

CS-GY 9223 I: Lecture 9

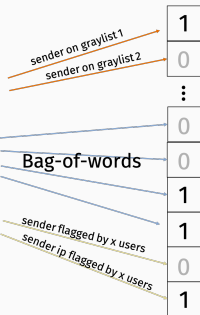
Low-rank approximation and singular value decomposition

NYU Tandon School of Engineering, Prof. Christopher Musco

Return to data compression:

```

MIME-Version: 1.0 Date: Mon, 7 Oct 2019
14:51:30 -0400 Message-ID: <CANV12i0ga=
39MAnnOpY24_juaxKQoua9NbiQC7Rf9Mda8@mail.g
ll.com> Subject: 92231 Reading Group, Meeting
2, tomorrow at 10am From: Christopher Musco
<cmusco@nyu.edu> To: aljeld@nyu.edu Content-
Type: multipart/alternative
boundary="000000000078ec240594568a53" --
0000000000078ec240594568a53 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
instantiation of the CS-OY 92231 reading
group. Nick Feng will be leading a discussion
about the paper Simple Analyses of the Sparse
Johnson-Lindenstrauss Transform
<http://arop.dagstuhl.de/opus/volltexte/2018
/8305/pdf/0802ca-005a-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* *401 578
2541* --0000000000078ec240594568a53 Content-
Transfer-Encoding: quoted-printable
    
```

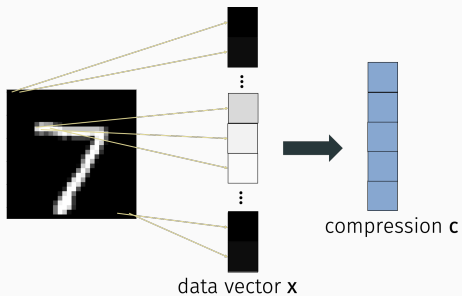


data vector x

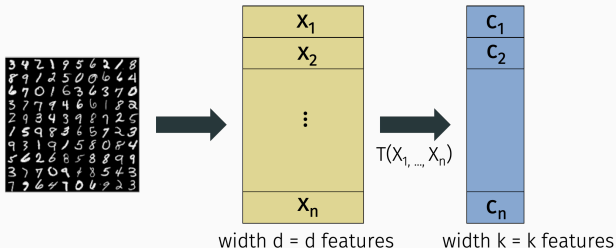


compression c

Return to data compression:



Main difference from randomized methods:



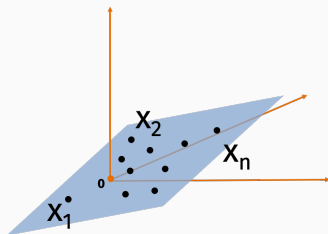
In this section, we will discuss data dependent data transformations. Johnson-Lindenstrauss, MinHash, SimHash were all data oblivious.

Advantages of data **independent** methods:

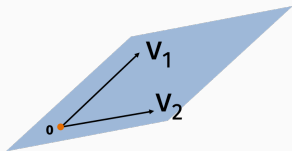
Advantages of data **dependent** methods:

LOW-RANK DATA

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



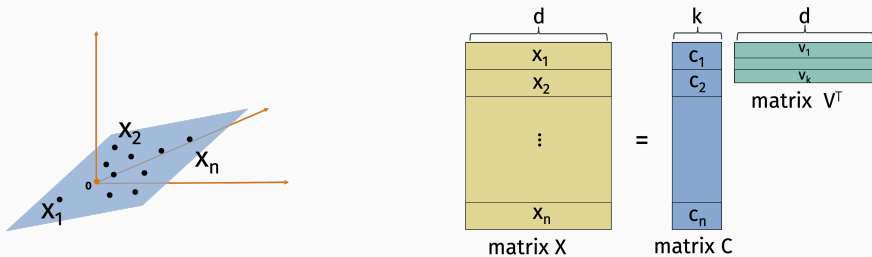
Let $\mathbf{v}_1, \dots, \mathbf{v}_k$ be orthogonal unit vectors spanning S .



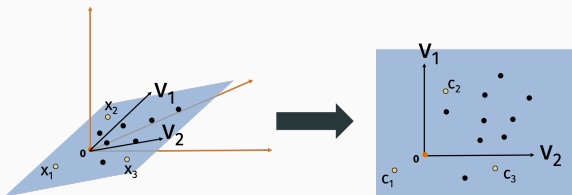
For all i , we can write:

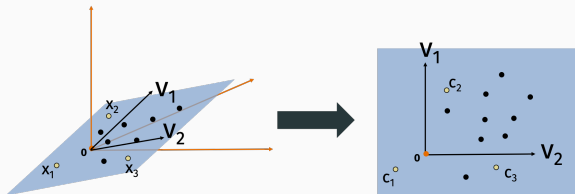
$$\mathbf{x}_i = c_{i,1}\mathbf{v}_1 + \dots + c_{i,k}\mathbf{v}_k.$$

LOW-RANK DATA



What are c_1, \dots, c_n ?

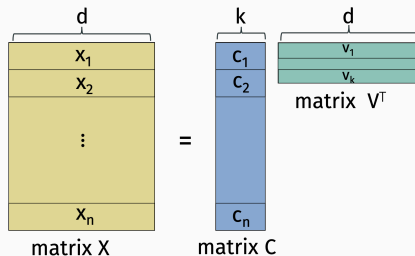




Lots of information preserved:

- $\|\mathbf{x}_i - \mathbf{x}_j\|_2 = \|\mathbf{c}_i - \mathbf{c}_j\|_2$ 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, $C = XV^T$:

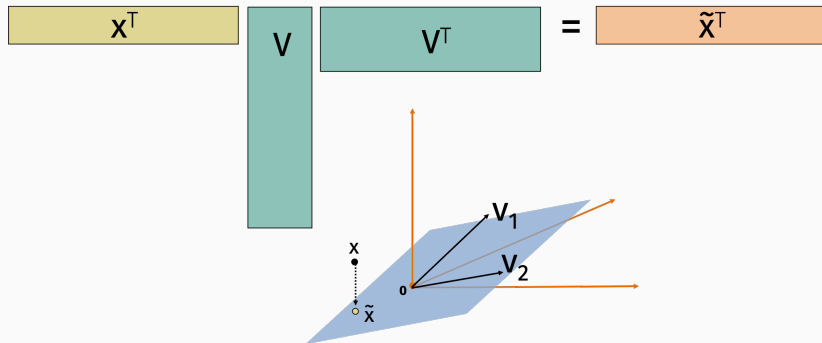
$$X = CV^T \Rightarrow XV = CV^T V$$

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

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

PROJECTION MATRICES

VV^T is a symmetric projection matrix.



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

LOW-RANK APPROXIMATION

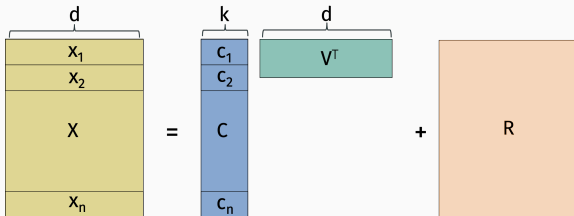
When X 's rows lie close to a k dimensional subspace, we can still approximate

$$X \approx XW^T.$$

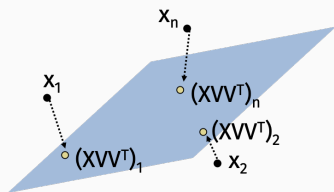
XW^T is a low-rank approximation for X .

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

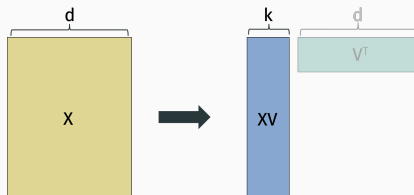
$$XW^T = \arg \min_C \|X - CV^T\|_F^2 = \sum_{i,j} (X_{i,j} - (CV^T)_{i,j})^2.$$



LOW-RANK APPROXIMATION



$$\|x_i - x_j\|_2 \approx \|(XV^T)_i - (XV^T)_j\|_2 = \|(XV^T)_i - (XV^T)_j\|_2$$

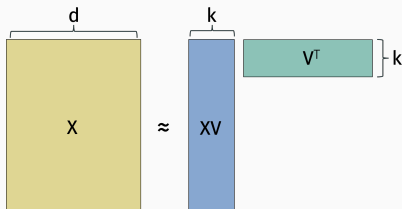


XV can be used as a compressed version of data matrix X .

WHY IS DATA APPROXIMATELY LOW-RANK?

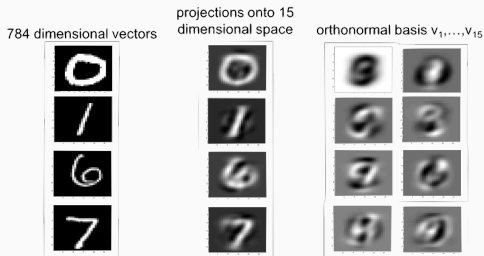
DUAL VIEW

Rows of \mathbf{X} (data points) are approximately spanned by k vectors. Columns of \mathbf{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 .



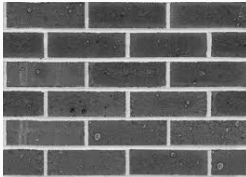
COLUMN REDUNDANCY

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

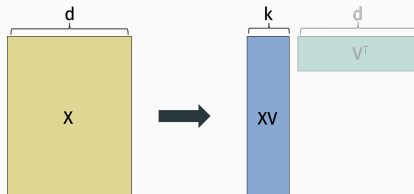
	bedrooms	bathrooms	sq.ft.	floors	list price	sale price
home 1	2	2	1800	2	200,000	195,000
home 2	4	2.5	2700	1	300,000	310,000
.
.
.
home n	5	3.5	3600	3	450,000	450,000

OTHER REASONS FOR LOW-RANK STRUCTURE

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



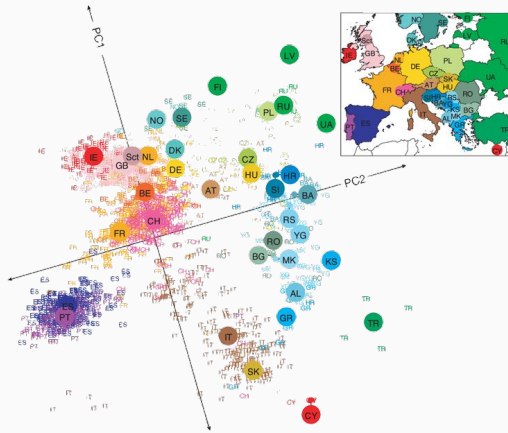
APPLICATIONS OF LOW-RANK APPROXIMATION



- $XV \cdot V^T$ 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 $k = 2, 3$.
- We will discuss many more next class.

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 \mathbf{V} ,

$$\|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2 = \min_{\mathbf{C}} \|\mathbf{X} - \mathbf{C}\mathbf{V}^T\|_F^2$$

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

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

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

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

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

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

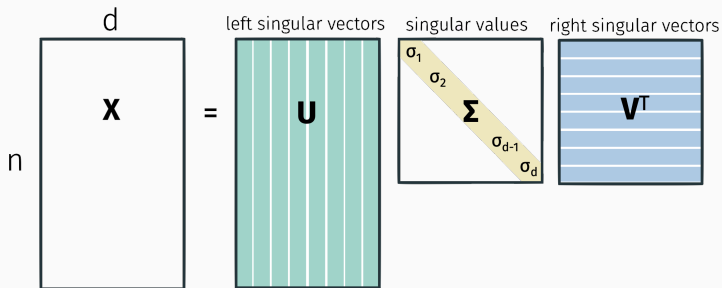
Choose \mathbf{v}_1 to be the top eigenvector of $\mathbf{X}^T \mathbf{X}$.

What about higher k ?

SINGULAR VALUE DECOMPOSITION

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

Any matrix X can be written:



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

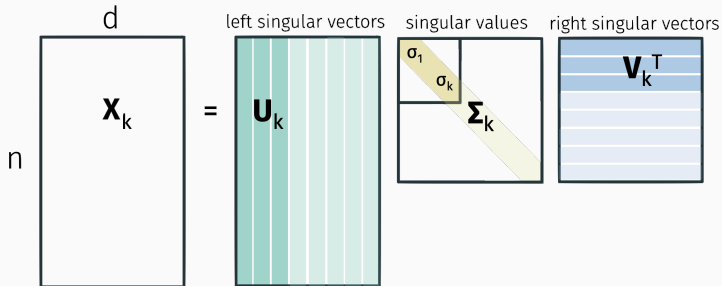
CONNECTION TO EIGENDECOMPOSITION

- \mathbf{U} contains the orthonormal eigenvectors of $\mathbf{X}\mathbf{X}^T$.
- \mathbf{V} contains the orthonormal eigenvectors of $\mathbf{X}^T\mathbf{X}$.
- $\sigma_i^2 = \lambda_i(\mathbf{X}\mathbf{X}^T) = \lambda_i(\mathbf{X}^T\mathbf{X})$

This can be checked directly:

SINGULAR VALUE DECOMPOSITION

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



$$\mathbf{X}_k = \mathbf{U}_k \mathbf{U}_k^T \mathbf{X}_k = \mathbf{X}_k \mathbf{V}_k \mathbf{V}_k^T$$

$$\mathbf{V}_k = \underset{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}}{\arg \min} \|\mathbf{X} - \mathbf{X} \mathbf{V} \mathbf{V}^T\|_F^2 = \underset{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}}{\arg \max} \|\mathbf{X} \mathbf{V} \mathbf{V}^T\|_F^2$$

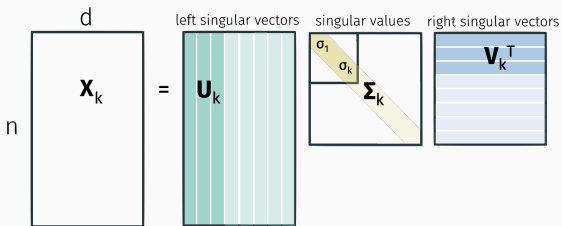
SINGULAR VALUE DECOMPOSITION

- \mathbf{V}_k 's columns are called the “top right singular vectors of \mathbf{X} ”
- \mathbf{U}_k 's columns are called the “top left singular vectors of \mathbf{X} ”
- $\sigma_1, \dots, \sigma_k$ are the “top singular values”. $\sigma_1, \dots, \sigma_d$ are sometimes called the “spectrum of \mathbf{X} ” (although this is more typically used to refer to eigenvalues).

Connection to **Principal Component Analysis**:

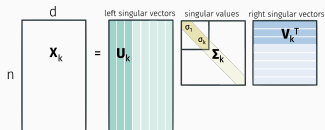
- Let $\bar{\mathbf{X}} = \mathbf{X} - \mathbf{1}\boldsymbol{\mu}^T$ where $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$. I.e. $\bar{\mathbf{X}}$ is obtained by mean centering \mathbf{X} 's rows.
- Let $\bar{\mathbf{U}}\bar{\boldsymbol{\Sigma}}\bar{\mathbf{V}}^T$ be the SVD of $\bar{\mathbf{X}}$.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



Observation 2: The optimal low-rank approximation error $E_k = \|\mathbf{X} - \mathbf{U}_k \mathbf{U}_k^T \mathbf{X}\|_F^2 = \|\mathbf{X} - \mathbf{X} \mathbf{V}_k \mathbf{V}_k^T\|_F^2$ can be written:

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

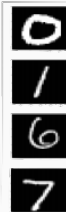
Observation 2: The optimal low-rank approximation error

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

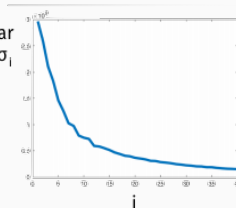
$$E_k = \sum_{i=k+1}^d \sigma_i^2.$$

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

784 dimensional vectors



singular
value σ_i



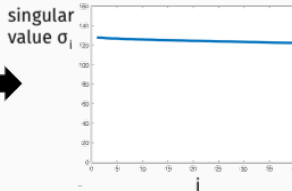
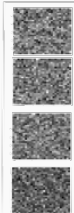
Observation 2: The optimal low-rank approximation error

$E_k = \|\mathbf{X} - \mathbf{U}_k \mathbf{U}_k^T \mathbf{X}\|_F^2 = \|\mathbf{X} - \mathbf{X} \mathbf{V}_k \mathbf{V}_k^T\|_F^2$ 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

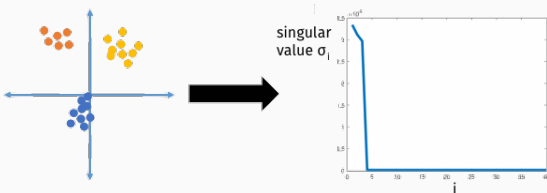


Observation 2: The optimal low-rank approximation error

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



Suffices to compute \mathbf{V} . Then $\mathbf{U}\mathbf{\Sigma} = \mathbf{X}\mathbf{V}$.

- Compute $\mathbf{X}^T\mathbf{X}$.
- Find eigendecomposition $\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = \mathbf{X}^T\mathbf{X}$.
- Compute $\mathbf{L} = \mathbf{X}\mathbf{V}$. Set $\sigma_j = \|\mathbf{L}_j\|_2$ and $\mathbf{U}_j = \mathbf{L}_j/\|\mathbf{L}_j\|_2$.

Total runtime \approx

COMPUTING THE SVD (FASTER)

- Use an iterative algorithm.
- 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.

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

Today: What about when $k = 1$?

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

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

Power method:

- Choose $\mathbf{z}^{(0)}$ randomly. E.g. $\mathbf{z}_0 \sim \mathcal{N}(0, 1)$.
- For $i = 1, \dots, T$
 - $\mathbf{z}^{(i)} = \mathbf{X}^T \cdot (\mathbf{X}\mathbf{z}^{(i-1)})$
 - $n_i = \|\mathbf{z}^{(i)}\|_2$
 - $\mathbf{z}^{(i)} = \mathbf{z}^{(i)} / n_i$

Return \mathbf{z}_T

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

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

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

Claim:

$$\mathbf{z}^{(1)} = \frac{1}{n_1} [c_1 \cdot \sigma_1^2 \mathbf{v}_1 + c_2 \cdot \sigma_2^2 \mathbf{v}_2 + \dots + c_d \cdot \sigma_d^2 \mathbf{v}_d]$$

Claim:

$$\mathbf{z}^{(T)} = \frac{1}{\prod_{i=1}^T n_i} [c_1 \cdot \sigma_1^{2T} \mathbf{v}_1 + c_2 \cdot \sigma_2^{2T} \mathbf{v}_2 + \dots + c_d \cdot \sigma_d^{2T} \mathbf{v}_d]$$

Theorem (Basis 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:

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

Total runtime:

First observation: For all i

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

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

Corollary:

$$\max_j \frac{c_j}{c_1} \leq O(d^3).$$

POWER METHOD FORMAL CONVERGENCE

$$\mathbf{z}^{(T)} = \frac{1}{\prod_{i=1}^T n_i} c_1 \cdot \sigma_1^{2T} \mathbf{v}_1 + \frac{1}{\prod_{i=1}^T n_i} c_2 \cdot \sigma_2^{2T} \mathbf{v}_2 + \dots + \frac{1}{\prod_{i=1}^T n_i} c_d \cdot \sigma_d^{2T} \mathbf{v}_d$$

POWER METHOD FORMAL CONVERGENCE

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

- $\alpha_1 \leq 1$.
- $\alpha_j^2 \leq (\epsilon/2d)^2$ for $j \geq 2$.
- $\alpha_1^2 \geq 1 - d \cdot (\epsilon/2d)^2 \implies \alpha_1 \geq 1 - \epsilon/2$.

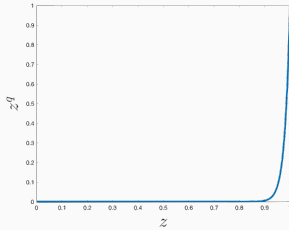
$$\|\mathbf{v}_1 - \mathbf{z}^{(T)}\|_2 \leq$$

Theorem (Basis Power Method Convergence)

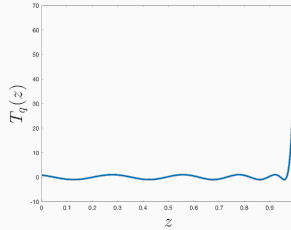
If Power Method is initialized with a random Gaussian vector then, with high probability, after $T = O\left(\frac{\log d/\epsilon}{\epsilon}\right)$ steps, we obtain a \mathbf{z} satisfying:

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

KRYLOV SUBSPACE METHODS



VS.



Lanczos method, Arnoldi method, etc. require $T = O\left(\frac{\log d/\epsilon}{\sqrt{\epsilon}}\right)$ steps for the same guarantee.

GENERALIZATIONS TO LARGE k

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

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

to obtain a nearly optimal low-rank approximation.