# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2023.
Lecture 21

## Logistics

See Piazza post about upcoming schedule information.

- No quiz due this week.

- Problem Set 4 is due 12/1.

- No class Thursday.

- Office hours next Monday at 10am in CS234.

- No class next Tuesday

- Class over Zoom next Thursday 11/30 at 10am. Office hours over Zoom at 9am. See Piazza for Zoom link.

- Second Linear Algebra Review Session on Monday 12/4 at 3pm in CS140.

## Summary

### Last Few Classes: Spectral Graph Partitioning

- Focus on separating graphs with small but relatively balanced cuts.
- Connection to second smallest eigenvector of graph Laplacian.
- Provable guarantees for stochastic block model.
- Expectation analysis in class. See slides for full analysis.

### This Class: Computing the SVD/eigendecomposition.

- Efficient algorithms for SVD/eigendecomposition.
- Iterative methods: power method, Krylov subspace methods.
- High level: a glimpse into fast methods for linear algebraic computation, which are workhorses behind data science.

Consider solving the optimization problem: $\min_{\{z \in \{-1,1\}^n: \text{ not all entries of } z \text{ are equal}\}} z^T L z$.

What is this optimization problem commonly known as?

- a. Computing the lowest eigenvalue of the graph Laplacian.
- b. Computing the second lowest eigenvalue of the graph Laplacian.
- c. Computing the minimum cut.
- d. Computing the smallest node degree.
- e. Computing the maximum eigenvalue of the graph Laplacian.

Check

We have talked about the eigendecomposition and SVD as ways to compress data, to embed entities like words and documents, to compress/cluster non-linearly separable data.

How efficient are these techniques? Can they be run on large datasets?

## Computing the SVD

**Basic Algorithm:** To compute the SVD of full-rank $X \in \mathbb{R}^{n \times d}$,
$X = U\Sigma V^T$:

- Compute $X^T X$ – $O(nd^2)$ runtime.
- Find eigendecomposition $X^T X = V\Lambda V^T$ – $O(d^3)$ runtime.
- Compute $L = XV$ – $O(nd^2)$ runtime. Note that $L = U\Sigma$.
- Set $\sigma_i = \|L_i\|_2$ and $U_i = L_i/\|L_i\|_2$. – $O(nd)$ runtime.

  **Total runtime:** $O(nd^2 + d^3) = O(nd^2)$ (assume w.l.o.g. $n \geq d$)

- If we have $n = 10$ million images with $200 \times 200 \times 3 = 120,000$ pixel values each, runtime is $1.5 \times 10^{17}$ operations!
- The worlds fastest super computers compute at $\approx 100$ petaFLOPS = $10^{17}$ FLOPS (floating point operations per second).
- This is a relatively easy task for them – but no one else.

To speed up SVD computation we will take advantage of the fact that we typically only care about computing the top (or bottom) $k$ singular vectors of a matrix $X \in \mathbb{R}^{n \times d}$ for $k \ll d$.

- Suffices to compute $V_k \in \mathbb{R}^{d \times k}$ and then compute $U_k \Sigma_k = X V_k$.

- Use an *iterative algorithm* to compute an *approximation* to the top $k$ singular vectors $V_k$ (the top $k$ eigenvectors of $X^T X$.)

- Runtime will be roughly $O(ndk)$ instead of $O(nd^2)$.

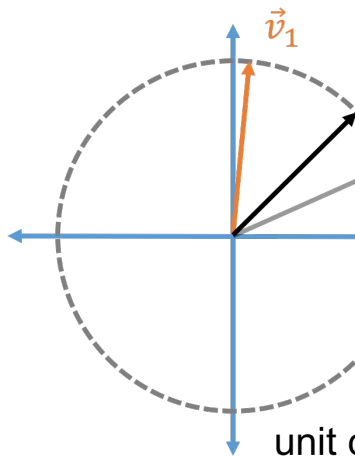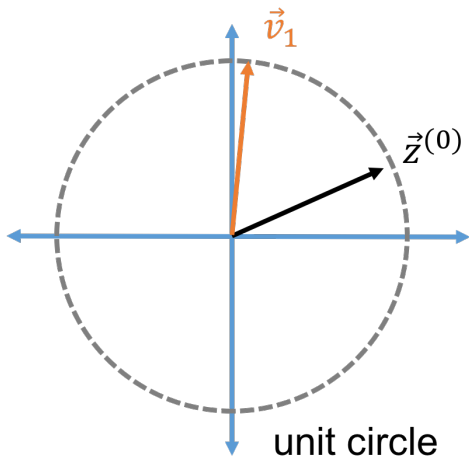Sparse (iterative) vs. Direct Method. `svd` vs. `svds`.

## Power Method

Power Method: The most fundamental iterative method for approximate SVD/eigendecomposition. Applies to computing $k = 1$ eigenvectors, but can be generalized to larger $k$.

Goal: Given symmetric $A \in \mathbb{R}^{d \times d}$, with eigendecomposition $A = V \Lambda V^T$, find $\vec{z} \approx \vec{v}_1$. I.e., the top eigenvector of $A$.

- **Initialize:** Choose $\vec{z}^{(0)}$ randomly. E.g. $\vec{z}^{(0)}(i) \sim \mathcal{N}(0, 1)$.

- For $i = 1, \ldots, t$
    - $\vec{z}^{(i)} := A \cdot \vec{z}^{(i-1)}$
    - $\vec{z}_i := \frac{\vec{z}^{(i)}}{\|\vec{z}^{(i)}\|_2}$

- Return $\vec{z}_t$

$\vec{v}_1$

$\vec{z}^{(0)}$

unit circle

$\vec{v}_1$

unit c

## Power Method Analysis

Power method:

- **Initialize:** Choose $\vec{z}^{(0)}$ randomly. E.g. $\vec{z}^{(0)}(i) \sim \mathcal{N}(0, 1)$.
- For $i = 1, \ldots, t$
    - $\vec{z}^{(i)} := A \cdot \vec{z}^{(i-1)}$
    - $\vec{z}_i := \frac{\vec{z}^{(i)}}{\|\vec{z}^{(i)}\|_2}$
- Return $\vec{z}_t$.

Theoretically equivalent to:

- For $i = 1, \ldots, t$
    - $\vec{z}^{(i)} := A \cdot \vec{z}^{(i-1)}$
- $\vec{z}_i := \frac{\vec{z}^{(i)}}{\|\vec{z}^{(i)}\|_2}$.
- Return $\vec{z}_t$.

## Power Method Analysis

Write $\vec{z}^{(0)}$ in **A**'s eigenvector basis:

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \ldots + c_d\vec{v}_d.$$

**Update step:** $\vec{z}^{(i)} = \mathbf{A} \cdot \vec{z}^{(i-1)} = \mathbf{V\Lambda V}^T \cdot \vec{z}^{(i-1)}$ (then normalize)

$$\mathbf{V}^T\vec{z}^{(0)} =$$

$$\mathbf{\Lambda V}^T\vec{z}^{(0)} =$$

$$\vec{z}^{(1)} = \mathbf{V\Lambda V}^T \cdot \vec{z}^{(0)} =$$

---

$\mathbf{A} \in \mathbb{R}^{d \times d}$: input matrix with eigendecomposition $\mathbf{A} = \mathbf{V\Lambda V}^T$. $\vec{v}_1$: top eigenvector, being computed, $\vec{z}^{(i)}$: iterate at step $i$, converging to $\vec{v}_1$.

## Power Method Analysis

**Claim 1:** Writing $\vec{z}^{(0)} = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \ldots + c_d \vec{v}_d$,

$$\vec{z}^{(1)} = c_1 \cdot \lambda_1 \vec{v}_1 + c_2 \cdot \lambda_2 \vec{v}_2 + \ldots + c_d \cdot \lambda_d \vec{v}_d.$$

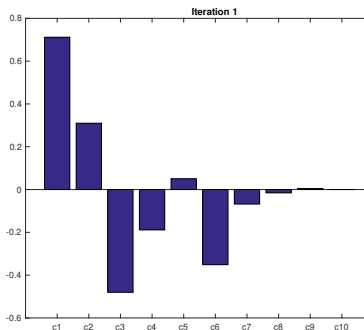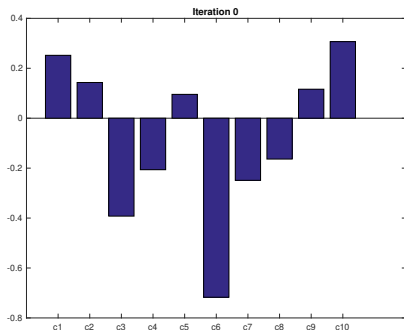$$\vec{z}^{(2)} = A\vec{z}^{(1)} = V\Lambda V^T \vec{z}^{(1)} =$$

**Claim 2:**

$$\vec{z}^{(t)} = c_1 \cdot \lambda_1^t \vec{v}_1 + c_2 \cdot \lambda_2^t \vec{v}_2 + \ldots + c_d \cdot \lambda_d^t \vec{v}_d.$$

---

$A \in \mathbb{R}^{d \times d}$: input matrix with eigendecomposition $A = V\Lambda V^T$. $\vec{v}_1$: top eigenvector, being computed, $\vec{z}^{(i)}$: iterate at step $i$, converging to $\vec{v}_1$.

After $t$ iterations, we have 'powered' up the eigenvalues, making the component in the direction of $v_1$ much larger, relative to the other components.

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \ldots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\lambda_1^t\vec{v}_1 + c_2\lambda_2^t\vec{v}_2 + \ldots + c_d\lambda_d^t\vec{v}_d$$
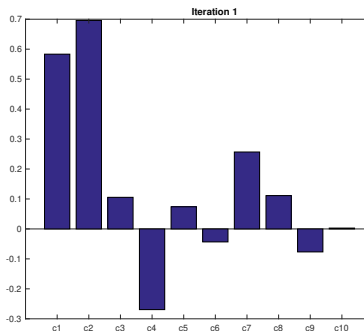


When will convergence be slow?

## Power Method Slow Convergence

**Slow Case: A** has eigenvalues: $\lambda_1 = 1, \lambda_2 = .99, \lambda_3 = .9, \lambda_4 = .8, \ldots$

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \ldots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\lambda_1^t\vec{v}_1 + c_2\lambda_2^t\vec{v}_2 + \ldots + c_d\lambda_d^t\vec{v}_d$$
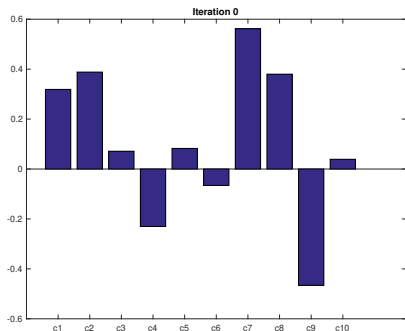
$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \ldots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\lambda_1^t\vec{v}_1 + c_2\lambda_2^t\vec{v}_2 + \ldots + c_d\lambda_d^t\vec{v}_d$

Write $|\lambda_2| = (1 - \gamma)|\lambda_1|$ for 'gap' $\gamma = \frac{|\lambda_1| - |\lambda_2|}{|\lambda_1|}$.

How many iterations $t$ does it take to have $|\lambda_2|^t \leq \delta \cdot |\lambda_1|^t$ for $\delta > 0$?

$$\begin{aligned} |\lambda_2|^t &= (1 - \gamma)^t \cdot |\lambda_1|^t \\ &= (1 - \gamma)^{1/\gamma \cdot \gamma t} \cdot |\lambda_1|^t \\ &\leq e^{-\gamma t} \cdot |\lambda_1|^t \end{aligned}$$

So it suffices to set $\gamma t = \ln(\delta/\gamma)$. Or $t = \frac{\ln(\delta/\gamma)}{\gamma}$.

How small must we set $\delta$ to ensure that $c_1\lambda_1^t$ dominates all other components and so $\vec{z}^{(t)}$ is very close to $\vec{v}_1$?

---

$\vec{v}_1$: top eigenvector, being computed, $\vec{z}^{(i)}$: iterate at step $i$, converging to $\vec{v}_1$. $\lambda_1, \lambda_2, \ldots \lambda_n$: eigenvalues of $\mathbf{A}$, $\gamma = \frac{|\lambda_1| - |\lambda_2|}{|\lambda_1|}$: eigengap controlling convergence rate

## Random Initialization

**Claim:** When $z^{(0)}$ is chosen with random Gaussian entries, writing $z^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \ldots + c_d\vec{v}_d$, with very high probability, for all $i$:

$$O(1/d^2) \leq |c_i| \leq O(\log d)$$

**Corollary:**

$$\max_j \left| \frac{c_j}{c_1} \right| \leq O(d^2 \log d).$$

$A \in \mathbb{R}^{d \times d}$: input matrix with eigendecomposition $A = V\Lambda V^T$. $\vec{v}_1$: top eigenvector, being computed, $\vec{z}^{(i)}$: iterate at step $i$, converging to $\vec{v}_1$.

# Random Initialization

**Claim 1:** When $z^{(0)}$ is chosen with random Gaussian entries, writing $z^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \ldots + c_d\vec{v}_d$, with very high probability, $\max_j \left| \frac{c_j}{c_1} \right| \leq O(d^2 \log d)$.

**Claim 2:** For gap $\gamma = \frac{|\lambda_1| - |\lambda_2|}{|\lambda_1|}$, and $t = \frac{\ln(1/\delta)}{\gamma}$, $\left| \frac{\lambda_i^t}{\lambda_1^t} \right| \leq \delta$ for all $i$.

$$\vec{z}^{(t)} := \frac{c_1\lambda_1^t\vec{v}_1 + \ldots + c_d\lambda_d^t\vec{v}_d}{\|c_1\lambda_1^t\vec{v}_1 + \ldots + c_d\lambda_d^t\vec{v}_d\|_2} \implies$$

$$\|\vec{z}^{(t)} - \vec{v}_1\|_2 \leq \left\| \frac{c_1\lambda_1^t\vec{v}_1 + \ldots + c_d\lambda_d^t\vec{v}_d}{\|c_1\lambda_1^t\vec{v}_1\|_2} - \vec{v}_1 \right\|_2$$

$$= \left\| \frac{c_2\lambda_2^t}{c_1\lambda_1^t}\vec{v}_2 + \ldots + \frac{c_d\lambda_d^t}{\lambda_1^t}\vec{v}_d \right\|_2 = \left| \frac{c_2\lambda_2^t}{c_1\lambda_1^t} \right| + \ldots + \left| \frac{c_d\lambda_d^t}{\lambda_1^t} \right| \leq \delta \cdot O(d^2 \log d) \cdot d.$$

Setting $\delta = O\left( \frac{\epsilon}{d^3 \log d} \right)$ gives $\|\vec{z}^{(t)} - \vec{v}_1\|_2 \leq \epsilon$.

> $A \in \mathbb{R}^{d \times d}$: input with eigenvalues $\lambda_1 \ldots, \lambda_d$ and eigenvectors $\vec{v}_1, \ldots, \vec{v}_d$. $\vec{z}^{(i)}$: iterate at step $i$. $c_1, \ldots, c_d$: coefficients of $\vec{z}^{(0)}$ in the eigenvector basis.

### Theorem (Basic Power Method Convergence)

*Let $\gamma = \frac{|\lambda_1| - |\lambda_2|}{|\lambda_1|}$ be the relative gap between the first and second eigenvalues. If Power Method is initialized with a random Gaussian vector $\vec{v}^{(0)}$ then, with high probability, after $t = O\left(\frac{\ln(d/\epsilon)}{\gamma}\right)$ steps:*

$$\|\vec{z}^{(t)} - \vec{v}_1\|_2 \leq \epsilon.$$

**Total runtime:** $O(t)$ matrix-vector multiplications. If $A = X^T X$:

$$O\left(\mathsf{nnz}(X) \cdot \frac{\ln(d/\epsilon)}{\gamma} \cdot \right) = O\left(nd \cdot \frac{\ln(d/\epsilon)}{\gamma}\right).$$

How is $\epsilon$ dependence?

How is $\gamma$ dependence?
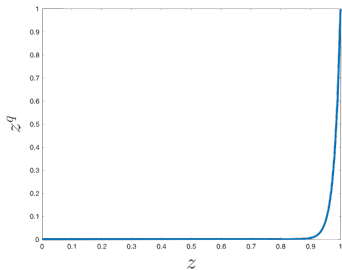
## krylov subspace methods

Krylov subspace methods (Lanczos method, Arnoldi method.)

- How `svds`/`eigs` are actually implemented. Only need $t = O\left(\frac{\ln(d/\epsilon)}{\sqrt{\gamma}}\right)$ steps for the same guarantee.
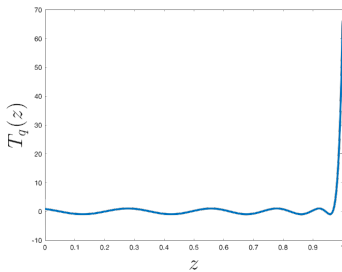
Main Idea: Need to separate $\lambda_1$ from $\lambda_i$ for $i \geq 2$.

- Power method: power up to $\lambda_1^t$ and $\lambda_i^t$.
- Krylov methods: apply a better degree $t$ polynomial $T_t(\cdot)$ to the eigenvalues to separate $T_t(\lambda_1)$ from $T_t(\lambda_i)$.
- Still requires just $t$ matrix vector multiplies. Why?

vs.

## generalizations to larger $k$

- Block Power Method (a.k.a. Simultaneous Iteration, Subspace Iteration, or Orthogonal Iteration)

- Block Krylov methods

$$\text{Runtime: } O\left(ndk \cdot \frac{\ln(d/\epsilon)}{\sqrt{\gamma}}\right)$$

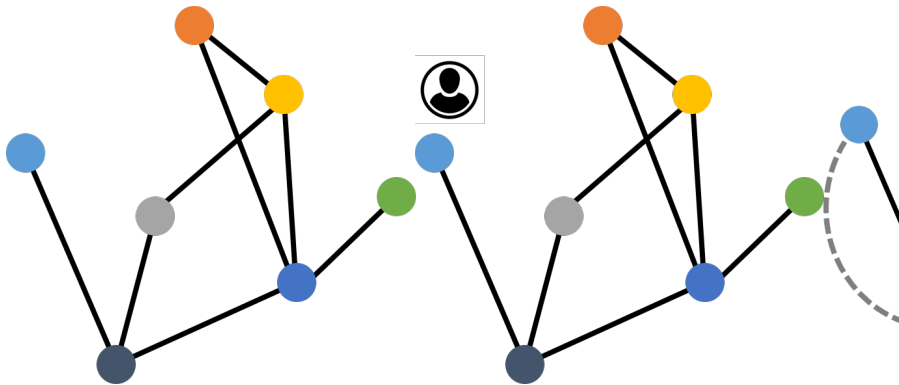to accurately compute the top $k$ singular vectors.

$$\text{'Gapless' Runtime: } O\left(ndk \cdot \frac{\ln(d/\epsilon)}{\sqrt{\epsilon}}\right)$$

if you just want a set of vectors that gives an $\epsilon$-optimal low-rank approximation when you project onto them.

# Connection Between Random Walks, Eigenvectors, and Power Method (Bonus Material)

Consider a random walk on a graph *G* with adjacency matrix **A**.



At each step, move to a random vertex, chosen uniformly at random from the neighbors of the current vertex.

## Connection to Random Walks

Let $\vec{p}^{(t)} \in \mathbb{R}^n$ have $i^{th}$ entry $\vec{p}_i^{(t)} = \text{Pr}$(walk at node i at step t).

- **Initialize:** $\vec{p}^{(0)} = [1, 0, 0, \ldots, 0]$.

- **Update:**

$$\text{Pr}(\text{walk at i at step t}) = \sum_{j \in neigh(i)} \text{Pr}(\text{walk at j at step t-1}) \cdot \frac{1}{degree(j)}$$
$$= \vec{z}^T \vec{p}^{(t-1)}$$

  where $\vec{z}(j) = \frac{1}{degree(j)}$ for all $j \in neigh(i)$, $\vec{z}(j) = 0$ for all $j \notin neigh(i)$.

- $\vec{z}$ is the $i^{th}$ row of the right normalized adjacency matrix $\mathbf{AD}^{-1}$.

- $\vec{p}^{(t)} = \mathbf{AD}^{-1}\vec{p}^{(t-1)} = \underbrace{\mathbf{AD}^{-1}\mathbf{AD}^{-1}\ldots\mathbf{AD}^{-1}}_{t \text{ times}}\vec{p}^{(0)}$

## Random Walking as Power Method

**Claim:** After $t$ steps, the probability that a random walk is at node $i$ is given by the $i^{th}$ entry of

$$\vec{p}^{(t)} = \underbrace{AD^{-1}AD^{-1}\ldots AD^{-1}}_{t \text{ times}}\vec{p}^{(0)}.$$

$$D^{-1/2}\vec{p}^{(t)} = \underbrace{(D^{-1/2}AD^{-1/2})(D^{-1/2}AD^{-1/2})\ldots(D^{-1/2}AD^{-1/2})}_{t \text{ times}}(D^{-1/2}\vec{p}^{(0)}).$$

- $D^{-1/2}\vec{p}^{(t)}$ is exactly what would obtained by applying $t/2$ iterations of power method to $D^{-1/2}\vec{p}^{(0)}$!

- Will converge to the top eigenvector of the normalized adjacency matrix $D^{-1/2}AD^{-1/2}$. Stationary distribution.

- Like the power method, the time a random walk takes to converge to its stationary distribution (mixing time) is dependent on the gap between the top two eigenvalues of $D^{-1/2}AD^{-1/2}$. The spectral gap.