

# Graph 3. Matrix Factorization via

---

## Eigendecomposition (ED) and Singular Value Decomposition (SVD)

2/27/2024

@Yiming Yang, S24 Lecture on ED  
& SVD

1

1

## Outline

---

- Motivation with PCA
- ED and SVD
- Convergence of HITS

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

2

2

## Why do we care about ED and SVD?

- We want to know why HITS converges and where does it converge.
- We want to understand how to visualize word embeddings in a 2D or 3D space (via PCA projection of high-dimensional vectors).
- We want to understand how to use a graph to propagate signals over nodes smoothly (later lectures on Laplacian Eigenmaps & Graph Convolution Networks)

2/27/2024

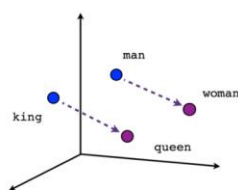
@Yiming Yang, S24 Lecture on ED &amp; SVD

3

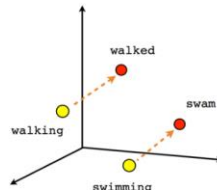
3

## Visualizing Word Embeddings in 2D or 3D

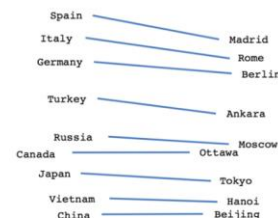
- Allowing us to see some interesting (analogical) patterns



Male-Female



Verb tense



Country-Capital

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

4

4

## Consider a word embedding matrix

$$X_{n \times m} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1m} \\ x_{21} & x_{22} & \cdots & x_{2m} \\ \cdots & \cdots & x_{ij} & \cdots \\ x_{n1} & x_{n2} & \cdots & x_{nm} \end{bmatrix} \quad \begin{array}{l} \text{n words} \\ \text{m features (factors)} \end{array}$$

- Each **row** is the embedding of a word.
- Each **column** is a latent feature of words.
- Each **cell** is the feature weight of the corresponding word.

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

5

5

## Principal Component Analysis (PCA)

[https://en.wikipedia.org/wiki/Principal\\_component\\_analysis](https://en.wikipedia.org/wiki/Principal_component_analysis)

- Projecting high-dimensional data (row vectors in  $X$ ) to  $k$ -dimensional

$$T_k = XV_k \quad (k < m)$$

$T_k \in \mathbb{R}^{n \times k}, \quad X \in \mathbb{R}^{n \times m}, \quad V_k \in \mathbb{R}^{m \times k},$

$V = (v_1, v_2, \dots, v_m)$  are the **eigenvectors** of  $X^T X$ .

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

6

6

## Steps and Intuition in PCA

[https://en.wikipedia.org/wiki/Principal\\_component\\_analysis](https://en.wikipedia.org/wiki/Principal_component_analysis)

- **Preprocessing:** Make each row of input matrix  $X^{(0)}$  to the centroid of all the row vectors

$$\mu := \frac{1}{n} \sum_{i=1}^n X_{i,:}^{(0)}, \quad X := X^{(0)} - \hat{1}\mu$$

- **Obtaining the top-k eigenvectors of matrix  $X^T X$**

$$V_k = (v_1, \dots, v_k), \quad k < m$$

- **Projecting** the row vectors of  $X$  onto the top-k eigenvectors

$$T_k = X V_k$$

- **Intuition:** Preserving most **variance in data with**

$$\text{Var}(v_1; X) \geq \text{var}(v_2; X) \geq \dots \geq \text{var}(v_m; X)$$

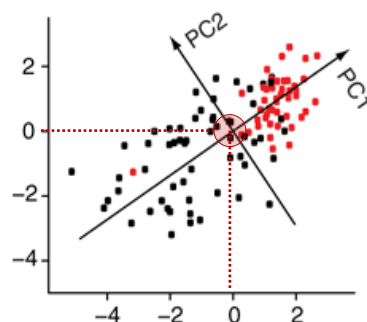
2/27/2024

@Yiming Yang, S24 Lecture on ED & SVD

7

7

## PCA of 2D Data



By rotating the orthogonal axes, we can see

- The 1<sup>st</sup> eigenvector (PC1= $v_1$ ) identifies the direction with the maximum variance in data.
- The 2<sup>nd</sup> eigenvector (PC2= $v_2$ ) identifies the direction with the secondary variance in data.

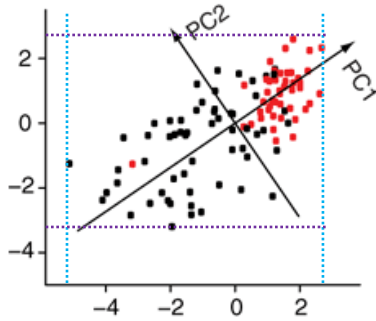
2/27/2024

@Yiming Yang, S24 Lecture on ED & SVD

8

8

## Variance in data on the original basis ( $X_1, X_2$ )



- Sample variance in direction  $X_1$

$$\text{var}(X_1) = \frac{1}{(n-1)} \sum_{i=1}^n x_{i1}^2$$

- Sample variance in direction  $X_2$

$$\text{var}(X_2) = \frac{1}{(n-1)} \sum_{i=1}^n x_{i2}^2$$

- Total variance in data

$$\text{var}(X) = \text{var}(X_1) + \text{var}(X_2)$$

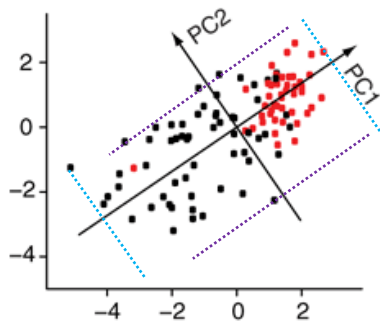
2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

9

9

## Variance on the new basis (PC1 and PC2)



- Sample variance in direction  $Z_1$  (PC1)

$$\text{var}(Z_1) = \frac{1}{(n-1)} \sum_{i=1}^n z_{i1}^2$$

- Sample variance in direction  $Z_2$  (PC2)

$$\text{var}(Z_2) = \frac{1}{(n-1)} \sum_{i=1}^n z_{i2}^2$$

- Total variance in data (**does not change**)

$$\text{var}(X) = \text{var}(Z_1) + \text{var}(Z_2)$$

- $Z_1$  is the direction preserving the maximal variance in data!

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

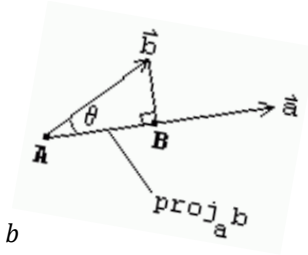
10

10

## Vector Projection

Given  $a, b \in R^m$

- $a \cdot b \triangleq \sum_{i=1}^n a_i b_i$
- $\cos(a, b) \triangleq \frac{a \cdot b}{\|a\| \times \|b\|}$
- Projection of vector  $b$  onto  $a$



$$\|b\| \cos(\theta) = \|b\| \frac{a \cdot b}{\|a\| \times \|b\|} = \frac{a \cdot b}{\|a\|} = \frac{a}{\|a\|} \cdot b$$

- Projection of any row vector  $x_i \in X$  onto unit vector  $v$  is  $(x_i \cdot v)$ .

10/8/2024

@Yiming Yang, Lecture on PCA, SVD and LSI

11

11

## First Component (PC1)

- Sample variance of  $X$  on any unit vector  $v$  is

$$\text{sample-var}(v; X) = \frac{1}{n-1} \sum_{i=1}^n (x_i \cdot v)^2$$

- To maximize the variance, we define PC1 as

$$\begin{aligned} v_1 &= \operatorname{argmax}_{\|v\|=1} \left\{ \sum_i (x_i \cdot v)^2 \right\} \\ &= \operatorname{argmax}_{\|v\|=1} \{ \|Xv\|^2 \} = \operatorname{argmax}_{\|v\|=1} \{ v^T X^T X v \} \end{aligned}$$

- Equivalently, we can solve  $(v_1, \lambda_1) = \operatorname{argmax}_{v, \lambda} \{ v^T X^T X v + \lambda(1 - v^T v) \}$
- Defining  $f(v, \lambda) = v^T X^T X v + \lambda(1 - v^T v)$ , we have the its mode(s) with

$$0 = \nabla_v f(v, \lambda) = 2X^T X v - 2\lambda v, \quad \text{i.e., } X^T X v = \lambda v$$

- Clearly,  $v$  and  $\lambda$  **must be** an eigenvector/eigenvalue of  $X^T X$ .

10/8/2024

@Yiming Yang, Lecture on PCA, SVD and LSI

12

12

## First and Other Components

- Necessary Condition
  - $v_1$  must be an eigenvector of  $X^T X$  and  $\lambda$  must be the corresponding eigenvalue.
- Sufficient condition
  - $v_1$  must be the eigenvector corresponding to  $\lambda_1 = \operatorname{argmax}_\lambda (v^T X^T X v)$
- How do we find  $v_1$ ?
  - Power Iteration with  $B = X^T X$  until convergence.
- How do we find next  $v_2, v_3, \dots$ ?

$$\hat{X}_k := X - \sum_{j=1}^{k-1} X v_j v_j^T, \quad v_k = \operatorname{argmax}_{\|v\|=1} \left\{ \frac{v^T \hat{X}_k^T \hat{X}_k v}{v^T v} \right\}$$

10/8/2024

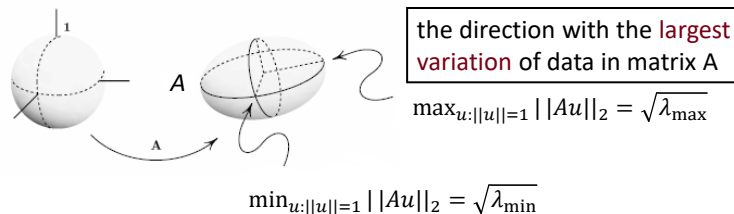
@Yiming Yang, Lecture on PCA, SVD and LSI

13

13

## Distortions by $Au$

(Carl Meyer, "Matrix Analysis and Applied Linear Algebra"  
<http://www.matrixanalysis.com>, Ch 5, p 281)



The eigenvalues of  $B = A^T A$  reflect *how much distortions may occur* when applying  $A$  to an arbitrary vector  $u$  on the surface of the unit ball.

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

14

14

## Distortion Factor: *Rayleigh Quotient*

$$\max_x \left\{ \frac{|Ax|}{|x|} \right\} = \sqrt{\lambda_1(B)} = \sigma_1(A)$$

- $\lambda_1$  is the largest **eigenvalue** of matrix  $B = A^T A$
- $\sigma_1$  is the largest **singular value** of matrix  $A$ .

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

15

15

## Outline

- Motivation with PCA
- ED and SVD
- Convergence of HITS

2/27/2024

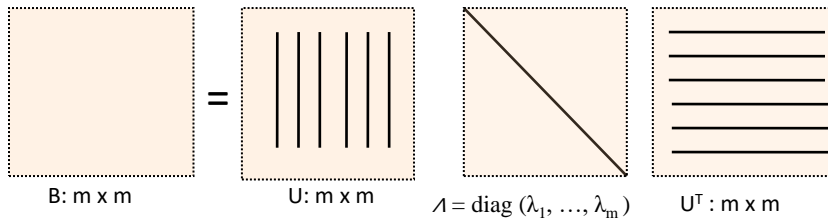
@Yiming Yang, S24 Lecture on ED &amp; SVD

16

16



## Eigendecomposition (Squared Matrix)



- Eigenvalues  $\lambda_1, \dots, \lambda_m$  are sorted by absolute value in decreasing order.
- Eigenvectors (the columns of matrix  $U$ ) are ordered accordingly.

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

17

17

## Properties of eigenvectors

- If  $u$  is an eigenvector, then  $v = cu$  is also an eigenvector for any  $c \in \mathbb{R}$ .

$$Bu = \lambda u \rightarrow B \underbrace{cu}_v = \lambda \underbrace{cu}_v$$

Thus, we only focus on **unit-length eigenvectors**.

- If  $u$  and  $v$  are both eigenvectors,  $u + v$  is also an eigenvector.

$$Bu = \lambda u, \quad Bv = \lambda v, \quad \rightarrow \quad B \underbrace{(u+v)}_w = \lambda \underbrace{(u+v)}_w$$

Thus, we focus on **linearly independent eigenvectors**, especially on **orthonormal** eigenvectors

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

18

18

## Properties of eigenvectors

- $B$  (  $m$ -by- $m$ ) has at most  $m$  linearly independent eigenvectors.
- The number of **distinct eigenvalues** may be equal or smaller than the number of the **linearly independent eigenvectors**.
- For example, identity matrix  $I$  has  $m$  independent eigenvectors (each column vector is orthogonal from other columns) but only one distinct eigenvalue  $\lambda = 1$ , as shown below:

$$Iv = v$$

$$\text{or} \quad 0 = \det(B - \lambda I) = (1 - \lambda)^m \det(I) \rightarrow \lambda = 1$$

- **We are interested in more general cases other than identity matrix.**

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

19

19

## Truncated Eigendecomposition

- Full dimension:  $B_{m \times m} = U_{m \times m} \Lambda_{m \times m} U_{m \times m}^T$
- Dimension reduced:

$$\begin{aligned} \hat{B}_{m \times m} &= \underbrace{U_{m \times k}}_{(u_1, u_2, \dots, u_k)} \underbrace{\Lambda_{k \times k}}_{\text{diag}\{\lambda_1, \dots, \lambda_k\}} U_{k \times m}^T \quad k < m \\ &= \sum_{i=1}^k \lambda_i u_i u_i^T \end{aligned}$$

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

20

20

## Singular Value Decomposition (SVD) of any matrix of rank $r \leq \min(n, m)$

$X: n \times m$        $U: n \times r$        $\Sigma: r \times r$        $V^T: r \times m$

$$X = U \Sigma V^T = (\vec{u}_1 \ \vec{u}_2 \ \dots \ \vec{u}_r) \begin{pmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \ddots \\ & & & \sigma_r \end{pmatrix} \begin{pmatrix} \vec{v}_1^T \\ \vec{v}_2^T \\ \vdots \\ \vec{v}_r^T \end{pmatrix} = \sum_{j=1}^r \sigma_j \vec{u}_j \vec{v}_j^T$$

$\Sigma = (\sigma_1, \sigma_2, \dots, \sigma_r)$  are the **singular values** ("spectrum");  
 $U = (u_1, u_2, \dots, u_r)$  are the **left singular vectors** (orthonormal);  
 $V = (v_1, v_2, \dots, v_r)$  are the **right singular vectors** (orthonormal).

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

21

21

21

## Truncated SVD (with $k < r$ )

$X: n \times m$        $U_k: n \times k$        $\Sigma: k \times k$        $V_k^T: k \times m$

$$X_k = U_k \Sigma_k V_k^T = (\vec{u}_1 \ \dots \ \vec{u}_k) \begin{pmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \ddots \\ & & & \sigma_k \end{pmatrix} \begin{pmatrix} \vec{v}_1^T \\ \vec{v}_2^T \\ \vdots \\ \vec{v}_k^T \end{pmatrix} = \sum_{j=1}^k \sigma_j \vec{u}_j \vec{v}_j^T$$

$(\sigma_1, \sigma_2, \dots, \sigma_k)$  are the **top-k singular values**;  
 $U_k$  and  $V_k$  contains the **top-k left/right singular vectors**, respectively.

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

22

22

22

## SVD and ED

- Rectangular matrix

$$A = U\Sigma V^T$$

- Squared matrices

$$B_a = A^T A = V\Sigma^T U^T U \Sigma V^T = V \Sigma^2 V^T = V \Lambda V^T$$

$$B_h = A A^T = U \Sigma V^T V \Sigma U^T = U \Sigma^2 U^T = U \Lambda U^T$$

$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_r), \lambda_i = \sigma_i^2 \forall i$$

- $U$  consists of the eigenvectors of  $B_h$ ;
- $V$  consists of the eigenvectors of  $B_a$ .

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

23

23

## Outline

- Motivation with PCA
- ED and SVD
- Convergence of HITS

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

24

24

## Recap the Statement about HITS' Convergence (in the lecture about HITS)

- [https://en.wikipedia.org/wiki/Power\\_iteration](https://en.wikipedia.org/wiki/Power_iteration)

"If we assume the matrix has an eigenvalue that is strictly greater in magnitude than its other eigenvalues and the starting vector has a nonzero component in the direction of an eigenvector associated with the dominant eigenvalue, then a subsequence converges to the eigenvector associated with the dominant eigenvalue."

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

25

25

## Convergence of Power Iteration in HITS (my sketch proof)

For precise proof see [https://en.wikipedia.org/wiki/Power\\_iteration](https://en.wikipedia.org/wiki/Power_iteration), which uses the Jordan Normal Decomposition instead of SVD.

$$B_a = A^T A \text{ where } A = U \Sigma V^T, \quad A \in \{0,1\}^{n \times n}, \quad U \in \mathbb{R}^{n \times r}, \quad V \in \mathbb{R}^{n \times r}, \quad r = \text{rank}(A)$$

$$B_a^k = (A^T A)^k = (V \Sigma U^T U \Sigma V^T)^k = V \Sigma^2 \dots \Sigma^2 V^T = V \Sigma^{2k} V^T$$

$$B_a^k z = V \Sigma^{2k} \underbrace{V^T z}_w = V \Sigma^{2k} w, \quad z \in \mathbb{R}^n, w \in \mathbb{R}^r$$

$$= [\vec{v}_1 \dots \vec{v}_r] \begin{bmatrix} \sigma_1^{2k} & & \\ & \ddots & \\ & & \sigma_r^{2k} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_r \end{bmatrix}$$

$$= w_1 \sigma_1^{2k} \vec{v}_1 + \dots + w_r \sigma_r^{2k} \vec{v}_r \quad \leftarrow \text{a linear combination of the eigenvectors}$$

$$= \sigma_1^{2k} \left( w_1 \vec{v}_1 + w_2 \left( \frac{\sigma_2}{\sigma_1} \right)^{2k} \vec{v}_2 + w_3 \left( \frac{\sigma_3}{\sigma_1} \right)^{2k} \vec{v}_3 \dots \right)$$

The 1<sup>st</sup> term dominates when k is large if w<sub>1</sub> is non-zero and if  $|\sigma_1| > |\sigma_2|$ .

2/27/2024

@Yiming Yang, S24 Lecture on ED &amp; SVD

26

26

26