CS-GY 6763: Lecture 11 Linear Programming, Singular Value Decomposition

NYU Tandon School of Engineering, Prof. Christopher Musco

Consider a convex function $f(\mathbf{x})$ be bounded between [-B, B] on a constraint set S.

Theorem (Dimension Dependent Convex Optimization) The Center-of-Gravity Method finds $\hat{\mathbf{x}}$ satisfying $f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x} \in S} f(\mathbf{x}) + \epsilon$ using $O(d \log(B/\epsilon))$ calls to a function and gradient oracle for convex f.

The center-of-gravity method is not computationally efficient, but inspired the polynomial time <u>ellipsoid method</u>.



Linear programs (LPs) are one of the most basic convex constrained, convex optimization problems:

Let $\mathbf{c} \in \mathbb{R}^d$, $\mathbf{b} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{n \times d}$ be fixed vectors that define the problem, and let \mathbf{x} be our variable parameter.

 $\min f(\mathbf{x}) = \mathbf{c}^{\mathsf{T}} \mathbf{x}$
subject to $\mathbf{A} \mathbf{x} \ge \mathbf{b}$.

Think about $Ax \ge b$ as a union of half-space constraints:

$$\{\mathbf{x} : \mathbf{a}_1^T \mathbf{x} \ge b_1\}$$
$$\{\mathbf{x} : \mathbf{a}_2^T \mathbf{x} \ge b_2\}$$
$$\vdots$$
$$\{\mathbf{x} : \mathbf{a}_n^T \mathbf{x} \ge b_n\}$$

- Classic optimization applications: industrial resource optimization problems were important original appications in the 70s.
- + Robust regression: $min_{x} \|Ax b\|_{1}$
- L1 constrained regression: $\min_{\mathbf{x}} ||\mathbf{x}||_1$ subject to $A\mathbf{x} = \mathbf{b}$. Lots of applications in sparse recovery/compressed sensing.
- $\cdot \ \text{Solve min}_x \, \|Ax b\|_\infty.$
- Polynomial time algorithms for Markov Decision Processes (reinforcement learning).
- Many combinatorial optimization problems can be solved via <u>LP relaxation</u>.

Theorem (Khachiyan, 1979)

Assume n = d. The ellipsoid method solves any linear program with L-bit integer valued constraints exactly in $O(n^4L)$ time.

A Soviet Discovery Rocks World of Mathematics

By MALCOLM W. BROWNE

A surprise discovery by an obscure Soviet mathematician has rocked the world of mathematics and computer analysis, and experts have begun exploring its practical applications.

Mathematicians describe the discoverv by L.G. Khachian as a method by which computers can find guaranteed solutions to a class of very difficult problems that have hitherto been tackled on a kind of hit-or-miss basis.

Apart from its profound theoretical interest, the discovery may be applicable sometimes involves so many steps that it

in weather prediction, complicated indus- could take billions of years to compute. trial processes, petroleum refining, the scheduling of workers at large factories. secret codes and many other things.

"I have been deluged with calls from virtually every department of government for an interpretation of the significance of this," a leading expert on computer methods, Dr. George B. Dantzig of Stanford University, said in an interview.

The solution of mathematical problems by computer must be broken down into a series of steps. One class of problem

The Russian discovery offers a way by which the number of steps in a solution can be dramatically reduced. It also offers the mathematician a way of learning quickly whether a problem has a solution or not, without having to complete the entire immense computation that may be required.

According to the American journal Sci-

Continued on Page A20, Column 3

ONLY \$10.00 A MONTH 24 Hr. Phone Answering Service, Totally New Concept" Increable!! 279-3870-ADVT.

Front page of New York Times, November 9, 1979.

Theorem (Karmarkar, 1984)

Assume n = d. The <u>interior point method</u> solves any linear program with L-bit integer valued constraints in $O(n^{3.5}L)$ time.



Front page of New York Times, November 19, 1984.

Lecture notes are posted on the website (optional reading).



Projected Gradient Descent Optimization Path

Lecture notes are posted on the website (optional reading).



Ideal Interior Point Optimization Path

Both results had a huge impact on the theory of optimization, although at the time neither the ellipsoid method or interior point method were faster than a heuristic known at the Simplex Method.

These days, improved interior point methods compete with and often outperform simplex.

Polynomial time linear programming algorithms have also had a huge impact of <u>combinatorial optimization</u>. They are often the work-horse behind approximation algorithms for NP-hard problems. Given a graph G with n nodes and edge set E. Each node is assigned a weight w_1, \ldots, w_n .



Goal: Select subset of nodes with minimum total weight that covers all edges.

NP-hard to solve exactly.



Given a graph G with n nodes and edge set E. Each node is assigned a weight w_1, \ldots, w_n .

Formally: Denote if node *i* is selected by assigning variable x_i to 0 or 1. Let $\mathbf{x} = [x_1, \dots, x_n]$.

$$\min_{\mathbf{x}} \sum_{i=1}^{n} x_i w_i \quad \text{subject to} \quad x_i \in \{0, 1\} \text{ for all } i$$
$$x_i + x_j \ge 1 \text{ for all } (i, j) \in E$$

We will use convex optimization give a 2-approximation in polynomial time.

Function to minimize is linear (so convex) but constraint set is not convex. Why?

High level approach:

- <u>Relax</u> to a problem with convex constraints.
- <u>Round</u> optimal solution of convex problem back to original constraint set.



High level approach:

- \cdot <u>Relax</u> to a problem with convex constraints.
- <u>Round</u> optimal solution of convex problem back to original constraint set.



High level approach:

- <u>Relax</u> to a problem with convex constraints.
- <u>Round</u> optimal solution of convex problem back to original constraint set.

Let $\overline{S} \supseteq S$ be the relaxed constraint set. Let $\mathbf{x}^* = \arg \min_{\mathbf{x} \in S} f(\mathbf{x})$ and let $\overline{\mathbf{x}}^* = \arg \min_{\mathbf{x} \in \overline{S}} f(\mathbf{x})$. We always have that:

$$f(\bar{\mathsf{X}}^*) \leq f(\mathsf{X}^*).$$

So typically the goal is to round \bar{x}^* to ${\cal S}$ in such a way that we don't increase the function value too much.

RELAXING VERTEX COVER

Vertex Cover:



$$x_i + x_j \ge 1$$
 for all $(i, j) \in E$

Relaxed Vertex Cover:



The second problem is a linear program! It can be solved in poly(n) time!

Any ideas on how to round this to a solution to the original problem? I.e., with constraints $x_i \in \{0, 1\}$ for all *i*.



Simply set all variable $x_i = 1$ of $\bar{x}_i^* \ge 1/2$ and $x_i = 0$ otherwise.



Observation 1: All edges remain covered. I.e., the constraint $x_i + x_j \ge 1$ for all $(i, j) \in E$ is not violated.

Observation 2: Let **x** be the rounded version of $\bar{\mathbf{x}}^*$. We have $f(\mathbf{x}) \leq 2 \cdot f(\bar{\mathbf{x}})$, and thus $f(\mathbf{x}) \leq 2 \cdot f(\mathbf{x}^*)$.

Proof:

So, a polynomial time algorithm for solving LPs immediately yields a 2-approximation algorithm for the NP-hard problem of vertex cover.

- Proven that it is NP-hard to do better than a 1.36 approximation in [Dinur, Safra, 2002].
- Recently improved to $\sqrt{2} \approx$ 1.41 in [Khot, Minzer, Safra 2018], which proved the 2-to-2 games conjecture.
- Widely believed that doing better than 2ϵ is NP-hard for any $\epsilon > 0$, and this is implied by Subhash Khot's Unique Games Conjecture.

There is a simpler greedy 2-approximation algorithm that doesn't use optimization at all!

BREAK

Next section of course: <u>Spectral methods</u> and <u>numerical linear</u> <u>algebra</u>.

Spectral methods generally refer to methods based on the "spectrum" of a matrix. I.e. on it's eigenvectors/eigenvalues and singular vectors/singular values. We will look at

- Applications to low-rank approximation and dimensionality reduction.
- Applications to graph problems.
- Fast algorithms for computing spectral information.

Reminder: A vector $\mathbf{v} \in \mathbb{R}^d$ is an <u>eigenvector</u> of a matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$, if there exists a scalar λ such that

 $Xv = \lambda v$

The scalar λ is called the <u>eigenvalue</u> associated with **v**.

Matrices can often be written completely in terms of their eigenvectors and eigenvalues. This is called eigendecomposition.

We will actually focus on a related tool called <u>singular value</u> decomposition.

If a <u>square</u> matrix has orthonormal rows, it also has orthonormal columns:

$$\mathsf{V}^{\mathsf{T}}\mathsf{V}=\mathsf{I}=\mathsf{V}\mathsf{V}^{\mathsf{T}}$$

$$\begin{bmatrix} -0.62 & 0.78 & -0.11 \\ -0.28 & -0.35 & -0.89 \\ -0.73 & -0.52 & 0.44 \end{bmatrix} \cdot \begin{bmatrix} -0.62 & -0.28 & -0.73 \\ 0.78 & -0.35 & -0.52 \\ -0.11 & -0.89 & 0.44 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

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

Same thing goes for Frobenius norm: for any matrix **X**, $\|\mathbf{V}\mathbf{X}\|_{F}^{2} = \|\mathbf{X}\|_{F}^{2}$ and $\|\mathbf{V}^{T}\mathbf{X}\|_{F}^{2} = \|\mathbf{X}\|_{F}^{2}$.

The same is not true for rectangular matrices.

$$\mathbf{V}^{\mathsf{T}} \quad \mathbf{V} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{V} \quad \mathbf{V}^{\mathsf{T}} = \begin{bmatrix} 5 & -1 & .7 & -2 \\ 1.6 & -44 & 4.2 & -1.5 \\ 7.8 & .42 & -5 & .67 \\ -2 & 2.0 & 1.1 & 8.0 \\ -1.5 & .55 & 3.2 & .5 \\ .67 & -2.8 & -2.4 & 1.6 \\ 9.0 & 8.7 & -7.7 & 7.8 \end{bmatrix}$$



For any **x**, $\|\mathbf{V}\mathbf{x}\|_{2}^{2} = \|\mathbf{x}\|_{2}^{2}$ but $\|\mathbf{V}^{\mathsf{T}}\mathbf{x}\|_{2}^{2} \neq \|\mathbf{x}\|_{2}^{2}$ in general.

Multiplying a vector by **V** with orthonormal columns <u>rotates</u> and/or reflects the vector.



Multiplying a vector by a rectangular matrix **V**^T with orthonormal rows <u>projects</u> the vector (representing it as coordinates in the lower dimensional space).



So we always have that $\|\mathbf{V}^{\mathsf{T}}\mathbf{x}\|_2 \leq \|\mathbf{x}\|_2$.

One of the most fundamental results in linear algebra. Any matrix **X** can be written:



Where $\mathbf{U}^{\mathsf{T}}\mathbf{U} = \mathbf{I}$, $\mathbf{V}^{\mathsf{T}}\mathbf{V} = \mathbf{I}$, and $\sigma_1 \geq \sigma_2 \geq \ldots \sigma_d \geq 0$.

Singular values are unique. Factors are not. E.g. would still get a valid SVD by multiplying both i^{th} column of V and U by -1.

Important <u>take away</u> from singular value decomposition.

Multiplying any vector **a** by a matrix **X** to form **Xa** can be viewed as a composition of 3 operations:

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

SINGULAR VALUE DECOMPOSITION: ROTATE/REFLECT



A square matrix has at most *d* linearly independent eigenvectors. If a matrix has a full set of *d* eigenvectors $\mathbf{v}_1, \ldots, \mathbf{v}_d$ with eigenvalues $\lambda_1, \ldots, \lambda_d$ it is called "diagonalizable" and can be written as:

$V\Lambda V^{-1}$.

V's columns are $\mathbf{v}_1, \ldots, \mathbf{v}_d$.

Singluar value decomposition

- Exists for all matrices, square or rectangular.
- Singular values are always positive.
- Factors **U** and **V** are orthogonal.

Eigendecomposition

- Exists for <u>some</u> square matrices.
- Eigenvalues can be positive, negative, or imaginary. Real if **X** is symmetric.
- Factor V is orthogonal if and only if X is symmetric.

CONNECTION TO EIGENDECOMPOSITION

- U contains the orthogonal eigenvectors of XX^{T} .
- V contains the orthogonal eigenvectors of $X^T X$.

•
$$\sigma_i^2 = \lambda_i (\mathbf{X}\mathbf{X}^T) = \lambda_i (\mathbf{X}^T\mathbf{X})$$

Lots of applications.

- Compute pseudoinverse $V\Sigma^{-1}U^{T}$.
- Read off condition number of **X**, σ_1^2/σ_d^2 .
- Compute matrix norms. E.g. $\|\mathbf{X}\|_2 = \sigma_1$, $\|\mathbf{X}\|_F = \sqrt{\sum_{i=1}^d \sigma_i^2}$.
- Compute matrix square root i.e. find a matrix B such that BB^T = X. Used e.g. in sampling from Gaussian with covariance X.
- Principal component analysis.

Killer app: Read off optimal <u>low-rank</u> approximation for X.
The column span of a matrix $X \in \mathbb{R}^{n \times d}$ is the set of all vectors that can be written as Xa for some a.

The dimension of the column span, D_c , is the maximum number of linear independent vectors in that set.

The row span of a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ is the set of all vectors that can be written as $\mathbf{b}^T \mathbf{X}$ for some **b**.

The dimension of the row span, D_r , is the maximum number of linear independent vectors in that set.

RANK

For a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ we have:

 $D_c \le d$ $D_r \le n$ $D_c = D_r.$

We call the value of $D_c = D_r$ the <u>rank</u> of **X**.

We always have that:

 $rank(A \cdot B \cdot C \cdot ...) \le min(rank(A), rank(B), rank(C), ...)$.

LOW-RANK APPROXIMATION

Approximate **X** as a rank *k* matrix:



Choose C and B to minimize:

$$\min_{B,C} \|X - CB\|$$

for some matrix norm. Common choice is $\|\mathbf{X} - \mathbf{CB}\|_{F}^{2}$.

APPLICATIONS OF LOW-RANK APPROXIMATION



- **CB** takes O(k(n + d)) space to store instead of O(nd).
 - Important in many applications, including e.g. LoRA: Low-Rank Adaptation of Large Language Models
 - Can be used to compress vector databases.
 - Many more applications.
- Many linear algebraic problems involving **CB** can be solved in $O(nk^2)$ instead of $O(nd^2)$ time.

Without loss of generality can assume that the right matrix is orthogonal. I.e. W^T with $W^TW = I$



Then we should choose left matrix **C** to minimize:

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

This is just *n* least squares regression problems!

$$\mathbf{c}_i = \underset{\mathbf{c}}{\operatorname{arg\,min}} \|\mathbf{W}\mathbf{c} - \mathbf{x}_i\|_2^2$$

$$\mathbf{c}_i = \mathbf{W}^T \mathbf{x}_i$$

 $\mathbf{C} = \mathbf{X} \mathbf{W}$

So our optimal low-rank approximation always has the form: $\mathbf{X} \approx \mathbf{X} \mathbf{W} \mathbf{W}^{T}$

WW^T is a symmetric projection matrix.



C = XW can be used as a meaningful compressed version of data matrix X. We have that:

$$\|\mathbf{x}_i - \mathbf{x}_j\|_2 \approx \|\mathbf{W}\mathbf{W}^T\mathbf{x}_i - \mathbf{W}\mathbf{W}^T\mathbf{x}_j\|_2 = \|\mathbf{c}_i - \mathbf{c}_i\|_2$$

So we expect that:

- · $\|\mathbf{x}_i\|_2 \approx \|\mathbf{c}_i\|_2$
- $\langle \mathbf{x}_i, \mathbf{x}_j \rangle \approx \langle \mathbf{c}_i, \mathbf{c}_j \rangle$
- etc.

How does this compare to Johnson-Lindenstrauss projection?

APPLICATIONS OF LOW-RANK APPROXIMATION

Also useful in:

• Data visualization when k = 2 or 3.



- Data denoising (e.g. distance triangulation).
- Feature selection.

PARTIAL SVD

Key result: Can find the best projection from the singular value decomposition. Note: $X_k = U_k \Sigma_k V_k^T = U_k U_k^T X = X V_k V_k^T$.



45

Goal: Minimize $\|\mathbf{X} - \mathbf{B}\|_{F}$.

Claim 1: Without loss of generality, can assume $B = UZV^T$ for some other rank *k* matrix Z.

Goal: Minimize $\|\mathbf{X} - \mathbf{B}\|_{F}$.

Claim 2: Should choose Z to be the best rank *k* approximation to Σ . (We will then show this equals Σ_k .)

OPTIMAL LOW-RANK APPROXIMATION



Claim 3:

$$\underset{W \in \mathbb{R}^{d \times k}}{\arg \min} \|\mathbf{X} - \mathbf{X} \mathbf{W} \mathbf{W}^{\mathsf{T}}\|_{F}^{2} = \underset{W \in \mathbb{R}^{d \times k}}{\arg \max} \|\mathbf{X} \mathbf{W} \mathbf{W}^{\mathsf{T}}\|_{F}^{2}$$

Follows from fact that for <u>all</u> orthogonal W:

$$\|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{W}^{\mathsf{T}}\|_{F}^{2} = \|\mathbf{X}\|_{F}^{2} - \|\mathbf{X}\mathbf{W}\mathbf{W}^{\mathsf{T}}\|_{F}^{2}$$

Claim 3:

$$\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{W}^{\mathsf{T}} \|_{\mathsf{F}}^2 = \|\mathbf{X}\|_{\mathsf{F}}^2 - \|\mathbf{X}\mathbf{W}\mathbf{W}^{\mathsf{T}}\|_{\mathsf{F}}^2$$



Final Step: Let $\mathbf{W}^* \in \mathbb{R}^{d \times k}$ contain the first *k* standard basis vectors. Then we claim that $\mathbf{W}^* = \arg \max_{\mathbf{W}} \|\mathbf{\Sigma} \mathbf{W} \mathbf{W}^T\|_F^2$.

USEFUL OBSERVATIONS



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

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

SPECTRAL PLOTS

Observation: The optimal low-rank approximation error $E_k = \|\mathbf{X} - \mathbf{X}_k\|_F^2 = \|\mathbf{X}\|_F^2 - \|\mathbf{X}_k\|_F^2$ 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:



SPECTRAL PLOTS

Observation: The optimal low-rank approximation error $E_k = \|\mathbf{X} - \mathbf{X}_k\|_F^2 = \|\mathbf{X}\|_F^2 - \|\mathbf{X}_k\|_F^2$ 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 right singular vectors **V**:

- Compute $\mathbf{X}^T \mathbf{X}$.
- Find eigendecomposition VAV^T = X^TX using e.g. QR algorithm.
- Compute $\mathbf{L} = \mathbf{XV}$. Set $\sigma_i = \|\mathbf{L}_i\|_2$ and $\mathbf{U}_i = \mathbf{L}_i / \|\mathbf{L}_i\|_2$.

Total runtime pprox

COMPUTING THE SVD (FASTER)

How to go faster?

- Compute <u>approximate</u> solution.
- Only compute top *k* singular vectors/values.
- <u>Iterative algorithms</u> achieve runtime $\approx O(ndk)$ vs. $O(nd^2)$ 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.

```
Today: Consider simlest case 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}^T.
```

Power method:

- Choose $z^{(0)}$ randomly. $z_0 \sim \mathcal{N}(0,1).$
- $\cdot \ z^{(0)} = z^{(0)} / \|z^{(0)}\|_2$
- For i = 1, ..., T
 - $\mathbf{z}^{(i)} = \mathbf{X}^{\mathsf{T}} \cdot (\mathbf{X} \mathbf{z}^{(i-1)})$

•
$$n_i = \|\mathbf{z}^{(i)}\|_2$$

•
$$\mathbf{z}^{(i)} = \mathbf{z}^{(i)}/n_i$$

Return **z**^(T)

POWER METHOD INTUITION



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:

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

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

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

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

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

$$\vdots$$

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

Note: $[c_1^{(i)}, \dots, c_d^{(i)}] = \mathbf{c}^{(i)} = \mathbf{V}^T \mathbf{z}^{(i)}$. Also: Since V is orthogonal and $\|\mathbf{z}^{(i)}\|_2 = 1$, $\|\mathbf{c}^{(i)}\|_2^2 = 1$.

ONE STEP ANALYSIS OF POWER METHOD

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

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

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

Equivalently: $\mathbf{c}^{(i)} = \frac{1}{n_i} \mathbf{\Sigma}^2 \mathbf{c}^{(i-1)}$.

Claim: After T updates:

$$\mathbf{z}^{(T)} = \frac{1}{\prod_{i=1}^{T} n_i} \left[c_1^{(0)} \sigma_1^{2T} \cdot \mathbf{v}_1 + c_2^{(0)} \sigma_2^{2T} \cdot \mathbf{v}_2 + \ldots + c_d^{(0)} \sigma_d^{2T} \cdot \mathbf{v}_d \right]$$

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 $j \neq 1$.

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 $\left|\frac{\alpha_i}{\alpha_1}\right| \leq \sqrt{\frac{\epsilon}{2d}}$ then we will have that $\|\mathbf{v}_1 - \mathbf{z}^{(T)}\|_2^2 \leq \epsilon$.

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

$$1 = \alpha_1^2 + \sum_{j=2}^d \alpha_d^2 \le \alpha_1^2 + \frac{\epsilon}{2}$$

$$\alpha_1^2 \ge 1 - \frac{\epsilon}{2}$$

$$|\alpha_1| \ge 1 - \frac{\epsilon}{2}$$

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

Let's see how many steps *T* it takes to ensure that $\left|\frac{\alpha_{j}}{\alpha_{1}}\right| \leq \sqrt{\frac{\epsilon}{2d}}$ where $\alpha_{j} = \frac{1}{\prod_{i=1}^{T} n_{i}} c_{j}^{(0)} \sigma_{j}^{2T}$. Answer will depend on $\gamma = \frac{\sigma_{1} - \sigma_{2}}{\sigma_{1}}$. Assumption: Starting coefficient on first eigenvector is not too small: $\left|c_{1}^{(0)}\right| \geq O\left(\frac{1}{\sqrt{d}}\right)$.

We will prove shortly that it holds with probability 99/100.

$$\frac{|\alpha_j|}{|\alpha_1|} = \frac{\sigma_j^{2^{\intercal}}}{\sigma_1^{2^{\intercal}}} \cdot \frac{|c_j^{(0)}|}{|c_1^{(0)}|} \le$$

Need to set T =

Need to prove: Starting coefficient on first eigenvector is not too small. I.e., with probability 99/100,

$$\left|c_{1}^{(0)}\right| \geq O\left(\frac{1}{\sqrt{d}}\right).$$

Prove using Gaussian <u>anti</u>-concentration. First use rotational invariance of Gaussian:

$$\mathbf{c}^{(0)} = \frac{\mathbf{V}^{\mathsf{T}} \mathbf{z}^{(0)}}{\|\mathbf{z}^{(0)}\|_2} = \frac{\mathbf{V}^{\mathsf{T}} \mathbf{z}^{(0)}}{\|\mathbf{V}^{\mathsf{T}} \mathbf{z}^{(0)}\|_2} \sim \frac{\mathbf{g}}{\|\mathbf{g}\|_2}$$

where $\mathbf{g} \sim \mathcal{N}(0, 1)^d$.

Need to show that with high probability, first entry of $\frac{g}{\|g\|_2} \ge c \cdot \frac{1}{\sqrt{d}}$.

Part 1: With super high probability (e.g. 99/100),

$\|\boldsymbol{g}\|_2^2 \leq$

Need to show that with high probability, the magnitude of the first entry of $\mathbf{g} \ge c$ for a constant c. Think e.g. c = 1/100.

Part 2: With probablility $1 - O(\alpha)$,

 $|g_1| \geq \alpha.$



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:

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

The method truly won't converge if γ is very small. Consider extreme case when $\gamma = 0$.

$$\mathbf{z}^{(T)} = \frac{1}{\prod_{i=1}^{T} n_i} \left[c_1^{(0)} \sigma_1^{2T} \cdot \mathbf{v}_1 + c_2^{(0)} \sigma_2^{2T} \cdot \mathbf{v}_2 + \ldots + c_d^{(0)} \sigma_d^{2T} \cdot \mathbf{v}_d \right]$$

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 **z** satisfying:

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

Intuition: For a good low-rank approximation, we don't actually need to converge to \mathbf{v}_1 if σ_1 and σ_2 are the same or very close. Would suffice to return either \mathbf{v}_1 or \mathbf{v}_2 , or some linear combination of the two.

• 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.
- $Z_0 = orth(G)$.
- For i = 1, ..., T
 - $Z^{(i)} = X^T \cdot (XZ^{(i-1)})$
 - · $Z^{(i)} = orth(Z^{(i)})$

Return **Z**^(T)

Guarantee: After
$$O\left(\frac{\log d/\epsilon}{\epsilon}\right)$$
 iterations:
 $\|\mathbf{X} - \mathbf{X}\mathbf{Z}\mathbf{Z}^{\mathsf{T}}\|_{F}^{2} \leq (1 + \epsilon)\|\mathbf{X} - \mathbf{X}\mathbf{V}_{\mathsf{k}}\mathbf{V}_{\mathsf{k}}^{\mathsf{T}}\|_{F}^{2}$.

Runtime: $O(nnz(X) \cdot k \cdot T) \leq O(ndk \cdot T)$.

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:

$$\|\mathbf{X} - \mathbf{X}\mathbf{Z}\mathbf{Z}^{\mathsf{T}}\|_{F}^{2} \leq (1+\epsilon)\|\mathbf{X} - \mathbf{X}\mathbf{V}_{\mathbf{k}}\mathbf{V}_{\mathbf{k}}^{\mathsf{T}}\|_{F}^{2}.$$

For a normalizing constant c, power method returns:

$$\mathsf{z}^{(q)} = c \cdot \left(\mathsf{X}^{\mathsf{T}}\mathsf{X}\right)^{q} \cdot \mathsf{g}$$

Along the way we computed:

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

 \mathcal{K} is called the <u>Krylov subspace of degree q</u>.

Idea behind Krlyov methods: Don't throw away everything before $(\mathbf{X}^T \mathbf{X})^q \cdot \mathbf{g}$.
Want to find **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 $\min_{v=Qw} \|\mathbf{X} \mathbf{X} v v^T\|_F^2$.
 - Find <u>best</u> vector in the Krylov subspace, instead of just using last vector.
 - Can be done in $O(ndk + dk^2)$ time.
 - What you're using when you run **svds** or **eigs** in MATLAB or Python.

For a degree *t* polynomial *p*, let $\mathbf{v}_p = \frac{p(\mathbf{X}^T \mathbf{X})\mathbf{g}}{\|p(\mathbf{X}^T \mathbf{X})\mathbf{g}\|_2}$. We always have that $\mathbf{v}_p \in \mathcal{K}_t$, the Krylov subspace contructed with *t* iterations. Power method returns:

$$\mathbf{v}_p$$
 where $p = x^q$ for $q = 2T$.

Lanczos method returns \mathbf{v}_{p^*} where:

$$p^* = \underset{\text{degree } t \ p}{\arg \min} \|\mathbf{X} - \mathbf{X} \mathbf{v}_p \mathbf{v}_p^T\|_F^2.$$

Claim: There is a $t = O\left(\sqrt{q \log \frac{1}{\Delta}}\right)$ degree polynomial \hat{p} approximating \mathbf{x}^q up to error Δ on $[0, \sigma_1^2]$.



$$\begin{aligned} \|\mathbf{X} - \mathbf{X}\mathbf{v}_{p^*}\mathbf{v}_{p^*}^T\|_F^2 &\leq \|\mathbf{X} - \mathbf{X}\mathbf{v}_{\hat{p}}\mathbf{v}_{\hat{p}}^T\|_F^2 \approx \|\mathbf{X} - \mathbf{X}\mathbf{v}_{x^q}\mathbf{v}_{x^q}^T\|_F^2 \approx \|\mathbf{X} - \mathbf{X}\mathbf{v}_{1}\mathbf{v}_{1}^T\|_F^2 \\ \text{Runtime: } O\left(\frac{\log(d/\epsilon)}{\sqrt{\epsilon}} \cdot \mathsf{nnz}(\mathbf{X})\right) \text{ vs. } O\left(\frac{\log(d/\epsilon)}{\epsilon} \cdot \mathsf{nnz}(\mathbf{X})\right) \end{aligned}$$

GENERALIZATIONS TO LARGER k

- Block Krylov methods
- Let $\mathbf{G} \in \mathbb{R}^{d \times k}$ be a random Gaussian matrix.

$$\boldsymbol{\cdot} \ \mathcal{K}_{q} = \left[\boldsymbol{G}, \left(\boldsymbol{X}^{T}\boldsymbol{X}\right) \cdot \boldsymbol{G}, \left(\boldsymbol{X}^{T}\boldsymbol{X}\right)^{2} \cdot \boldsymbol{G}, \ldots, \left(\boldsymbol{X}^{T}\boldsymbol{X}\right)^{q} \cdot \boldsymbol{G}\right]$$

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