tolerance $\epsilon > 0$, the **numerical rank** is the number of singular values greater than ϵ

Linear maps

The SVD decomposes the action of a matrix $A \in \mathbb{R}^{m \times n}$ on a vector $\vec{x} \in \mathbb{R}^n$ into:

1. Rotation

$$V^T \vec{x} = \sum_{i=1}^n \langle \vec{v}_i, \vec{x} \rangle \vec{e}_i$$

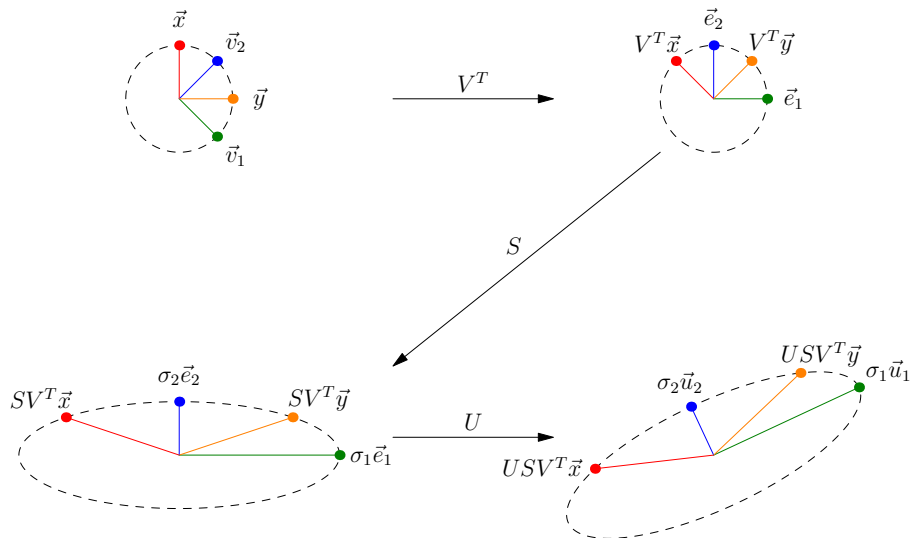
2. Scaling

$$SV^T \vec{x} = \sum_{i=1}^n \sigma_i \langle \vec{v}_i, \vec{x} \rangle \vec{e}_i$$

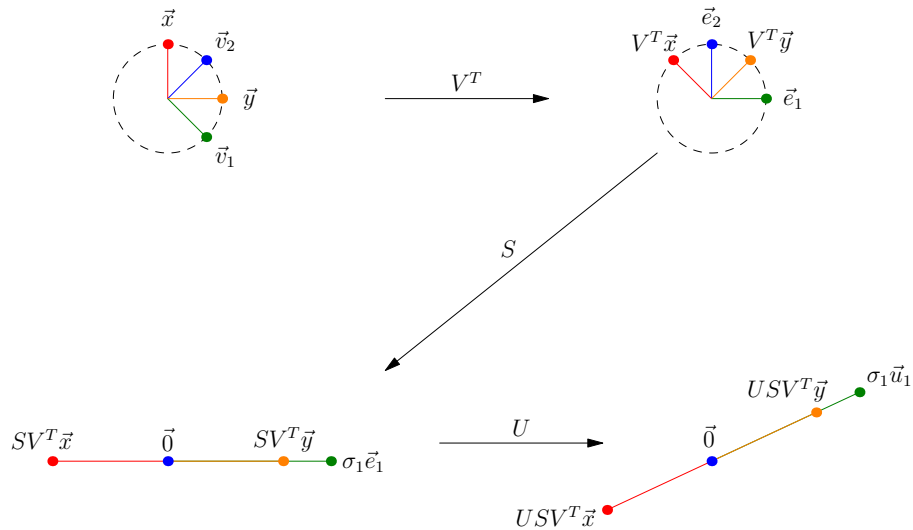
3. Rotation

$$USV^T \vec{x} = \sum_{i=1}^n \sigma_i \langle \vec{v}_i, \vec{x} \rangle \vec{u}_i$$

Linear maps



Linear maps



Singular values

For any orthogonal matrices $\tilde{U} \in \mathbb{R}^{m \times m}$ and $\tilde{V} \in \mathbb{R}^{n \times n}$ the singular values of $\tilde{U}A$ and $A\tilde{V}$ are **the same** as those of A

Proof:

$$\begin{aligned}\bar{U} &:= \tilde{U}U \\ \bar{V}^T &:= V^T\tilde{V}\end{aligned}$$

are orthogonal, so $\bar{U}SV^T$ and $US\bar{V}^T$ are valid SVDs for $\tilde{U}A$ and $A\tilde{V}$

Singular values

The singular values satisfy

$$\begin{aligned}\sigma_1 &= \max_{\{\|\vec{x}\|_2=1 \mid \vec{x} \in \mathbb{R}^n\}} \|A\vec{x}\|_2 \\ &= \max_{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m\}} \|A^T \vec{y}\|_2 \\ \sigma_i &= \max_{\{\|\vec{x}\|_2=1 \mid \vec{x} \in \mathbb{R}^n, \vec{x} \perp \vec{u}_1, \dots, \vec{u}_{i-1}\}} \|A\vec{x}\|_2 \\ &= \max_{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m, \vec{y} \perp \vec{v}_1, \dots, \vec{v}_{i-1}\}} \|A^T \vec{y}\|_2, \quad 2 \leq i \leq \min\{m, n\}\end{aligned}$$

Proof

Consider $\vec{x} \in \mathbb{R}^n$ such that $\|\vec{x}\|_2 = 1$ and for a fixed $1 \leq i \leq n$

$$\vec{x} \perp \vec{v}_1, \dots, \vec{v}_{i-1}$$

We decompose \vec{x} into

$$\vec{x} = \sum_{j=i}^n \alpha_j \vec{v}_j + \mathcal{P}_{\text{row}(A)^\perp} \vec{x}$$

where $1 = \|\vec{x}\|_2^2 \geq \sum_{j=i}^n \alpha_j^2$

Proof

$$\|A\vec{x}\|_2^2$$

Proof

$$\|A\vec{x}\|_2^2 = \left\langle \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k, \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k \right\rangle$$

Proof

$$\begin{aligned}\|A\vec{x}\|_2^2 &= \left\langle \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k, \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k \right\rangle \\ &= \sum_{k=1}^n \sigma_k^2 \langle \vec{v}_k, \vec{x} \rangle^2 \quad \text{because } \vec{u}_1, \dots, \vec{u}_n \text{ are orthonormal}\end{aligned}$$

Proof

$$\begin{aligned}\|A\vec{x}\|_2^2 &= \left\langle \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k, \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k \right\rangle \\ &= \sum_{k=1}^n \sigma_k^2 \langle \vec{v}_k, \vec{x} \rangle^2 \quad \text{because } \vec{u}_1, \dots, \vec{u}_n \text{ are orthonormal} \\ &= \sum_{k=1}^n \sigma_k^2 \left\langle \vec{v}_k, \sum_{j=1}^n \alpha_j \vec{v}_j + \mathcal{P}_{\text{row}(A)^\perp} \vec{x} \right\rangle^2\end{aligned}$$

Proof

$$\begin{aligned}\|A\vec{x}\|_2^2 &= \left\langle \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k, \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k \right\rangle \\ &= \sum_{k=1}^n \sigma_k^2 \langle \vec{v}_k, \vec{x} \rangle^2 \quad \text{because } \vec{u}_1, \dots, \vec{u}_n \text{ are orthonormal} \\ &= \sum_{k=1}^n \sigma_k^2 \left\langle \vec{v}_k, \sum_{j=i}^n \alpha_j \vec{v}_j + \mathcal{P}_{\text{row}(A)^\perp} \vec{x} \right\rangle^2 \\ &= \sum_{j=i}^n \sigma_j^2 \alpha_j^2 \quad \text{because } \vec{v}_1, \dots, \vec{v}_n \text{ are orthonormal}\end{aligned}$$

Proof

$$\begin{aligned}\|A\vec{x}\|_2^2 &= \left\langle \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k, \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k \right\rangle \\ &= \sum_{k=1}^n \sigma_k^2 \langle \vec{v}_k, \vec{x} \rangle^2 \quad \text{because } \vec{u}_1, \dots, \vec{u}_n \text{ are orthonormal} \\ &= \sum_{k=1}^n \sigma_k^2 \left\langle \vec{v}_k, \sum_{j=i}^n \alpha_j \vec{v}_j + \mathcal{P}_{\text{row}(A)^\perp} \vec{x} \right\rangle^2 \\ &= \sum_{j=i}^n \sigma_j^2 \alpha_j^2 \quad \text{because } \vec{v}_1, \dots, \vec{v}_n \text{ are orthonormal} \\ &\leq \sigma_i^2 \sum_{j=i}^n \alpha_j^2 \quad \text{because } \sigma_i \geq \sigma_{i+1} \geq \dots \geq \sigma_n\end{aligned}$$

Proof

$$\begin{aligned}\|A\vec{x}\|_2^2 &= \left\langle \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k, \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{x} \rangle \vec{u}_k \right\rangle \\ &= \sum_{k=1}^n \sigma_k^2 \langle \vec{v}_k, \vec{x} \rangle^2 \quad \text{because } \vec{u}_1, \dots, \vec{u}_n \text{ are orthonormal} \\ &= \sum_{k=1}^n \sigma_k^2 \left\langle \vec{v}_k, \sum_{j=i}^n \alpha_j \vec{v}_j + \mathcal{P}_{\text{row}(A)^\perp} \vec{x} \right\rangle^2 \\ &= \sum_{j=i}^n \sigma_j^2 \alpha_j^2 \quad \text{because } \vec{v}_1, \dots, \vec{v}_n \text{ are orthonormal} \\ &\leq \sigma_i^2 \sum_{j=i}^n \alpha_j^2 \quad \text{because } \sigma_i \geq \sigma_{i+1} \geq \dots \geq \sigma_n \\ &\leq \sigma_i^2\end{aligned}$$

Singular vectors

The right singular vectors satisfy

$$\begin{aligned}\vec{v}_1 &= \arg \max_{\{\|\vec{x}\|_2=1 \mid \vec{x} \in \mathbb{R}^n\}} \|A\vec{x}\|_2 \\ \vec{v}_i &= \arg \max_{\{\|\vec{x}\|_2=1 \mid \vec{x} \in \mathbb{R}^n, \vec{x} \perp \vec{v}_1, \dots, \vec{v}_{i-1}\}} \|A\vec{x}\|_2, \quad 2 \leq i \leq m\end{aligned}$$

and the left singular vectors satisfy

$$\begin{aligned}\vec{u}_1 &= \arg \max_{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m\}} \|A^T \vec{y}\|_2 \\ \vec{u}_i &= \arg \max_{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m, \vec{y} \perp \vec{u}_1, \dots, \vec{u}_{i-1}\}} \|A^T \vec{y}\|_2, \quad 2 \leq i \leq n\end{aligned}$$

Proof

\vec{v}_i achieves the maximum

$$\|A\vec{v}_i\|_2^2$$

Same proof for \vec{u}_i replacing A by A^T

Proof

\vec{v}_i achieves the maximum

$$\|A\vec{v}_i\|_2^2 = \left\langle \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{v}_i \rangle \vec{u}_k, \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{v}_i \rangle \vec{u}_k \right\rangle$$

Same proof for \vec{u}_i replacing A by A^T

Proof

\vec{v}_i achieves the maximum

$$\begin{aligned}\|A\vec{v}_i\|_2^2 &= \left\langle \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{v}_i \rangle \vec{u}_k, \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{v}_i \rangle \vec{u}_k \right\rangle \\ &= \sum_{k=1}^n \sigma_k^2 \langle \vec{v}_k, \vec{v}_i \rangle^2\end{aligned}$$

Same proof for \vec{u}_i replacing A by A^T

Proof

\vec{v}_i achieves the maximum

$$\begin{aligned}\|A\vec{v}_i\|_2^2 &= \left\langle \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{v}_i \rangle \vec{u}_k, \sum_{k=1}^n \sigma_k \langle \vec{v}_k, \vec{v}_i \rangle \vec{u}_k \right\rangle \\ &= \sum_{k=1}^n \sigma_k^2 \langle \vec{v}_k, \vec{v}_i \rangle^2 \\ &= \sigma_i^2\end{aligned}$$

Same proof for \vec{u}_i replacing A by A^T

Optimal subspace for orthogonal projection

Given a set of vectors $\vec{a}_1, \vec{a}_2, \dots, \vec{a}_n$ and a fixed dimension $k \leq n$, the SVD of

$$A := [\vec{a}_1 \quad \vec{a}_2 \quad \cdots \quad \vec{a}_n] \in \mathbb{R}^{m \times n}$$

provides the k -dimensional subspace that captures the most energy

$$\sum_{i=1}^n \|\mathcal{P}_{\text{span}(\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k)} \vec{a}_i\|_2^2 \geq \sum_{i=1}^n \|\mathcal{P}_{\mathcal{S}} \vec{a}_i\|_2^2$$

for any subspace \mathcal{S} of dimension k

Proof

Because $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k$ are orthonormal

$$\sum_{i=1}^n \left\| \mathcal{P}_{\text{span}(\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k)} \vec{a}_i \right\|_2^2$$

Proof

Because $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k$ are orthonormal

$$\sum_{i=1}^n \left\| \mathcal{P}_{\text{span}(\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k)} \vec{a}_i \right\|_2^2 = \sum_{i=1}^n \sum_{j=1}^k \langle \vec{u}_j, \vec{a}_i \rangle^2$$

Proof

Because $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k$ are orthonormal

$$\begin{aligned}\sum_{i=1}^n \left\| \mathcal{P}_{\text{span}(\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k)} \vec{a}_i \right\|_2^2 &= \sum_{i=1}^n \sum_{j=1}^k \langle \vec{u}_j, \vec{a}_i \rangle^2 \\ &= \sum_{j=1}^k \left\| A^T \vec{u}_j \right\|_2^2\end{aligned}$$

Proof

Because $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k$ are orthonormal

$$\begin{aligned}\sum_{i=1}^n \left\| \mathcal{P}_{\text{span}(\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k)} \vec{a}_i \right\|_2^2 &= \sum_{i=1}^n \sum_{j=1}^k \langle \vec{u}_j, \vec{a}_i \rangle^2 \\ &= \sum_{j=1}^k \left\| A^T \vec{u}_j \right\|_2^2\end{aligned}$$

Induction on k

Proof

Because $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k$ are orthonormal

$$\begin{aligned}\sum_{i=1}^n \left\| \mathcal{P}_{\text{span}(\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k)} \vec{a}_i \right\|_2^2 &= \sum_{i=1}^n \sum_{j=1}^k \langle \vec{u}_j, \vec{a}_i \rangle^2 \\ &= \sum_{j=1}^k \left\| A^T \vec{u}_j \right\|_2^2\end{aligned}$$

Induction on k

The base case $k = 1$ follows from

$$\vec{u}_1 = \underset{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m\}}{\text{arg max}} \left\| A^T \vec{y} \right\|_2$$

Proof

Let \mathcal{S} be a subspace of dimension k

$\mathcal{S} \cap \text{span}(\vec{u}_1, \dots, \vec{u}_{k-1})^\perp$ contains a nonzero vector \vec{b}

If $\dim(\mathcal{V})$ has dimension n , $\mathcal{S}_1, \mathcal{S}_2 \subseteq \mathcal{V}$ and $\dim(\mathcal{S}_1) + \dim(\mathcal{S}_2) > n$,
then $\dim(\mathcal{S}_1 \cap \mathcal{S}_2) \geq 1$

Proof

Let \mathcal{S} be a subspace of dimension k

$\mathcal{S} \cap \text{span}(\vec{u}_1, \dots, \vec{u}_{k-1})^\perp$ contains a nonzero vector \vec{b}

If $\dim(\mathcal{V})$ has dimension n , $\mathcal{S}_1, \mathcal{S}_2 \subseteq \mathcal{V}$ and $\dim(\mathcal{S}_1) + \dim(\mathcal{S}_2) > n$,
then $\dim(\mathcal{S}_1 \cap \mathcal{S}_2) \geq 1$

There exists an orthonormal basis $\vec{b}_1, \vec{b}_2, \dots, \vec{b}_k$ for \mathcal{S} such that $\vec{b}_k := \vec{b}$
is **orthogonal** to $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_{k-1}$

Induction hypothesis

$$\begin{aligned}\sum_{i=1}^{k-1} \left\| A^T \vec{u}_i \right\|_2^2 &= \sum_{i=1}^n \left\| \mathcal{P}_{\text{span}(\vec{u}_1, \vec{u}_2, \dots, \vec{u}_{k-1})} \vec{a}_i \right\|_2^2 \\ &\geq \sum_{i=1}^n \left\| \mathcal{P}_{\text{span}(\vec{b}_1, \vec{b}_2, \dots, \vec{b}_{k-1})} \vec{a}_i \right\|_2^2 \\ &= \sum_{i=1}^{k-1} \left\| A^T \vec{b}_i \right\|_2^2\end{aligned}$$

Proof

Recall that

$$\vec{u}_k = \underset{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m, \vec{y} \perp \vec{u}_1, \dots, \vec{u}_{k-1}\}}{\text{arg max}} \left\| A^T \vec{y} \right\|_2$$

which implies

$$\left\| A^T \vec{u}_k \right\|_2^2 \geq \left\| A^T \vec{b}_k \right\|_2^2$$

Proof

Recall that

$$\vec{u}_k = \arg \max_{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m, \vec{y} \perp \vec{u}_1, \dots, \vec{u}_{k-1}\}} \left\| A^T \vec{y} \right\|_2$$

which implies

$$\left\| A^T \vec{u}_k \right\|_2^2 \geq \left\| A^T \vec{b}_k \right\|_2^2$$

We conclude

$$\begin{aligned} \sum_{i=1}^n \left\| \mathcal{P}_{\text{span}(\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k)} \vec{a}_i \right\|_2^2 &= \sum_{i=1}^k \left\| A^T \vec{u}_i \right\|_2^2 \\ &\geq \sum_{i=1}^k \left\| A^T \vec{b}_i \right\|_2^2 \\ &= \sum_{i=1}^n \left\| \mathcal{P}_S \vec{a}_i \right\|_2^2 \end{aligned}$$

Best rank- k approximation

Let USV^T be the SVD of a matrix $A \in \mathbb{R}^{m \times n}$

The truncated SVD $U_{:,1:k} S_{1:k,1:k} V_{:,1:k}^T$ is the **best rank- k approximation**

$$U_{:,1:k} S_{1:k,1:k} V_{:,1:k}^T = \arg \min_{\{\tilde{A} \mid \text{rank}(\tilde{A})=k\}} \left\| A - \tilde{A} \right\|_F$$

Proof

Let \tilde{A} be an arbitrary matrix in $\mathbb{R}^{m \times n}$ with $\text{rank}(\tilde{A}) = k$

Let $\tilde{U} \in \mathbb{R}^{m \times k}$ be a matrix with orthonormal columns such that $\text{col}(\tilde{U}) = \text{col}(\tilde{A})$

$$\begin{aligned} \left\| U_{:,1:k} U_{:,1:k}^T A \right\|_F^2 &= \sum_{i=1}^n \left\| \mathcal{P}_{\text{col}(U_{:,1:k})} \vec{a}_i \right\|_2^2 \\ &\geq \sum_{i=1}^n \left\| \mathcal{P}_{\text{col}(\tilde{U})} \vec{a}_i \right\|_2^2 \\ &= \left\| \tilde{U} \tilde{U}^T A \right\|_F^2 \end{aligned}$$

Orthogonal column spaces

If the column spaces of $A, B \in \mathbb{R}^{m \times n}$ are orthogonal then

$$\|A + B\|_F^2 = \|A\|_F^2 + \|B\|_F^2$$

Proof

$\text{col}(A - \tilde{U}\tilde{U}^T A)$ is orthogonal to $\text{col}(\tilde{A}) = \text{col}(\tilde{U})$

$$\|A - \tilde{A}\|_F^2 = \|A - \tilde{U}\tilde{U}^T A\|_F^2 + \|\tilde{A} - \tilde{U}\tilde{U}^T A\|_F^2$$

Proof

$\text{col}(A - \tilde{U}\tilde{U}^T A)$ is orthogonal to $\text{col}(\tilde{A}) = \text{col}(\tilde{U})$

$$\begin{aligned}\|A - \tilde{A}\|_F^2 &= \|A - \tilde{U}\tilde{U}^T A\|_F^2 + \|\tilde{A} - \tilde{U}\tilde{U}^T A\|_F^2 \\ &\geq \|A - \tilde{U}\tilde{U}^T A\|_F^2\end{aligned}$$

Proof

$\text{col}(A - \tilde{U}\tilde{U}^T A)$ is orthogonal to $\text{col}(\tilde{A}) = \text{col}(\tilde{U})$

$$\begin{aligned}\|A - \tilde{A}\|_F^2 &= \|A - \tilde{U}\tilde{U}^T A\|_F^2 + \|\tilde{A} - \tilde{U}\tilde{U}^T A\|_F^2 \\ &\geq \|A - \tilde{U}\tilde{U}^T A\|_F^2 \\ &= \|A\|_F^2 - \|\tilde{U}\tilde{U}^T A\|_F^2\end{aligned}$$

Proof

$\text{col}(A - \tilde{U}\tilde{U}^T A)$ is orthogonal to $\text{col}(\tilde{A}) = \text{col}(\tilde{U})$

$$\begin{aligned}\|A - \tilde{A}\|_F^2 &= \|A - \tilde{U}\tilde{U}^T A\|_F^2 + \|\tilde{A} - \tilde{U}\tilde{U}^T A\|_F^2 \\ &\geq \|A - \tilde{U}\tilde{U}^T A\|_F^2 \\ &= \|A\|_F^2 - \|\tilde{U}\tilde{U}^T A\|_F^2 \\ &\geq \|A\|_F^2 - \|U_{:,1:k} U_{:,1:k}^T A\|_F^2\end{aligned}$$

Proof

$\text{col}(A - \tilde{U}\tilde{U}^T A)$ is orthogonal to $\text{col}(\tilde{A}) = \text{col}(\tilde{U})$

$$\begin{aligned}\|A - \tilde{A}\|_F^2 &= \|A - \tilde{U}\tilde{U}^T A\|_F^2 + \|\tilde{A} - \tilde{U}\tilde{U}^T A\|_F^2 \\ &\geq \|A - \tilde{U}\tilde{U}^T A\|_F^2 \\ &= \|A\|_F^2 - \|\tilde{U}\tilde{U}^T A\|_F^2 \\ &\geq \|A\|_F^2 - \|U_{:,1:k} U_{:,1:k}^T A\|_F^2 \\ &= \|A - U_{:,1:k} U_{:,1:k}^T A\|_F^2\end{aligned}$$

Reminder

For any pair of $m \times n$ matrices A and B

$$\operatorname{tr} \left(B^T A \right) := \operatorname{tr} \left(A B^T \right)$$

Frobenius norm

For any matrix $A \in \mathbb{R}^{m \times n}$, with singular values $\sigma_1, \dots, \sigma_{\min\{m,n\}}$

$$\|A\|_F = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2}.$$

Frobenius norm

For any matrix $A \in \mathbb{R}^{m \times n}$, with singular values $\sigma_1, \dots, \sigma_{\min\{m,n\}}$

$$\|A\|_F = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2}.$$

Proof:

$$\|A\|_F^2 = \text{tr}(A^T A)$$

Frobenius norm

For any matrix $A \in \mathbb{R}^{m \times n}$, with singular values $\sigma_1, \dots, \sigma_{\min\{m,n\}}$

$$\|A\|_F = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2}.$$

Proof:

$$\begin{aligned}\|A\|_F^2 &= \text{tr}(A^T A) \\ &= \text{tr}(V S U^T U S V^T)\end{aligned}$$

Frobenius norm

For any matrix $A \in \mathbb{R}^{m \times n}$, with singular values $\sigma_1, \dots, \sigma_{\min\{m,n\}}$

$$\|A\|_F = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2}.$$

Proof:

$$\begin{aligned}\|A\|_F^2 &= \text{tr}(A^T A) \\ &= \text{tr}(V S U^T U S V^T) \\ &= \text{tr}(V S S V^T) \quad \text{because } U^T U = I\end{aligned}$$

Frobenius norm

For any matrix $A \in \mathbb{R}^{m \times n}$, with singular values $\sigma_1, \dots, \sigma_{\min\{m,n\}}$

$$\|A\|_F = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2}.$$

Proof:

$$\begin{aligned}\|A\|_F^2 &= \text{tr}(A^T A) \\ &= \text{tr}(V S U^T U S V^T) \\ &= \text{tr}(V S S V^T) \quad \text{because } U^T U = I \\ &= \text{tr}(V^T V S S)\end{aligned}$$

Frobenius norm

For any matrix $A \in \mathbb{R}^{m \times n}$, with singular values $\sigma_1, \dots, \sigma_{\min\{m,n\}}$

$$\|A\|_F = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2}.$$

Proof:

$$\begin{aligned}\|A\|_F^2 &= \text{tr}(A^T A) \\ &= \text{tr}(V S U^T U S V^T) \\ &= \text{tr}(V S S V^T) \quad \text{because } U^T U = I \\ &= \text{tr}(V^T V S S) \\ &= \text{tr}(S S) \quad \text{because } V^T V = I\end{aligned}$$

Operator norm

The operator norm of a linear map and the corresponding matrix $A \in \mathbb{R}^{m \times n}$ is

$$\begin{aligned}\|A\| &:= \max_{\{\|\vec{x}\|_2=1 \mid \vec{x} \in \mathbb{R}^n\}} \|A\vec{x}\|_2 \\ &= \sigma_1\end{aligned}$$

Nuclear norm

The nuclear norm of a matrix $A \in \mathbb{R}^{m \times n}$ with singular values $\sigma_1, \dots, \sigma_{\min\{m,n\}}$ is

$$\|A\|_* := \sum_{i=1}^{\min\{m,n\}} \sigma_i$$

Multiplication by orthogonal matrices

For any orthogonal matrices $\tilde{U} \in \mathbb{R}^{m \times m}$ and $\tilde{V} \in \mathbb{R}^{n \times n}$ the singular values of $\tilde{U}A$ and $A\tilde{V}$ are the same as those of A

Consequence:

The operator, Frobenius and nuclear norm of $\tilde{U}A$ and $A\tilde{V}$ are **the same** as those of A

Hölder's inequality for matrices

For any matrix $A \in \mathbb{R}^{m \times n}$,

$$\|A\|_* = \sup_{\{\|B\| \leq 1 \mid B \in \mathbb{R}^{m \times n}\}} \langle A, B \rangle.$$

Consequence: nuclear norm satisfies **triangle inequality**

$$\|A + B\|_*$$

Hölder's inequality for matrices

For any matrix $A \in \mathbb{R}^{m \times n}$,

$$\|A\|_* = \sup_{\{\|B\| \leq 1 \mid B \in \mathbb{R}^{m \times n}\}} \langle A, B \rangle.$$

Consequence: nuclear norm satisfies **triangle inequality**

$$\|A + B\|_* = \sup_{\{\|C\| \leq 1 \mid C \in \mathbb{R}^{m \times n}\}} \langle A + B, C \rangle$$

Hölder's inequality for matrices

For any matrix $A \in \mathbb{R}^{m \times n}$,

$$\|A\|_* = \sup_{\{\|B\| \leq 1 \mid B \in \mathbb{R}^{m \times n}\}} \langle A, B \rangle.$$

Consequence: nuclear norm satisfies **triangle inequality**

$$\begin{aligned} \|A + B\|_* &= \sup_{\{\|C\| \leq 1 \mid C \in \mathbb{R}^{m \times n}\}} \langle A + B, C \rangle \\ &\leq \sup_{\{\|C\| \leq 1 \mid C \in \mathbb{R}^{m \times n}\}} \langle A, C \rangle + \sup_{\{\|D\| \leq 1 \mid D \in \mathbb{R}^{m \times n}\}} \langle B, D \rangle \end{aligned}$$

Hölder's inequality for matrices

For any matrix $A \in \mathbb{R}^{m \times n}$,

$$\|A\|_* = \sup_{\{\|B\| \leq 1 \mid B \in \mathbb{R}^{m \times n}\}} \langle A, B \rangle.$$

Consequence: nuclear norm satisfies **triangle inequality**

$$\begin{aligned} \|A + B\|_* &= \sup_{\{\|C\| \leq 1 \mid C \in \mathbb{R}^{m \times n}\}} \langle A + B, C \rangle \\ &\leq \sup_{\{\|C\| \leq 1 \mid C \in \mathbb{R}^{m \times n}\}} \langle A, C \rangle + \sup_{\{\|D\| \leq 1 \mid D \in \mathbb{R}^{m \times n}\}} \langle B, D \rangle \\ &= \|A\|_* + \|B\|_* \end{aligned}$$

Proof

The proof relies on the following lemma:

For any $Q \in \mathbb{R}^{n \times n}$

$$\max_{1 \leq i \leq n} |Q_{ii}| \leq \|Q\|$$

Proof:

Since $\|\vec{e}_i\|_2 = 1$,

$$\begin{aligned} \max_{1 \leq i \leq n} |Q_{ii}| &\leq \max_{1 \leq i \leq n} \sqrt{\sum_{j=1}^n Q_{ji}^2} \\ &= \max_{1 \leq i \leq n} \|Q \vec{e}_i\|_2 \\ &\leq \|Q\| \end{aligned}$$

Proof

Let $A := USV^T$,

$$\sup_{\|B\| \leq 1} \text{Trace}(A^T B)$$

Proof

Let $A := USV^T$,

$$\sup_{\|B\| \leq 1} \text{Trace}(A^T B) = \sup_{\|B\| \leq 1} \text{Trace}(V S U^T B)$$

Proof

Let $A := USV^T$,

$$\begin{aligned}\sup_{\|B\| \leq 1} \text{Trace}(A^T B) &= \sup_{\|B\| \leq 1} \text{Trace}(V S U^T B) \\ &= \sup_{\|B\| \leq 1} \text{Trace}(S B U^T V)\end{aligned}$$

Proof

Let $A := USV^T$,

$$\begin{aligned}\sup_{\|B\| \leq 1} \text{Trace}(A^T B) &= \sup_{\|B\| \leq 1} \text{Trace}(V S U^T B) \\ &= \sup_{\|B\| \leq 1} \text{Trace}(S B U^T V) \\ &\leq \sup_{\|M\| \leq 1} \text{Trace}(S M) \quad \text{since } \|B\| = \|B U^T V\| \\ &\leq \sup_{\max_{1 \leq i \leq n} |M_{ii}| \leq 1} \text{Trace}(S M) \quad \text{by the lemma}\end{aligned}$$

Proof

Let $A := USV^T$,

$$\begin{aligned} \sup_{\|B\| \leq 1} \text{Trace}(A^T B) &= \sup_{\|B\| \leq 1} \text{Trace}(V S U^T B) \\ &= \sup_{\|B\| \leq 1} \text{Trace}(S B U^T V) \\ &\leq \sup_{\|M\| \leq 1} \text{Trace}(S M) \quad \text{since } \|B\| = \|B U^T V\| \\ &\leq \sup_{\max_{1 \leq i \leq n} |M_{ii}| \leq 1} \text{Trace}(S M) \quad \text{by the lemma} \\ &\leq \sup_{\max_{1 \leq i \leq n} |M_{ii}| \leq 1} \sum_{i=1}^n M_{ii} \sigma_i \end{aligned}$$

Proof

Let $A := USV^T$,

$$\begin{aligned}\sup_{\|B\| \leq 1} \text{Trace}(A^T B) &= \sup_{\|B\| \leq 1} \text{Trace}(V S U^T B) \\ &= \sup_{\|B\| \leq 1} \text{Trace}(S B U^T V) \\ &\leq \sup_{\|M\| \leq 1} \text{Trace}(S M) \quad \text{since } \|B\| = \|B U^T V\| \\ &\leq \sup_{\max_{1 \leq i \leq n} |M_{ii}| \leq 1} \text{Trace}(S M) \quad \text{by the lemma} \\ &\leq \sup_{\max_{1 \leq i \leq n} |M_{ii}| \leq 1} \sum_{i=1}^n M_{ii} \sigma_i \\ &\leq \sum_{i=1}^n \sigma_i\end{aligned}$$

Proof

Let $A := USV^T$,

$$\begin{aligned} \sup_{\|B\| \leq 1} \text{Trace}(A^T B) &= \sup_{\|B\| \leq 1} \text{Trace}(V S U^T B) \\ &= \sup_{\|B\| \leq 1} \text{Trace}(S B U^T V) \\ &\leq \sup_{\|M\| \leq 1} \text{Trace}(S M) \quad \text{since } \|B\| = \|B U^T V\| \\ &\leq \sup_{\max_{1 \leq i \leq n} |M_{ii}| \leq 1} \text{Trace}(S M) \quad \text{by the lemma} \\ &\leq \sup_{\max_{1 \leq i \leq n} |M_{ii}| \leq 1} \sum_{i=1}^n M_{ii} \sigma_i \\ &\leq \sum_{i=1}^n \sigma_i \\ &= \|A\|_* \end{aligned}$$

Proof

To complete the proof, we show that the equality holds

UV^T has operator norm equal to one

$$\langle A, UV^T \rangle$$

Proof

To complete the proof, we show that the equality holds

UV^T has operator norm equal to one

$$\langle A, UV^T \rangle = \text{tr} (A^T UV^T)$$

Proof

To complete the proof, we show that the equality holds

UV^T has operator norm equal to one

$$\begin{aligned}\langle A, UV^T \rangle &= \text{tr} \left(A^T UV^T \right) \\ &= \text{tr} \left(V S U^T UV^T \right)\end{aligned}$$

Proof

To complete the proof, we show that the equality holds

UV^T has operator norm equal to one

$$\begin{aligned}\langle A, UV^T \rangle &= \text{tr} \left(A^T UV^T \right) \\ &= \text{tr} \left(V S U^T UV^T \right) \\ &= \text{tr} \left(V^T V S \right)\end{aligned}$$

Proof

To complete the proof, we show that the equality holds

UV^T has operator norm equal to one

$$\begin{aligned}\langle A, UV^T \rangle &= \text{tr} (A^T UV^T) \\ &= \text{tr} (V S U^T UV^T) \\ &= \text{tr} (V^T V S) \\ &= \text{tr} (S)\end{aligned}$$

Proof

To complete the proof, we show that the equality holds

UV^T has operator norm equal to one

$$\begin{aligned}\langle A, UV^T \rangle &= \text{tr} (A^T UV^T) \\ &= \text{tr} (V S U^T UV^T) \\ &= \text{tr} (V^T V S) \\ &= \text{tr} (S) \\ &= \|A\|_*\end{aligned}$$

Basic properties

Singular value decomposition

Denoising

Collaborative filtering

Principal component analysis

Probabilistic interpretation

Dimensionality reduction

Eigendecomposition

Denosing correlated signals

Aim: Estimating n m -dimensional signals $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbb{R}^m$ from

$$\vec{y}_i = \vec{x}_i + \vec{z}_i, \quad 1 \leq i \leq n$$

Assumption 1: Signals are similar and approximately span an **unknown** low-dimensional subspace

Assumption 2: Noisy perturbations are independent / uncorrelated

Denoising correlated signals

Aim: Estimating n m -dimensional signals $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbb{R}^m$ from

$$\vec{y}_i = \vec{x}_i + \vec{z}_i, \quad 1 \leq i \leq n$$

Assumption 1: Signals are similar and approximately span an **unknown** low-dimensional subspace

Assumption 2: Noisy perturbations are independent / uncorrelated

Denoising correlated signals

By the assumptions

$$X := [\vec{x}_1 \quad \vec{x}_2 \quad \cdots \quad \vec{x}_n]$$

is approximately low rank, whereas

$$Z := [\vec{z}_1 \quad \vec{z}_2 \quad \cdots \quad \vec{z}_n]$$

is full rank

If Z is not too large, low-rank approximation to

$$\begin{aligned} Y &:= [\vec{y}_1 \quad \vec{y}_2 \quad \cdots \quad \vec{y}_n] \\ &= X + Z \end{aligned}$$

should correspond mostly to X

Denoising via SVD truncation

1. Stack the vectors as the columns of a matrix $Y \in \mathbb{R}^{m \times n}$
2. Compute the SVD of $Y = USV^T$
3. Truncate the SVD to produce the low-rank estimate L

$$L := U_{:,1:k} S_{1:k,1:k} V_{:,1:k}^T,$$

for a fixed value of $k \leq \min \{m, n\}$

Important decision

What rank k to choose?

- ▶ *Large k*
- ▶ *Small k*

Important decision

What rank k to choose?

- ▶ *Large k* will approximate signals well but not suppress noise
- ▶ *Small k*

Important decision

What rank k to choose?

- ▶ *Large* k will approximate signals well but not suppress noise
- ▶ *Small* k will suppress noise but may not approximate signals well

Denosing of digit images

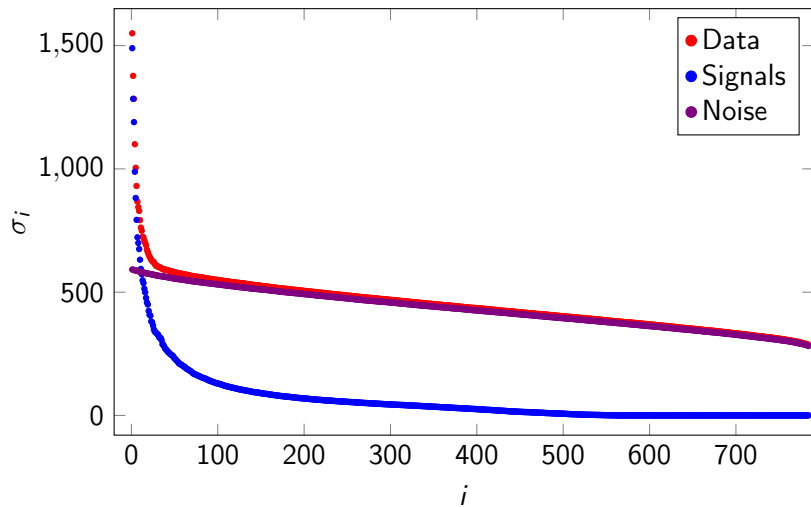
MNIST data

Signals: 6131 28×28 images of the number 3

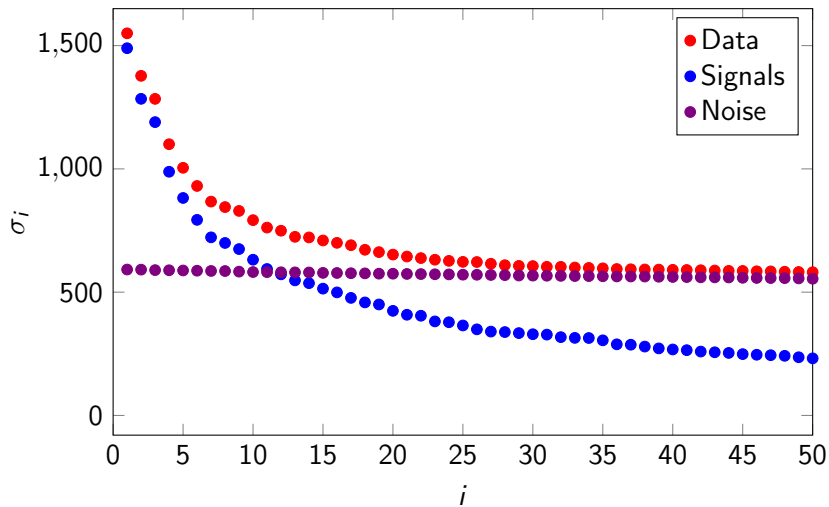
Noise: iid Gaussian so that SNR is 0.5 in ℓ_2 norm

More noise than signal!

Denoising of digit images

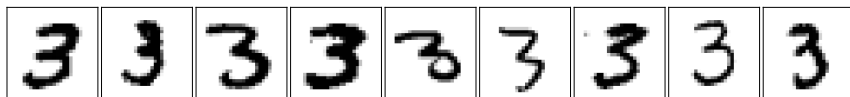


Denoising of digit images



$k = 40$

Signals



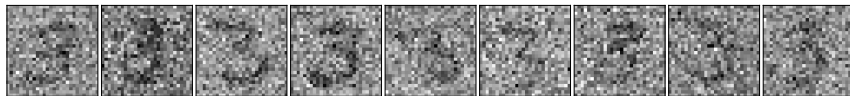
Rank
40
approx



Estimate

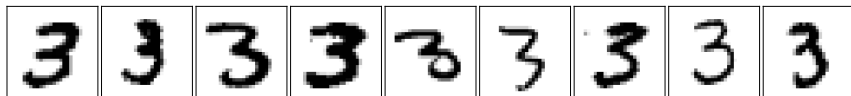


Data

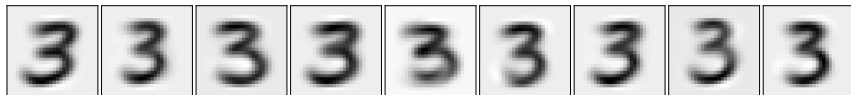


$k = 10$

Signals



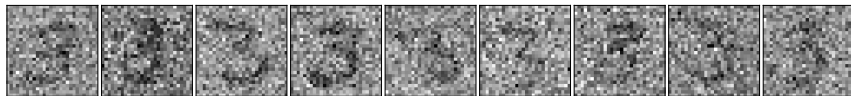
Rank
40
approx



Estimate



Data



Basic properties

Singular value decomposition

Denoising

Collaborative filtering

Principal component analysis

Probabilistic interpretation

Dimensionality reduction

Eigendecomposition

Collaborative filtering

$$A := \begin{pmatrix} 1 & 1 & 5 & 4 \\ 2 & 1 & 4 & 5 \\ 4 & 5 & 2 & 1 \\ 5 & 4 & 2 & 1 \\ 4 & 5 & 1 & 2 \\ 1 & 2 & 5 & 5 \end{pmatrix} \begin{array}{l} \text{The Dark Knight} \\ \text{Spiderman 3} \\ \text{Love Actually} \\ \text{Bridget Jones's Diary} \\ \text{Pretty Woman} \\ \text{Superman 2} \end{array}$$

	Bob	Molly	Mary	Larry	
	1	1	5	4	The Dark Knight
	2	1	4	5	Spiderman 3
	4	5	2	1	Love Actually
	5	4	2	1	Bridget Jones's Diary
	4	5	1	2	Pretty Woman
	1	2	5	5	Superman 2

Intuition

Some people have similar tastes and hence produce similar ratings

Some movies are similar and hence elicit similar reactions

This tends to induce **low-rank** structure in the matrix of ratings

SVD

$$A - \mu \vec{1} \vec{1}^T = U \Sigma V^T = U \begin{bmatrix} 7.79 & 0 & 0 & 0 \\ 0 & 1.62 & 0 & 0 \\ 0 & 0 & 1.55 & 0 \\ 0 & 0 & 0 & 0.62 \end{bmatrix} V^T$$

$$\mu := \frac{1}{n} \sum_{i=1}^m \sum_{j=1}^n A_{ij}$$

Rank 1 model

$$\bar{A} + \sigma_1 \vec{u}_1 \vec{v}_1^T = \begin{pmatrix} \text{Bob} & \text{Molly} & \text{Mary} & \text{Larry} \\ 1.34 (1) & 1.19 (1) & 4.66 (5) & 4.81 (4) \\ 1.55 (2) & 1.42 (1) & 4.45 (4) & 4.58 (5) \\ 4.45 (4) & 4.58 (5) & 1.55 (2) & 1.42 (1) \\ 4.43 (5) & 4.56 (4) & 1.57 (2) & 1.44 (1) \\ 4.43 (4) & 4.56 (5) & 1.57 (1) & 1.44 (2) \\ 1.34 (1) & 1.19 (2) & 4.66 (5) & 4.81 (5) \end{pmatrix} \begin{matrix} \text{The Dark Knight} \\ \text{Spiderman 3} \\ \text{Love Actually} \\ \text{B.J.'s Diary} \\ \text{Pretty Woman} \\ \text{Superman 2} \end{matrix}$$

First left singular vector

$$\vec{u}_1 = \begin{pmatrix} \text{D. Knight} & \text{Sp. 3} & \text{Love Act.} & \text{B.J.'s Diary} & \text{P. Woman} & \text{Sup. 2} \\ -0.45 & -0.39 & 0.39 & 0.39 & 0.39 & -0.45 \end{pmatrix}$$

Coefficients cluster movies into action (+) and romantic (-)

First right singular vector

$$\vec{v}_1 = \begin{matrix} & \text{Bob} & \text{Molly} & \text{Mary} & \text{Larry} \\ & (0.48 & 0.52 & -0.48 & -0.52) \end{matrix}$$

Coefficients cluster people into action (-) and romantic (+)

Generalization

Each rating is a sum of k terms

$$\text{rating}(\text{movie } i, \text{user } j) = \sum_{l=1}^k \sigma_l \vec{u}_l [i] \vec{v}_l [j]$$

Singular vectors cluster users and movies in different ways

Singular values weight the importance of the different factors.

Basic properties

Singular value decomposition

Denoising

Collaborative filtering

Principal component analysis

Probabilistic interpretation

Dimensionality reduction

Eigendecomposition

Sample covariance matrix

The sample covariance matrix of $\{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\} \in \mathbb{R}^m$

$$\Sigma(\vec{x}_1, \dots, \vec{x}_n) := \frac{1}{n-1} \sum_{i=1}^n (\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n)) (\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n))^T$$

$$\text{av}(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n) := \frac{1}{n} \sum_{i=1}^n \vec{x}_i$$

$$\Sigma(\vec{x}_1, \dots, \vec{x}_n)_{ij} = \begin{cases} \text{var}(\vec{x}_1[i], \dots, \vec{x}_n[i]) & \text{if } i = j, \\ \text{cov}((\vec{x}_1[i], \vec{x}_1[j]), \dots, (\vec{x}_n[i], \vec{x}_n[j])) & \text{if } i \neq j \end{cases}$$

Variation in a certain direction

For a unit vector $\vec{d} \in \mathbb{R}^m$

$$\text{var} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right)$$

Variation in a certain direction

For a unit vector $\vec{d} \in \mathbb{R}^m$

$$\begin{aligned} & \text{var} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right) \\ &= \frac{1}{n-1} \sum_{i=1}^n \left(\vec{d}^T \vec{x}_i - \text{av} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right) \right)^2 \end{aligned}$$

Variation in a certain direction

For a unit vector $\vec{d} \in \mathbb{R}^m$

$$\begin{aligned} & \text{var} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right) \\ &= \frac{1}{n-1} \sum_{i=1}^n \left(\vec{d}^T \vec{x}_i - \text{av} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right) \right)^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n \left(\vec{d}^T \left(\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n) \right) \right)^2 \end{aligned}$$

Variation in a certain direction

For a unit vector $\vec{d} \in \mathbb{R}^m$

$$\begin{aligned} & \text{var} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right) \\ &= \frac{1}{n-1} \sum_{i=1}^n \left(\vec{d}^T \vec{x}_i - \text{av} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right) \right)^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n \left(\vec{d}^T \left(\vec{x}_i - \text{av} \left(\vec{x}_1, \dots, \vec{x}_n \right) \right) \right)^2 \\ &= \vec{d}^T \left(\frac{1}{n-1} \sum_{i=1}^n \left(\vec{x}_i - \text{av} \left(\vec{x}_1, \dots, \vec{x}_n \right) \right) \left(\vec{x}_i - \text{av} \left(\vec{x}_1, \dots, \vec{x}_n \right) \right)^T \right) \vec{d} \end{aligned}$$

Variation in a certain direction

For a unit vector $\vec{d} \in \mathbb{R}^m$

$$\begin{aligned} & \text{var} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right) \\ &= \frac{1}{n-1} \sum_{i=1}^n \left(\vec{d}^T \vec{x}_i - \text{av} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right) \right)^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n \left(\vec{d}^T \left(\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n) \right) \right)^2 \\ &= \vec{d}^T \left(\frac{1}{n-1} \sum_{i=1}^n \left(\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n) \right) \left(\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n) \right)^T \right) \vec{d} \\ &= \vec{d}^T \Sigma(\vec{x}_1, \dots, \vec{x}_n) \vec{d} \end{aligned}$$

Covariance matrix captures variance **in every direction!**

Principal component analysis

Given n data vectors $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbb{R}^d$,

1. Center the data,

$$\vec{c}_i = \vec{x}_i - \text{av}(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n), \quad 1 \leq i \leq n$$

2. Group the centered data as columns of a matrix

$$C = [\vec{c}_1 \quad \vec{c}_2 \quad \cdots \quad \vec{c}_n].$$

3. Compute the SVD of C

The left singular vectors are the **principal directions**

The **principal values** are the coefficients of the centered vectors in the basis of principal directions.

Directions of maximum variance

The principal directions satisfy

$$\vec{u}_1 = \arg \max_{\{\|\vec{d}\|_2=1 \mid \vec{d} \in \mathbb{R}^n\}} \text{var} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right)$$

$$\vec{u}_i = \arg \max_{\{\|\vec{d}\|_2=1 \mid \vec{d} \in \mathbb{R}^n, \vec{d} \perp \vec{u}_1, \dots, \vec{u}_{i-1}\}} \text{var} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right), \quad 2 \leq i \leq k$$

Directions of maximum variance

The associated singular values satisfy

$$\frac{\sigma_1}{\sqrt{n-1}} = \max_{\{\|\vec{d}\|_2=1 \mid \vec{d} \in \mathbb{R}^n\}} \text{std} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right)$$

$$\frac{\sigma_i}{\sqrt{n-1}} = \max_{\{\|\vec{d}\|_2=1 \mid \vec{d} \in \mathbb{R}^n, \vec{d} \perp \vec{u}_1, \dots, \vec{u}_{i-1}\}} \text{std} \left(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n \right), \quad 2 \leq i \leq k$$

Proof

$$\begin{aligned}\Sigma(\vec{x}_1, \dots, \vec{x}_n) &= \frac{1}{n-1} \sum_{i=1}^n (\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n)) (\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n))^T \\ &= \frac{1}{n-1} CC^T\end{aligned}$$

Proof

$$\begin{aligned}\Sigma(\vec{x}_1, \dots, \vec{x}_n) &= \frac{1}{n-1} \sum_{i=1}^n (\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n)) (\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n))^T \\ &= \frac{1}{n-1} CC^T\end{aligned}$$

For any vector \vec{d}

$$\text{var}(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n) = \vec{d}^T \Sigma(\vec{x}_1, \dots, \vec{x}_n) \vec{d}$$

Proof

$$\begin{aligned}\Sigma(\vec{x}_1, \dots, \vec{x}_n) &= \frac{1}{n-1} \sum_{i=1}^n (\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n)) (\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n))^T \\ &= \frac{1}{n-1} CC^T\end{aligned}$$

For any vector \vec{d}

$$\begin{aligned}\text{var}(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n) &= \vec{d}^T \Sigma(\vec{x}_1, \dots, \vec{x}_n) \vec{d} \\ &= \frac{1}{n-1} \vec{d}^T CC^T \vec{d}\end{aligned}$$

Proof

$$\begin{aligned}\Sigma(\vec{x}_1, \dots, \vec{x}_n) &= \frac{1}{n-1} \sum_{i=1}^n (\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n)) (\vec{x}_i - \text{av}(\vec{x}_1, \dots, \vec{x}_n))^T \\ &= \frac{1}{n-1} CC^T\end{aligned}$$

For any vector \vec{d}

$$\begin{aligned}\text{var}(\vec{d}^T \vec{x}_1, \dots, \vec{d}^T \vec{x}_n) &= \vec{d}^T \Sigma(\vec{x}_1, \dots, \vec{x}_n) \vec{d} \\ &= \frac{1}{n-1} \vec{d}^T CC^T \vec{d} \\ &= \frac{1}{n-1} \left\| C^T \vec{d} \right\|_2^2\end{aligned}$$

Singular values

The singular values of A satisfy

$$\begin{aligned}\sigma_1 &= \max_{\{\|\vec{x}\|_2=1 \mid \vec{x} \in \mathbb{R}^n\}} \|A\vec{x}\|_2 \\ &= \max_{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m\}} \|A^T \vec{y}\|_2 \\ \sigma_i &= \max_{\{\|\vec{x}\|_2=1 \mid \vec{x} \in \mathbb{R}^n, \vec{x} \perp \vec{u}_1, \dots, \vec{u}_{i-1}\}} \|A\vec{x}\|_2 \\ &= \max_{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m, \vec{y} \perp \vec{v}_1, \dots, \vec{v}_{i-1}\}} \|A^T \vec{y}\|_2, \quad 2 \leq i \leq \min\{m, n\}\end{aligned}$$

Singular vectors

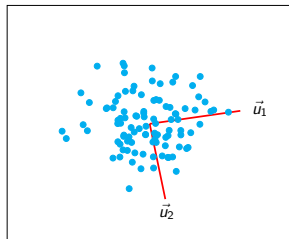
The left singular vectors of A satisfy

$$\vec{u}_1 = \arg \max_{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m\}} \left\| A^T \vec{y} \right\|_2$$

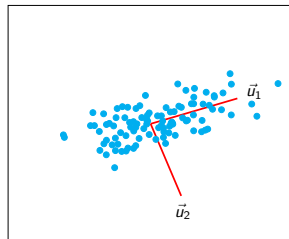
$$\vec{u}_i = \arg \max_{\{\|\vec{y}\|_2=1 \mid \vec{y} \in \mathbb{R}^m, \vec{y} \perp \vec{v}_1, \dots, \vec{v}_{i-1}\}} \left\| A^T \vec{y} \right\|_2, \quad 2 \leq i \leq n$$

PCA in 2D

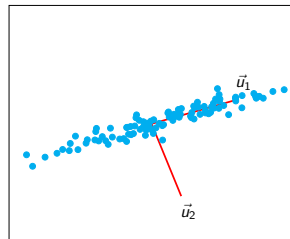
$$\begin{aligned}\sigma_1/\sqrt{n-1} &= 0.705, \\ \sigma_2/\sqrt{n-1} &= 0.690\end{aligned}$$



$$\begin{aligned}\sigma_1/\sqrt{n-1} &= 0.983, \\ \sigma_2/\sqrt{n-1} &= 0.356\end{aligned}$$



$$\begin{aligned}\sigma_1/\sqrt{n-1} &= 1.349, \\ \sigma_2/\sqrt{n-1} &= 0.144\end{aligned}$$



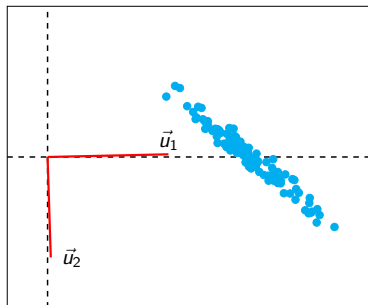
Centering

$$\sigma_1/\sqrt{n-1} = 5.077$$

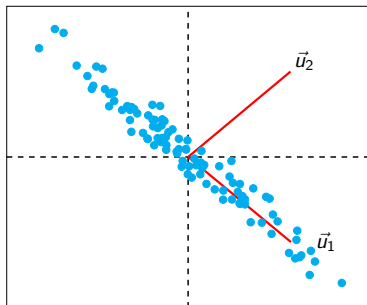
$$\sigma_2/\sqrt{n-1} = 0.889$$

$$\sigma_1/\sqrt{n-1} = 1.261$$

$$\sigma_2/\sqrt{n-1} = 0.139$$



Uncentered data



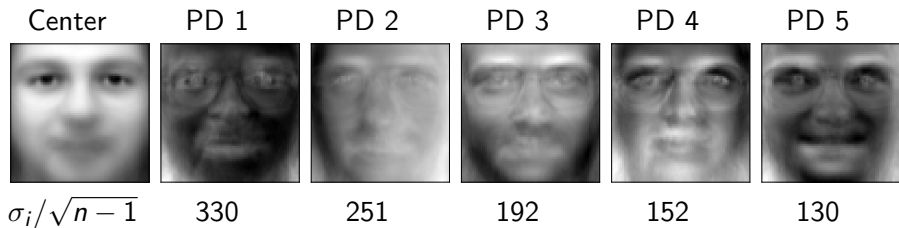
Centered data

PCA of faces

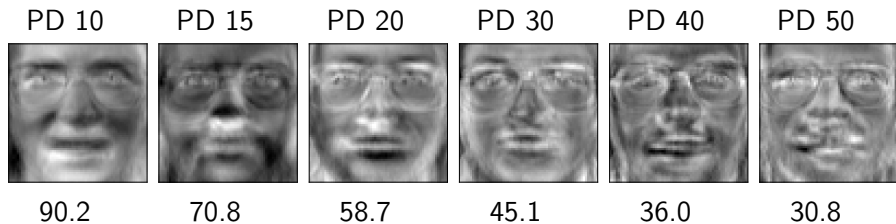
Data set of 400 64×64 images from 40 subjects (10 per subject)

Each face is vectorized and interpreted as a vector in \mathbb{R}^{4096}

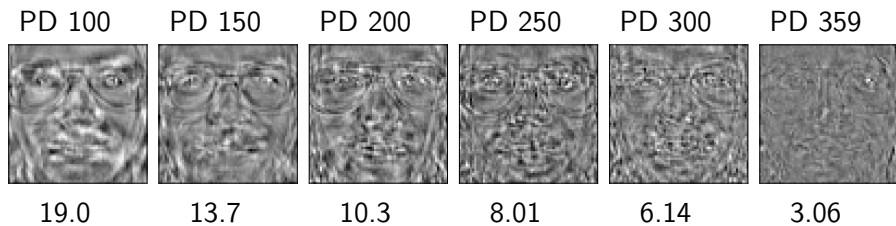
PCA of faces



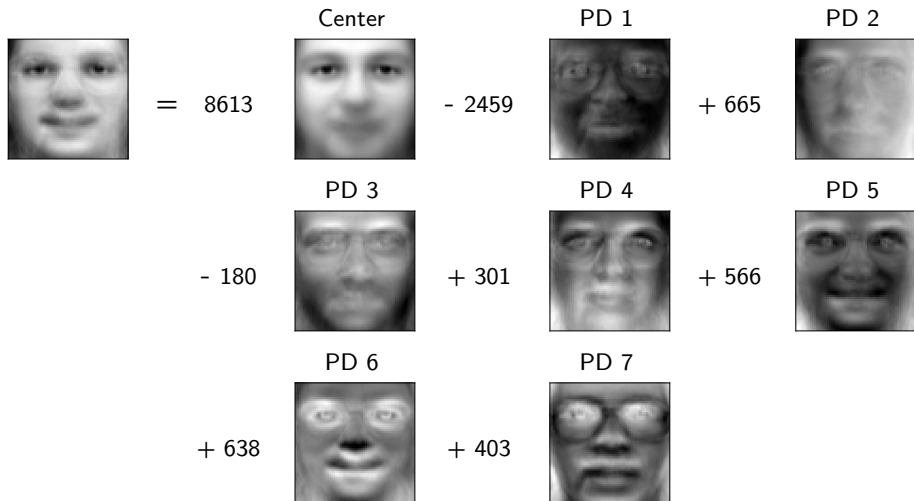
PCA of faces



PCA of faces



Projection onto first 7 principal directions



Projection onto first k principal directions

Signal



5 PDs



10 PDs



20 PDs



30 PDs



50 PDs



100 PDs



150 PDs



200 PDs



250 PDs



300 PDs



359 PDs



Basic properties

Singular value decomposition

Denoising

Collaborative filtering

Principal component analysis

Probabilistic interpretation

Dimensionality reduction

Eigendecomposition

Covariance matrix

The covariance matrix of a random vector \vec{x} is defined as

$$\begin{aligned}\Sigma_{\vec{x}} &:= \begin{bmatrix} \text{Var}(\vec{x}[1]) & \text{Cov}(\vec{x}[1], \vec{x}[2]) & \cdots & \text{Cov}(\vec{x}[1], \vec{x}[n]) \\ \text{Cov}(\vec{x}[2], \vec{x}[1]) & \text{Var}(\vec{x}[2]) & \cdots & \text{Cov}(\vec{x}[2], \vec{x}[n]) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(\vec{x}[n], \vec{x}[1]) & \text{Cov}(\vec{x}[n], \vec{x}[2]) & \cdots & \text{Var}(\vec{x}[n]) \end{bmatrix} \\ &= \text{E}(\vec{x}\vec{x}^T) - \text{E}(\vec{x})\text{E}(\vec{x})^T\end{aligned}$$

If the covariance matrix is **diagonal**, the entries are uncorrelated

Covariance matrix after a linear transformation

Let $\Sigma \in \mathbb{R}^{n \times n}$ be the covariance matrix of \vec{x} . For any matrix $A \in \mathbb{R}^{m \times n}$

$$\Sigma_{A\vec{x}+\vec{b}} = A\Sigma_{\vec{x}}A^T$$

Proof:

$$\Sigma_{A\vec{x}}$$

Covariance matrix after a linear transformation

Let $\Sigma \in \mathbb{R}^{n \times n}$ be the covariance matrix of \vec{x} . For any matrix $A \in \mathbb{R}^{m \times n}$

$$\Sigma_{A\vec{x}+\vec{b}} = A\Sigma\vec{x}A^T$$

Proof:

$$\Sigma_{A\vec{x}} = E\left((A\vec{x})(A\vec{x})^T\right) - E(A\vec{x})E(A\vec{x})^T$$

Covariance matrix after a linear transformation

Let $\Sigma \in \mathbb{R}^{n \times n}$ be the covariance matrix of \vec{x} . For any matrix $A \in \mathbb{R}^{m \times n}$

$$\Sigma_{A\vec{x}+\vec{b}} = A\Sigma_{\vec{x}}A^T$$

Proof:

$$\begin{aligned}\Sigma_{A\vec{x}} &= E\left((A\vec{x})(A\vec{x})^T\right) - E(A\vec{x})E(A\vec{x})^T \\ &= A\left(E(\vec{x}\vec{x}^T) - E(\vec{x})E(\vec{x})^T\right)A^T\end{aligned}$$

Covariance matrix after a linear transformation

Let $\Sigma \in \mathbb{R}^{n \times n}$ be the covariance matrix of \vec{x} . For any matrix $A \in \mathbb{R}^{m \times n}$

$$\Sigma_{A\vec{x}+\vec{b}} = A\Sigma_{\vec{x}}A^T$$

Proof:

$$\begin{aligned}\Sigma_{A\vec{x}} &= E\left((A\vec{x})(A\vec{x})^T\right) - E(A\vec{x})E(A\vec{x})^T \\ &= A\left(E(\vec{x}\vec{x}^T) - E(\vec{x})E(\vec{x})^T\right)A^T \\ &= A\Sigma_{\vec{x}}A^T\end{aligned}$$

Variance in a fixed direction

The variance of \vec{x} in the direction of a unit-norm vector \vec{v} equals

$$\text{Var} \left(\vec{v}^T \vec{x} \right) = \vec{v}^T \Sigma_{\vec{x}} \vec{v}$$

SVD of covariance matrix

$$\begin{aligned}\Sigma_{\vec{x}} &= U\Lambda U^T \\ &= [\vec{u}_1 \quad \vec{u}_2 \quad \cdots \quad \vec{u}_n] \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ & & \cdots & \\ 0 & 0 & \cdots & \sigma_n \end{bmatrix} [\vec{u}_1 \quad \vec{u}_2 \quad \cdots \quad \vec{u}_n]^T\end{aligned}$$

Directions of maximum variance

The SVD of the covariance matrix $\Sigma_{\vec{x}}$ of a random vector \vec{x} satisfies

$$\sigma_1 = \max_{\|\vec{v}\|_2=1} \text{Var} \left(\vec{v}^T \vec{x} \right)$$

$$\vec{u}_1 = \arg \max_{\|\vec{v}\|_2=1} \text{Var} \left(\vec{v}^T \vec{x} \right)$$

$$\sigma_k = \max_{\|\vec{v}\|_2=1, \vec{v} \perp \vec{u}_1, \dots, \vec{u}_{k-1}} \text{Var} \left(\vec{v}^T \vec{x} \right)$$

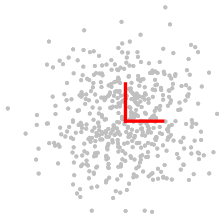
$$\vec{u}_k = \arg \max_{\|\vec{v}\|_2=1, \vec{v} \perp \vec{u}_1, \dots, \vec{u}_{k-1}} \text{Var} \left(\vec{v}^T \vec{x} \right)$$

Directions of maximum variance

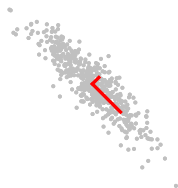
$$\sqrt{\sigma_1} = 1.22, \sqrt{\sigma_2} = 0.71$$



$$\sqrt{\sigma_1} = 1, \sqrt{\sigma_2} = 1$$

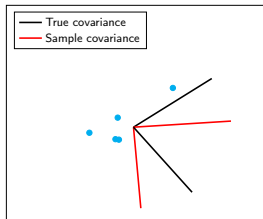


$$\sqrt{\sigma_1} = 1.38, \sqrt{\sigma_2} = 0.32$$

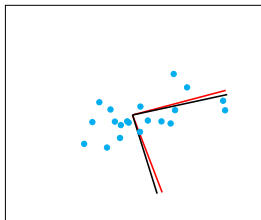


Probabilistic interpretation of PCA

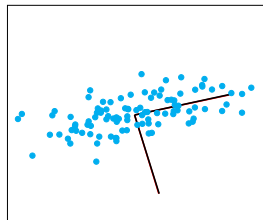
$n = 5$



$n = 20$



$n = 100$



Basic properties

Singular value decomposition

Denoising

Collaborative filtering

Principal component analysis

Probabilistic interpretation

Dimensionality reduction

Eigendecomposition

Dimensionality reduction

Data with a large number of features can be difficult to analyze or process

Dimensionality reduction is a useful preprocessing step

If data modeled are vectors in \mathbb{R}^m we can reduce the dimension by **projecting** onto \mathbb{R}^k , where $k < m$

For **orthogonal** projections, the new representation is $\langle \vec{b}_1, \vec{x} \rangle, \langle \vec{b}_2, \vec{x} \rangle, \dots, \langle \vec{b}_k, \vec{x} \rangle$ for a basis $\vec{b}_1, \dots, \vec{b}_k$ of the subspace that we project on

Problem: How do we choose the subspace?

Optimal subspace for orthogonal projection

Given a set of vectors $\vec{a}_1, \vec{a}_2, \dots, \vec{a}_n$ and a fixed dimension $k \leq n$, the SVD of

$$A := [\vec{a}_1 \quad \vec{a}_2 \quad \cdots \quad \vec{a}_n] \in \mathbb{R}^{m \times n}$$

provides the k -dimensional subspace that captures the most energy

$$\sum_{i=1}^n \|\mathcal{P}_{\text{span}(\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k)} \vec{a}_i\|_2^2 \geq \sum_{i=1}^n \|\mathcal{P}_{\mathcal{S}} \vec{a}_i\|_2^2$$

for any subspace \mathcal{S} of dimension k

Nearest neighbors in principal-component space

Nearest neighbors classification (Algorithm 4.2 in Lecture Notes 1) computes n distances in \mathbb{R}^m for each new example

Cost: $\mathcal{O}(nmp)$ for p examples

Idea: Project onto first k main principal directions beforehand

Cost:

- ▶ $\mathcal{O}(m^2n)$, if $m < n$, to compute principal dimensions
- ▶ kmn operations to project training set
- ▶ kmp operations to project test set
- ▶ knp to perform nearest-neighbor classification

Faster if $p > m$

Face recognition

Training set: 360 64×64 images from 40 different subjects (9 each)

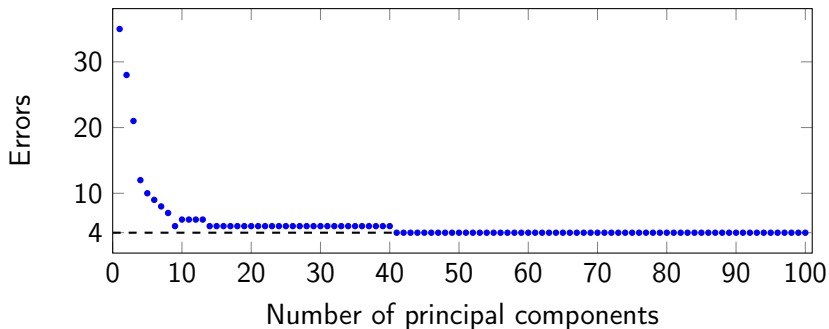
Test set: 1 new image from each subject

We model each image as a vector in \mathbb{R}^{4096} ($m = 4096$)

To classify we:

1. Project onto first k principal directions
2. Apply nearest-neighbor classification using the ℓ_2 -norm distance in \mathbb{R}^k

Performance



Nearest neighbor in \mathbb{R}^{41}

Test image



Projection



Closest projection



Corresponding image



Dimensionality reduction for visualization

Motivation: Visualize high-dimensional features projected onto 2D or 3D

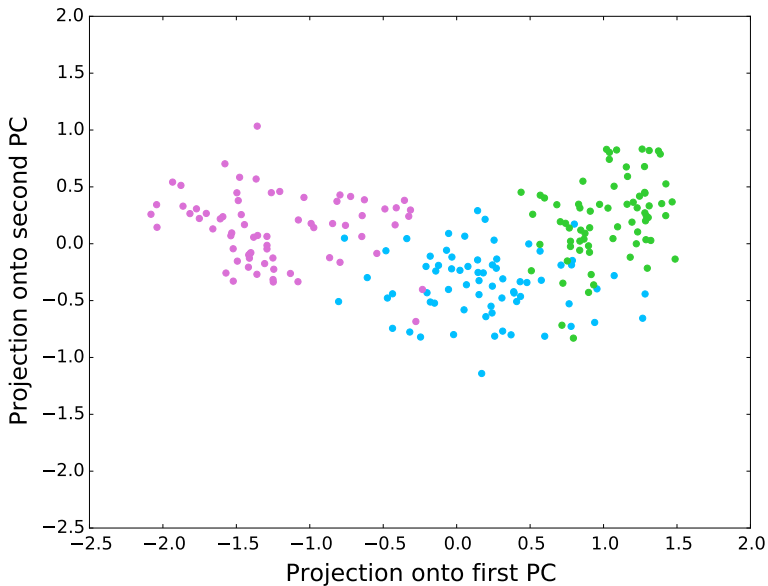
Example:

Seeds from three different varieties of wheat: Kama, Rosa and Canadian

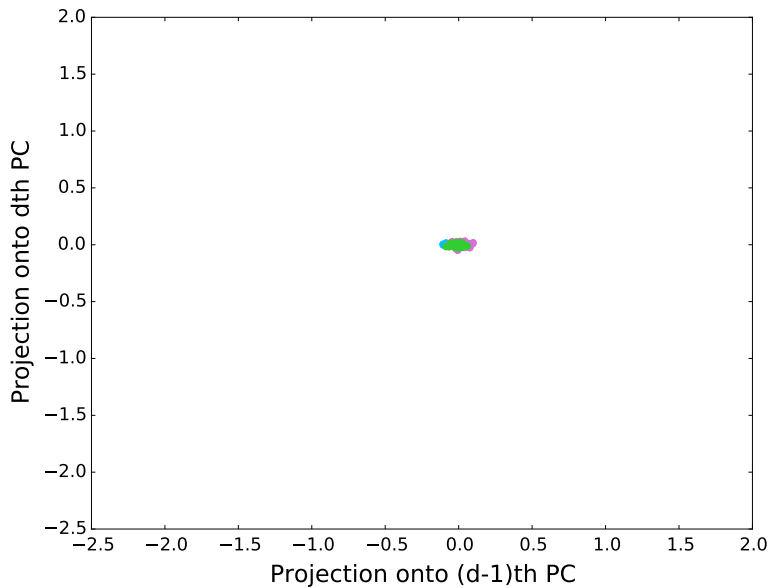
Features:

- ▶ Area
- ▶ Perimeter
- ▶ Compactness
- ▶ Length of kernel
- ▶ Width of kernel
- ▶ Asymmetry coefficient
- ▶ Length of kernel groove

Projection onto two first PCs



Projection onto two last PCs



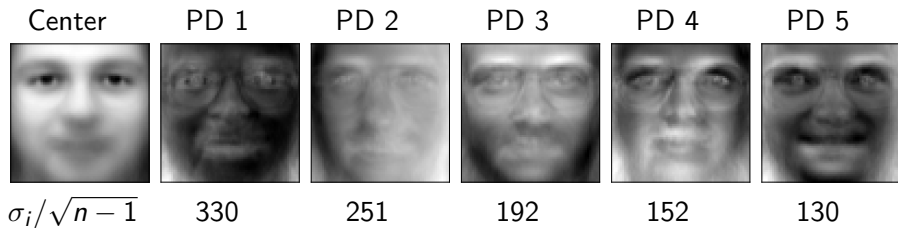
Whitening

Motivation: Dominating principal directions are not necessarily the most informative

Principal directions corresponding to small singular values may contain information that is *drowned* by main directions

Intuitively, **linear skew** obscures useful structure

PCA of faces



Whitening

Given $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbb{R}^d$ we:

1. Center the data,

$$\vec{c}_i = \vec{x}_i - \text{av}(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n), \quad 1 \leq i \leq n$$

2. Group the centered data as columns of a matrix

$$C = [\vec{c}_1 \quad \vec{c}_2 \quad \dots \quad \vec{c}_n]$$

3. Compute the SVD of $C = USV^T$

4. Whiten by applying the linear map $US^{-1}U^T$

$$\vec{w}_i := US^{-1}U^T \vec{c}_i$$

Covariance of whitened data

Matrix of whitened vectors

$$W = US^{-1}U^T C$$

The covariance matrix of the whitened data is

$$\Sigma(\vec{c}_1, \dots, \vec{c}_n)$$

Covariance of whitened data

Matrix of whitened vectors

$$W = US^{-1}U^T C$$

The covariance matrix of the whitened data is

$$\Sigma(\vec{c}_1, \dots, \vec{c}_n) = \frac{1}{n-1} WW^T$$

Covariance of whitened data

Matrix of whitened vectors

$$W = US^{-1}U^T C$$

The covariance matrix of the whitened data is

$$\begin{aligned}\Sigma(\vec{c}_1, \dots, \vec{c}_n) &= \frac{1}{n-1} WW^T \\ &= \frac{1}{n-1} US^{-1}U^T CC^T US^{-1}U^T\end{aligned}$$

Covariance of whitened data

Matrix of whitened vectors

$$W = US^{-1}U^T C$$

The covariance matrix of the whitened data is

$$\begin{aligned}\Sigma(\vec{c}_1, \dots, \vec{c}_n) &= \frac{1}{n-1} WW^T \\ &= \frac{1}{n-1} US^{-1}U^T CC^T US^{-1}U^T \\ &= \frac{1}{n-1} US^{-1}U^T USV^T VSU^T US^{-1}U^T\end{aligned}$$

Covariance of whitened data

Matrix of whitened vectors

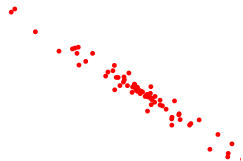
$$W = US^{-1}U^T C$$

The covariance matrix of the whitened data is

$$\begin{aligned}\Sigma(\vec{c}_1, \dots, \vec{c}_n) &= \frac{1}{n-1} WW^T \\ &= \frac{1}{n-1} US^{-1}U^T CC^T US^{-1}U^T \\ &= \frac{1}{n-1} US^{-1}U^T USV^T VSU^T US^{-1}U^T \\ &= \frac{1}{n-1} I\end{aligned}$$

Whitening

$\vec{x}_1, \dots, \vec{x}_n$



$U^T \vec{x}_1, \dots, U^T \vec{x}_n$



$S^{-1} U^T \vec{x}_1, \dots, S^{-1} U^T \vec{x}_n$



Whitened faces

\vec{x}



$US^{-1}U^T\vec{x}$



Basic properties

Singular value decomposition

Denoising

Collaborative filtering

Principal component analysis

Probabilistic interpretation

Dimensionality reduction

Eigendecomposition

Eigenvectors

An **eigenvector** \vec{q} of a square matrix $A \in \mathbb{R}^{n \times n}$ satisfies

$$A\vec{q} = \lambda\vec{q}$$

for a scalar λ which is the corresponding **eigenvalue**

Even if A is real, its eigenvectors and eigenvalues can be complex

Eigendecomposition

If a matrix has n linearly independent eigenvectors then it is **diagonalizable**

Let $\vec{q}_1, \dots, \vec{q}_n$ be lin. indep. eigenvectors of $A \in \mathbb{R}^{n \times n}$ with eigenvalues $\lambda_1, \dots, \lambda_n$

$$A = [\vec{q}_1 \quad \vec{q}_2 \quad \cdots \quad \vec{q}_n] \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ & & \cdots & \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix} [\vec{q}_1 \quad \vec{q}_2 \quad \cdots \quad \vec{q}_n]^{-1}$$
$$= Q \Lambda Q^{-1}$$

Proof

AQ

Proof

$$AQ = [A\vec{q}_1 \quad A\vec{q}_2 \quad \cdots \quad A\vec{q}_n]$$

Proof

$$\begin{aligned}AQ &= [A\vec{q}_1 \quad A\vec{q}_2 \quad \cdots \quad A\vec{q}_n] \\ &= [\lambda_1\vec{q}_1 \quad \lambda_2\vec{q}_2 \quad \cdots \quad \lambda_n\vec{q}_n]\end{aligned}$$

Proof

$$\begin{aligned}AQ &= [A\vec{q}_1 \quad A\vec{q}_2 \quad \cdots \quad A\vec{q}_n] \\ &= [\lambda_1\vec{q}_1 \quad \lambda_2\vec{q}_2 \quad \cdots \quad \lambda_n\vec{q}_n] \\ &= Q\Lambda\end{aligned}$$

Not all matrices have an eigendecomposition

$$A := \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$$

Assume A has an eigenvector \vec{q} associated to an eigenvalue λ

$$\begin{bmatrix} \vec{q}[2] \\ 0 \end{bmatrix}$$

Not all matrices have an eigendecomposition

$$A := \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$$

Assume A has an eigenvector \vec{q} associated to an eigenvalue λ

$$\begin{bmatrix} \vec{q}[2] \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \vec{q}[1] \\ \vec{q}[2] \end{bmatrix}$$

Not all matrices have an eigendecomposition

$$A := \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$$

Assume A has an eigenvector \vec{q} associated to an eigenvalue λ

$$\begin{aligned} \begin{bmatrix} \vec{q}[2] \\ 0 \end{bmatrix} &= \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \vec{q}[1] \\ \vec{q}[2] \end{bmatrix} \\ &= \begin{bmatrix} \lambda \vec{q}[1] \\ \lambda \vec{q}[2] \end{bmatrix} \end{aligned}$$

Spectral theorem for symmetric matrices

If $A \in \mathbb{R}^n$ is symmetric, then it has an eigendecomposition of the form

$$A = U\Lambda U^T$$

where the matrix of eigenvectors U is an orthogonal matrix

Eigendecomposition vs SVD

Symmetric matrices also have singular value decomposition

$$A = USV^T$$

Left singular vectors = eigenvectors

Singular values = magnitude of eigenvalues

$$\vec{v}_i = \vec{u}_i \text{ if } \lambda_i \geq 0$$

$$\vec{v}_i = -\vec{u}_i \text{ if } \lambda_i < 0$$

Application of eigendecomposition

Fast computation of

$$AA \cdots A\vec{x} = A^k \vec{x}$$

$$A^k$$

Application of eigendecomposition

Fast computation of

$$AA \cdots A\vec{x} = A^k \vec{x}$$

$$A^k = Q\Lambda Q^{-1}Q\Lambda Q^{-1} \cdots Q\Lambda Q^{-1}$$

Application of eigendecomposition

Fast computation of

$$AA \cdots A \vec{x} = A^k \vec{x}$$

$$\begin{aligned} A^k &= Q \Lambda Q^{-1} Q \Lambda Q^{-1} \cdots Q \Lambda Q^{-1} \\ &= Q \Lambda^k Q^{-1} \end{aligned}$$

Application of eigendecomposition

Fast computation of

$$AA \cdots A \vec{x} = A^k \vec{x}$$

$$\begin{aligned} A^k &= Q \Lambda Q^{-1} Q \Lambda Q^{-1} \cdots Q \Lambda Q^{-1} \\ &= Q \Lambda^k Q^{-1} \\ &= Q \begin{bmatrix} \lambda_1^k & 0 & \cdots & 0 \\ 0 & \lambda_2^k & \cdots & 0 \\ & & \cdots & \\ 0 & 0 & \cdots & \lambda_n^k \end{bmatrix} Q^{-1} \end{aligned}$$

Application of eigendecomposition

For any vector \vec{x}

$$A^k \vec{x}$$

Application of eigendecomposition

For any vector \vec{x}

$$A^k \vec{x} = \sum_{i=1}^n \alpha_i A^k \vec{q}_i$$

Application of eigendecomposition

For any vector \vec{x}

$$\begin{aligned} A^k \vec{x} &= \sum_{i=1}^n \alpha_i A^k \vec{q}_i \\ &= \sum_{i=1}^n \alpha_i \lambda_i^k \vec{q}_i \end{aligned}$$

Application of eigendecomposition

For any vector \vec{x}

$$\begin{aligned} A^k \vec{x} &= \sum_{i=1}^n \alpha_i A^k \vec{q}_i \\ &= \sum_{i=1}^n \alpha_i \lambda_i^k \vec{q}_i \end{aligned}$$

If $|\lambda_1| > |\lambda_2| \geq \dots$ then \vec{q}_1 will eventually dominate

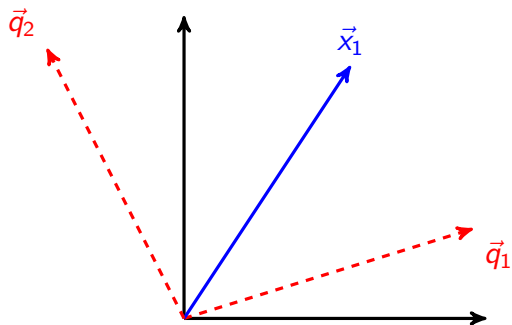
Power method

Set $\vec{x}_1 := \vec{x} / \|\vec{x}\|_2$, where \vec{x} is a randomly generated.

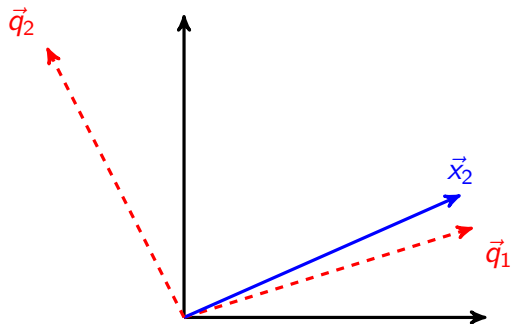
For $i = 1, 2, 3, \dots$, compute

$$\vec{x}_i := \frac{A\vec{x}_{i-1}}{\|A\vec{x}_{i-1}\|_2}$$

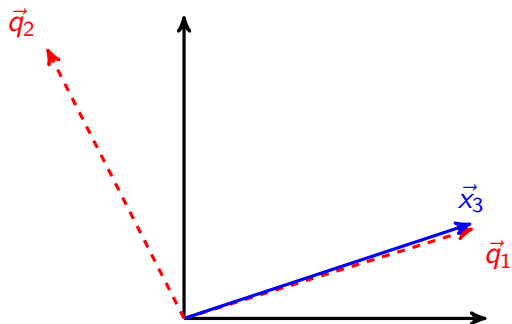
Power method



Power method



Power method



Deer and wolfs

Model for deer d_{n+1} and wolf w_{n+1} population in year $n + 1$

$$\begin{aligned}d_{n+1} &= \frac{5}{4}d_n - \frac{3}{4}w_n \\w_{n+1} &= \frac{1}{4}d_n + \frac{1}{4}w_n \quad n = 0, 1, 2, \dots\end{aligned}$$

Deer and wolfs

$$\begin{bmatrix} d_n \\ w_n \end{bmatrix}$$

Deer and wolfs

$$\begin{bmatrix} d_n \\ w_n \end{bmatrix} = Q \Lambda^n Q^{-1} \begin{bmatrix} d_0 \\ w_0 \end{bmatrix}$$

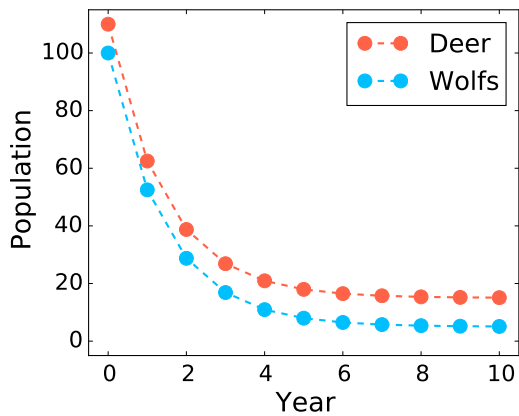
Deer and wolfs

$$\begin{aligned} \begin{bmatrix} d_n \\ w_n \end{bmatrix} &= Q \Lambda^n Q^{-1} \begin{bmatrix} d_0 \\ w_0 \end{bmatrix} \\ &= \begin{bmatrix} 3 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 0.5^n \end{bmatrix} \begin{bmatrix} 0.5 & -0.5 \\ -0.5 & 1.5 \end{bmatrix} \begin{bmatrix} d_0 \\ w_0 \end{bmatrix} \end{aligned}$$

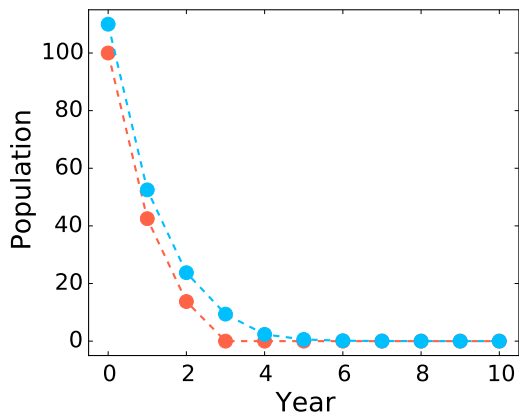
Deer and wolfs

$$\begin{aligned}\begin{bmatrix} d_n \\ w_n \end{bmatrix} &= Q\Lambda^n Q^{-1} \begin{bmatrix} d_0 \\ w_0 \end{bmatrix} \\ &= \begin{bmatrix} 3 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 0.5^n \end{bmatrix} \begin{bmatrix} 0.5 & -0.5 \\ -0.5 & 1.5 \end{bmatrix} \begin{bmatrix} d_0 \\ w_0 \end{bmatrix} \\ &= \frac{d_0 - w_0}{2} \begin{bmatrix} 3 \\ 1 \end{bmatrix} + \frac{3w_0 - d_0}{8^n} \begin{bmatrix} 1 \\ 1 \end{bmatrix}\end{aligned}$$

Deer and wolfs



Deer and wolfs



Deer and wolfs

