

# COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

---

Cameron Musco

University of Massachusetts Amherst. Fall 2019.

Lecture 17

- Problem Set 3 was released on Saturday. Due next Friday 11/15 at 8pm.
- I will hold office hours after class until 12:30pm today.

## Last Class:

- Finished up spectral clustering and stochastic block model.
- Started discussion of efficient algorithms for SVD/eigendecomposition.

## This Class:

- Finish efficient algorithms for SVD/eigendecomposition.
- Iterative methods: power method, Krylov subspace methods.
- ~~Start optimization unit. Gradient Descent.~~

$$O(nd^2)$$

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  $\underline{V}_k \in \mathbb{R}^{d \times k}$  and then compute  $U_k \Sigma_k = X V_k$ .  

$$X = U \Sigma V^T \quad X V_k = U \Sigma V^T V_k = U \Sigma_k$$
- Use an *iterative algorithm* to compute an *approximation* to the top  $k$  singular vectors  $V_k$ .
- Runtime will be roughly  $O(ndk)$  instead of  $O(nd^2)$ .

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  $\mathbf{X} \in \mathbb{R}^{n \times k}$  for  $k \ll d$ .

- Suffices to compute  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$  and then compute  $\mathbf{U}_k \mathbf{\Sigma}_k = \mathbf{XV}_k$ .
- Use an *iterative algorithm* to compute an *approximation* to the top  $k$  singular vectors  $\mathbf{V}_k$ .
- Runtime will be roughly  $O(ndk)$  instead of  $O(nd^2)$ .

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

**Power Method:** The most fundamental iterative method for approximate SVD. Applies to computing  $k = 1$  singular vectors, but can easily be generalized to larger  $k$ .

**Power Method:** The most fundamental iterative method for approximate SVD. Applies to computing  $k = 1$  singular vectors, but can easily be generalized to larger  $k$ .

**Goal:** Given  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ , find  $\vec{z} \approx \vec{v}_1$ .

**Power Method:** The most fundamental iterative method for approximate SVD. Applies to computing  $k = 1$  singular vectors, but can easily be generalized to larger  $k$ .

**Goal:** Given  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}$ , find  $\vec{z} \approx \vec{v}_1$ .

- **Initialize:** Choose  $\vec{z}^{(0)}$  randomly. E.g.  $\vec{z}^{(0)}(i) \sim \mathcal{N}(0, 1)$ .  $\in \mathbb{R}^d$
- For  $i = 1, \dots, t$ 
  - $\vec{z}^{(i)} = (\mathbf{X}^T \mathbf{X}) \cdot \vec{z}^{(i-1)}$
  - $n_i = \|\vec{z}^{(i)}\|_2$
  - $\vec{z}^{(i)} = \vec{z}^{(i)} / n_i$

Return  $\vec{z}_t$



**Power Method:** The most fundamental iterative method for approximate SVD. Applies to computing  $k = 1$  singular vectors, but can easily be generalized to larger  $k$ .

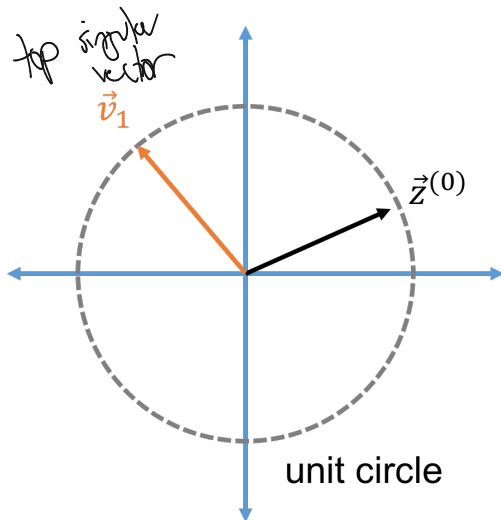
**Goal:** Given  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}$ , find  $\vec{z} \approx \vec{v}_1$ .

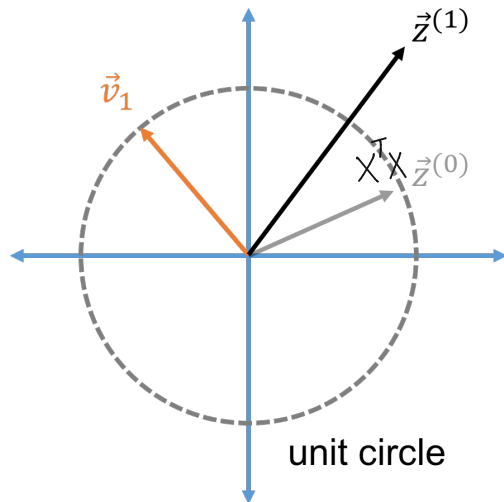
- **Initialize:** Choose  $\vec{z}^{(0)}$  randomly. E.g.  $\vec{z}^{(0)}(i) \sim \mathcal{N}(0, 1)$ .
- For  $i = 1, \dots, t$ 
  - $\vec{z}^{(i)} = (\mathbf{X}^T \mathbf{X}) \cdot \vec{z}^{(i-1)}$  Runtime:  $2 \cdot nd$   $\mathbf{X}^T \mathbf{X} = O(nd^2)$
  - $n_i = \|\vec{z}^{(i)}\|_2$  Runtime:  $d$
  - $\vec{z}^{(i)} = \vec{z}^{(i)} / n_i$  Runtime:  $d$   $O(nd)$

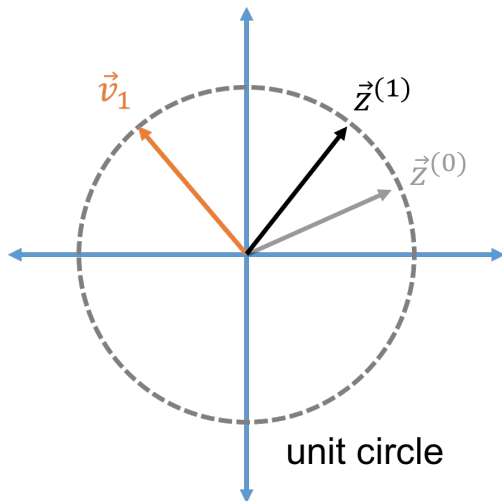
Return  $\vec{z}_t$

Total Runtime:  $O(ndt) \approx O(nd)$

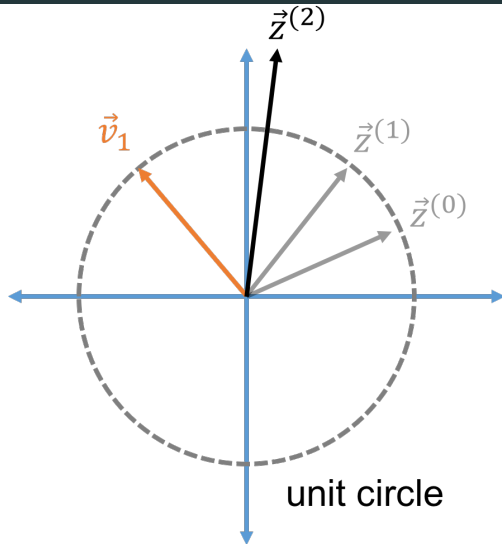
# POWER METHOD

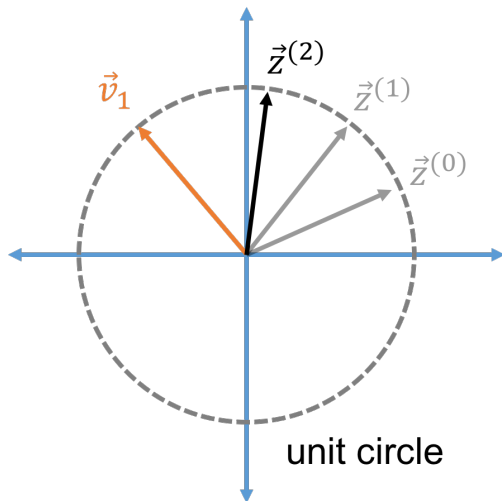




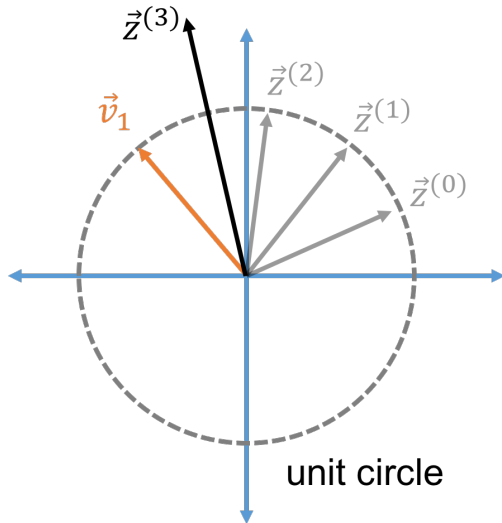


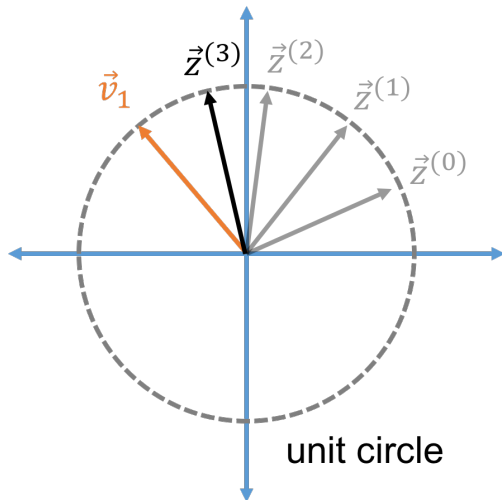
# POWER METHOD



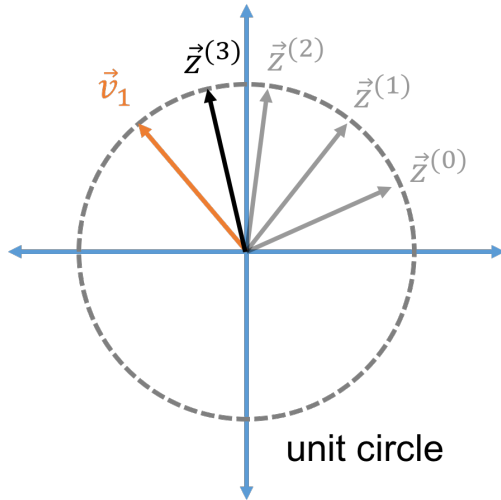


# POWER METHOD









Why is it converging towards  $\vec{v}_1$ ?

*random*

Write  $\vec{z}^{(0)}$  in the right singular vector basis:

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

$X \in \mathbb{R}^{n \times d}$ : input matrix with SVD  $X = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

## POWER METHOD INTUITION

Write  $\vec{z}^{(0)}$  in the right singular vector basis:

$$X = U \Sigma V^T \quad \vec{z}^{(0)} = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_d \vec{v}_d.$$

Update step:  $\vec{z}^{(i)} = \frac{1}{\|X \vec{z}^{(i-1)}\|} X \vec{z}^{(i-1)} = \frac{1}{\|X \vec{z}^{(i-1)}\|} V \Sigma^2 V^T \cdot \vec{z}^{(i-1)}$  (then normalize)

$X \in \mathbb{R}^{n \times d}$ : input matrix with SVD  $X = U \Sigma V^T$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

# POWER METHOD INTUITION

Write  $\vec{z}^{(0)}$  in the right singular vector basis:

$$* \quad \vec{z}^{(0)} = \underbrace{c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_d \vec{v}_d}$$

Update step:  $\vec{z}^{(i)} = \mathbf{X}\mathbf{X}^T \cdot \vec{z}^{(i-1)} = \mathbf{V}\mathbf{\Sigma}^2\mathbf{V}^T \cdot \vec{z}^{(i-1)}$  (then normalize)

$$\begin{aligned} \underline{\mathbf{V}^T \vec{z}^{(0)}} &= \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_d \end{bmatrix} \\ \mathbf{V}^T \vec{z}^{(0)} &= \mathbf{V}^T (c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots) \\ \mathbf{V}^T c_1 \vec{v}_1 &= c_1 \mathbf{V}^T \vec{v}_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \\ &= c_1 \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + c_2 \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \end{aligned}$$

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : input matrix with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

## POWER METHOD INTUITION

Write  $\vec{z}^{(0)}$  in the right singular vector basis:

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

Update step:  $\vec{z}^{(i)} = \mathbf{X}\mathbf{X}^T \cdot \vec{z}^{(i-1)} = \mathbf{V}(\boldsymbol{\Sigma}^2 \mathbf{V}^T \cdot \vec{z}^{(i-1)})$  (then normalize)

$$\mathbf{V}^T \vec{z}^{(0)} = \begin{bmatrix} c_1 \\ \vdots \\ c_d \end{bmatrix}$$

$$\boldsymbol{\Sigma}^2 \mathbf{V}^T \vec{z}^{(0)} = \begin{bmatrix} \sigma_1^2 \cdot c_1 \\ \sigma_2^2 \cdot c_2 \\ \vdots \\ \sigma_d^2 \cdot c_d \end{bmatrix}$$

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : input matrix with SVD  $\mathbf{X} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

# POWER METHOD INTUITION

Write  $\vec{z}^{(0)}$  in the right singular vector basis:

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

Update step:  $\vec{z}^{(i)} = \frac{X^T X}{\|X^T X \vec{z}^{(i-1)}\|} \vec{z}^{(i-1)} = \mathbf{V} \Sigma^2 \mathbf{V}^T \cdot \vec{z}^{(i-1)}$  (then normalize)

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

$$\Sigma^2 \mathbf{V}^T \vec{z}^{(0)} = \begin{bmatrix} \sigma_1^2 c_1 \\ \vdots \\ \sigma_d^2 c_d \end{bmatrix}$$

$$\vec{z}^{(1)} = \mathbf{V} (\Sigma^2 \mathbf{V}^T \cdot \vec{z}^{(0)}) = \mathbf{V} \begin{bmatrix} \sigma_1^2 c_1 \\ \vdots \\ \sigma_d^2 c_d \end{bmatrix} = \sigma_1^2 c_1 \vec{v}_1 + \sigma_2^2 c_2 \vec{v}_2 + \dots + \sigma_d^2 c_d \vec{v}_d$$

$X \in \mathbb{R}^{n \times d}$ : input matrix with SVD  $X = \mathbf{U} \Sigma \mathbf{V}^T$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

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

$$\vec{z}^{(1)} = c_1 \cdot \sigma_1^2 \vec{v}_1 + c_2 \cdot \sigma_2^2 \vec{v}_2 + \dots + c_d \cdot \sigma_d^2 \vec{v}_d.$$

$X \in \mathbb{R}^{n \times d}$ : input matrix with SVD  $X = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

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

$$\vec{z}^{(1)} = c_1 \cdot \sigma_1^2 \vec{v}_1 + c_2 \cdot \sigma_2^2 \vec{v}_2 + \dots + c_d \cdot \sigma_d^2 \vec{v}_d.$$

$$\vec{z}^{(2)} = \mathbf{X}^T \mathbf{X} \vec{z}^{(1)} = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^T \vec{z}^{(1)} = c_1 \sigma_1^4 \vec{v}_1 + \dots + c_d \sigma_d^4 \vec{v}_d$$

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : input matrix with SVD  $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .



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

$$\vec{z}^{(1)} = c_1 \cdot \underline{\sigma_1^2} \vec{v}_1 + c_2 \cdot \underline{\sigma_2^2} \vec{v}_2 + \dots + c_d \cdot \underline{\sigma_d^2} \vec{v}_d.$$

$$\vec{z}^{(2)} = \mathbf{X}^T \mathbf{X} \vec{z}^{(1)} = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^T \vec{z}^{(1)} =$$

Claim 2:

$$\vec{z}^{(t)} = \underbrace{c_1 \cdot \sigma_1^{2t} \vec{v}_1 + c_2 \cdot \sigma_2^{2t} \vec{v}_2 + \dots + c_d \cdot \sigma_d^{2t} \vec{v}_d}_{\|\vec{z}^{(t)}\|}$$

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : input matrix with SVD  $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

## POWER METHOD CONVERGENCE

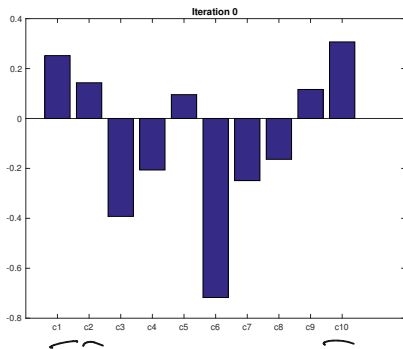
After  $t$  iterations, we have 'powered' up the singular values, making the component in the direction of  $v_1$  much larger, relative to the other components.

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$

# POWER METHOD CONVERGENCE

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

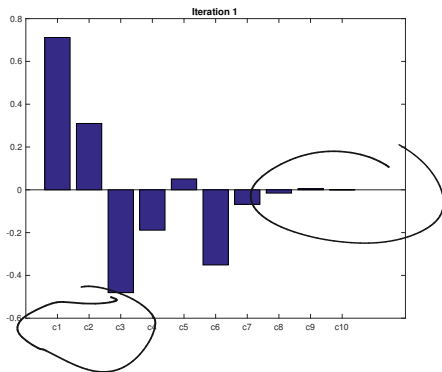
$$\boxed{\vec{z}^{(0)}} = \underbrace{c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_d \vec{v}_d}_{\substack{c_1^+ \\ \underbrace{\hspace{1cm}} \\ c_2^+ \\ \underbrace{\hspace{1cm}} \\ \dots \\ c_d^+}} \implies \vec{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d$$



# POWER METHOD CONVERGENCE

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

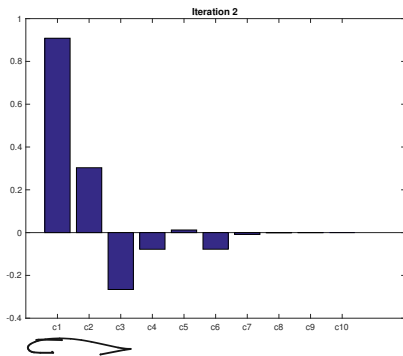
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



## POWER METHOD CONVERGENCE

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

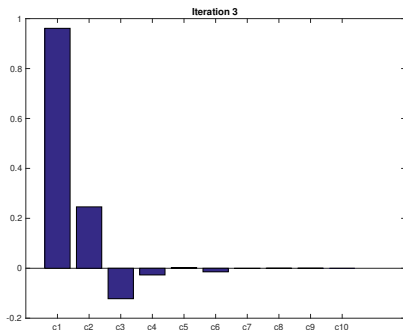
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



## POWER METHOD CONVERGENCE

After  $t$  iterations, we have ‘powered’ up the singular values, making the component in the direction of  $v_1$  much larger, relative to the other components.

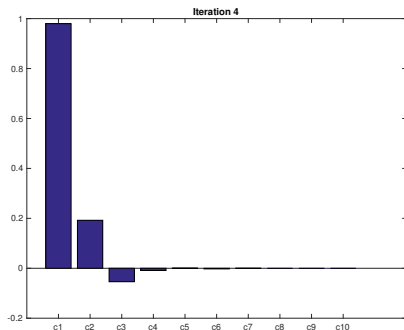
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



## POWER METHOD CONVERGENCE

After  $t$  iterations, we have ‘powered’ up the singular values, making the component in the direction of  $v_1$  much larger, relative to the other components.

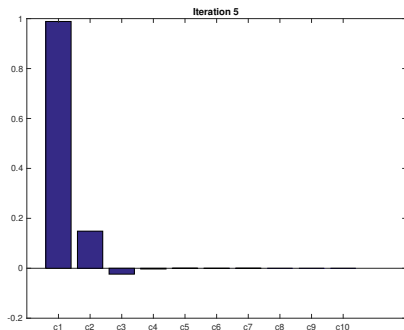
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



## POWER METHOD CONVERGENCE

After  $t$  iterations, we have ‘powered’ up the singular values, making the component in the direction of  $v_1$  much larger, relative to the other components.

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$

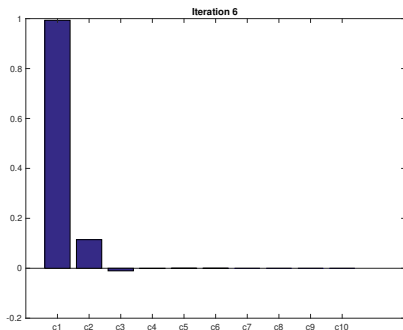




## POWER METHOD CONVERGENCE

After  $t$  iterations, we have ‘powered’ up the singular values, making the component in the direction of  $v_1$  much larger, relative to the other components.

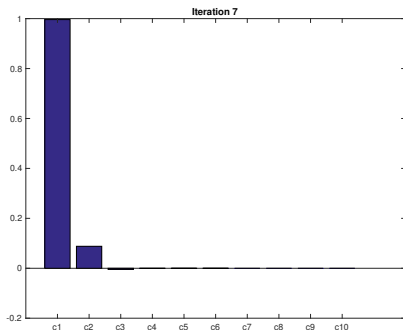
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



## POWER METHOD CONVERGENCE

After  $t$  iterations, we have ‘powered’ up the singular values, making the component in the direction of  $v_1$  much larger, relative to the other components.

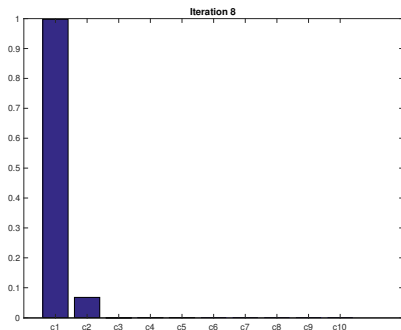
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



## POWER METHOD CONVERGENCE

After  $t$  iterations, we have 'powered' up the singular values, making the component in the direction of  $\vec{v}_1$  much larger, relative to the other components.

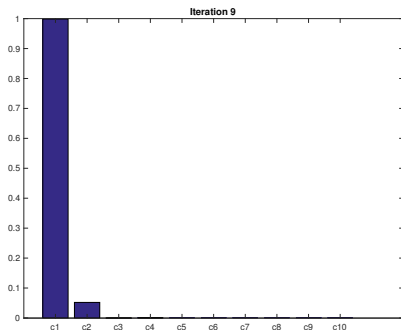
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



## POWER METHOD CONVERGENCE

After  $t$  iterations, we have ‘powered’ up the singular values, making the component in the direction of  $v_1$  much larger, relative to the other components.

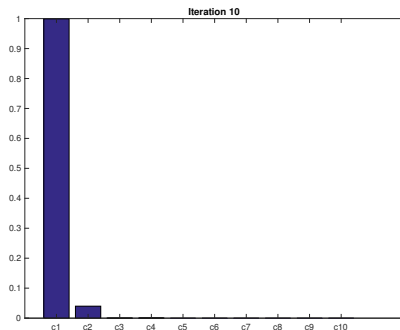
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



## POWER METHOD CONVERGENCE

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

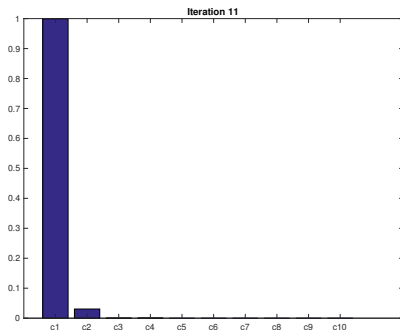
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



## POWER METHOD CONVERGENCE

After  $t$  iterations, we have ‘powered’ up the singular values, making the component in the direction of  $\vec{v}_1$  much larger, relative to the other components.

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



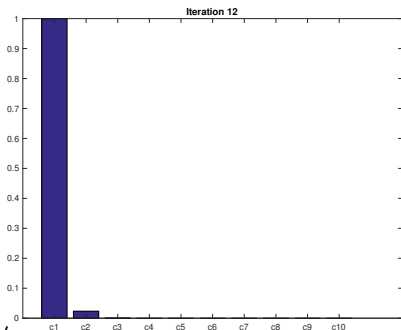
## POWER METHOD CONVERGENCE

After  $t$  iterations, we have 'powered' up the singular values, making the component in the direction of  $v_1$  much larger, relative to the other components.  $c_i$  positive or negative

$$\vec{z}^{(0)} = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_d \vec{v}_d \implies \vec{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d$$

$$\vec{z}^{(t)} \approx \vec{v}_1 \cdot \frac{c_1 \sigma_1^{2t}}{\sigma_1^{2t}}$$

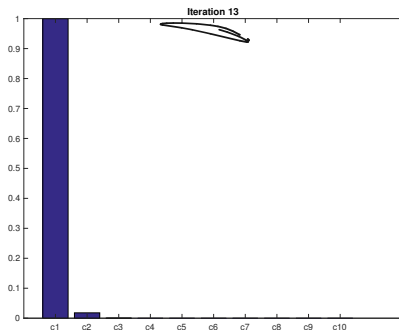
$$\vec{z}^{(0)} = 1 \cdot \vec{v}_1 + 0 \cdot \vec{v}_2 + \dots + 0 \cdot \vec{v}_d$$



# POWER METHOD CONVERGENCE

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

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



When will convergence be slow?



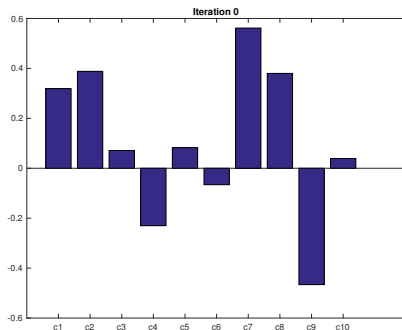
Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$

# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

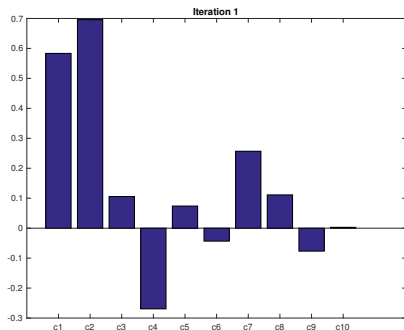
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

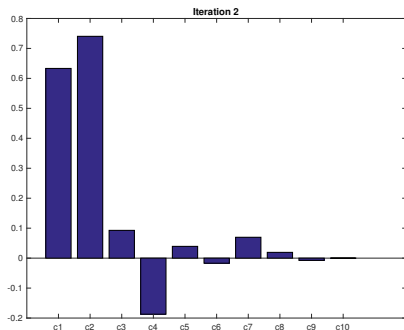
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

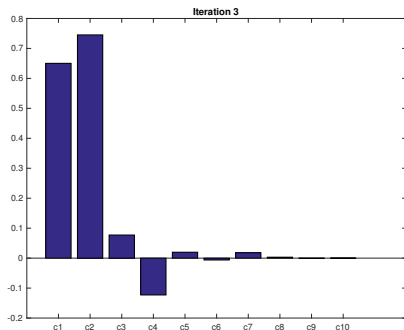
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

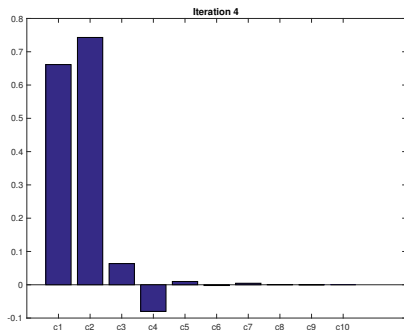
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

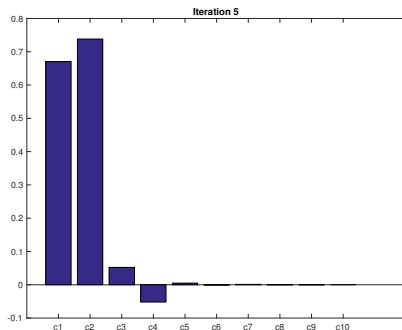
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

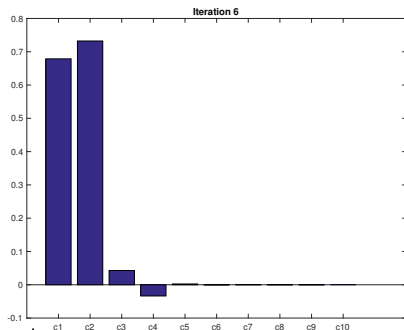
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

$$\vec{z}^{(0)} = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_d \vec{v}_d \implies \vec{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d$$



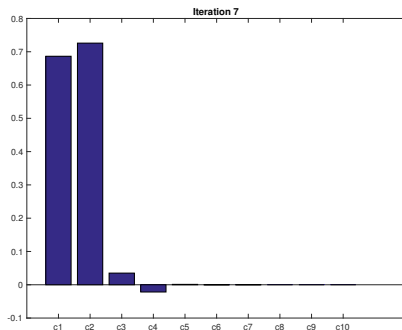
$c_i \vec{z}^{(t)}$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

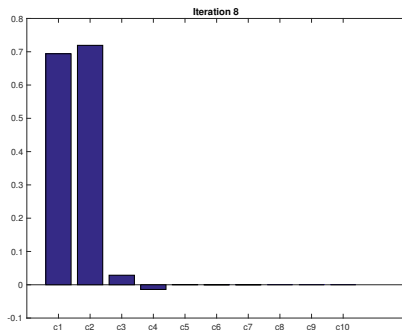
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

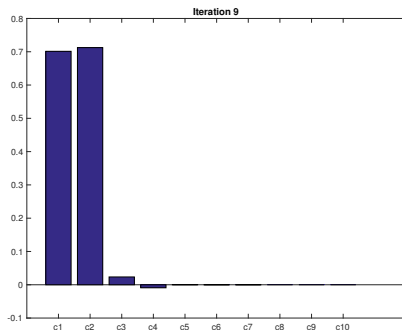
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

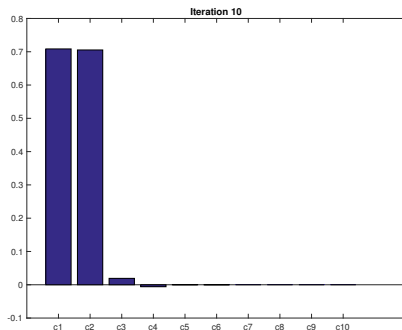
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

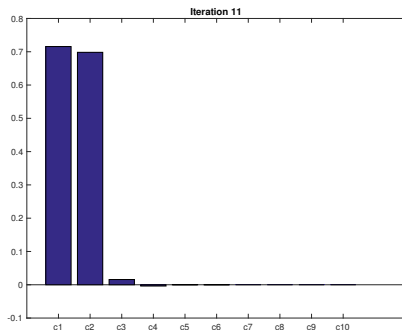
$$\vec{z}^{(0)} = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_d \vec{v}_d \implies \vec{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

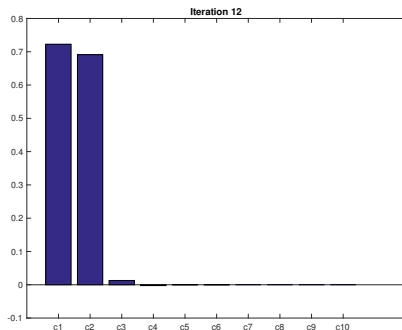
$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$



# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$

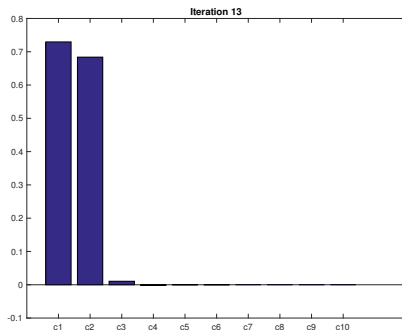


# POWER METHOD SLOW CONVERGENCE

Slow Case:  $X$  has singular values:  $\sigma_1 = 1, \sigma_2 = .99, \sigma_3 = .9, \sigma_4 = .8, \dots$

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$

$z^t \neq v_1$



# POWER METHOD CONVERGENCE RATE

$\gamma > 0$

"Rule of 70"

$$\vec{z}^{(0)} = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_d \vec{v}_d \implies \vec{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d$$

Write  $\sigma_2 = (1 - \gamma) \sigma_1$  for 'gap'  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ . How many iterations  $t$  does it take to have  $\sigma_2^{2t} \leq \frac{1}{2} \cdot \sigma_1^{2t}$ ?

$$\sigma_2^{2t} \leq \frac{1}{2} \sigma_1^{2t} \rightarrow \left(\frac{\sigma_2}{\sigma_1}\right)^{2t} \leq \frac{1}{2}$$

$$(1 - \gamma)^{2t} \leq \frac{1}{2}$$

$$2t = \log_{1-\gamma} (1/2)$$

$$(1 - \gamma)^{1/\gamma} = 1/e$$

$$t = O(1/\gamma) \rightarrow (1 - \gamma)^{2t} \leq \frac{1}{2}$$

$X \in \mathbb{R}^{n \times d}$ : matrix with SVD  $X = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .



## POWER METHOD CONVERGENCE RATE

$$\vec{z}^{(0)} = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_d \vec{v}_d \implies \vec{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d$$

Write  $\sigma_2 = (1 - \gamma)\sigma_1$  for 'gap'  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ . How many iterations  $t$  does it take to have  $\sigma_2^{2t} \leq \frac{1}{2} \cdot \sigma_1^{2t}$ ?  $O(1/\gamma)$ .  $O(\log_{1-\gamma} 1/2)$

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : matrix with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

## POWER METHOD CONVERGENCE RATE

$$\vec{z}^{(0)} = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_d \vec{v}_d \implies \vec{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d$$

Write  $\sigma_2 = (1 - \gamma)\sigma_1$  for 'gap'  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ . How many iterations  $t$  does it take to have  $\sigma_2^{2t} \leq \frac{1}{2} \cdot \sigma_1^{2t}$ ?  $O(1/\gamma)$ .

How many iterations  $t$  does it take to have  $\sigma_2^{2t} \leq \delta \cdot \sigma_1^{2t}$ ?  $\log_{1-\gamma} \delta$

$$O(1/\gamma)$$

$$\sigma_2^{2t} \leq \frac{1}{2} \sigma_1^{2t}$$

$$\sigma_2^{2t \cdot m} \leq \frac{1}{2^m} \cdot \sigma_1^{2t} \quad m = \log(1/\delta)$$

$$O(1/\gamma) \cdot O(\log(1/\delta)) = \frac{\log(1/\delta)}{\gamma}$$

$X \in \mathbb{R}^{n \times d}$ : matrix with SVD  $X = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

## POWER METHOD CONVERGENCE RATE

$$\vec{z}^{(0)} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_d\vec{v}_d \implies \vec{z}^{(t)} = c_1\sigma_1^{2t}\vec{v}_1 + c_2\sigma_2^{2t}\vec{v}_2 + \dots + c_d\sigma_d^{2t}\vec{v}_d$$

Write  $\sigma_2 = (1 - \gamma)\sigma_1$  for 'gap'  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ . How many iterations  $t$  does it take to have  $\sigma_2^{2t} \leq \frac{1}{2} \cdot \sigma_1^{2t}$ ?  $\mathcal{O}(1/\gamma)$ .

How many iterations  $t$  does it take to have  $\sigma_2^{2t} \leq \delta \cdot \sigma_1^{2t}$ ?  $\mathcal{O}\left(\frac{\log(1/\delta)}{\gamma}\right)$ .

$X \in \mathbb{R}^{n \times d}$ : matrix with SVD  $X = \mathbf{U}\Sigma\mathbf{V}^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

## POWER METHOD CONVERGENCE RATE

$$\vec{z}^{(0)} = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_d \vec{v}_d \implies \vec{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d$$

Write  $\sigma_2 = (1 - \gamma)\sigma_1$  for 'gap'  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ . How many iterations  $t$  does it take to have  $\sigma_2^{2t} \leq \frac{1}{2} \cdot \sigma_1^{2t}$ ?  $\mathcal{O}(1/\gamma)$ .

How many iterations  $t$  does it take to have  $\sigma_2^{2t} \leq \delta \cdot \sigma_1^{2t}$ ?  $\mathcal{O}\left(\frac{\log(1/\delta)}{\gamma}\right)$ .

How small must we set  $\delta$  to ensure that  $c_1 \sigma_1^{2t}$  dominates all other components and so  $\vec{z}^{(t)}$  is very close to  $\vec{v}_1$ ?

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : matrix with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

# RANDOM INITIALIZATION

**Claim:** When  $z^{(0)}$  is chosen with random Gaussian entries, writing  $z^{(0)} = c_1 v_1 + c_2 v_2 + \dots + c_d v_d$ , with very high probability, for all  $i$ :

$$\left[ \frac{O(1/d^2)}{\sqrt{\log d}} \leq c_i \leq O(\sqrt{\log d}) \right]$$

How is  $c_i$  distributed?

$$V^T z^{(0)} = \begin{bmatrix} c_1 \\ \vdots \\ c_d \end{bmatrix}$$

$$c_i = \langle v_i, z^{(0)} \rangle = \sum_j v_i(j) \cdot z^{(0)}(j)$$

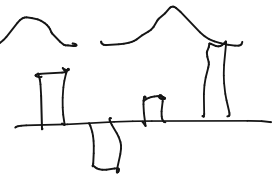
$\swarrow$   $\searrow$   
 $N(0,1)$

$$N(0, v_i(j)^2)$$

$$c_i \sim N(0, 1)$$

"rotation invariance of Gaussian"

$$c_i \sim N(0, \sum_j v_i(j)^2) = N(0, \|v_i\|^2)$$



$$e^{-x^2} \quad e^{-\frac{x^2}{2}} = \frac{1}{\sqrt{2}}$$

$X \in \mathbb{R}^{n \times d}$ : matrix with SVD  $X = U \Sigma V^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_i$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_i$ .

## RANDOM INITIALIZATION

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

$$\log d \geq 1 \quad \underline{O(1/d^2)} \leq |c_i| \leq \underline{O(\log d)}$$

Corollary:

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

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : matrix with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_1$ : top right singular vector, 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  $z^{(0)} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \dots + c_d \mathbf{v}_d$ , with very high probability,  $\max_j \frac{c_j}{c_1} \leq O(d^2 \log d)$ .

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : matrix with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_1$ : top right singular vector, 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  $z^{(0)} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \dots + c_d \mathbf{v}_d$ , with very high probability,  $\max_j \frac{c_j}{c_1} \leq O(d^2 \log d)$ .

**Claim 2:** For gap  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ , after  $t = O\left(\frac{\log(1/\delta)}{\gamma}\right)$  iterations:

$$\bar{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d \propto \underbrace{c_1 \vec{v}_1} + \underbrace{c_2 \delta \vec{v}_2} + \dots + \underbrace{c_d \delta \vec{v}_d}$$

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : matrix with SVD  $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\bar{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .



## RANDOM INITIALIZATION

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

**Claim 2:** For gap  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ , after  $t = O\left(\frac{\log(1/\delta)}{\gamma}\right)$  iterations:

$$\vec{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d \propto c_1 \vec{v}_1 + c_2 \delta \vec{v}_2 + \dots + c_d \delta \vec{v}_d$$

If we set  $\delta = O\left(\frac{\epsilon}{d^2 \log d}\right)$  by Claim 1 will have:

$$\frac{\epsilon}{\sqrt{d^2 \log d}} \cdot d^2 \log d = \frac{\epsilon}{d}$$

$$\vec{z}^{(t)} \propto \vec{v}_1 + \frac{\epsilon}{d} (\vec{v}_2 + \dots + \vec{v}_d).$$

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : matrix with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

# RANDOM INITIALIZATION

$$c_1 = 0$$

**Claim 1:** When  $z^{(0)}$  is chosen with random Gaussian entries, writing  $z^{(0)} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \dots + c_d \mathbf{v}_d$ , with very high probability,

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

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

**Claim 2:** For gap  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$ , after  $t = O\left(\frac{\log(1/\delta)}{\gamma}\right)$  iterations:

$$\vec{z}^{(t)} = c_1 \sigma_1^{2t} \vec{v}_1 + c_2 \sigma_2^{2t} \vec{v}_2 + \dots + c_d \sigma_d^{2t} \vec{v}_d \propto c_1 \vec{v}_1 + c_2 \delta \vec{v}_2 + \dots + c_d \delta \vec{v}_d$$

If we set  $\delta = O\left(\frac{\epsilon}{d^3 \log d}\right)$  by Claim 1 will have:

$$\vec{z}^{(t)} \propto \vec{v}_1 + \frac{\epsilon}{d} (\vec{v}_2 + \dots + \vec{v}_d)$$

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

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : matrix with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Singular values  $\sigma_1, \sigma_2, \dots, \sigma_d$ .  $\vec{v}_1$ : top right singular vector, being computed,  $\vec{z}^{(i)}$ : iterate at step  $i$ , converging to  $\vec{v}_1$ .

## Theorem (Basic Power Method Convergence)

Let  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$  be the relative gap between the first and second largest singular values. If Power Method is initialized with a random Gaussian vector  $\vec{v}^{(0)}$  then, with high probability, after  $t = O\left(\frac{\log d/\epsilon}{\gamma}\right)$  steps:

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

### Theorem (Basic Power Method Convergence)

Let  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$  be the relative gap between the first and second largest singular values. If Power Method is initialized with a random Gaussian vector  $\vec{v}^{(0)}$  then, with high probability, after  $t = O\left(\frac{\log d/\epsilon}{\gamma}\right)$  steps:

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

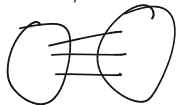
**Total runtime:**  $O(t)$  matrix-vector multiplications.

$$O\left(\underline{\text{nnz}(\mathbf{X})} \cdot \frac{\log(d/\epsilon)}{\gamma}\right) = O\left(nd \cdot \frac{\log(d/\epsilon)}{\gamma}\right) \cdot O(n^2)$$

# POWER METHOD THEOREM

## Theorem (Basic Power Method Convergence)

Let  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$  be the relative gap between the first and second largest singular values. If Power Method is initialized with a random Gaussian vector  $\vec{v}^{(0)}$  then, with high probability, after  $t = O\left(\frac{\log d/\epsilon}{\gamma}\right)$  steps:



$$\|\vec{z}^{(t)} - \vec{v}_1\|_2 \leq \epsilon.$$
$$\|\vec{z}^{(i)} - \vec{z}^{(i-1)}\|$$

Total runtime:  $O(t)$  matrix-vector multiplications.

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

$$\epsilon \rightarrow \frac{\epsilon}{2}$$

$$\log(1/\epsilon) \rightarrow \log(1/\epsilon) + \log(2)$$

How is  $\epsilon$  dependence?

"linearly convergent"

How is  $\gamma$  dependence?

$$\sigma_1 > \sigma_2$$

Krylov subspace methods (Lanczos method, Arnoldi method.)

- How `svds/eigs` are actually implemented. Only need  $t = O\left(\frac{\log d/\epsilon}{\sqrt{\gamma}}\right)$  steps for the same guarantee.

Krylov subspace methods (Lanczos method, Arnoldi method.)

- How **svds/eigs** are actually implemented. Only need  $t = O\left(\frac{\log d/\epsilon}{\sqrt{\gamma}}\right)$  steps for the same guarantee.

**Main Idea:** Need to separate  $\sigma_1$  from  $\sigma_i$  for  $i \geq 2$ .

Krylov subspace methods (Lanczos method, Arnoldi method.)

- How **svds/eigs** are actually implemented. Only need  $t = O\left(\frac{\log d/\epsilon}{\sqrt{\gamma}}\right)$  steps for the same guarantee.

**Main Idea:** Need to separate  $\sigma_1$  from  $\sigma_i$  for  $i \geq 2$ .

- Power method: power up to  $\sigma_1^{2 \cdot t}$  and  $\sigma_i^{2 \cdot t}$ .



## Krylov subspace methods (Lanczos method, Arnoldi method.)

- How **svds/eigs** are actually implemented. Only need  $t = O\left(\frac{\log d/\epsilon}{\sqrt{\gamma}}\right)$  steps for the same guarantee.

**Main Idea:** Need to separate  $\sigma_1$  from  $\sigma_i$  for  $i \geq 2$ .

- Power method: power up to  $\sigma_1^{2 \cdot t}$  and  $\sigma_i^{2 \cdot t}$ .  $\rho(x) = x^t \rightarrow \rho(\sigma_1^2) \gg \rho(\sigma_i^2)$
- Krylov methods: apply a better degree  $t$  polynomial  $T_t(\sigma_1^2)$  and  $T_t(\sigma_i^2)$ .

## Krylov subspace methods (Lanczos method, Arnoldi method.)

- How `svds/eigs` are actually implemented. Only need  $t = O\left(\frac{\log d/\epsilon}{\sqrt{\gamma}}\right)$  steps for the same guarantee.

**Main Idea:** Need to separate  $\sigma_1$  from  $\sigma_i$  for  $i \geq 2$ .

- Power method: power up to  $\sigma_1^{2 \cdot t}$  and  $\sigma_i^{2 \cdot t}$ .
- Krylov methods: apply a **better** degree  $t$  polynomial  $T_t(\sigma_1^2)$  and  $T_t(\sigma_i^2)$ .
- Still requires just  $2t$  matrix vector multiplies. **Why?**

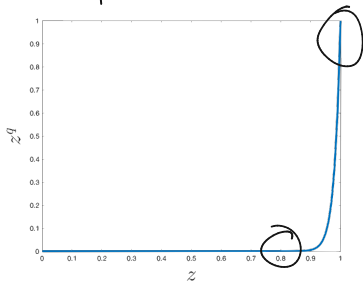
$$X^T X v^0 \quad (X^T X) v_1 \quad (X^T X)^2 v_2$$

$$X^T X v^0 \quad \underline{(X^T X)^2 v_0} \quad (X^T X)^3 v_0$$

$$\left. \begin{aligned} & a_1 X^T X + a_2 (X^T X)^2 + \dots \end{aligned} \right) v_0$$

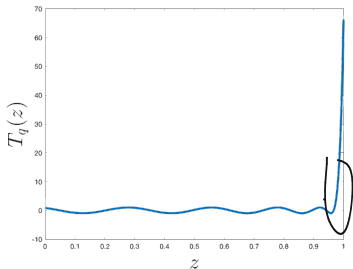
$$\begin{matrix} X^T X X^T X v_0 \\ \text{only} \end{matrix}$$

$$P(x) = x^t$$



$$\frac{1}{\sqrt{y}} \approx \frac{1}{\sqrt{y}}$$

rather bad



vs.

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 aka Simultaneous Iteration aka Subspace Iteration aka Orthogonal Iteration
- Block Krylov methods

**Runtime:**  $O\left(\underbrace{ndk}_{\substack{\hookrightarrow \sigma_k \\ \sigma_{k+1}}} \cdot \frac{\log d/\epsilon}{\sqrt{\gamma}}\right)$

to accurately compute the top  $k$  singular vectors.

- Block Power Method aka Simultaneous Iteration aka Subspace Iteration aka Orthogonal Iteration
- Block Krylov methods

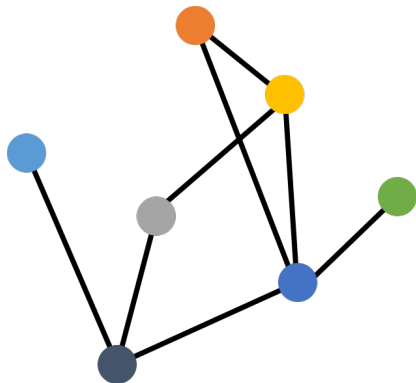
$$\text{Runtime: } O\left(ndk \cdot \frac{\log d/\epsilon}{\sqrt{\gamma}}\right)$$

to accurately compute the top  $k$  singular vectors.

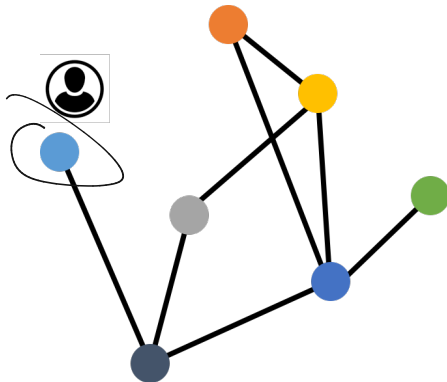
$$\text{'Gapless' Runtime: } O\left(ndk \cdot \frac{\log d/\epsilon}{\sqrt{\epsilon}}\right)$$

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

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

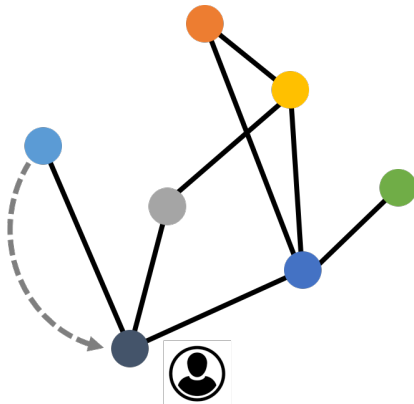


Consider a random walk on a graph  $G$  with adjacency matrix  $A$ .



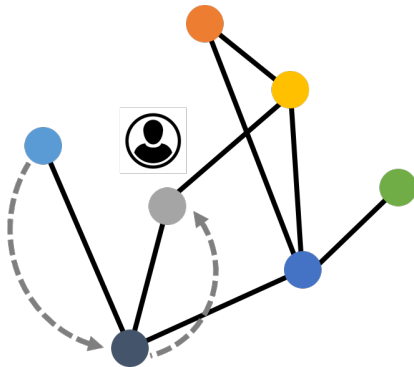
At each step, move to a random vertex, chosen uniformly at random from the neighbors of the current vertex.

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  $i^{\text{th}}$  entry  $\vec{p}_i^{(t)} = \Pr(\text{walk at node } i \text{ at step } t)$ .

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

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

## CONNECTION TO RANDOM WALKS

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

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

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

$$z_j \vec{p}_i^{(t-1)}$$

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

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

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

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

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

## CONNECTION TO RANDOM WALKS

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

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

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

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

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

$$\begin{aligned} \mathcal{P}^{(t)} &= \mathbf{AD}^{-1} \mathcal{P}^{(t-1)} \\ \mathcal{P}^{(t)} &= \underbrace{\mathbf{AD}^{-1} \mathbf{AD}^{-1} \dots}_{t \text{ times}} \mathcal{P}^{(0)} \end{aligned}$$

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

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

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

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

- $\vec{z}$  is just the  $i^{\text{th}}$  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  $i^{\text{th}}$  entry  $\vec{p}_i^{(t)} = \Pr(\text{walk at node } i \text{ at step } t)$ .

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

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

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

- $\vec{z}$  is just the  $i^{\text{th}}$  row of the right normalized adjacency matrix  $\mathbf{AD}^{-1}$ .
- $\vec{p}^{(t)} = \mathbf{AD}^{-1} \vec{p}^{(t-1)} = \underbrace{\mathbf{AD}^{-1} \mathbf{AD}^{-1} \dots \mathbf{AD}^{-1}}_{t \text{ times}} \vec{p}^{(0)}$

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

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

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

$$\vec{p}^{(t)} = \underbrace{\mathbf{A}\mathbf{D}^{-1}\mathbf{A}\mathbf{D}^{-1} \dots \mathbf{A}\mathbf{D}^{-1}}_{t \text{ times}} \vec{p}^{(0)}.$$

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

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

$$\vec{p}^{(t)} = \underbrace{\mathbf{A}\mathbf{D}^{-1}\mathbf{A}\mathbf{D}^{-1} \dots \mathbf{A}\mathbf{D}^{-1}}_{t \text{ times}} \vec{p}^{(0)}.$$

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

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

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

$$\vec{p}^{(t)} = \underbrace{\mathbf{A}\mathbf{D}^{-1}\mathbf{A}\mathbf{D}^{-1} \dots \mathbf{A}\mathbf{D}^{-1}}_{t \text{ times}} \vec{p}^{(0)}.$$

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

- $\mathbf{D}^{-1/2} \vec{p}^{(t)}$  is exactly what would be obtained by applying  $t/2$  iterations of power method to  $\mathbf{D}^{-1/2} \vec{p}^{(0)}$ !
- Will converge to the top singular vector (eigenvector) of the normalized adjacency matrix  $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ . **Stationary distribution.**

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

$$\vec{p}^{(t)} = \underbrace{\mathbf{A}\mathbf{D}^{-1}\mathbf{A}\mathbf{D}^{-1} \dots \mathbf{A}\mathbf{D}^{-1}}_{t \text{ times}} \vec{p}^{(0)}.$$

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

- $\mathbf{D}^{-1/2} \vec{p}^{(t)}$  is exactly what would be obtained by applying  $t/2$  iterations of power method to  $\mathbf{D}^{-1/2} \vec{p}^{(0)}$ !
- Will converge to the top singular vector (eigenvector) of the normalized adjacency matrix  $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ . **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  $\mathbf{A} \mathbf{D}^{-1}$ . The **spectral gap.**