

## CS-GY 6763: Lecture 9

# Low-rank approximation and singular value decomposition

---

NYU Tandon School of Engineering, Prof. Christopher Musco

- Reading group tomorrow at 9:30am. **Pat and Hogyong** are presenting on a method for the Frequent Items problem that improve on the CountMin method we learned in class in many scenarios.
- Midterms will be returned at the end of class.
- Problem Set 2 is being graded.

# SPECTRAL METHODS

Return to data compression:

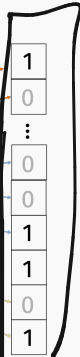
```
HTTP-Version: 1.0 Date: Mon, 7 Oct 2018
11:10 -0400 Message-Id: <CAWV1stQp-a-
jMLAnn0fy2j_jeaX5QwueRbIqcrRyNDa8aail.gn
il.com> Subject: 92231 Reading Group, Meeting
2, tomorrow at 10am From: Christopher Musco
<cmusco@nyu.edu> To: alj1d8@nyu.edu Content-
Type: multipart/alternative;
boundary="000000000078ec240594568a53" --
000000000078ec240594568a53 Content-Type:
text/plain charset="UTF-8" I hope everyone
had a good weekend! Tomorrow at 10am in 370
Jay St. #1114* we will meet for the second
installation of the CS-QY 92231 reading
group. Nick Feng will be leading a discussion
about the paper Simple Analyses of the Sparse
Johnson-Lindenstrauss Transform
<http://drops.dagstuhl.de/opus/volltexte/2018
/8305/pdf/OASIca-2018-2018-15.pdf>. Please
read the abstract and introduction before the
meeting. Best, - CM *Christopher Musco,
Assistant Professor* *New York University,
Tandon School of Engineering* *4011 578
2541* --000000000078ec240594568a53 Content-
Type: text/html charset="UTF-8" Content-
Transfer-Encoding: quoted-printable
```

sender on graylist 1  
sender on graylist 2

Bag-of-words

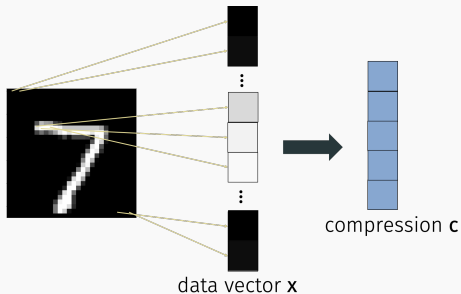
sender flagged by x users  
sender ip flagged by x users

data vector x



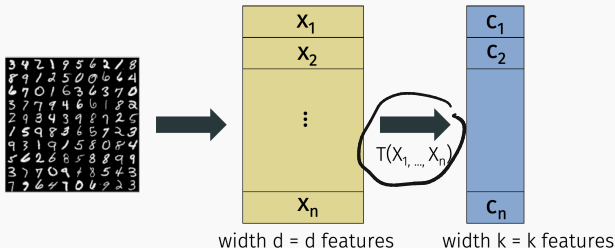
compression c

Return to data compression:



Main difference from randomized methods:

$$\underline{\underline{X_1 = \prod X_1}}$$



In this section, we will discuss data dependent transformations. Johnson-Lindenstrauss, MinHash, SimHash were all data oblivious.

Advantages of data **independent** methods:

- computationally more "light weight"
- don't require "full view" of data
  - streaming
  - distributed

Advantages of data **dependent** methods:

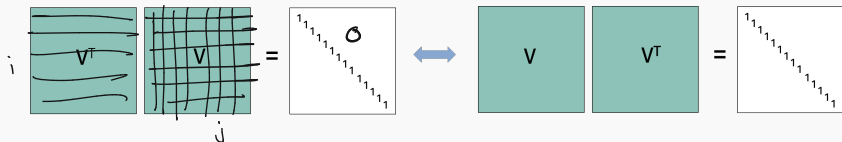
- better compression when there's underlying structure in the data

# LINEAR ALGEBRA REMINDER

If a square matrix has orthonormal rows, it also have orthonormal columns:

$$\|A\|_F^2 = \sum_{i,j} A_{ij}^2 = \sum_i \|\mathbf{a}_i\|_2^2$$

↑  
column in A



$$\underline{V^T V} = \underline{I} = \underline{W W^T}$$

Implies that for any vector  $\mathbf{x}$ ,  $\|\mathbf{V}\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$  and  $\|\mathbf{V}^T\mathbf{x}\|_2^2$ .

$$\mathbf{x}^T \mathbf{V}^T \mathbf{V} \mathbf{x} = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2$$

$$\mathbf{x}^T \mathbf{V} \mathbf{V}^T \mathbf{x} = \mathbf{x}^T \mathbf{I} \mathbf{x} = \|\mathbf{x}\|_2^2$$

Same thing goes for Frobenius norm: for any matrix  $\mathbf{X}$ ,

$$\underline{\|\mathbf{V}\mathbf{X}\|_F^2} = \underline{\|\mathbf{X}\|_F^2} \text{ and } \underline{\|\mathbf{V}^T\mathbf{X}\|_F^2} = \underline{\|\mathbf{X}\|_F^2}.$$

# LINEAR ALGEBRA REMINDER

The same is not true for rectangular matrices:

$V^T V = I$       but       $V V^T = P$

*projection matrix*

.5	-1	.7	-2
1/6	-.44	4.2	-1.5
7.8	.42	-.5	.67
-2	2.0	1.1	8.0
-1.5	.55	3.2	.5
.67	-2.8	-2.4	1.6
9.0	8.7	-7.7	7.8

For any  $x$ ,  $\|Vx\|_2^2 = \|x\|_2^2$  but  $\|V^T x\|_2^2 \neq \|x\|_2^2$  in general.

$$\| \cdot \| = \| \cdot \|$$

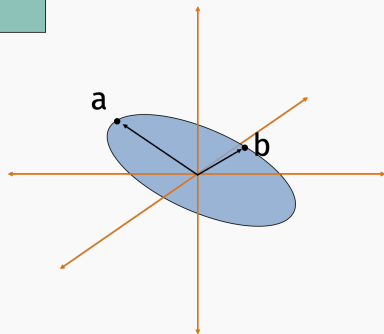
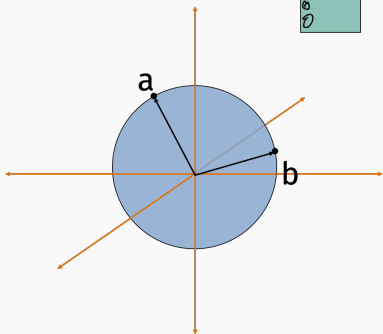
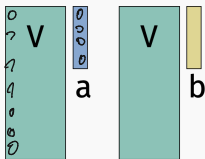


# LINEAR ALGEBRA REMINDER

Multiplying a vector by  $V$  with orthonormal columns rotates and/or reflects the vector.

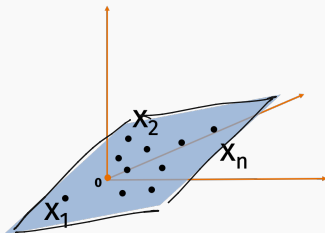
$$\|V(a-b)\|_2 = \|a-b\|_2$$

$$\|V(a+b)\|_2 = \|a+b\|_2$$

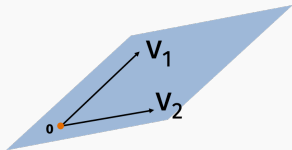


## LOW-RANK DATA

Suppose  $\underline{x}_1, \dots, \underline{x}_n \in \mathbb{R}^d$  lie on a low-dimensional subspace  $S$  through the origin. I.e. our data set is **rank  $k$**  for  $k < d$ .



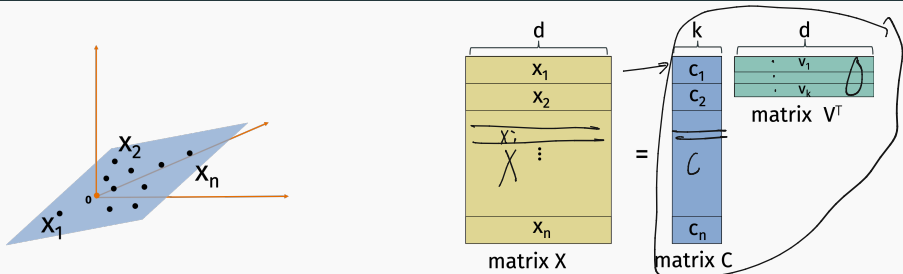
Let  $\mathbf{v}_1, \dots, \mathbf{v}_k$  be orthogonal unit vectors spanning  $S$ .



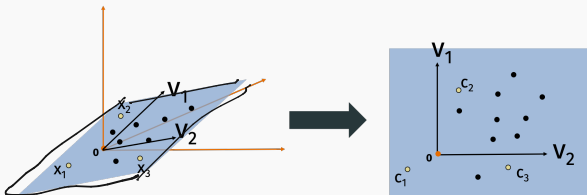
For all  $i$ , we can write:

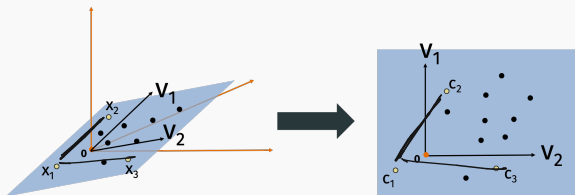
$$\mathbf{x}_i = \underline{c_{i,1}} \mathbf{v}_1 + \dots + \underline{c_{i,k}} \mathbf{v}_k.$$

# LOW-RANK DATA



What are  $c_1, \dots, c_n$ ?

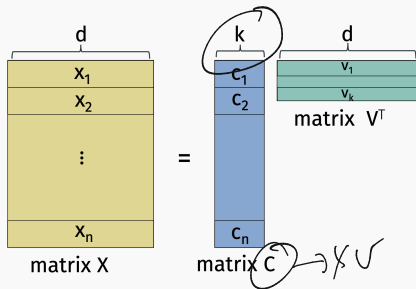




Lots of information preserved:

- $\|\mathbf{x}_i - \mathbf{x}_j\|_2 = \|\mathbf{c}_i - \mathbf{c}_j\|_2$  for all  $i, j$ .
- $\mathbf{x}_i^T \mathbf{x}_j = \mathbf{c}_i^T \mathbf{c}_j$  for all  $i, j$ .
- Norms preserved, linear separability preserved,  $\min \|\mathbf{X}\mathbf{y} - \mathbf{b}\| = \min \|\mathbf{C}\mathbf{z} - \mathbf{b}\|$ , etc., etc.

# LOW-RANK DATA



Formally,  $C = XV$ .

$$\underline{X} = \underline{C}V^T \Rightarrow \underline{XV} = \underline{C}V^T\underline{V} \quad \begin{matrix} \nearrow I \\ XV = C \end{matrix}$$

Since  $V$ 's columns are an orthonormal basis,  $V^T V = I$ .

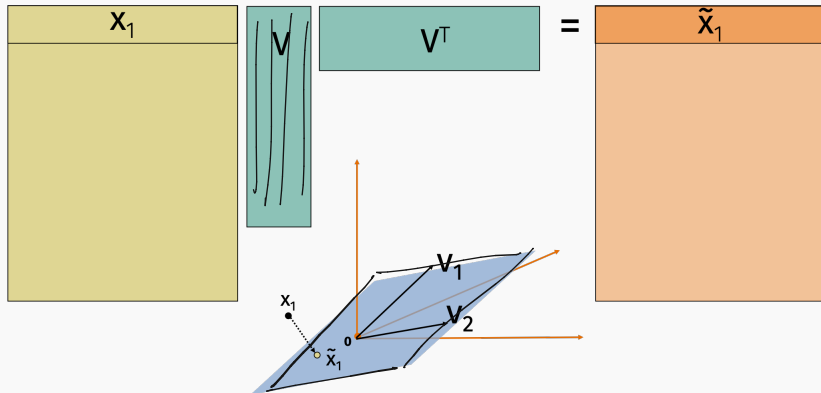
$$\text{So } X = XVV^T.$$

for some  $V$  with  $k$  columns

# PROJECTION MATRICES

$VV^T$  is a symmetric projection matrix.

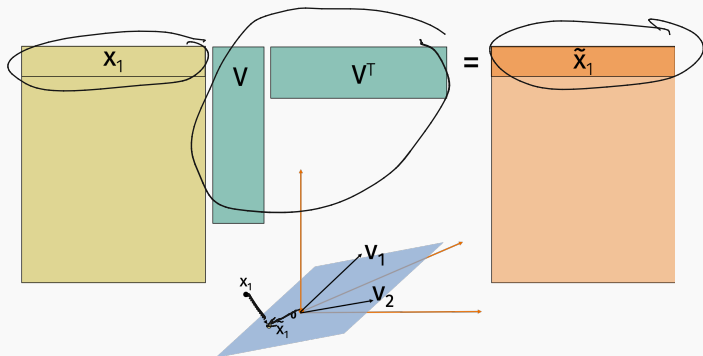
$$V^T V = I$$



When all data points already lie in the subspace spanned by  $V$ 's columns, projection doesn't do anything. So  $X = \underline{XVV^T}$ .

# PROJECTION MATRICES

$VV^T$  is a symmetric projection matrix.



$\tilde{x}_1 = \underline{x_1^T V V^T}$  is the projection of  $x_1^T$  onto the subspace.

By pythagorean theorem,  $\|x_1^T - x_1^T V V^T\|_2^2 = \|x_1^T\|_2^2 - \|x_1^T V V^T\|_2^2$  and  
 by apply to all rows,  $\|X - \underline{X V V^T}\|_F^2 = \underline{\|X\|_F^2} - \underline{\|X V V^T\|_F^2}$ .

## LOW-RANK APPROXIMATION

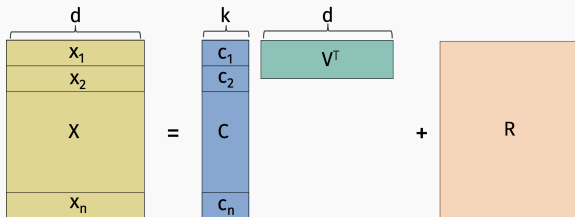
When  $X$ 's rows lie close to a  $k$  dimensional subspace, we can still approximate

$$\underline{X} \approx \underline{XV}^T.$$

$XV^T$  is a low-rank approximation for  $X$ .

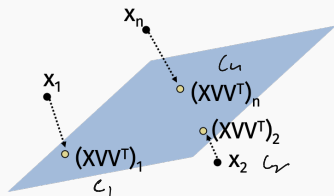
For a given subspace  $\mathcal{V}$  spanned by the columns in  $V$ ,

$$XV^T = \arg \min_C \|X - CV^T\|_F^2 = \sum_{i,j} (X_{i,j} - (CV^T)_{i,j})^2.$$

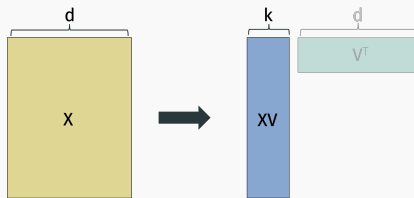




## LOW-RANK APPROXIMATION

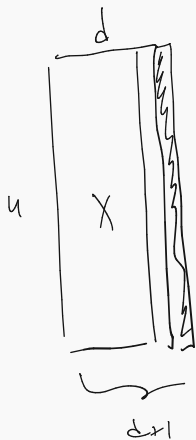


$$\|x_i - x_j\|_2 \approx \|x_i^T \mathbf{V}\mathbf{V}^T - x_j^T \mathbf{V}\mathbf{V}^T\|_2 = \|x_i^T \mathbf{V} - x_j^T \mathbf{V}\|_2$$



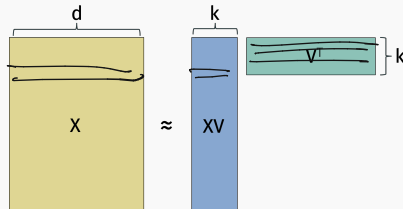
$XV$  can be used as a compressed version of data matrix  $X$ .

# WHY IS DATA APPROXIMATELY LOW-RANK?



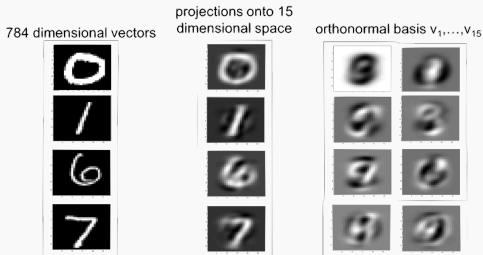
## DUAL VIEW

Rows of  $X$  (data points) are approximately spanned by  $k$  vectors. Columns of  $X$  (data features) are approximately spanned by  $k$  vectors.



## ROW REDUNDANCY

If a data set only had  $k$  unique data points, it would be exactly rank  $k$ . If it has  $k$  “clusters” of data points (e.g. the 10 digits) it’s often very close to rank  $k$ .



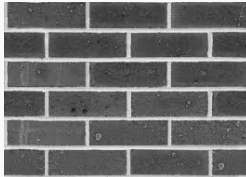
## COLUMN REDUNDANCY

Colinearity/correlation of data features leads to a low-rank data matrix.

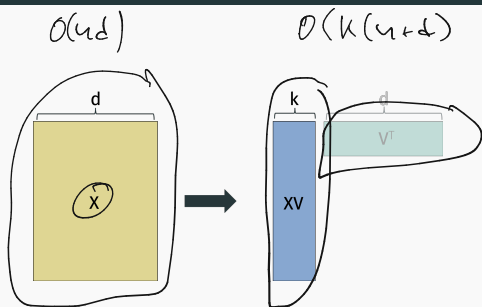
	bedrooms	bathrooms	sq.ft.	floors	list price	sale price
home 1	2	2	1800	2	200,000	195,000
home 2	4	2.5	2700	1	300,000	310,000
.	.	.	.	.	.	.
.	.	.	.	.	.	.
.	.	.	.	.	.	.
home n	5	3.5	3600	3	450,000	450,000

## OTHER REASONS FOR LOW-RANK STRUCTURE

When encoded as a matrix, which image has lower approximate rank?



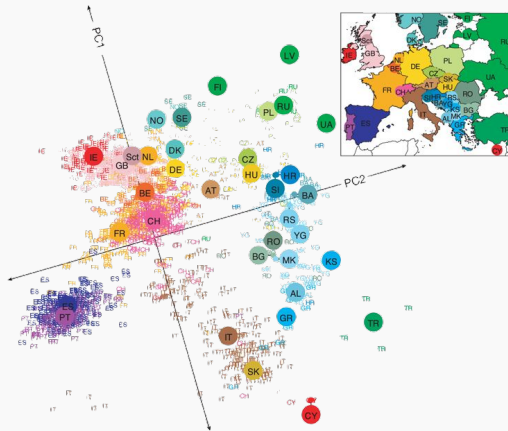
## APPLICATIONS OF LOW-RANK APPROXIMATION



- $XV \cdot V^T$  takes  $O(k(n + d))$  space to store instead of  $O(nd)$ .
- Regression problems involving  $XV \cdot V^T$  can be solved in  $O(nk^2)$  instead of  $O(nd^2)$  time.
- $XV$  can be used for visualization when  $k = 2, 3$ .

# APPLICATIONS OF LOW-RANK APPROXIMATION

“Genes Mirror Geography Within Europe” – Nature, 2008.



Each data vector  $\mathbf{x}_i$  contains genetic information for one person in Europe. Set  $k = 2$  and plot  $(XV)_i$  for each  $i$  on a 2-d plane. Color points by what country they are from.



## COMPUTATIONAL QUESTION

Given a subspace  $\mathcal{V}$  spanned by the  $k$  columns in  $V$ ,

$$\|X - XVV^T\|_F^2 \neq \min \|X - CV^T\|_F^2$$

$$\begin{aligned} & \|Ax\|_2 \\ \max_{x'} & \frac{\|Ax\|_2}{\|x\|_2} \end{aligned}$$

We want to find the best  $V \in \mathbb{R}^{d \times k}$ :

$$\min_{\text{orthonormal } V \in \mathbb{R}^{d \times k}} \|X - XVV^T\|_F^2 \quad (1)$$

Note that  $\|X - XVV^T\|_F^2 = \|X\|_F^2 - \|XVV^T\|_F^2$  for all orthonormal  $V$  (since  $VV^T$  is a projection). Equivalent form:

$$\max_{\text{orthonormal } V \in \mathbb{R}^{d \times k}} \|XVV^T\|_F^2 = \|XV\|_F^2 \quad (2)$$

If  $k = 1$ , want to find a single vector  $\mathbf{v}_1$  which maximizes:

$$\|\mathbf{X}\mathbf{v}_1\mathbf{v}_1^T\|_F^2 = \|\mathbf{X}\mathbf{v}_1\|_F^2 = \|\mathbf{X}\mathbf{v}_1\|_2^2 = \mathbf{v}_1^T \mathbf{X}^T \mathbf{X} \mathbf{v}_1.$$

Choose  $\mathbf{v}_1$  to be the top eigenvector of  $\mathbf{X}^T \mathbf{X}$ .

What about higher  $k$ ?

$$\begin{aligned} \max_{\mathbf{V}} & \mathbf{V}_1^T \mathbf{X}^T \mathbf{X} \mathbf{V}_1 \\ \mathbf{V}_1 & : \|\mathbf{V}_1\|_2 = 1 \end{aligned}$$

# SINGULAR VALUE DECOMPOSITION

$$\text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB)$$

One-stop shop for computing optimal low-rank approximations.

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

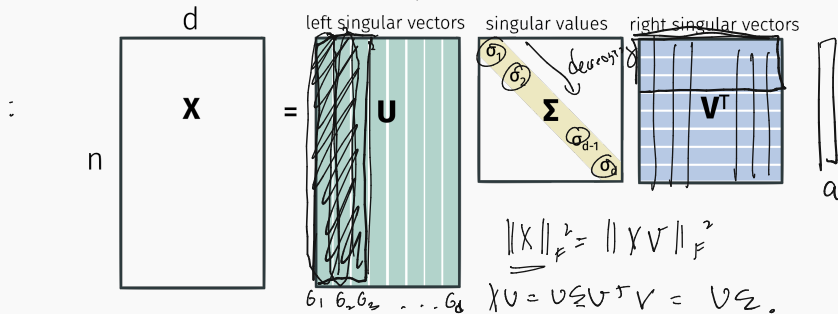
$$X^T X = V \Sigma U^T U \Sigma V^T = V \Sigma^2 V^T$$

Any matrix X can be written:

$$\text{tr}(U \Sigma^2 V^T V) = \text{tr}(\Sigma^2 V^T V) = \text{tr}(\Sigma^2)$$

$$\|X\|_F^2 = \text{tr}(X^T X)$$

$$\|X\|_2^2 = \lambda_{\max}(X^T X)$$



Where  $U^T U = I$ ,  $V^T V = I$ , and  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d \geq 0$ .

Note that  $\sum_{i=1}^d \sigma_i^2 = \|X\|_F^2$ .

## CONNECTION TO EIGENDECOMPOSITION

$$v_1, \dots, v_d$$

- $V_k$ 's columns are called the "top right singular vectors of  $X$ "
- $U_k$ 's columns are called the "top left singular vectors of  $X$ "
- $\sigma_1, \dots, \sigma_r$  are the "top singular values".  $\sigma_1, \dots, \sigma_d$  are sometimes called the "spectrum of  $X$ " (although this is more typically used to refer to eigenvalues).
- $U$  contains the orthonormal eigenvectors of  $XX^T$ .
- $V$  contains the orthonormal eigenvectors of  $X^T X$ .
- $\sigma_i^2 = \lambda_i(XX^T) = \lambda_i(X^T X)$

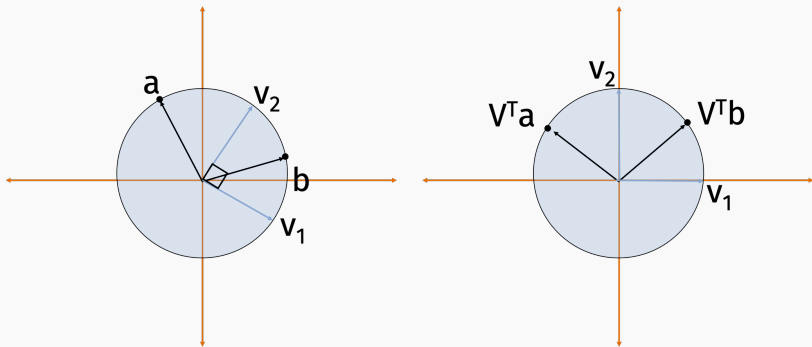
**Exercise:** Check this can be checked directly.

Important take away from singular value decomposition.

Multiplying any vector  $\mathbf{a}$  by a matrix  $\mathbf{X}$  to form  $\mathbf{Xa}$  can be viewed as a composition of 3 operations:

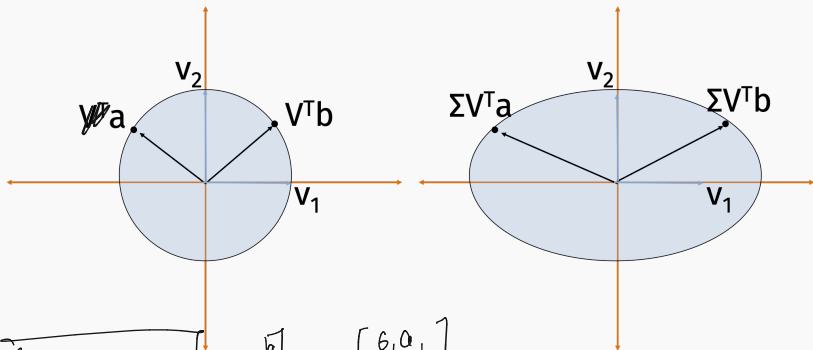
1. Rotate/reflect the vector (multiplication by  $\mathbf{V}^T$ ).
2. Scale the coordinates (multiplication by  $\mathbf{\Sigma}$ ).
3. Rotate/reflect the vector again (multiplication by  $\mathbf{U}$ ).

# SINGULAR VALUE DECOMPOSITION: ROTATE/REFLECT



# SINGULAR VALUE DECOMPOSITION: STRETCH

$\sigma_i > 0$



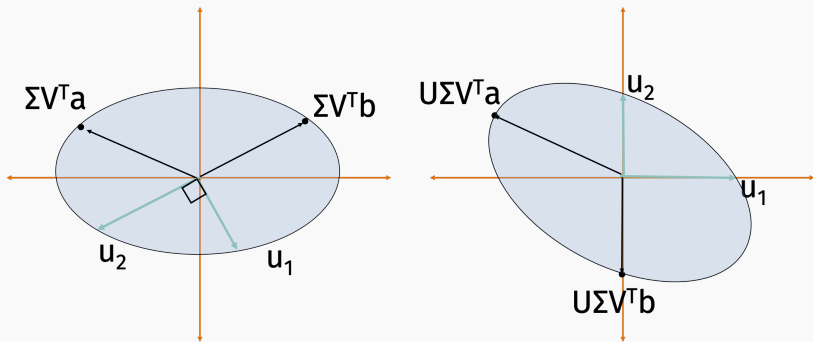
$$\begin{bmatrix} \sigma_1 & & & 0 \\ & \sigma_2 & & \\ & & \dots & \\ 0 & & & \sigma_d \end{bmatrix}$$

$\Sigma$

$$\begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_d \end{bmatrix} = \begin{bmatrix} \sigma_1 \alpha_1 \\ \vdots \\ \sigma_d \alpha_d \end{bmatrix}$$

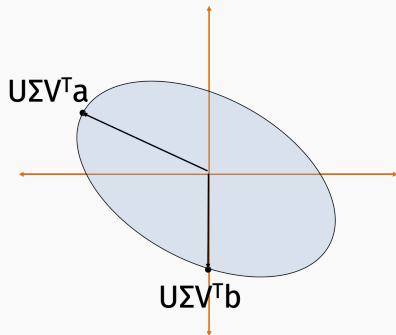
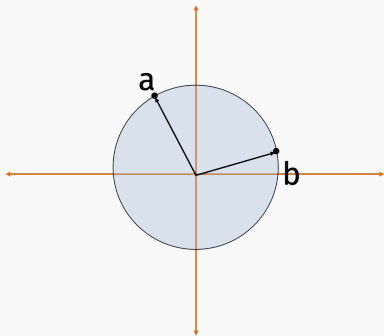
$\alpha$

# SINGULAR VALUE DECOMPOSITION: ROTATE/REFLECT



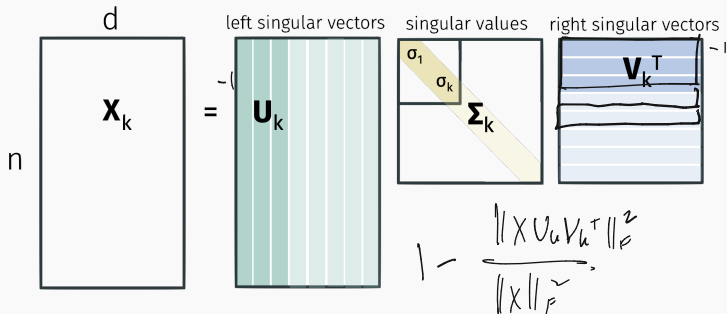


# SINGULAR VALUE DECOMPOSITION



# SINGULAR VALUE DECOMPOSITION

Can read off optimal low-rank approximations from the SVD:



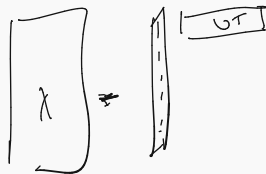
$$X_k = U_k \Sigma_k V_k^T = U_k U_k^T X = X V_k V_k^T$$

$$\boxed{V_k} = \arg \min_{\text{orthonormal } V \in \mathbb{R}^{d \times k}} \|X - X V V^T\|_F^2 = \arg \max_{\text{orthonormal } V \in \mathbb{R}^{d \times k}} \|X V V^T\|_F^2$$

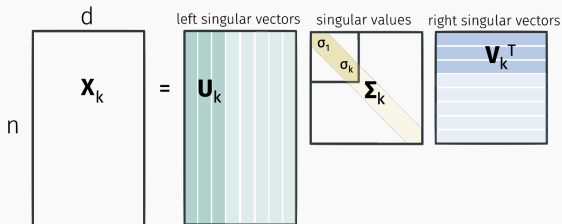
# SINGULAR VALUE DECOMPOSITION

Connection to **Principal Component Analysis**:

- Let  $\bar{X} = X - \mathbf{1}\mu^T$  where  $\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ . I.e.  $\bar{X}$  is obtained by mean centering  $X$ 's rows.
- Let  $\bar{U}\bar{\Sigma}\bar{V}^T$  be the SVD of  $\bar{X}$ .  $\bar{U}$ 's first columns are the "top principal components" of  $X$ .  $\bar{V}$ 's first columns are the "weight vectors" for these principal components.



## USEFUL OBSERVATIONS

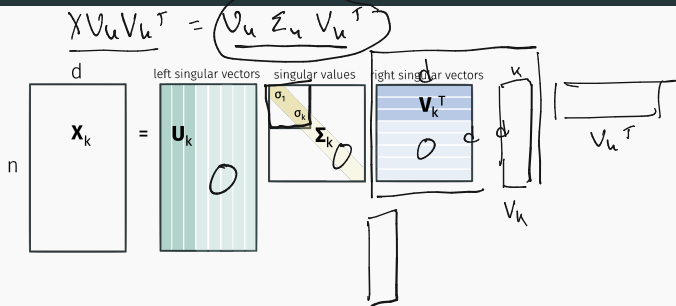


~~Observation 1: The optimal compression  $\mathbf{XV}_k$  has orthogonal columns.~~

# USEFUL OBSERVATIONS

$$V^T V_k = \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ & & 1 & 0 \\ & & & 0 \end{bmatrix}$$

$$\begin{bmatrix} V_k^T \\ E \end{bmatrix} \begin{bmatrix} V_k \\ \end{bmatrix}$$



**Observation 2:** The optimal low-rank approximation error

$E_k = \|X - X V_k V_k^T\|_F^2 = \|X\|_F^2 - \|X V_k V_k^T\|_F^2$  can be written:

$$\sum_{i=1}^d G_i^2$$

$$E_k = \sum_{i=k+1}^d \sigma_i^2.$$

$$\sum_{i=1}^k G_i^2$$

$$X V_k V_k^T = U \underbrace{\Sigma V^T V_k V_k^T}_{}$$

## SPECTRAL PLOTS

**Observation 2:** The optimal low-rank approximation error

$E_k = \|\mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^T\|_F^2 = \|\mathbf{X}\|_F^2 - \|\mathbf{X}\mathbf{V}_k\mathbf{V}_k^T\|_F^2$  can be written:

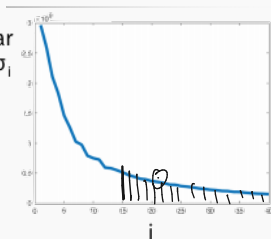
$$E_k = \sum_{i=k+1}^d \sigma_i^2.$$

Can immediately get a sense of “how low-rank” a matrix is from it’s spectrum:

784 dimensional vectors



singular  
value  $\sigma_i$



## SPECTRAL PLOTS

**Observation 2:** The optimal low-rank approximation error

$E_k = \|\mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^T\|_F^2 = \|\mathbf{X}\|_F^2 - \|\mathbf{X}\mathbf{V}_k\mathbf{V}_k^T\|_F^2$  can be written:

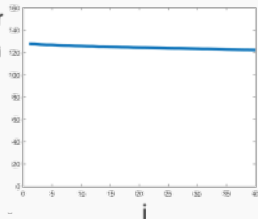
$$E_k = \sum_{i=k+1}^d \sigma_i^2.$$

Can immediately get a sense of “how low-rank” a matrix is from it’s spectrum:

784 dimensional vectors



singular  
value  $\sigma_i$

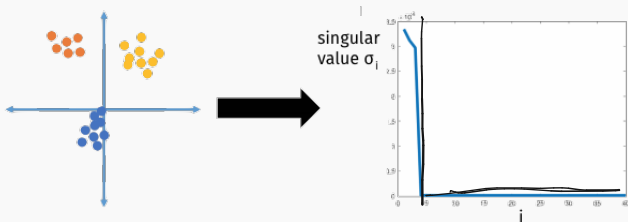


**Observation 2:** The optimal low-rank approximation error

$E_k = \|\mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^T\|_F^2 = \|\mathbf{X}\|_F^2 - \|\mathbf{X}\mathbf{V}_k\mathbf{V}_k^T\|_F^2$  can be written:

$$E_k = \sum_{i=k+1}^d \sigma_i^2.$$

Can immediately get a sense of “how low-rank” a matrix is from it’s spectrum:





## COMPUTING THE SVD

Suffices to compute right singular vectors  $V$ :  $\rightarrow O(d^3)$

- Compute  $X^T X$ .  $d \times d \rightarrow O(d^2)$
- Find eigendecomposition  $V \Lambda V^T = X^T X$ .
- Compute  $L = XV$ . Set  $\sigma_i = \|L_i\|_2$  and  $U_i = L_i / \|L_i\|_2$ .

Total runtime  $\approx \underline{O(d^2)} + O(d^3)$

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

$$\underline{XV} = U \Sigma$$

## COMPUTING THE SVD (FASTER)

- Compute approximate solution.
- Only compute top  $k$  singular vectors/values. Runtime will depend on  $k$ . When  $k = d$  we can't do any better than classical algorithms based on eigendecomposition.
- Iterative algorithms achieve runtime  $\approx \underline{O(ndk)}$  vs.  $\underline{O(nd^2)}$  time.
  - Krylov subspace methods like the Lanczos method are most commonly used in practice.
  - **Power method** is the simplest Krylov subspace method, and still works very well.

**What we won't discuss today:** sketching methods and stochastic methods (which are faster in some settings).

# POWER METHOD

Today: What about when  $k = 1$ ?

Goal: Find some  $\underline{z} \approx \underline{v}_1$

Input:  $\underline{X} \in \mathbb{R}^{n \times d}$  with SVD  $\underline{U}\underline{\Sigma}\underline{V}^T$ .

$$\underline{X}^T \underline{X} = \underline{V} \underline{\Sigma} \underline{U}^T \underline{U} \underline{\Sigma} \underline{V}^T = \underline{V} \underline{\Sigma}^2 \underline{V}^T$$

$$\underline{V} \underline{\Sigma}^2 \underline{V}^T \underline{z} \quad \underline{V}^T \underline{z} = \begin{bmatrix} c_1 \\ \vdots \\ c_d \end{bmatrix}$$

$$\underline{z}^{(1)} = c_1 \underline{v}_1 + c_2 \underline{v}_2 + \dots + c_d \underline{v}_d$$

Power method:

$$\underline{\Sigma}^2 \underline{V}^T \underline{z} = \begin{bmatrix} \sigma_1^2 c_1 \\ \sigma_2^2 c_2 \\ \vdots \\ \sigma_d^2 c_d \end{bmatrix} \underline{z}^{(i+1)} = \underline{c}_1 \underline{\sigma}_1^2 \underline{v}_1 + \dots + \underline{c}_d \underline{\sigma}_d^2 \underline{v}_d$$

- Choose  $\underline{z}^{(0)}$  randomly. E.g.  $\underline{z}_0 \sim \mathcal{N}(0, 1)$ .

- $\underline{z}^{(0)} = \underline{z}^{(0)} / \|\underline{z}^{(0)}\|_2$

- For  $i = 1, \dots, T$

- $\underline{z}^{(i)} = \underline{X}^T \cdot (\underline{X} \underline{z}^{(i-1)})$

- $n_i = \|\underline{z}^{(i)}\|_2$

- $\underline{z}^{(i)} = \underline{z}^{(i)} / n_i$

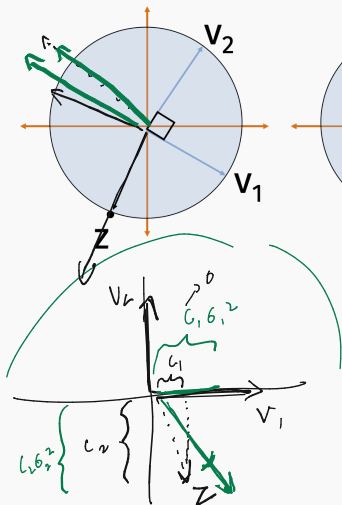
Return  $\underline{z}^{(T)}$

$$\underline{z}^{(i)} = \frac{1}{\prod_{j=1}^i n_j} (\underline{X}^T \underline{X}) (\underline{X}^T \underline{X}) \dots (\underline{X}^T \underline{X}) \underline{z}^{(0)}$$

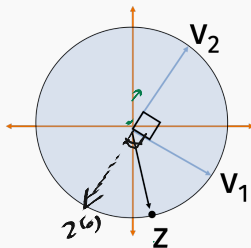
$$(\underline{X}^T \underline{X})^i \underline{z}^{(0)}$$

# POWER METHOD INTUITION

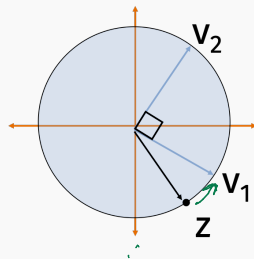
0 iterations



1 iterations



2 iterations



## POWER METHOD FORMAL CONVERGENCE

### Theorem (Basic Power Method Convergence)

Let  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$  be parameter capturing the “gap” between the first and second largest singular values. If Power Method is initialized with a random Gaussian vector then, with high probability, after  $T = O\left(\frac{\log(d/\epsilon)}{\gamma}\right)$  steps, we have either:

$$\|v_1 - z^{(T)}\|_2 \leq \epsilon \quad \text{or} \quad \|v_1 - (-z^{(T)})\|_2 \leq \epsilon.$$

**Total runtime:**  $O\left(nd \cdot \frac{\log d/\epsilon}{\gamma} \cdot k\right)$

$$O(nd^2)$$

## ONE STEP ANALYSIS OF POWER METHOD

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

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

$$\mathbf{z}^{(1)} = c_1^{(1)}\mathbf{v}_1 + c_2^{(1)}\mathbf{v}_2 + \dots + c_d^{(1)}\mathbf{v}_d$$

$\vdots$

$$\mathbf{z}^{(i)} = c_1^{(i)}\mathbf{v}_1 + c_2^{(i)}\mathbf{v}_2 + \dots + c_d^{(i)}\mathbf{v}_d$$

Note:  $[c_1^{(i)}, \dots, c_d^{(i)}] = \mathbf{c}^{(i)} = \underline{\mathbf{V}^T \mathbf{z}^{(i)}}$ .

Also:  $\sum_{j=1}^d \underbrace{(c_j^{(i)})^2}_{\leftarrow} \neq 1.$

## ONE STEP ANALYSIS OF POWER METHOD

Claim: After update  $\underline{\underline{z}}^{(i)} = \frac{1}{n_i} X^T X \underline{\underline{z}}^{(i-1)}$ ,

$$c_j^{(i)} = \frac{1}{n_i} \sigma_j^2 c_j^{(i-1)}$$

$$\underline{\underline{z}}^{(i)} = \frac{1}{n_i} \left[ c_1^{(i-1)} \sigma_1^2 \underline{\underline{v}}_1 + c_2^{(i-1)} \sigma_2^2 \underline{\underline{v}}_2 + \dots + c_d^{(i-1)} \sigma_d^2 \underline{\underline{v}}_d \right]$$

# MULTI-STEP ANALYSIS OF POWER METHOD

Claim: After  $T$  updates:

$$z^{(T)} = \frac{1}{\prod_{i=1}^T n_i} \left[ \overset{\text{near } 1}{c_1^{(0)} \sigma_1^{2T}} \cdot v_1 + \overset{\text{near } 0}{c_2^{(0)} \sigma_2^{2T}} \cdot v_2 + \dots + \overset{\text{near } 0}{c_d^{(0)} \sigma_d^{2T}} \cdot v_d \right]$$

$\approx v_1$

$$[1 \cdot v_1 + 0 \cdot v_2 + \dots + \dots + 0 \cdot v_d]$$

Let  $\alpha_j = \frac{1}{\prod_{i=1}^T n_i} c_j^{(0)} \sigma_j^{2T}$ . **Goal:** Show that  $\alpha_j \ll \alpha_1$  for all  $j \neq 1$ .

$$\alpha_1 \approx 1$$



## POWER METHOD FORMAL CONVERGENCE

Since  $\mathbf{z}^{(T)}$  is a unit vector,  $\sum_{i=1}^d \alpha_i^2 = 1$ . So  $\alpha_1 \leq 1$ .

If we can prove that  $\frac{\alpha_j}{\alpha_1} \leq \sqrt{\frac{\epsilon}{d}}$  then:

$$\alpha_j^2 \leq \alpha_1^2 \cdot \frac{\epsilon}{d}$$

$$1 = \alpha_1^2 + \sum_{j=2}^d \alpha_j^2 \leq \alpha_1^2 + \epsilon$$

$$\alpha_1^2 \geq 1 - \epsilon$$

$$\underline{|\alpha_1|} \geq 1 - \epsilon$$

$$\|\mathbf{v}_1 - \mathbf{z}^{(T)}\|_2 = 2 - 2\langle \mathbf{v}_1, \mathbf{z}^{(T)} \rangle \leq 2\epsilon$$

# POWER METHOD FORMAL CONVERGENCE

Lets prove that  $\frac{\alpha_j}{\alpha_1} \leq \sqrt{\frac{\epsilon}{d}}$  where  $\alpha_j = \frac{1}{\prod_{i=1}^T n_i} c_j^{(0)} \sigma_j^{2T}$

**First observation:** Starting coefficients are all roughly equal.

For all  $j$

$$O(1/d^3) \leq c_j^{(0)} \leq 1$$

with probability  $1 - \frac{1}{d}$ . This is a very loose bound, but it's all that we will need. **Prove using Gaussian concentration.**

$$\frac{\alpha_j}{\alpha_1} = \frac{\sigma_j^{2T} c_j^{(0)}}{\sigma_1^{2T} c_1^{(0)}} \leq \frac{\sigma_j^{2T}}{\sigma_1^{2T}} \cdot \frac{1}{O(1/d^3)} \cdot \frac{\sigma_j^{2T}}{\sigma_1^{2T}} \cdot O(d^3)$$

$$y = \frac{\sigma_1 - \sigma_2}{\sigma_1}$$

Need  $T =$

$$y = 1 - \frac{\sigma_2}{\sigma_1} \quad \frac{\sigma_2}{\sigma_1} = 1 - y$$

$$= \left(\frac{\sigma_j}{\sigma_1}\right)^{2T} \cdot O(d^3)$$

$$\leq (1-y)^{2T} \cdot O(d^3)$$

$$\leq \frac{\epsilon^c}{d^c} \text{ if } T = O\left(\frac{\log(d/\epsilon)}{y}\right)$$

### Theorem (Gapless Power Method Convergence)

If Power Method is initialized with a random Gaussian vector then, with high probability, after  $T = O\left(\frac{\log d/\epsilon}{\epsilon}\right)$  steps, we obtain a  $\mathbf{z}$  satisfying:

$$\|X - X\mathbf{z}\mathbf{z}^T\|_F^2 \leq (1 + \epsilon)\|X - X\mathbf{v}_1\mathbf{v}_1^T\|_F^2$$

## GENERALIZATIONS TO LARGER $k$

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

### Power method:

- Choose  $\mathbf{G} \in \mathbb{R}^{d \times k}$  be a random Gaussian matrix.
- $\mathbf{Z}_0 = \text{orth}(\mathbf{G})$ .
- For  $i = 1, \dots, T$ 
  - $\mathbf{Z}^{(i)} = \mathbf{X}^T \cdot (\mathbf{X}\mathbf{Z}^{(i-1)})$
  - $\mathbf{Z}^{(i)} = \text{orth}(\mathbf{Z}^{(i)})$

Return  $\mathbf{Z}^{(T)}$

$$\begin{array}{c} \lambda^T \mathbf{X} \mathbf{Z}^{(i-1)} \\ \underbrace{\hspace{10em}} \\ \text{rank} \\ \underbrace{\hspace{10em}} \\ O(\epsilon dk) \end{array}$$

**Runtime:**  $O\left(\frac{\log d/\epsilon}{\epsilon}\right)$  iterations to obtain a nearly optimal low-rank approximation:

$$\|\mathbf{X} - \mathbf{X}\mathbf{Z}\mathbf{Z}^T\|_F^2 \leq (1 + \epsilon) \|\mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^T\|_F^2.$$

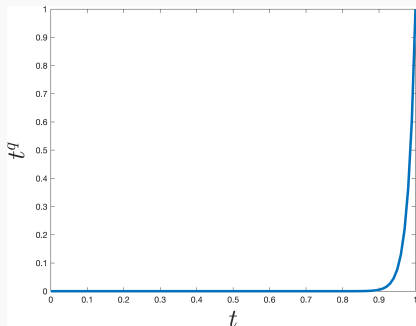
Possible to “accelerate” these methods.

**Convergence Guarantee:**  $T = O\left(\frac{\log d/\epsilon}{\sqrt{\epsilon}}\right)$  iterations to obtain a nearly optimal low-rank approximation:

$$\|\mathbf{A} - \mathbf{AZZ}^T\|_F^2 \leq (1 + \epsilon)\|\mathbf{A} - \mathbf{AV}_k\mathbf{V}_k^T\|_F^2.$$

**Runtime:**  $O(\text{nnz}(\mathbf{X}) \cdot k \cdot T) \leq O(ndk \cdot T)$ .

$$\mathbf{z}^{(q)} = c \cdot \underbrace{(\mathbf{X}^T \mathbf{X})^q}_{\text{circled}} \cdot \mathbf{g}$$



$$\begin{aligned} & \underline{\underline{P(\mathbf{X}^T \mathbf{X})}} \cdot \mathbf{g} \\ & (c_0 \mathbf{g} + c_1 \mathbf{X}^T \mathbf{X} \mathbf{g} \\ & + c_2 (\mathbf{X}^T \mathbf{X})^2 \mathbf{g} \\ & + \dots \\ & c_r (\mathbf{X}^T \mathbf{X})^r \mathbf{g}) \end{aligned}$$

$$\mathbf{z}^{(q)} = c \cdot \left[ c_1 \cdot \sigma_1^{2q} \mathbf{v}_1 + c_2 \cdot \sigma_2^{2q} \mathbf{v}_2 + \dots + c_n \cdot \sigma_n^{2q} \mathbf{v}_n \right]$$

$$z^{(q)} = \underline{c} \cdot (X^T X)^q \cdot \mathbf{g}$$

Along the way we computed:

$$\mathcal{K}_q = \left[ \mathbf{g}, (X^T X) \cdot \mathbf{g}, (X^T X)^2 \cdot \mathbf{g}, \dots, (X^T X)^q \cdot \mathbf{g} \right]$$

$\mathcal{K}$  is called the Krylov subspace of degree  $q$ .

**Idea behind Krylov methods:** Don't throw away everything before  $(X^T X)^q \cdot \mathbf{g}$ . What you're using when you run `svds` or `eigs` in MATLAB or Python.

Want to find  $\mathbf{v}$ , which minimizes  $\|\mathbf{X} - \mathbf{X}\mathbf{v}\mathbf{v}^T\|_F^2$ .

### Lanczos method:

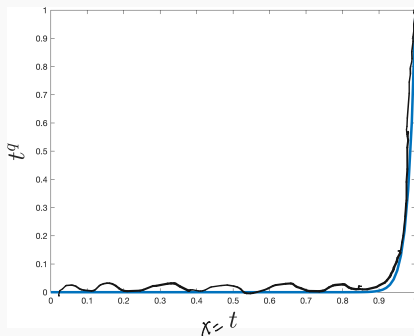
- Let  $\mathbf{Q} \in \mathbb{R}^{d \times k}$  be an orthonormal span for the vectors in  $\mathcal{K}$ .
- Solve  $\min_{\mathbf{v}=\mathbf{Q}\mathbf{w}} \|\mathbf{X} - \mathbf{X}\mathbf{v}\mathbf{v}^T\|_F^2$ .
  - Find best vector in the Krylov subspace, instead of just using last vector.
  - Can be done in  $O(\text{nnz}(\mathbf{X}) \cdot k + dk^2)$  time.



# LANCZOS METHOD ANALYSIS

**Claim:** There is an  $O\left(\sqrt{q \log \frac{1}{\epsilon}}\right)$  degree polynomial  $\hat{p}$  approximating  $\underline{\underline{x}}^q$  up to error  $\epsilon \sigma_1^2$  on  $[0, \sigma_1^2]$ .

$$6, 2 = 1$$



$$\|X - Xv_{p^*}v_{p^*}^T\|_F^2 \leq \|X - Xv_{\hat{p}}v_{\hat{p}}^T\|_F^2 \approx \|X - Xv_{x^q}v_{x^q}^T\|_F^2 \approx \|X - Xv_1v_1^T\|_F^2$$

**Runtime:**  $O\left(\frac{\log(d/\epsilon)}{\sqrt{\gamma}} \cdot \text{nnz}(X)\right)$  vs.  $O\left(\frac{\log(d/\epsilon)}{\gamma} \cdot \text{nnz}(X)\right)$

## POWER METHOD – NO GAP DEPENDENCE

Convergence is slow when  $\gamma = \frac{\sigma_1 - \sigma_2}{\sigma_1}$  is small.  $\mathbf{z}^{(q)}$  has large components of both  $\mathbf{v}_1$  and  $\mathbf{v}_2$ . But in this case:

$$\|\mathbf{X} - \mathbf{X}\mathbf{v}_1\mathbf{v}_1^T\|_F^2 = \sum_{i \neq 1} \sigma_i^2 \approx \sum_{i \neq 2} \sigma_i^2 = \sigma_i^2 \|\mathbf{X} - \mathbf{X}\mathbf{v}_2\mathbf{v}_2^T\|_F^2.$$

So we don't care! Either  $\mathbf{v}_1$  or  $\mathbf{v}_2$  give good rank-1 approximations.

**Claim:** To achieve

$$\|\mathbf{X} - \mathbf{X}\mathbf{z}\mathbf{z}^T\|_F^2 \leq (1 + \epsilon) \|\mathbf{X} - \mathbf{X}\mathbf{v}_1\mathbf{v}_1^T\|_F^2$$

we need  $O\left(\frac{\log(d/\epsilon)}{\epsilon}\right)$  power method iterations or  $O\left(\frac{\log(d/\epsilon)}{\sqrt{\epsilon}}\right)$  Lanczos iterations.

- Block Krylov methods
- Let  $\mathbf{G} \in \mathbb{R}^{d \times k}$  be a random Gaussian matrix.
- $\mathcal{K}_q = \left[ \mathbf{G}, (\mathbf{X}^T \mathbf{X}) \cdot \mathbf{G}, (\mathbf{X}^T \mathbf{X})^2 \cdot \mathbf{G}, \dots, (\mathbf{X}^T \mathbf{X})^q \cdot \mathbf{G} \right]$

**Runtime:**  $O\left(\text{nnz}(\mathbf{X}) \cdot k \cdot \frac{\log d/\epsilon}{\sqrt{\epsilon}}\right)$  to obtain a nearly optimal low-rank approximation.