## COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2021.
Lecture 22

- Problem Set 4 due December 1.
- No quiz this week.
- We're going to start on optimization after break. And just cover a bit less material.

Last Class:

- Efficient algorithms for SVD/eigendecomposition.
- Start on iterative methods: intuition behind the power method.

This Class:

- Finish power method analysis.
- Krylov subspace methods.
- Connections to random walks and Markov chains.

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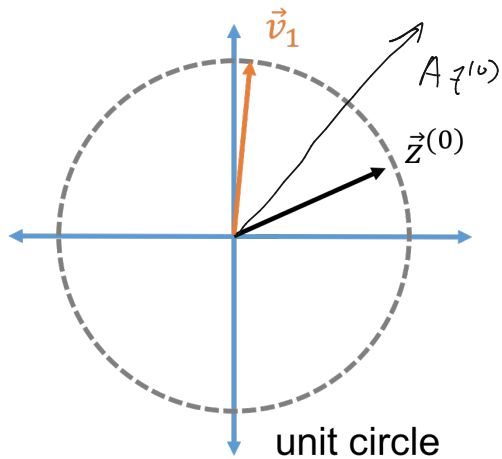
$X^T X$

Goal: Given symmetric $A \in \mathbb{R}^{d \times d}$, with eigendecomposition $A = V \Lambda V^T$, find $\vec{z} \approx \vec{v}_1$ – the top eigenvector of $A$.
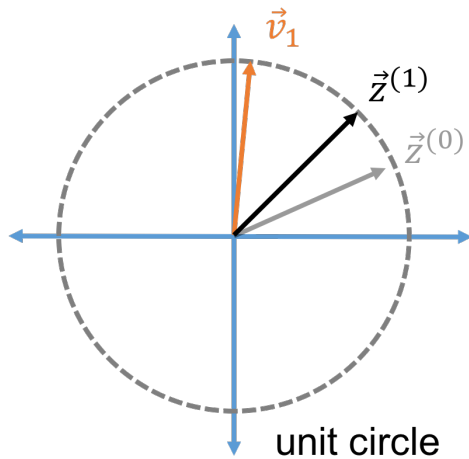
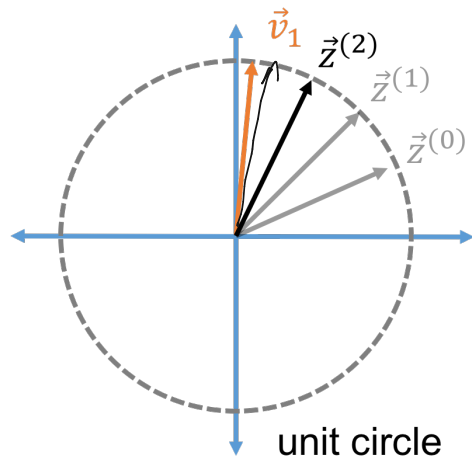- Initialize: Choose $\vec{z}^{(0)}$ randomly. E.g. $\vec{z}^{(0)}(i) \sim \mathcal{N}(0, 1)$.
- For $i = 1, \ldots, t$
  - $\vec{z}^{(i)} := A \cdot \vec{z}^{(i-1)}$
  - $\vec{z}^{(i)} := \frac{\vec{z}^{(i)}}{\|\vec{z}^{(i)}\|_2}$

  Return $\vec{z}^{(i)}$

compute $\vec{z}^{(t)} = A^t \vec{z}^{(0)}$

equivalent code $z^{(t)} = \frac{z^{(t)}}{\|z^{(t)}\|_2}$

3

unit circle

unit circle

unit circle

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.

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

$\lambda_1^t > \lambda_i^t$ for all $i > 1$.



Iteration 0
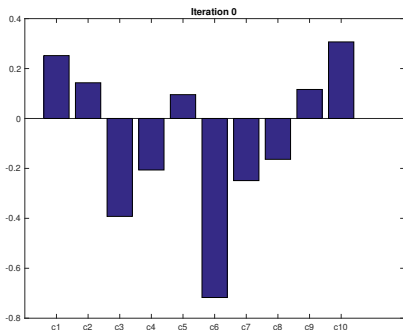
5

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.

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

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.

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



5

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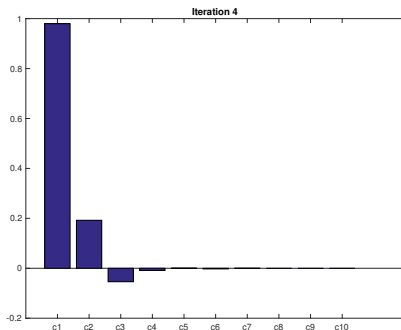
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \ldots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\lambda_1^t\vec{v}_1 + c_2\lambda_2^t\vec{v}_2 + \ldots + c_d\lambda_d^t\vec{v}_d$$



5

## POWER METHOD CONVERGENCE

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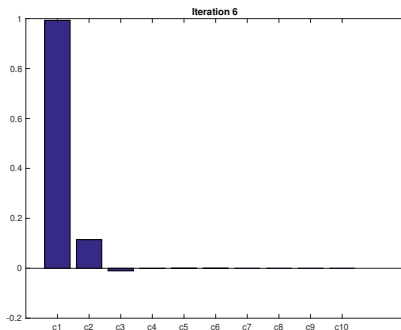
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \ldots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\lambda_1^t\vec{v}_1 + c_2\lambda_2^t\vec{v}_2 + \ldots + c_d\lambda_d^t\vec{v}_d$$

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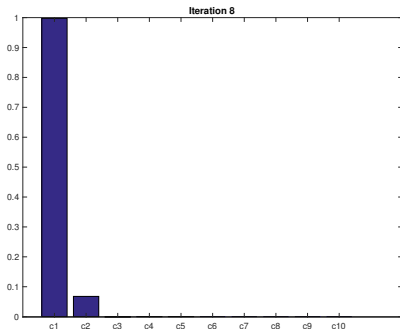
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \ldots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\lambda_1^t\vec{v}_1 + c_2\lambda_2^t\vec{v}_2 + \ldots + c_d\lambda_d^t\vec{v}_d$$



$\lambda_2 \simeq \lambda_1$
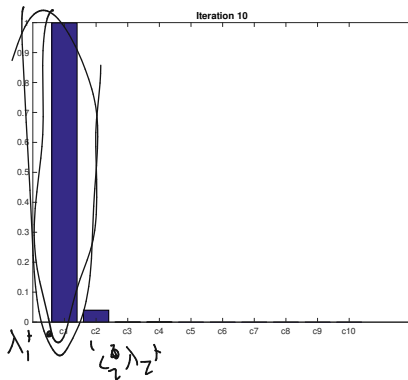
$\lambda_1^t$    $\lambda_2^t$

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

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 $|\lambda_2|^t \leq \frac{1}{e} \cdot |\lambda_1|^t$?

$$\left((1-\gamma)|\lambda_1|\right)^t \leq \frac{1}{e}|\lambda_1|^t$$

$$(1-\gamma)^t < \frac{1}{e} \qquad t > 1/\gamma.$$

$$\lim_{x \to 0} (1-x)^{1/x} = \frac{1}{e}$$

$$\lim_{x \to 0} (1+x)^{1/x} = e$$

$$1-x \approx \frac{1}{1+x}$$

$$.99 \approx \frac{1}{1.01}$$

Example

$$\gamma = 1/3$$

$$\lambda_2 = \frac{2}{3}\lambda_1$$

$$t = 3$$

$$\lambda_2^3 = \frac{8}{27}\lambda_1^3 < \frac{1}{3}\lambda_1^3 \leq \frac{1}{e}\lambda_1^3$$

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

6

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

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 $|\lambda_2|^t \leq \frac{1}{e} \cdot |\lambda_1|^t$? $1/\gamma$.

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

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

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 $|\lambda_2|^t \leq \frac{1}{e} \cdot |\lambda_1|^t$?  **$1/\gamma$.**

How many iterations $t$ does it take to have $|\lambda_2|^t \leq \underline{\delta} \cdot |\lambda_1|^t$?

$$t = O\left(\frac{\ln(1/\delta)}{\gamma}\right)$$

$$\ln_{1-\gamma}(\delta)$$

$$|\lambda_2|^{1/\gamma} \leq \frac{1}{e}|\lambda_1|^{1/\gamma}$$

$$\left(|\lambda_2|^{1/\gamma}\right)^{\ln(1/\delta)} \leq \frac{1}{e^{\ln(1/\delta)}} \cdot \left(|\lambda_1|^{1/\gamma}\right)^{\ln(1/\delta)}$$

$$|\lambda_2|^{\frac{\ln(1/\delta)}{\gamma}} \leq \delta|\lambda_1|^{\frac{\ln(1/\delta)}{\gamma}}$$

---

$\mathbf{A} \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

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

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 $|\lambda_2|^t \leq \frac{1}{e} \cdot |\lambda_1|^t$?  $1/\gamma$.

How many iterations $t$ does it take to have $|\lambda_2|^t \leq \delta \cdot |\lambda_1|^t$?  $\frac{\ln(1/\delta)}{\gamma}$.

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

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

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 $|\lambda_2|^t \leq \frac{1}{e} \cdot |\lambda_1|^t$?  $1/\gamma$.

How many iterations $t$ does it take to have $|\lambda_2|^t \leq \delta \cdot |\lambda_1|^t$?  $\frac{\ln(1/\delta)}{\gamma}$.

Will have for all $i > 1$, $|\lambda_i|^t \leq |\lambda_2|^t \leq \delta \cdot |\lambda_1|^t$.

---

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

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

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 $|\lambda_2|^t \leq \frac{1}{e} \cdot |\lambda_1|^t$?  $1/\gamma$.

How many iterations $t$ does it take to have $|\lambda_2|^t \leq \delta \cdot |\lambda_1|^t$?  $\frac{\ln(1/\delta)}{\gamma}$.

Will have for all $i > 1$, $|\lambda_i|^t \leq |\lambda_2|^t \leq \delta \cdot |\lambda_1|^t$.

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 $\vec{v}_1$?

---

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

**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 $i$:

$c_1 = \langle \vec{v}_1, z^{(0)} \rangle$     $O(1/d^2) \le |c_i| \le O(\log d)$

**Corollary:** $\mathcal{N}(0,1)$

$$\max_j \left| \frac{c_j}{c_1} \right| \le O(d^2 \log d).$$     $\frac{\log d}{1/d^2}$

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

**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 $i$.

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

**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 $i$.

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

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

8

**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 $i$.

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

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



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

**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 $i$.

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

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

$$c_1 \lambda_1^t \vec{v}_1 = \vec{v}_1$$

$$\|c_1 \lambda_1^t\|_2 = c_1 \lambda_1^t$$

$$= \left\| \frac{c_2 \lambda_2^t}{c_1 \lambda_1^t} \vec{v}_2 + \ldots + \frac{c_d \lambda_d^t}{c_1 \lambda_1^t} \vec{v}_d \right\|_2$$

---

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

8

**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 $i$.

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

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

$$= \left\| \frac{c_2\lambda_2^t}{c_1\lambda_1^t}\vec{v}_2 + \ldots + \frac{c_d\lambda_d^t}{\lambda_1^t}\vec{v}_d \right\|_2 = \left| \frac{c_2\lambda_2^t}{c_1\lambda_1^t} \right| + \ldots + \left| \frac{c_d\lambda_d^t}{\lambda_1^t} \right|$$

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

8

**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 $i$.

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

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

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

$< \alpha(d^2 \log d)$

$\leq \delta$

by random    because $t = \frac{\ln(1/\delta)}{\gamma}$

initialization

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

8

**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 $i$.

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

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

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

$$t = O\left(\frac{\log(d/\epsilon)}{\gamma}\right)$$

$$\log(1/\delta) = \log\left(\frac{d^3 \log d}{\epsilon}\right)$$

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

$$= O\left(\log\left(\frac{d/\epsilon}{\epsilon}\right)\right)$$

$A \in \mathbb{R}^{d \times d}$: 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 $\vec{v}_1$.

### 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:*

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

### 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:*

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

with A.

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

$$O\left(\underbrace{nnz(X)} \cdot \underbrace{\frac{\ln(d/\epsilon)}{\gamma}}\right) = O\left(\underbrace{nd} \cdot \underbrace{\frac{\ln(d/\epsilon)}{\gamma}}\right). \qquad O(nd^2)$$

Full SVD

### 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 $A = X^T X$:

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

How is $\epsilon$ dependence? — super good

How is $\gamma$ dependence? — not great

9

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.

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 $i \geq 2$.

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 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 $\underbrace{\lambda_1^t}$ and $\underbrace{\lambda_i^t}$.     $T(x) = x^t$
- Krylov methods: apply a better degree $t$ polynomial $T_t(\cdot)$ to the eigenvalues to separate $T_t(\lambda_1)$ from $T_t(\lambda_i)$.

10

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 $\underline{T_t(\cdot)}$ to the eigenvalues to separate $T_t(\lambda_1)$ from $T_t(\lambda_i)$.
- Still requires just $t$ matrix vector multiplies. Why?

$$\left[\ \underbrace{A^t z^{(0)}}\ \ldots\cdots\ \underbrace{A^2 z^{(0)}}\ \ \underbrace{A z^{(0)}}\ \ z^{(0)} \right.$$

$$\underbrace{T_t(A)}\ \underbrace{z^0} = c_t A^t z^{(0)} + c_{t-1} A^{t-1} z^0 + \ldots c_1 A z^{(0)} + c_0 z^{(0)}$$

10

Optimal 'jump' polynomial in general is given by a degree $t$ Chebyshev polynomial. Krylov methods find a polynomial tuned to the input matrix that does at least as well.

Block Power Method (a.k.a. Simultaneous Iteration,
Subspace Iteration, or Orthogonal Iteration)
Standard Krylov methods (i.e., svds/eigs)
· Block Krylov methods

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

to accurately compute the top $k$ singular vectors.

- Block Power Method (a.k.a. Simultaneous Iteration, Subspace Iteration, or Orthogonal Iteration)
- Standard Krylov methods (i.e., `svds`/`eigs`)
- 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

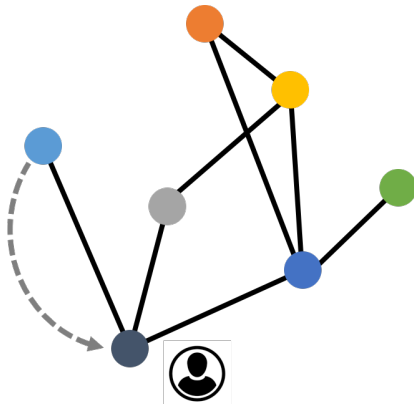Consider a random walk on a graph $G$ with adjacency matrix $\mathbf{A}$.

Consider a random walk on a graph $G$ with adjacency matrix $\mathbf{A}$.
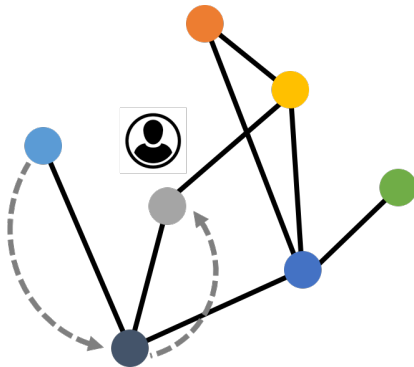


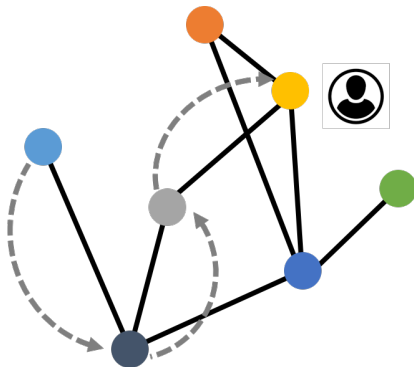At each step, move to a random vertex, chosen uniformly at random from the neighbors of the current vertex.

Consider a random walk on a graph $G$ with adjacency matrix $\mathbf{A}$.

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

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

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

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

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

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

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

- Update:

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

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

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

- Update:

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

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

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

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

- Update:

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

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 $\mathbf{AD}^{-1}$.

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

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

- Update:

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

  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 $\mathbf{AD}^{-1}$.

- $\vec{p}^{(t)} = \mathbf{AD}^{-1}\vec{p}^{(t-1)}$

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

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

· **Update:**

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

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 $\mathbf{AD}^{-1}$.

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

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

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

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

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

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

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 $D^{-1/2}\vec{p}^{(0)}$!

16

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.

**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.

16

A small spectral gap for $D^{-1/2}AD^{-1/2}$ corresponds to a small second smallest eigenvalue for the normalized Laplacian $D^{-1/2}LD^{-1/2}$. Why?

Why does this make sense intuitively given what we know about the second smallest eigenvalue of the Laplacian?