CS-GY 6763: Lecture 10 Singular value decomposition, low-rank approximation, Krylov subspace methods

NYU Tandon School of Engineering, Prof. Christopher Musco

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 <u>not true</u> 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 $\|\mathbf{V}^T\mathbf{x}\|_2^2 \neq \|\mathbf{x}\|_2^2$ in general.

Multiplying a vector by ${\bf V}$ with orthonormal columns $\underline{rotates}$ and/or reflects the vector.



Multiplying a vector by a rectangular matrix \mathbf{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 $\mathbf{U}^T \mathbf{U} = \mathbf{I}$, $\mathbf{V}^T \mathbf{V} = \mathbf{I}$, and $\sigma_1 \ge \sigma_2 \ge \ldots \sigma_d \ge 0$.

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

Important take away 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



SINGULAR VALUE DECOMPOSITION: STRETCH



SINGULAR VALUE DECOMPOSITION: ROTATE/REFLECT



SINGULAR VALUE DECOMPOSITION



Recall that an eigenvalue of a <u>square</u> matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$ is any vector \mathbf{v} such that $\mathbf{X}\mathbf{v} = \lambda\mathbf{v}$. A 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 $\lambda_1, \ldots, \lambda_n$ it is called "diagonalizable" and can be written as:

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

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 or negative.
- 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 low-rank approximations for X.

LOW-RANK APPROXIMATION

Approximate **X** as the product of two rank *k* matrices:



Typically 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).
- Regression problems involving **CB** can be solved in $O(nk^2)$ instead of $O(nd^2)$ time.
- Will see a bunch more in a minute.

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



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

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

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.



LOW-RANK APPROXIMATION





C = XW can be used as a compressed version of data matrix X.

Let C = XW. We have that:

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

Similarly, we expect that:

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

How does this compare to Johnson-Lindenstrauss projection?

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



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*.



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

APPLICATIONS OF LOW-RANK APPROXIMATION

Fact that $\|\mathbf{x}_i - \mathbf{x}_j\|_2 \approx \|\mathbf{x}_i^T \mathbf{W} \mathbf{W}^T - \mathbf{x}_j^T \mathbf{W} \mathbf{W}^T\|_2 = \|\mathbf{c}_i - \mathbf{c}_i\|_2$ leads to lots of applications.

- Data compression. E.g. used in state-of-the-art data dependence methods for nearest neighbor search.
- Data visualization when k = 2 or 3.





• Entity embeddings (next lecture).

APPLICATIONS OF LOW-RANK APPROXIMATION

• Reduced order modeling for solving physical equations.



- · Constructing preconditioners in optimization.
- · Many more.

Can find the best projection from the singular value decomposition.



Claim:
$$X_k = U_k \Sigma_k V_k^T = X V_k V_k^T$$
.

OPTIMALITY OF SVD



OPTIMALITY OF SVD

Claim 2:

$$\underset{\text{rank }k \ B}{\arg\min} \|\boldsymbol{\Sigma} \boldsymbol{V}^{T} - \boldsymbol{B}\|_{F}^{2} = \underset{\text{rank }k \ B}{\arg\min} \|\boldsymbol{V}\boldsymbol{\Sigma} - \boldsymbol{B}^{T}\|_{F}^{2}$$

Claim 3:

$$\underset{\text{rank }k \text{ B}}{\arg\min} \| \mathbf{V} \mathbf{\Sigma} - \mathbf{B}^T \|_F^2 = \underset{\text{rank }k \text{ B}}{\arg\min} \| \mathbf{\Sigma} - \mathbf{V}^T \mathbf{B}^T \|_F^2$$

Chose \mathbf{B}^T so that $\mathbf{V}^T \mathbf{B}^T = \mathbf{\Sigma}_k$.

USEFUL OBSERVATIONS



Observation 1:

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

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

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

USEFUL OBSERVATIONS

Claim:

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

USEFUL OBSERVATIONS



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:

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

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:



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:



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:



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)

- Compute <u>approximate</u> 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.
- <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.

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

•
$$n_i = \|\mathbf{z}^{(i)}\|_2$$

•
$$z^{(i)} = z^{(i)}/n_i$$

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:

$$\mathbf{z}^{(0)} = c_1^{(0)} \mathbf{v}_1 + c_2^{(0)} \mathbf{v}_2 + \ldots + c_d^{(0)} \mathbf{v}_d$$
$$\mathbf{z}^{(1)} = c_1^{(1)} \mathbf{v}_1 + c_2^{(1)} \mathbf{v}_2 + \ldots + c_d^{(1)} \mathbf{v}_d$$
$$\vdots$$
$$\mathbf{z}^{(i)} = c_1^{(i)} \mathbf{v}_1 + c_2^{(i)} \mathbf{v}_2 + \ldots + c_d^{(i)} \mathbf{v}_d$$

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

Also: $\|\mathbf{c}^{(i)}\|_2^2 = \sum_{j=1}^d (c_j^{(i)})^2 = 1.$

ONE STEP ANALYSIS OF POWER METHOD

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

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

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

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 $|\alpha_1| \le 1$. If we can prove that $\left|\frac{\alpha_i}{\alpha_1}\right| \le \sqrt{\frac{\epsilon}{d}}$ then:

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

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

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

$$|\alpha_1| \ge 1 - \epsilon$$

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

Lets proves that $\left|\frac{\alpha_{j}}{\alpha_{1}}\right| \leq \sqrt{\frac{\epsilon}{d}}$ where $\alpha_{j} = \frac{1}{\prod_{i=1}^{T} n_{i}} c_{j}^{(0)} \sigma_{j}^{2T}$ **Assumption:** Starting coefficients are all <u>roughly</u> equal. For all j $O(1/d^{1.5}) \leq \left|c_{j}^{(0)}\right| \leq 1$. This is a very loose bound, but it's all that we will need. We will prove shortly that it holds with probability 99/100.

$$\frac{|\alpha_j|}{|\alpha_1|} = \frac{\sigma_j^{27}}{\sigma_1^{27}} \cdot \frac{|c_j^{(0)}|}{|c_1^{(0)}|} \le$$

Need T =

Need to prove: Starting coefficients are all roughly equal.

For all
$$j$$
 $O(1/d^{1.5}) \le |c_j^{(0)}| \le 1$

with probability 99/100. Prove using Gaussian (anti)-concentration.

Right hand side is immediate from fact that $\sum_{j} (c_{j}^{(0)})^{2} = 1$.

To show the left hand side we 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, every entry of $\frac{g}{\|g\|_2} \ge c \cdot \frac{1}{d^{1.5}}.$

Part 1: With probablility 999/100,

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

Need to show that with high probability, the magnitude of every entry of $\mathbf{g} \geq c \cdot \frac{1}{d}$.

Part 2: With probablility 1 - c/d,

$$\min_i |g_i| \ge O\left(\frac{c}{d}\right).$$



Applying union bound completes the result.

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
 - $\cdot \ \mathsf{Z}^{(i)} = \mathsf{X}^{\mathsf{T}} \cdot (\mathsf{X} \mathsf{Z}^{(i-1)})$
 - · $Z^{(i)} = orth(Z^{(i)})$

Return **Z**^(T)

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

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

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:

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

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

KRYLOV SUBSPACE METHODS



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

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

Along the way we computed:

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

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

Idea behind Krlyov methods: Don't throw away everything before $(X^TX)^q \cdot g$. What you're using when you run svds or eigs in MATLAB or Python.

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.

For a degree *t* polynomial *p*, let $\mathbf{v}_p = \frac{p(\mathbf{X}^T \mathbf{X})\mathbf{g}}{\|p(\mathbf{X}^T \mathbf{X})\mathbf{g}\|_2}$. Power method returns:

 V_{X^t} .

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 $\Delta \sigma_1^2$ on $[0, \sigma_1^2]$.



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

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

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

So we don't care! Either v_1 or v_2 give good rank-1 approximations.

Claim: To achieve

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

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.

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.