

CS 58500 – Theoretical Computer Science Toolkit

Lecture 13 (03/03) Spectral Methods

https://ruizhezhang.com/course_spring_2026.html



Today's Lecture

- Introduction to Spectral Method
- SVD and Best-Fit Subspace
- k -Means Clustering
- Complexity of SVD

Introduction to Spectral Methods: Motivation

Datasets are the foundation of progress in AI

For text:

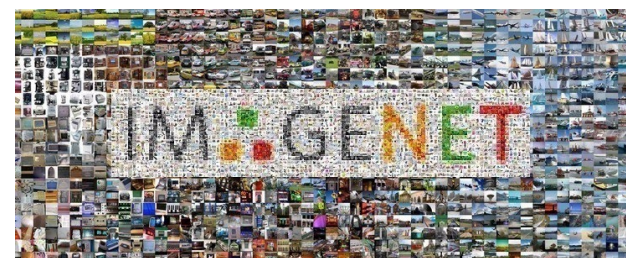
- GPT-1 (2018): **3 B** Tokens
- GPT-2 (2019): **30 B** Tokens
- GPT-3 (2020): **300 B** Tokens
- GPT-4 (2023): **13k B** Tokens
- GPT-5 (2025): **70k B** Tokens (?)



20000x growth in 7 years

For images:

- ImageNet (2009): **1 Million** images
- LAION-5B (2022): **5 Billion** Images



5000x growth in 5 years

Introduction to Spectral Methods: Motivation

Challenges in modern data science

- Enormous data
- Curse of dimensionality
- Imperfect data (noisy, messy, missing features, ...)

*“Spectral method refers to a collection of algorithms built upon the **eigenvectors** (resp. **singular vectors**) and **eigenvalues** (resp. **singular values**) of some **properly designed matrices** generated from data”*

(Spectral methods for data science: A statistical perspective, Chen et al. '21)

Introduction to Spectral Methods: Applications

Clustering

- Community detection in networks
- Joint image alignment
- Ranking

Learning Hidden Structures

- Synchronization in cryo-EM
- Tensor estimation

Numerical Linear Algebra

- Dimensionality reduction
- low-rank matrix estimation

Econometric and financial modeling

Spectral graph theory

- Random walk
- Electrical flows and effective resistance
- Expander graph
- Graph sparsification
- Graph partitioning

Sum-of-squares (SoS) / “Finding a needle in a haystack”

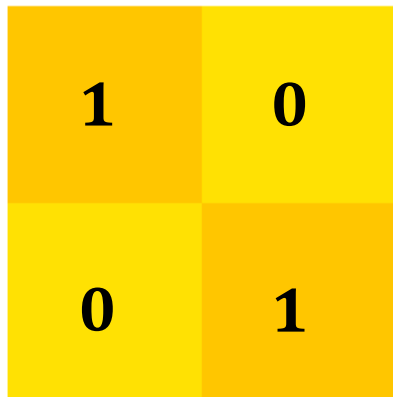
- Planted sparse vector
- Tensor PCA
- Refuting random CSP

Introduction to Spectral Methods: Examples

Clustering: grouping the elements based on their mutual similarities

- Suppose n people are divided into two groups, with the first $n/2$ people in group A, and the remaining $n/2$ people in group B
- We observe pairwise similarity measurements based on their group memberships:

$$A_{ij} = \begin{cases} 1 & \text{if } i, j \text{ are in the same group} \\ 0 & \text{otherwise} \end{cases}$$

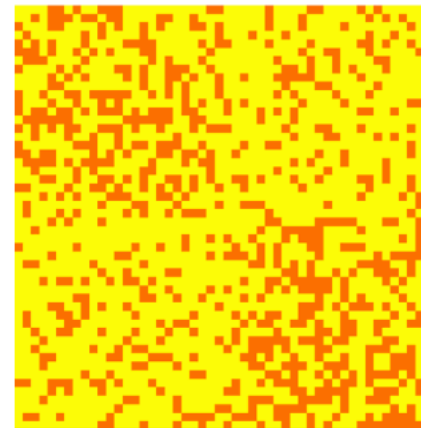


Ideal A

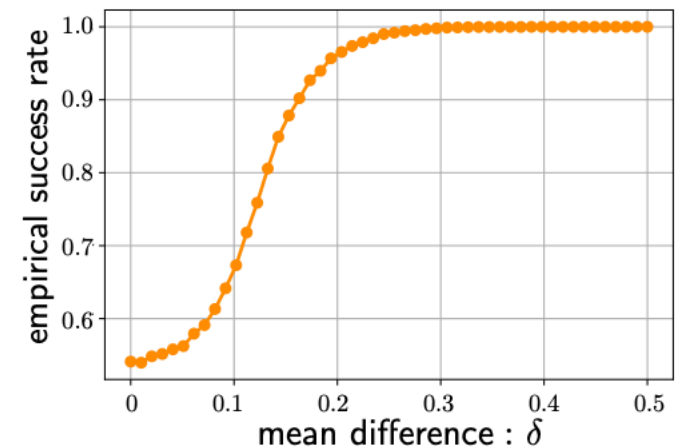
$$\mathbf{v}_1 = \mathbf{1}_n$$
$$\mathbf{v}_2 = \begin{bmatrix} +\mathbf{1}_{\frac{n}{2}} \\ -\mathbf{1}_{\frac{n}{2}} \end{bmatrix}$$



$$\tilde{A}_{ij} \sim \text{Bernoulli}\left(\frac{1 \pm \delta}{2}\right)$$



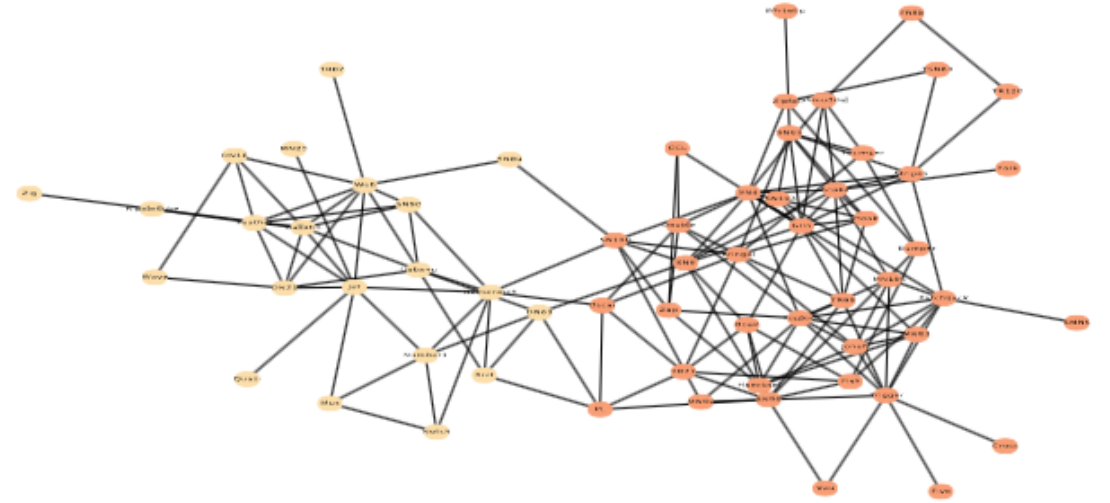
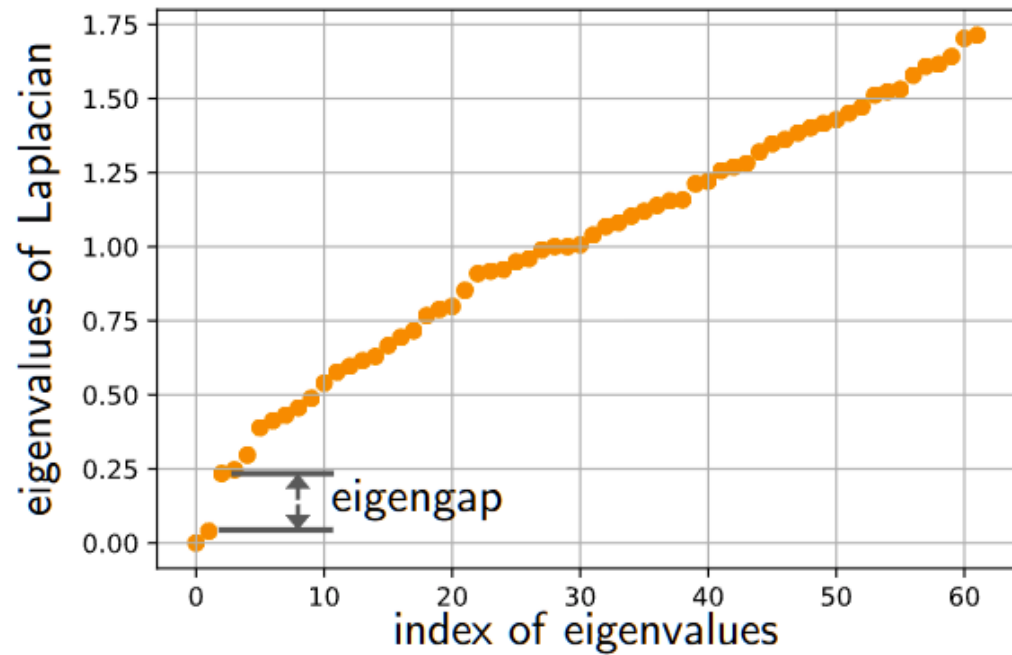
Noisy \tilde{A}



Introduction to Spectral Methods: Examples

Clustering: grouping the elements based on their mutual similarities

- Real-world data: two communities of dolphins in Doubtful Sound, New Zealand



Introduction to Spectral Methods: Examples

Principal component analysis (PCA): identifying a rank- r subspace to fit the data samples $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$

- Consider the sample covariance matrix:

$$\mathbf{M} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \in \mathbb{R}^{d \times d}$$

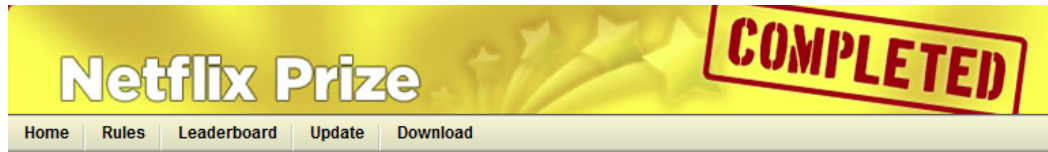
- If $\mathbf{x}_1, \dots, \mathbf{x}_n$ are approximately lie in a rank- r subspace \mathbf{U}_*
- Then we can take the rank- r leading eigenspace of \mathbf{M} , if the signal-to-noise ratio (SNR) is sufficiently high
- Real-world example: [Eigenface](#)



Introduction to Spectral Methods: Examples

Matrix completion: recovering a low-rank matrix from sub-sampled or incomplete observations

- Netflix prize



Leaderboard

Showing Test Score. [Click here to show quiz score](#)

Display top leaders.

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
Grand Prize - RMSE = 0.8567 - Winning Team: BellKor's Pragmatic Chaos				
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8582	9.90	2009-07-10 21:24:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries!	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43

	movie 1	movie 2	movie 3	movie 4	movie 5	movie 6	movie 7	movie 8	movie 9	movie 10	..	movie 17770
user 1			1		2							3
user 2		2		3	3			4				
user 3							5	3		4		
user 4	2				3			2				2
user 5		4				5			3			4
user 6			2									
user 7			2					4	2	3		
user 8	3	4				4						
user 9									3			
user 10			1		2							2
...												
user 480189		4			3			3				

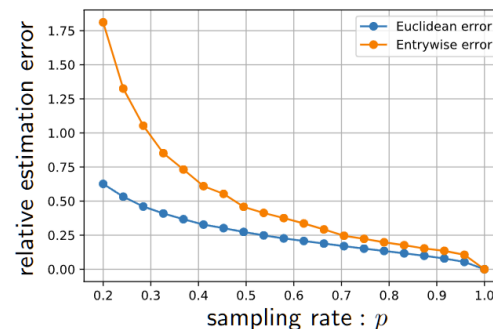
Introduction to Spectral Methods: Examples

Matrix completion: recovering a low-rank matrix from sub-sampled or incomplete observations

- Let $\mathbf{M}^* \in \mathbb{R}^{n \times n}$ be the ground-truth low-rank matrix
- Suppose the entries of \mathbf{M}^* are randomly sampled with probability p independently
- An unbiased estimator for \mathbf{M}^* is:

$$\hat{M}_{ij} := \begin{cases} \frac{1}{p} M_{ij}^* & \text{if } (i, j) \text{ is observed} \\ 0 & \text{otherwise} \end{cases}$$

- Then, we compute the best rank- r approximation of $\hat{\mathbf{M}}$ by the rank- r singular value decomposition



Introduction to Spectral Methods: Examples

Ranking based on pairwise comparison

- People struggle to rank many items at once, but find it much easier to make pairwise comparisons
- **Bradley-Terry model**: suppose there are latent scores $\{w_i^*\}_{i \in [n]}$ for each item
- The outcome of the comparison is generated as:

$$\Pr[i \succ j] = \frac{w_i^*}{w_i^* + w_j^*} \quad \forall i, j \in [n]$$

- Let $u_i^* := \log w_i^*$ for $i \in [n]$. We also have

$$\Pr[i \succ j] = \frac{\exp(u_i^*)}{\exp(u_i^*) + \exp(u_j^*)} = \sigma(u_i^* - u_j^*) \quad \forall i, j \in [n]$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$ is the logistic function

- Very popular recently due to applications in **aligning LLMs with human preference** (RLHF, DPO...)

Introduction to Spectral Methods: Examples

Ranking based on pairwise comparison

- **Bradley-Terry model:** suppose there are latent scores $\{w_i^*\}_{i \in [n]}$ for each item
- The outcome of the comparison is generated as:

$$\Pr[i \succ j] = \frac{w_i^*}{w_i^* + w_j^*} \quad \forall i, j \in [n]$$

- We can recover the global ranking by considering a probability transition matrix \mathbf{P}^* of a Markov chain:

$$P_{ij}^* = \begin{cases} \frac{1}{n} \frac{w_j^*}{w_i^* + w_j^*} & \text{if } i \neq j \\ 1 - \sum_{k \neq i} P_{ik}^* & \text{if } i = j \end{cases}$$

- Every entry of \mathbf{P}^* is non-negative, and every row sums to one, i.e., \mathbf{P}^* is a **stochastic matrix**

Introduction to Spectral Methods: Examples

Ranking based on pairwise comparison

- **Bradley-Terry model:** suppose there are latent scores $\{w_i^*\}_{i \in [n]}$ for each item
- The outcome of the comparison is generated as:

$$\Pr[i \succ j] = \frac{w_i^*}{w_i^* + w_j^*} \quad \forall i, j \in [n]$$

- We can recover the global ranking by considering a probability transition matrix \mathbf{P}^* of a Markov chain:

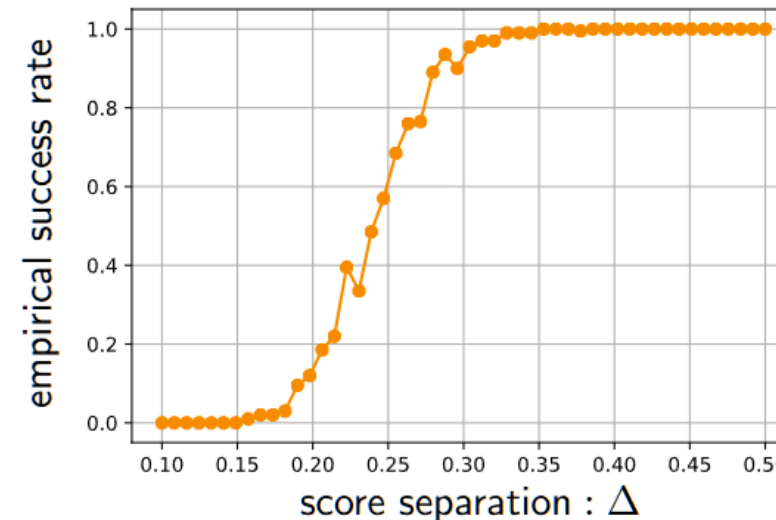
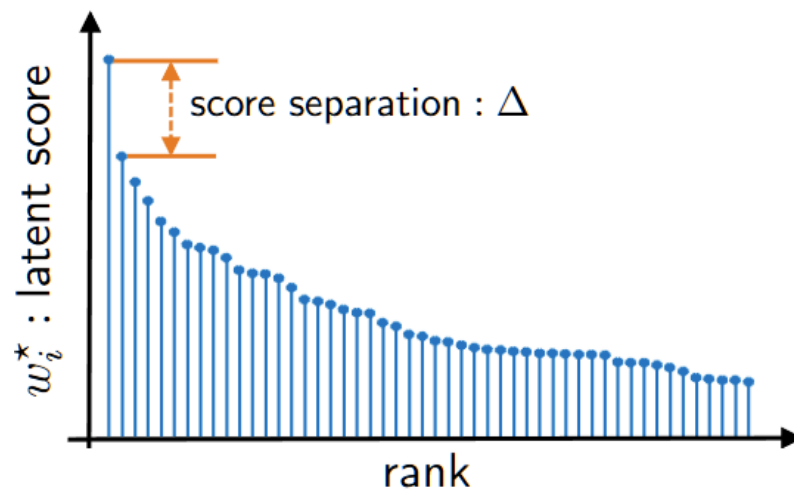
$$P_{ij}^* = \begin{cases} \frac{1}{n} \frac{w_j^*}{w_i^* + w_j^*} & \text{if } i \neq j \\ 1 - \sum_{k \neq i} P_{ik}^* & \text{if } i = j \end{cases}$$

- $\mathbf{w}^{*\top} \mathbf{P}^* = \mathbf{w}^{*\top}$, i.e., the latent score vector is a left eigenvector with eigenvalue one

Introduction to Spectral Methods: Examples

Ranking based on pairwise comparison

- Spectral ranking algorithm:
 - Use pairwise comparison to form an unbiased estimator of P^*
 - Compute its leading left eigenvector and report the ranking according to this eigenvector



- Google's PageRank algorithm uses a similar idea

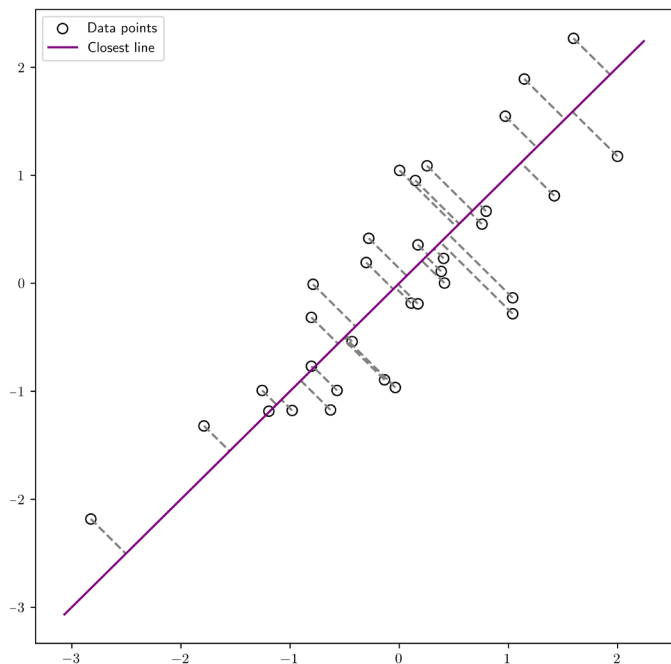
Introduction to Spectral Method: Unified Framework

1. Identify a key matrix M^* whose eigenvectors or singular vectors contains the target information (but M^* is usually inaccessible)
 2. Construct a surrogate matrix M of M^* using the data samples in hand, and compute the corresponding eigenvectors or singular vectors of this surrogate matrix
- Sample complexity?
 - Time complexity?
 - Noise stability?
 - Beyond worst-case (e.g., average-case, smoothed analysis)?

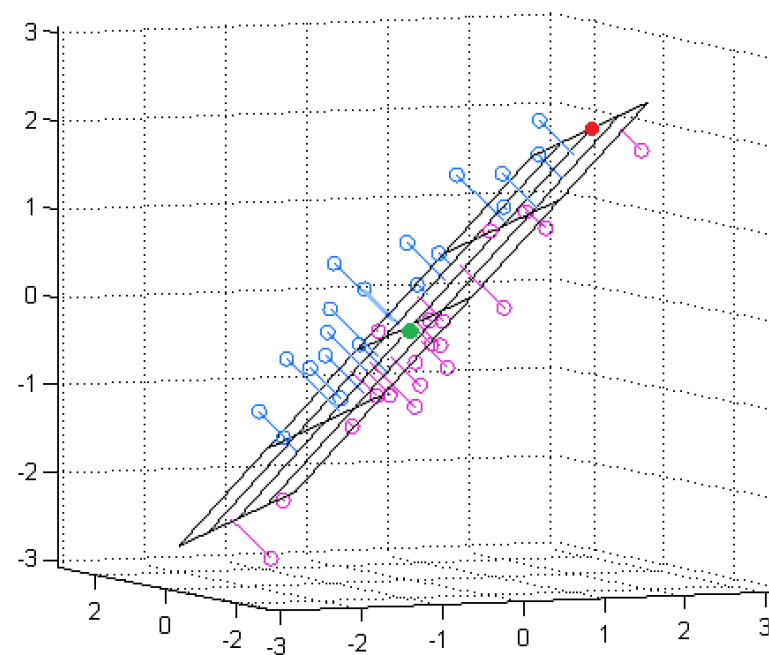
Today's Lecture

- Introduction to Spectral Method
- **SVD and Best-Fit Subspace**
- k -Means Clustering
- Complexity of SVD

SVD and Best-Fit Subspace



Linear regression



Best-fit subspace

Linear Algebra Self-Test

- Invertible matrix
- Diagonalizable matrix
- Normal matrix
- Determinant
- Similarity transformation
- Transpose, Hermitian transpose
- Pseudoinverse
- Orthogonal matrix, Unitary matrix
- Permutation matrix
- Vandermonde matrix
- Hankle matrix, Toeplitz matrix
- Eigendecomposition
- Singular value decomposition
- QR decomposition/Gram-Schmidt
- Jordan decomposition
- Polar decomposition
- Operator norm
- Frobenius norm
- Schatten p -norms
- Rank
- Kernel, Range

Singular Value Decomposition

- For an $n \times n$ matrix \mathbf{A} , an eigenvalue λ with corresponding eigenvector \mathbf{v} satisfies

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$

- If $\det(\mathbf{A}) \neq 0$, then \mathbf{A} has n nonzero eigenvalues and n corresponding eigenvectors

- For an $m \times n$ matrix \mathbf{A} , a **singular value** σ and corresponding **singular vectors** $\mathbf{u} \in \mathbb{R}^m, \mathbf{v} \in \mathbb{R}^n$ satisfies

$$\mathbf{A}\mathbf{v} = \sigma\mathbf{u}, \quad \mathbf{u}^\top\mathbf{A} = \sigma\mathbf{v}^\top$$

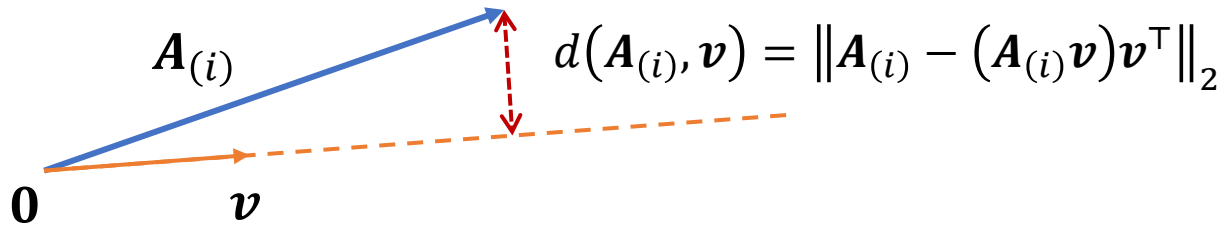
- You can observe that u and v are **equal length**:

$$\mathbf{u}^\top\mathbf{A}\mathbf{v} = \sigma\mathbf{u}^\top\mathbf{u} = \sigma\|\mathbf{u}\|_2^2 = \sigma\mathbf{v}^\top\mathbf{v} = \sigma\|\mathbf{v}\|_2^2$$

- We may always assume that $\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1$

Proposition. The vector $\mathbf{v}_1 = \arg \max_{\|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2^2$ is a singular vector, and $\|\mathbf{A}\mathbf{v}_1\|_2$ is the largest singular value

Best-Fit Line



$$\max_{\|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2^2 = \|\mathbf{A}\|_F^2 - \min_{\|\mathbf{v}\|_2=1} \sum_{i \in [m]} d(\mathbf{A}_{(i)}, \mathbf{v})^2$$

- If we take the data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^n$ to be the rows of A
- Then the top singular vector $\mathbf{v}_1 = \arg \max_{\|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2^2$ defines the direction of the line through the origin that best fits the data points in the mean squared error (MSE) sense

Singular Value Decomposition

Lemma. The (right) singular vectors of A are the eigenvectors of $A^T A$

Proof.

- Let \mathbf{v} be a singular vector of A with singular value σ . Then, we have

$$A^T A \mathbf{v} = A^T \sigma \mathbf{u} = \sigma (\mathbf{u}^T A)^T = \sigma^2 \mathbf{v}$$

- Conversely, let \mathbf{v} be an eigenvector of $A^T A$ with eigenvalue λ . Then, we have

$$\mathbf{v}^T A^T A \mathbf{v} = \lambda \|\mathbf{v}\|_2^2 \implies \mathbf{u} := \frac{A \mathbf{v}}{\sqrt{\lambda}} \text{ and } \sigma = \sqrt{\lambda}$$

- The proof also implies that the singular values $\sigma \geq 0$



Singular Value Decomposition

Proposition. The vector $\mathbf{v}_1 = \arg \max_{\|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2^2$ is a singular vector, and $\|\mathbf{A}\mathbf{v}_1\|_2$ is the largest singular value

Proof.

- Define the **Lagrangian** $\mathcal{L}(\mathbf{v}, \lambda) := \|\mathbf{A}\mathbf{v}\|_2^2 - \lambda(\|\mathbf{v}\|_2^2 - 1) = \mathbf{v}^\top \mathbf{A}^\top \mathbf{A} \mathbf{v} - \lambda(\mathbf{v}^\top \mathbf{v} - 1)$
- Then, we have
$$\nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v}, \lambda) = 2(\mathbf{A}^\top \mathbf{A})\mathbf{v} - 2\lambda\mathbf{v} = 0 \quad \implies \quad (\mathbf{A}^\top \mathbf{A})\mathbf{v} = \lambda\mathbf{v}$$
- That is, every eigenvector of $\mathbf{A}^\top \mathbf{A}$ is a stationary point for this constrained optimization problem
- Thus, \mathbf{v}_1 is the top eigenvector of $\mathbf{A}^\top \mathbf{A}$, which is also the top singular vector of \mathbf{A} by the previous lemma



Singular Value Decomposition

Theorem. Define the k -dimensional subspace V_k as the span of the following k vectors:

$$\mathbf{v}_1 := \arg \max_{\|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2^2$$

$$\mathbf{v}_2 := \arg \max_{\|\mathbf{v}\|_2=1, \mathbf{v} \perp \mathbf{v}_1} \|\mathbf{A}\mathbf{v}\|_2^2$$

\vdots

$$\mathbf{v}_k := \arg \max_{\|\mathbf{v}\|_2=1, \mathbf{v} \perp \mathbf{v}_1, \dots, \mathbf{v}_{k-1}} \|\mathbf{A}\mathbf{v}\|_2^2$$

Then, V_k is optimal in the sense that

$$V_k = \arg \min_{\dim(V)=k} \sum_{i \in [m]} d(\mathbf{A}_{(i)}, V)^2$$

Moreover, $\mathbf{v}_1, \dots, \mathbf{v}_n$ are all singular vectors with corresponding singular values

$$\sigma_1 = \|\mathbf{A}\mathbf{v}_1\|_2 \geq \sigma_2 = \|\mathbf{A}\mathbf{v}_2\|_2 \geq \dots \geq \sigma_n = \|\mathbf{A}\mathbf{v}_n\|_2$$

And $\mathbf{A} = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ with $\mathbf{u}_i := \sigma_i^{-1} \mathbf{A}\mathbf{v}_i$ is a **singular value decomposition (SVD)** of \mathbf{A}

Singular Value Decomposition

Proof.

- We prove the optimality of V_k by induction. $k = 1$ is trivial. Suppose that V_{k-1} is optimal
- Suppose V'_k is the optimal k -dimension subspace.
- Let $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^n$ be a set of orthonormal basis for V'_k such that $\mathbf{w}_k \in V_{k-1}^\perp$

- Then, by the optimality of V_{k-1} , we have

$$\|\mathbf{A}\mathbf{w}_1\|_2^2 + \dots + \|\mathbf{A}\mathbf{w}_{k-1}\|_2^2 + \|\mathbf{A}\mathbf{w}_k\|_2^2 \leq \|\mathbf{A}\mathbf{v}_1\|_2^2 + \dots + \|\mathbf{A}\mathbf{v}_{k-1}\|_2^2 + \|\mathbf{A}\mathbf{w}_k\|_2^2$$

- Thus, $V'_k = \text{span}\{V_{k-1} \cup \{\mathbf{w}_k\}\}$
- By the optimality of V'_k , $\mathbf{w}_k = \arg \max_{\|\mathbf{v}\|_2=1, \mathbf{v} \perp \mathbf{v}_1, \dots, \mathbf{v}_{k-1}} \|\mathbf{A}\mathbf{v}\|_2^2$
- Hence, V_k is also optimal

Singular Value Decomposition

Proof of 'Moreover'.

- We can generalize the previous proposition to $\mathbf{v}_2, \dots, \mathbf{v}_n$.
- Consider $\mathbf{v}_2 = \arg \max_{\|\mathbf{v}\|_2=1, \mathbf{v} \perp \mathbf{v}_1} \|\mathbf{A}\mathbf{v}\|_2^2$. Define the Lagrangian
$$\mathcal{L}(\mathbf{v}, \lambda, \mu) := \|\mathbf{A}\mathbf{v}\|_2^2 - \lambda(\|\mathbf{v}\|_2^2 - 1) - \mu \mathbf{v}_1^\top \mathbf{v}$$
- $\nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v}, \lambda, \mu) = 2(\mathbf{A}^\top \mathbf{A})\mathbf{v} - 2\lambda\mathbf{v} - \mu\mathbf{v}_1 = 0$
- Multiplying \mathbf{v}_1^\top on the left:
$$2\mathbf{v}_1^\top (\mathbf{A}^\top \mathbf{A})\mathbf{v} - 2\lambda\mathbf{v}_1^\top \mathbf{v} - \mu\mathbf{v}_1^\top \mathbf{v}_1 = 2\sigma_1^2 \mathbf{v}_1^\top \mathbf{v} - 2\lambda\mathbf{v}_1^\top \mathbf{v} - \mu = -\mu = 0$$
- Thus, $(\mathbf{A}^\top \mathbf{A})\mathbf{v} = \lambda\mathbf{v}$, i.e., \mathbf{v}_2 is an eigenvector of $\mathbf{A}^\top \mathbf{A}$, which is also a singular vector of \mathbf{A}
- For the decomposition $\mathbf{A} = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$, notice that $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is an orthonormal basis of \mathbb{R}^n , and

$$\left(\sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^\top \right) \mathbf{v}_k = \sigma_k \mathbf{u}_k = \|\mathbf{A}\mathbf{v}_k\|_2 \frac{\mathbf{A}\mathbf{v}_k}{\|\mathbf{A}\mathbf{v}_k\|_2} = \mathbf{A}\mathbf{v}_k \quad \forall k \in [n]$$



Singular Value Decomposition

In general, for $\mathbf{A} = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ with

- $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$
- $\{\mathbf{u}_i\}$ and $\{\mathbf{v}_i\}$ are orthonormal

Then we say it is a singular value decomposition of \mathbf{A}

- $V_k = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ gives the best-fit subspace for $\mathbf{A}_{(1)}, \dots, \mathbf{A}_{(m)} \in \mathbb{R}^n$

Theorem (Best rank- k approximation). Among all rank k matrices $\tilde{\mathbf{A}}$, the matrix $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ is the one which minimizes $\|\mathbf{A} - \tilde{\mathbf{A}}\|_F^2$. Moreover, $\|\mathbf{A} - \mathbf{A}_k\|_F^2 = \sum_{i=k+1}^n \sigma_i^2$

Singular Value Decomposition

Theorem (Best rank- k approximation). Among all rank k matrices $\tilde{\mathbf{A}}$, the matrix $\mathbf{A} = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ is the one which minimizes $\|\mathbf{A} - \tilde{\mathbf{A}}\|_F^2$. Moreover, $\|\mathbf{A} - \mathbf{A}_k\|_F^2 = \sum_{i=k+1}^n \sigma_i^2$

Proof.

- Note that

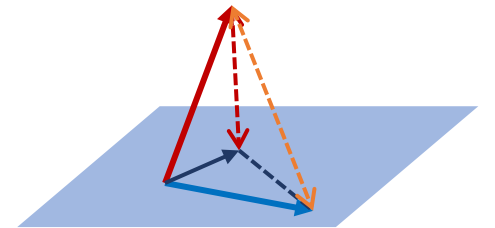
$$\|\mathbf{A} - \tilde{\mathbf{A}}\|_F^2 = \sum_{i=1}^m \|\mathbf{A}_{(i)} - \tilde{\mathbf{A}}_{(i)}\|_2^2$$

- Since $\text{rank}(\tilde{\mathbf{A}}) = k$, let V be the k -dimensional subspace spanned by the rows of $\tilde{\mathbf{A}}$

- $\|\mathbf{A}_{(i)} - \tilde{\mathbf{A}}_{(i)}\|_2^2 \geq \|\mathbf{A}_{(i)} - \mathcal{P}_V(\mathbf{A}_{(i)})\|_2^2 = d(\mathbf{A}_{(i)}, V)^2$

- Then, we have

$$\|\mathbf{A} - \tilde{\mathbf{A}}\|_F^2 = \sum_{i=1}^m d(\mathbf{A}_{(i)}, V)^2 \geq \|\mathbf{A} - \mathbf{A}_k\|_F^2$$



Singular Value Decomposition

Matrix norm and singular values

- Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with singular values $\sigma_1 \geq \dots \geq \sigma_n$
- The **spectral norm** $\|\mathbf{A}\| = \sigma_1$, i.e., the ℓ_∞ -norm of the singular values:

$$\|\mathbf{A}\| = \sqrt{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})} = \sigma_1 = \sup_{\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n: \|\mathbf{x}\|_2=1, \|\mathbf{y}\|_2=1} \mathbf{x}^\top \mathbf{A} \mathbf{y}$$

- The **Frobenius norm** is the ℓ_2 -norm of the singular values:

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |A_{ij}|^2} = \sqrt{\text{tr}[\mathbf{A}^\top \mathbf{A}]} = \sqrt{\sum_{i=1}^n \sigma_i^2}$$

- The **Schatten p -norm** is defined to be the ℓ_p -norm of the singular values

Today's Lecture

- Introduction to Spectral Method
- SVD and Best-Fit Subspace
- *k*-Means Clustering
- Complexity of SVD

k -Means Clustering

Given m points $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathbb{R}^n$, find k points $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_k\} \subset \mathbb{R}^n$ such that

$$f_{\mathcal{X}}(\mathcal{C}) := \sum_{i=1}^m \text{dist}(\mathbf{x}_i, \mathcal{C})^2$$

is minimized

- **NP**-hard even for $k = 2$
- The optimal solution is obtained by k clusters of $\mathcal{X} = S_1 \uplus \dots \uplus S_k$, and the cluster center \mathbf{c}_j is the **centroid** of the points in S_j

➤ For any set $\mathcal{S} = \{\mathbf{y}_1, \dots, \mathbf{y}_r\}$, let $\bar{\mathbf{y}} := \frac{1}{r}(\mathbf{y}_1 + \dots + \mathbf{y}_r)$. Then, for any $\mathbf{c} \in \mathbb{R}^n$, we have

$$\sum_{i=1}^r \|\mathbf{y}_i - \mathbf{c}\|_2^2 = \sum_{i=1}^r \|\mathbf{y}_i - \bar{\mathbf{y}}\|_2^2 + r\|\bar{\mathbf{y}} - \mathbf{c}\|_2^2 \geq \sum_{i=1}^r \|\mathbf{y}_i - \bar{\mathbf{y}}\|_2^2$$

k -Means Clustering

Let $g_x(V)$ be the sum of squared distance to a k -dimensional subspace V :

$$g_x(V) := \sum_{i=1}^m \text{dist}(\mathbf{x}_i, V)^2$$

$g_x(V)$ can be minimized using SVD. Moreover, it provides a **lower bound** for k -means clustering:

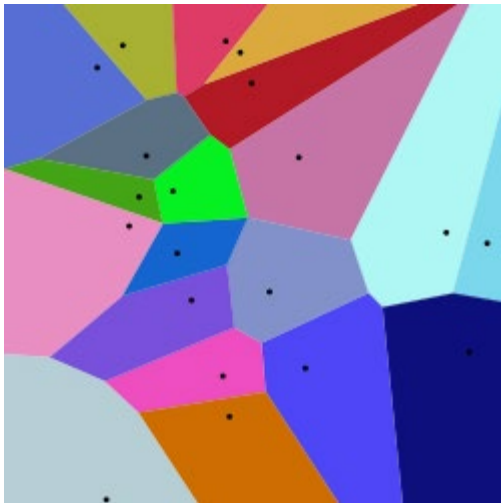
$$f_x(\mathcal{C}) \geq g_x(\text{span}(\mathcal{C}))$$

- The best-fit subspace problem can be regarded as a **relaxation** of the k -means problem
- We can design a **2-approximation** algorithm for k -means by:
 1. Identifying a best-fit subspace to reduce the dimensionality
 2. Exactly solving the k -means problem in this lower-dimensional subspace

k -Means Clustering

Theorem (Exact k -means). The k -means problem can be solved in $\mathcal{O}(m^{k^2 d/2})$ time when the input dataset $\mathcal{X} \subset \mathbb{R}^d$ of size m

- The most brute-force method is to enumerate all possible ways to form k clusters, and the complexity is $\mathcal{O}(k^m m d) = \exp(\mathcal{O}(m \log k)) \gg \exp(\mathcal{O}(k^2 d \log m))$
- k centers \Leftrightarrow Voronoi diagram



\mathbb{R}^d is partitioned into k cells:

$$R_i := \left\{ \mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{c}_i\|_2 \leq \|\mathbf{x} - \mathbf{c}_j\|_2 \quad \forall j \neq i \right\}$$

- There are at most $\binom{k}{2}$ faces (boundary hyperplanes)
- We can enumerate these $\mathcal{O}(k^2)$ hyperplanes

k -Means Clustering

Since only the cluster assignments of the m points matter, we may freely translate or rotate any boundary hyperplane as long as it does **not cross** any data point

Algorithm (Voronoi- k -means).

1. Enumerate the sets of t hyperplanes such that **each hyperplane contains exactly d data points**, where $k \leq t \leq k^2/2$, and the total number of such sets is:

$$\sum_{t=k}^{k^2/2} \binom{m}{d}^t = \mathcal{O}(m^{dk^2/2})$$

2. Check for each set whether the t hyperplanes form k cells
3. If so, assign the boundary data points to the left or the right cells
#assignments = 2^{td}
4. Compute the centroid of the data points in each cell as the centers

k -Means Clustering

Algorithm (2-approximation for k -means).

1. Compute the rank- k best-fit subspace V by SVD
 2. Let $\mathbf{x}_1^V, \mathbf{x}_2^V, \dots, \mathbf{x}_m^V$ be the projection of the dataset to V
 3. Output $\mathcal{C}^V \leftarrow \text{Voronoi-}k\text{-means}(\{\mathbf{x}_1^V, \mathbf{x}_2^V, \dots, \mathbf{x}_m^V\})$
- Let OPT be the optimal value of the k -means problem. We have

$$\text{OPT} \geq \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{x}_i^V\|_2^2$$

- We also have

$$\text{OPT} = \sum_{i=1}^m d(\mathbf{x}_i, \mathcal{C}_*)^2 \geq \sum_{i=1}^m d(\mathbf{x}_i^V, \mathcal{C}_*^V)^2 \geq \sum_{i=1}^m d(\mathbf{x}_i^V, \mathcal{C}^V)^2$$

k-Means Clustering

- Let OPT be the optimal value of the *k*-means problem. We have

$$\text{OPT} \geq \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{x}_i^V\|_2^2$$

- We also have

$$\text{OPT} = \sum_{i=1}^m d(\mathbf{x}_i, \mathcal{C}_\star)^2 \geq \sum_{i=1}^m d(\mathbf{x}_i^V, \mathcal{C}_\star^V)^2 \geq \sum_{i=1}^m d(\mathbf{x}_i^V, \mathcal{C}^V)^2$$

- Summing them together, we get

$$2\text{OPT} \geq \sum_{i=1}^m \left(\|\mathbf{x}_i - \mathbf{x}_i^V\|_2^2 + d(\mathbf{x}_i^V, \mathcal{C}^V)^2 \right) = \sum_{i=1}^m d(\mathbf{x}_i, \mathcal{C}^V)^2 =: \text{ALG}$$

- Thus, the algorithm can produce a 2-approximate solution in time

$$\mathcal{O}(mn^2 + m^{k^3/2})$$

Today's Lecture

- Introduction to Spectral Method
- SVD and Best-Fit Subspace
- k -Means Clustering
- Complexity of SVD

Complexity of SVD

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

Practical version

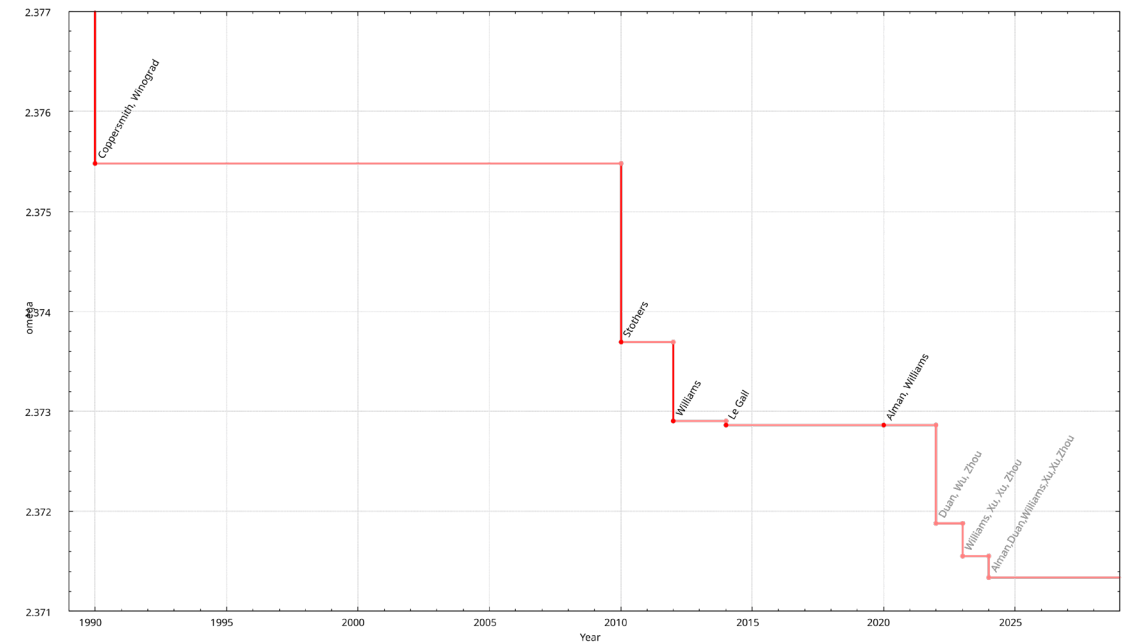
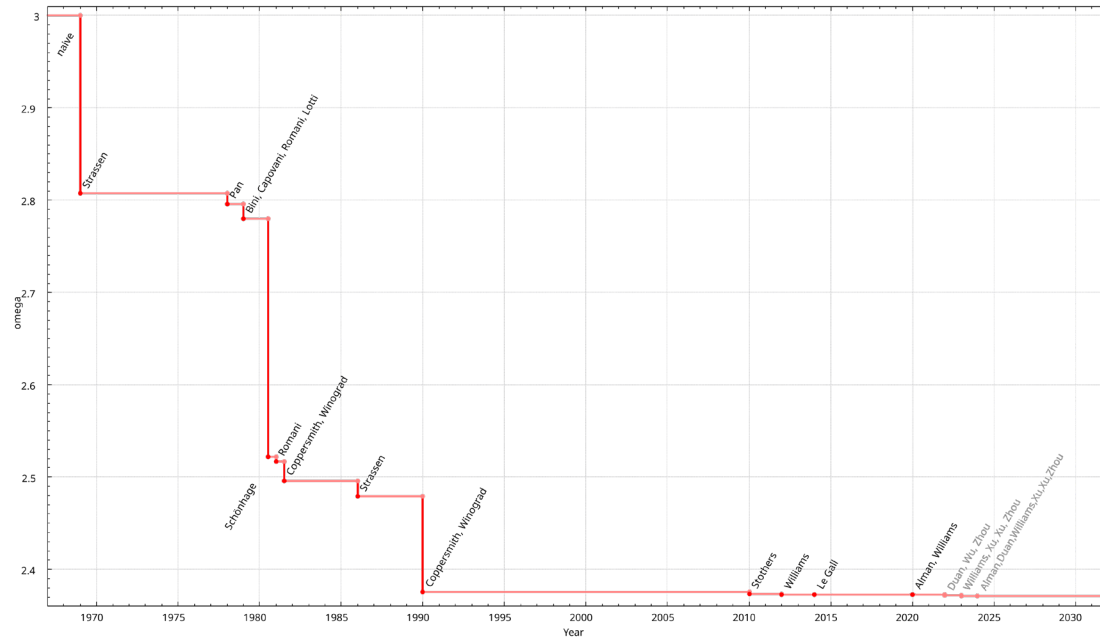
- Complexity: $\mathcal{O}(mn^2 + n^3)$
- $\mathbf{A} \rightarrow \mathbf{A}^\top \mathbf{A} \rightarrow$ eigendecomposition \rightarrow SVD
- LAPACK (used by Matlab and Numpy): highly optimized version of this approach

Complexity of SVD

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

Theoretical (TCS) version

- Complexity: $\mathcal{O}^*(mn^{\omega-1})$ where $\omega \approx 2.371339$ is the fast matrix multiplication
- See (Kacham-Woodruff '24, Appendix A) and (Banks-Garza-Vargas-Kulkarni-Srivastava '23)



Complexity of SVD

Let $A \in \mathbb{R}^{m \times n}$ and $m \geq n$. Suppose we only want to approximate the top- k singular values and singular vectors (e.g., low-rank approximation or PCA)

LazySVD

- $A \rightarrow$ top singular vector $v_1 \rightarrow$ left-project $(I - vv^T)A \rightarrow$ repeat k times
- Complexity: $\mathcal{O}\left(\frac{k \cdot \text{nnz}(A) + k^2 n}{\sqrt{\epsilon}}\right)$ (Allen-Zhu-Li '17, Kacham-Woodruff '24)

How to compute the top-singular vector?

- Allen-Zhu-Li's algorithms use **convex optimizers** including accelerated gradient descent (AGD) and accelerated stochastic variance reduced gradient (SVRG)
- We'll see a simpler approach: **the Power Method**

Detour: Computational Models

In optimization and numerical linear algebra, we are dealing with **real number operations**

There are various computational models:

- **Exact real arithmetic:** infinite precision
- **Variable precision rational arithmetic:** rationals are stored exactly as numerators and denominators
- **Finite precision arithmetic:** real numbers are rounded to a fixed number of bits which may depend on the input size and accuracy

The third model is more **realistic** than the first, and potentially more **efficient** in terms of the bit complexity than the second

Detour: Finite Precision Arithmetic

- Numbers are stored and manipulated approximately up to some **machine precision**

$$\varepsilon := \varepsilon_{\text{mach}}(n, \delta) > 0$$

which depends on the instance size n and the desired accuracy δ

- Every $x \in \mathbb{R}$ is stored as **$\text{fl}(x) = (1 + \Delta)x$** for some adversarially chosen $\Delta \in \mathbb{R}$ such that $|\Delta| \leq \varepsilon$

- For each arithmetic operator $\circ \in \{+, -, \times, \div\}$, it holds that

$$\text{fl}(x \circ y) = (1 + \Delta)(x \circ y), \quad |\Delta| \leq \varepsilon$$

- $\text{fl}(\sqrt{x}) = (1 + \Delta)\sqrt{x}$ for $|\Delta| \leq \varepsilon$

- The bit lengths of numbers stored in this form is **fixed** at $\log_2(1/\varepsilon)$

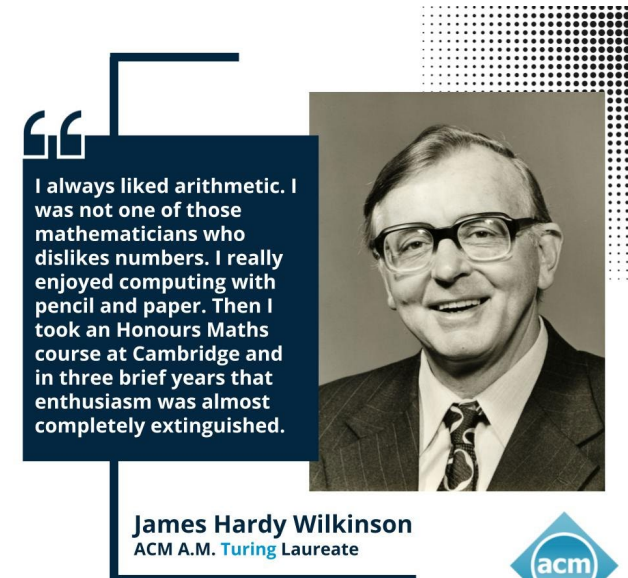
- The **bit complexity** of an algorithm is the number of arithmetic operations times $\tilde{O}(\log_2(1/\varepsilon))$, the running time of standard floating-point arithmetic

Detour: Finite Precision Arithmetic

- An iterative algorithm that can be implemented in finite precision (i.e., bit length $\text{polylog}(n/\delta)$) is called **numerically stable**
- Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a function. Its **condition number** is defined as the smallest $\kappa_f \geq 0$ such that

$$\underbrace{\frac{\|f(x + \delta x) - f(x)\|_2}{\|f(x)\|_2}}_{\text{Forward error}} \leq \kappa_f \cdot \underbrace{\frac{\|\delta x\|_2}{\|x\|_2}}_{\text{Backward error}} + \mathcal{O}\left(\left(\frac{\|\delta x\|_2}{\|x\|_2}\right)^2\right)$$

- Backward stability is the gold standard for algorithms
 - Suppose an algorithm outputs y
 - $\delta x := \arg \min\{\|\Delta\|_2 : f(x + \Delta) = y\}$
 - Backward error $:= \frac{\|\delta x\|_2}{\|x\|_2} \leq \text{poly}(n) \cdot \varepsilon + \mathcal{O}(\varepsilon^2)$, then **backward stable**



Detour: Finite Precision Arithmetic

Examples

- Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, there is an algorithm that outputs $\mathbf{C} \in \mathbb{R}^{n \times n}$ such that

$$\|\mathbf{C} - \mathbf{AB}\| \leq \text{poly}(n) \cdot \varepsilon \cdot \|\mathbf{A}\| \|\mathbf{B}\|$$

on a floating-point machine with precision ε , in $\mathcal{O}^*(n^\omega)$ arithmetic operations ([Demmel-Dumitriu-Holtz-Kleinberg '07](#))

- Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be an invertible matrix with **condition number** $\kappa(\mathbf{A}) := \frac{\sigma_1}{\sigma_n}$. There is an algorithm that outputs $\mathbf{C} \in \mathbb{R}^{n \times n}$ such that

$$\|\mathbf{C} - \mathbf{A}^{-1}\| \leq \kappa(\mathbf{A})^{\text{poly log } n} \cdot \varepsilon \cdot \|\mathbf{A}^{-1}\|$$

on a floating-point machine with precision ε , in $\mathcal{O}^*(n^\omega)$ arithmetic operations ([Demmel-Dumitriu-Holtz '07](#))

Complexity of SVD: Power Method

Let's first assume that $A \in \mathbb{R}^{n \times n}$ is square symmetric and has the same left- and right-singular vectors:

$$A = \sum_{i=1}^r \sigma_i \mathbf{v}_i \mathbf{v}_i^\top$$

- Squaring it gives

$$A^2 = \left(\sum_{i=1}^r \sigma_i \mathbf{v}_i \mathbf{v}_i^\top \right) \left(\sum_{j=1}^r \sigma_j \mathbf{v}_j \mathbf{v}_j^\top \right) = \sum_{i,j=1}^r \sigma_i \sigma_j \mathbf{v}_i (\mathbf{v}_i^\top \mathbf{v}_j) \mathbf{v}_j^\top = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^\top$$

- Multiplying k times gives

$$A^k = \sum_{i=1}^r \sigma_i^k \mathbf{v}_i \mathbf{v}_i^\top, \quad \frac{A^k}{\|A^k\|_F} = \sum_{i=1}^r \frac{\sigma_i^k}{(\sigma_1^{2k} + \dots + \sigma_r^{2k})^{1/2}} \mathbf{v}_i \mathbf{v}_i^\top \rightarrow \mathbf{v}_1 \mathbf{v}_1^\top \text{ if } \sigma_1 > \sigma_2$$

Complexity of SVD: Power Method

In general, let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be rectangular. Then, we can apply the power method to $\mathbf{B} = \mathbf{A}^\top \mathbf{A}$:

$$\mathbf{B}^k = \left(\left(\sum_{i=1}^r \sigma_i \mathbf{v}_i \mathbf{u}_i^\top \right) \left(\sum_{j=1}^r \sigma_j \mathbf{u}_j \mathbf{v}_j^\top \right) \right)^k = \left(\sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^\top \right)^k = \sum_{i=1}^r \sigma_i^{2k} \mathbf{v}_i \mathbf{v}_i^\top$$

- As k increases, σ_i / σ_1 goes to zero, and $\mathbf{B}^k \approx \sigma_1^{2k} \mathbf{v}_1 \mathbf{v}_1^\top$
- However, matrix multiplication is expensive
- We can pick a vector \mathbf{x} and compute $\mathbf{x}^{(k)} := \mathbf{B}^k \mathbf{x} = \mathbf{A}^\top \mathbf{A} \mathbf{x}^{(k-1)}$
- Each iteration involves two Mat-Vec products, which take $\mathcal{O}(\text{nnz}(\mathbf{A}))$ time
- $\mathbf{x}^{(k)} \approx \sigma_1^{2k} \mathbf{v}_1 \langle \mathbf{x}, \mathbf{v}_1 \rangle$. If $\langle \mathbf{x}, \mathbf{v}_1 \rangle \neq 0$, we can recover \mathbf{v}_1 from $\mathbf{x}^{(k)}$ by normalization

Complexity of SVD: Power Method

Lemma. Let $\mathbf{x} = (x_1, x_2, \dots, x_n)$ be a unit d -dimensional vector picked at random from \mathbb{S}^{n-1} . The probability that $|x_1| \geq \frac{1}{20\sqrt{n}}$ is at least 0.9

Proof.

- Let $\alpha := 1/(20\sqrt{n})$
- We first show that for a uniformly random $\mathbf{v} \sim B_2^n$, $\Pr[|v_1| \geq \alpha] \geq 0.9$
- Then, by taking $\mathbf{x} := \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$, \mathbf{x} is a uniformly random unit vector and $x_1 \geq v_1$

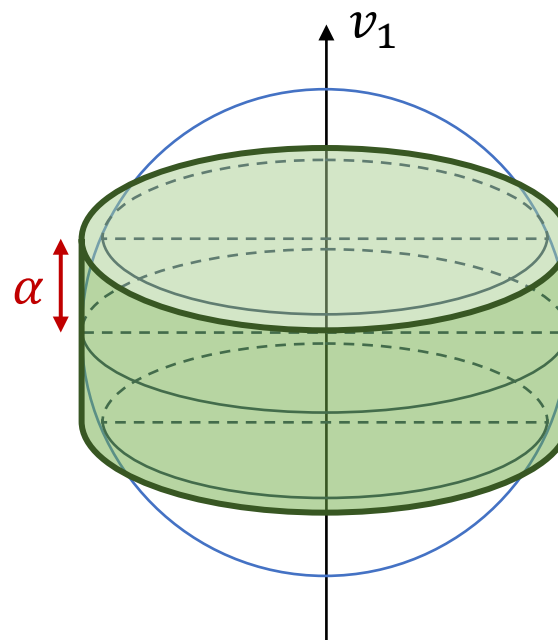
Complexity of SVD: Power Method

Lemma. Let $\mathbf{x} = (x_1, x_2, \dots, x_n)$ be a unit d -dimensional vector picked at random from \mathbb{S}^{n-1} . The probability that $|x_1| \geq \frac{1}{20\sqrt{n}}$ is at least 0.9

Proof.

- Notice that

$$\begin{aligned}\Pr[|v_1| \leq \alpha] &\leq \frac{2\alpha \cdot \text{Vol}_{n-1}(B_2^{n-1})}{\text{Vol}_n(B_2^n)} \\ &\sim 2\alpha \left(\frac{2\pi e}{n-1}\right)^{(n-1)/2} \left(\frac{n}{2\pi e}\right)^{n/2} \\ &= \frac{\alpha\sqrt{2n}}{\sqrt{\pi e}} \left(1 + \frac{1}{n-1}\right)^{(n-1)/2} \\ &\sim \frac{\alpha\sqrt{2n}}{\sqrt{\pi}} = \frac{\sqrt{2n}}{20\sqrt{n\pi}} \approx 0.04\end{aligned}$$



$$\Pr_{v \sim B_2^n}[|v_1| \leq \alpha]$$



Complexity of SVD: Power Method

Theorem. Let A be an $m \times n$ matrix and \mathbf{x} a random unit length vector. Let V be the space spanned by the left singular vectors of A corresponding to singular values greater than $(1 - \epsilon)\sigma_1$.

Let $k = \Omega\left(\frac{\ln(n/\epsilon)}{\epsilon}\right)$. Let \mathbf{w} be unit vector after k iterations of the power method:

$$\mathbf{w} = \frac{\mathbf{x}^{(k)}}{\|\mathbf{x}^{(k)}\|_2} = \frac{(AA^\top)^k \mathbf{x}}{\|(AA^\top)^k \mathbf{x}\|_2}$$

Then, the probability that $d(\mathbf{w}, V) \geq \epsilon$ is at most $1/10$

- It does not rely on any **explicit** singular value **gap assumption**
- If we know that $\sigma_1 - \sigma_2 \geq \Delta$, then we can take $\epsilon = \Delta/2$ and recover v_1 with ϵ accuracy in **$\tilde{O}(\text{nnz}(A)\Delta^{-1})$ time**

Complexity of SVD: Power Method

Proof.

- Let $\mathbf{A} = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ be the full SVD
- Expand \mathbf{x} in the right singular vector basis: $\mathbf{x} = \sum_{i=1}^n c_i \mathbf{v}_i$ where $c_i \in \mathbb{R}$ and $\sqrt{c_1^2 + \dots + c_n^2} = 1$
- Then, we have $\mathbf{x}^{(k)} = \sum_{i=1}^n \sigma_i^{2k} \mathbf{v}_i \langle \mathbf{x}, \mathbf{v}_i \rangle = \sum_{i=1}^n \sigma_i^{2k} c_i \mathbf{v}_i$
- Since \mathbf{x} is a random unit vector, and $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is a fixed orthonormal basis of \mathbb{R}^n
- We get that $\mathbf{c} := (c_1, \dots, c_n)$ is a random unit vector
- The [lemma](#) implies that $|c_1| \geq \frac{1}{20\sqrt{n}}$ with probability ≥ 0.9
- Suppose $\sigma_1 \geq \dots \geq \sigma_s \geq (1 - \epsilon)\sigma_1 > \sigma_{s+1} \geq \dots \geq \sigma_n$, then we have

$$\|\mathbf{x}^{(k)}\|_2^2 = \sum_{i=1}^n \sigma_i^{4k} c_i^2 \geq \sigma_1^{4k} c_1^2 \geq \frac{\sigma_1^{4k}}{400n}$$

Complexity of SVD: Power Method

Proof.

- Suppose $\sigma_1 \geq \dots \geq \sigma_s \geq (1 - \epsilon)\sigma_1 > \sigma_{s+1} \geq \dots \geq \sigma_n$, then we have

$$\|\mathbf{x}^{(k)}\|_2^2 = \sum_{i=1}^n \sigma_i^{4k} c_i^2 \geq \sigma_1^{4k} c_1^2 \geq \frac{\sigma_1^{4k}}{400n}$$

- Thus, for $V = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_s\}$, we have

$$d(\mathbf{x}^{(k)}, V)^2 = \sum_{i=s+1}^n \sigma_i^{4k} c_i^2 \leq (1 - \epsilon)^{4k} \sigma_1^{4k} \sum_{i=s+1}^n c_i^2 \leq (1 - \epsilon)^{4k} \sigma_1^{4k}$$

- Hence,

$$d(\mathbf{w}, V) = \frac{d(\mathbf{x}^{(k)}, V)}{\|\mathbf{x}^{(k)}\|_2} \leq \frac{(1 - \epsilon)^{2k} \sigma_1^{2k}}{\sigma_1^{2k} / 20\sqrt{n}} = 20(1 - \epsilon)^{2k} \sqrt{n} \leq \epsilon$$

if we take $k = \Omega(\ln(n/\epsilon)/\epsilon)$

