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

NYU Tandon School of Engineering, Prof. Christopher Musco

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



$$V^T V = I = VV^T$$

$$(V^T)^T U^T = I$$

$v_1^T v_j \text{ for } j \neq i$

$\text{this} = 0$

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

Implies that for any vector $\mathbf{x}$, $\underline{\underline{\|\mathbf{V}\mathbf{x}\|_2^2}} = \|\mathbf{x}\|_2^2$ and $\underline{\underline{\|\mathbf{V}^T\mathbf{x}\|_2^2}}$. $\begin{bmatrix} x_1, \ldots x_k \end{bmatrix}$

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

$$\|M\|_F^2 = \sum_{ij} M_{ij}^2$$

$$\| Vx \|_2^2 = \|x\|_2^2$$

$$= (Vx)^T Vx = x^T \underbrace{U^T V}_{=I} x = x^T x = \|x\|_2^2$$

$$VX = \begin{bmatrix} Vx_1, \ldots Vx_k \end{bmatrix}$$

$$\| VX \|_F^2 = \sum_{i=1}^{k} \| Vx_i \|_2^2 = \sum_{i=1}^{u} \|x_i\|_2^2 = \|X\|_F^2$$

The same is <u>not true</u> for rectangular matrices:



$$V^T V = I \qquad \text{but} \qquad VV^T \neq I$$

For any $\mathbf{x}$, $\|V\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$ <u>but</u> $\|V^T\mathbf{x}\|_2^2 \neq \|\mathbf{x}\|_2^2$ in general.

$$x^T V^T V x = x^T x = \|x\|_2^2 \qquad\qquad \|Vx\|_F^2 = \|x\|_F^2$$
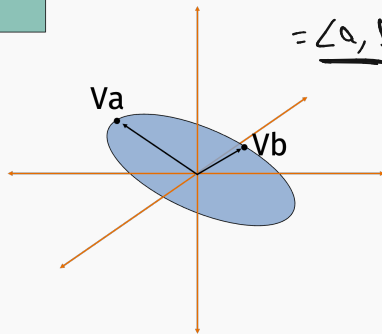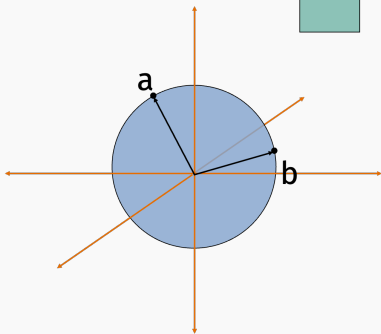
4

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



$$\|Va\|_2^2 = \|a\|_2^2$$

$$\langle Va, Vb \rangle$$
$$= a^\top V^\top V b = a^\top b$$
$$= \underline{\langle a, b \rangle}$$

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 $\|\mathbf{V}^T\mathbf{x}\|_2 \leq \|\mathbf{x}\|_2$.

*One of the most fundamental results in linear algebra.*

Underlined: *Any* matrix $X$ can be written:

Handwritten annotations: *Reduced SVD*, *Economy SVD*, $U$, $\Sigma$, $V^T$, $D$



| d | left singular vectors | singular values | right singular vectors |

Handwritten: $U^T u_i = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

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

$\sigma_1$, $\sigma_2$, $\sigma_{d-1}$, $\sigma_d$

$u_i$, $n$

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

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

Important take away from singular value decomposition.

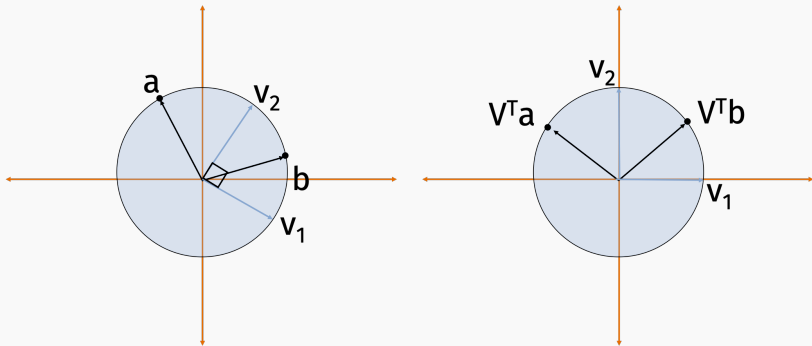Multiplying any vector **a** by a matrix **X** to form **Xa** can be viewed as a composition of 3 operations:

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

$$X_a = U\Sigma V^T a$$
$$= U(\Sigma(V^T a))$$

Recall that an eigenvalue of a <u>square</u> matrix $X \in \mathbb{R}^{d \times d}$ is any vector $v$ such that $\underline{Xv} = \underline{\lambda v}$. A matrix has at most $d$ linearly independent eigenvectors. If a matrix has a full set of $d$ eigenvectors $\underline{v_1, \ldots, v_d}$ with eigenvalues $\lambda_1, \ldots, \lambda_n$ it is called "diagonalizable" and can be written as:

$\rightarrow$ scaler

$X v_i = \lambda_i v_i$

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

$V^{-1} \neq V^T$

$\begin{vmatrix} | & & \\ & | & \\ & & | \end{vmatrix}$

$\overline{\begin{vmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{vmatrix}}$

### 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 <u>X</u> is <u>symmetric</u>.

$$V^{-1} = V^{T} \qquad V \wedge V^{T}$$

- $U$ contains the orthogonal eigenvectors of $XX^T$.
- $V$ contains the orthogonal eigenvectors of $X^TX$.
- $\sigma_i^2 = \lambda_i(XX^T) = \lambda_i(X^TX)$

$$\sigma_1$$
$$\ddots$$
$$\sigma_n$$

$$\Sigma \cdot \Sigma = \begin{matrix} \sigma_1^2 \\ & \ddots \\ & & \sigma_n^2 \end{matrix}$$

$$X = U\Sigma U^T \qquad X^T = V\Sigma U^T \qquad XX^T = U\Sigma\underbrace{V^T V}_{I}\Sigma U^T$$

$$\underline{XX^T} = U\Sigma^2 U^T$$

$$XX^T u_1 = U\Sigma^2 U^T u_1$$
$$= \lambda_1 u_1 \qquad \underbrace{\phantom{U\Sigma^2 U^T u_1}}_{e_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}}$$

$$\begin{bmatrix} \sigma_1^2 \\ & \ddots \\ & & \sigma_n^2 \end{bmatrix}$$

$$= V\Sigma^2 e_1$$

$$= U\begin{bmatrix} \sigma_1^2 \\ 0 \\ 0 \end{bmatrix} = \underline{\sigma_1^2 \cdot u_1}$$

15

Lots of applications.

$$\min_{a} \|Xa - b\|_2^2$$

- Compute pseudoinverse $\underline{V}\underline{\Sigma}^{-1}\underline{U}^T$.

$$\max \frac{\|Ax\|_2}{\|x\|_2}$$

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

Killer app: Read off optimal low-rank approximations for $\underline{X}$.

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



$$(n \times k)(k \times d)$$
$$= (n \times d)$$

Typically choose C and B to minimize:

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

for some matrix norm. Common choice is $\|X - CB\|_F^2$.

17

matrix X ≈ matrix C · matrix B

- **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 $W^T W = I$
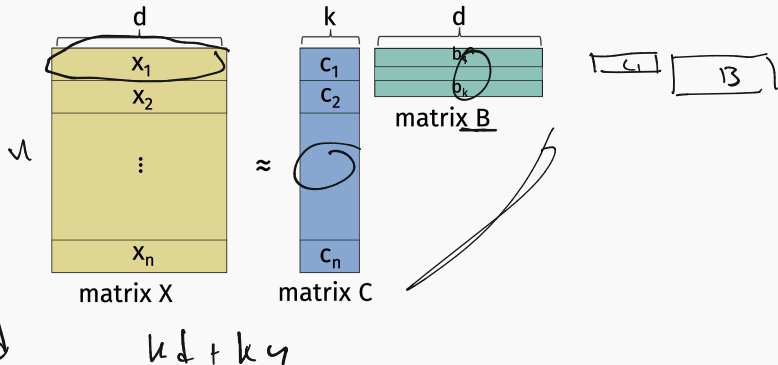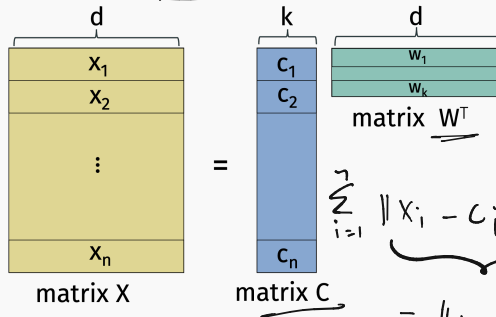
$X_1 \approx C_1 B$
$\downarrow$
$z_1 \omega^T = C_1 B$



d
$$\begin{array}{|c|}\hline x_1 \\ \hline x_2 \\ \hline \vdots \\ \hline x_n \\ \hline \end{array}$$
matrix X

=

k
$$\begin{array}{|c|}\hline c_1 \\ \hline c_2 \\ \hline \\ \hline c_n \\ \hline \end{array}$$
matrix C

d
$$\begin{array}{|c|}\hline w_1 \\ \hline w_k \\ \hline \end{array}$$
matrix $W^T$

$\sum_{i=1}^{n} \| x_i - c_i \omega^T \|_2^2$

$= \| x_i - \omega c_i \|_2^2$

Then we should choose C to minimize:

$$\min_{C, \omega} \| X - CW^T \|_F^2$$

This is just *n* least squares regression problems!

19

$$\min \|Az - b\|_2^2 \qquad z^* = (A^TA)^{-1}A^Tb$$

$$2A^T(Az - b) = 0$$

$$c_i = \arg\min_c \|Wc - x_i\|_2^2$$

$$(\omega^T\omega)^{-1}\omega^T x_i$$

$$I$$

$$c_i = W^Tx_i$$
$$C = XW$$

So our optimal low-rank approximation always has the form:

$$X \approx XWW^T$$

$$\min_\omega \|X - X\omega\omega^T\|_F$$

20

$WW^T$ is a symmetric <u>projection matrix</u>.

$X \omega \omega^T \longrightarrow$ Projection

$P \cdot P = P$



$\tilde{X} = X \omega \omega^T$

$\rightarrow$ row vector

$X_1 \, \omega \omega^T$

$= \omega \omega^T X_1 \mapsto$ column vector

$\omega \omega^T \underbrace{\omega \omega^T X}_{I}$

$= \omega \omega^T X$

21

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

Let $C = XW$. We have that:

$$\tilde{x}_i \quad \tilde{x}_j$$

$$\|x_i - x_j\|_2 \approx \|x_i^T WW^T - x_j^T WW^T\|_2 = \|c_i - c_j\|_2$$

Similarly, we expect that:

$$\|x_i - x_j\|_2 \approx \|c_i - c_j\|_2$$

- $\|x_i\|_2 \approx \|c_i\|_2$
- $\langle x_i, x_j \rangle \approx \langle c_i, c_j \rangle$

$$c_i, c_j \in \left[\begin{array}{c}\end{array}\right] k \qquad x_i, x_j \in \left[\begin{array}{c}\end{array}\right] d$$

- 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*.



784 dimensional vectors    projections onto 15 dimensional space    orthonormal basis $v_1,\ldots,v_{15}$

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 |

Fact that $\|x_i - x_j\|_2 \approx \|x_i^T WW^T - x_j^T WW^T\|_2 = \|c_i - c_j\|_2$ 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).

- Reduced order modeling for solving physical equations.



- Constructing preconditioners in optimization.
- Many more.

$$X_u = U_u U_u^\top X$$

Can find the best projection from the singular value decomposition.

$$X_u = X V_u V_u^\top$$

$$O(ud^2)$$

$$O(udk)$$



d

left singular vectors    singular values    right singular vectors

$\mathbf{X}_k$    =    $\mathbf{U}_k$    $\sigma_1$   $\sigma_k$   $\mathbf{\Sigma}_k$    $\mathbf{V}_k^\top$

n

rank

$$\mathbf{V}_k = \operatorname*{arg\,min}_{\text{orthogonal } \mathbf{W} \in \mathbb{R}^{d \times k}} \|X - XWW^T\|_F^2$$

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

$$X U_k V_k^T = U_u \Sigma_u V_u^T$$

$$X V_u = U_u \Sigma_u$$

$$U \Sigma V^T V_k = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 0 & c \\ 0 & 0 & c & c \\ 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 6_1 & \cdots & \\ & & 6_u \\ & 6 & \end{bmatrix}$$

Claim 1:

$$U \; \overline{\Sigma}_u \; \overline{U}_u^T = U_u \Sigma_u V_u^{-T}$$

$$\arg\min_{\text{rank } k \; \mathbf{B}} \|\mathbf{X} - \mathbf{B}\|_F^2 = U \; \arg\min_{\text{rank } k \; \mathbf{B}} \|\mathbf{\Sigma V}^T - \mathbf{B}\|_F^2$$



$\mathbf{\Sigma V^T}$

$U \Sigma U^T$

Rotated by U
on the left

Claim 2:

$$B^T = \overline{V_k} \,\overline{\Sigma_k}$$

$$B = \Sigma_k V_k^T$$



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

Claim 3:

$$\to \|V^T(V\Sigma - B^T)\|_F^2$$

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

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

$$\|X - XGG^\top\|_F \geq \|X - XU_uU_u^\top\|_F$$



d

left singular vectors    singular values    right singular vectors

$$n \quad X_k \quad = \quad U_k \quad \begin{matrix} \sigma_1 \\ \quad \sigma_k \end{matrix} \Sigma_k \quad V_k^\top$$

Observation 1:

$$\arg\min_{W\in\mathbb{R}^{d\times k}} \|X - XWW^T\|_F^2 = \left( \arg\max_{W\in\mathbb{R}^{d\times k}} \|XWW^T\|_F^2 \right)$$

Follows from fact that for all orthogonal W:

$$\|X - XWW^T\|_F^2 = \|X\|_F^2 - \|XWW^T\|_F^2$$

Pythagorean theorem

33

Claim:

$$\sum_{i=1}^{q} \| x_i - x_i \omega \omega^T \|_2^2$$

$$\| X - XWW^T \|_F^2 = \| X \|_F^2 - \| XWW^T \|_F^2$$



$$\| x_1 \|_2^2 = \| x_1 - x_1 \omega \omega^T \|_2^2 + \| x_1 \omega \omega^T \|_2^2$$

$$\| x_1 - x_1 \omega \omega^T \|_2^2$$
$$= \| x_1 \|_2^2 - \| x_1 \omega \omega^T \|_2^2$$

$\mathbf{x_1}$

$\mathbf{x_1 WW^T}$

$x_1 - x_1 \omega \omega^T$

d

left singular vectors   singular values   right singular vectors

$X_k$ = $U_k$     $V_k^T$

n

$\sigma_1$
$\sigma_k$

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

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

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

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

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

Observation 2: The optimal low-rank approximation error
$E_k = \|X - XV_kV_k^T\|_F^2 = \|X\|_F^2 - \|XV_kV_k^T\|_F^2$ 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 $\sigma_i$

i

Observation 2: The optimal low-rank approximation error
$E_k = \|X - XV_k V_k^T\|_F^2 = \|X\|_F^2 - \|XV_k V_k^T\|_F^2$ 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:



37

**Observation 2:** The optimal low-rank approximation error
$E_k = \|X - XV_kV_k^T\|_F^2 = \|X\|_F^2 - \|XV_kV_k^T\|_F^2$ can be written:
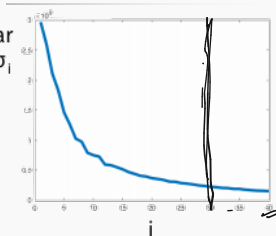
$$E_k = \sum_{i=k+1}^{d} \sigma_i^2.$$

Can immediately get a sense of "how low-rank" a matrix is from it's spectrum:

Suffices to compute right singular vectors $V$:

$X = U \Sigma V^T$

$X = U \Sigma$

- Compute $X^T X$.
- Find eigendecomposition $V \Lambda V^T = X^T X$ using e.g. QR algorithm.
- Compute $L = XV$. Set $\sigma_i = \|L_i\|_2$ and $U_i = L_i/\|L_i\|_2$.

$X \to n \times d$

Total runtime $\approx O(nd^2)$ $+ d^3 \log(\log(1/\epsilon))$

- Compute <u>approximate</u> solution.
- Only compute <u>top $k$ singular vectors/values</u>. 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 $k = 1$?

$$\| z - \boxed{v_1} \|_2^2 \le \varepsilon$$

**Goal:** Find some $z \approx v_1$.

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

$$O(nd) \qquad O(nd)$$

$$X^T(X z^{(i-1)})$$

**Power method:**

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

$$z = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

2nd of

$$d$$

$$n \quad \boxed{\phantom{xxxx}} \quad O(ndT)$$

41

0 iterations    1 iterations    2 iterations

$$Z = V_1$$

$$\frac{X^T X Z}{\|(X^T X Z)\|_2} \quad : \quad \frac{\lambda_1 (k^t x) \cdot V_1}{\| X^T X Z \|_\sim}$$

$$\left( C \cdot V_1 \right)$$

### Theorem (Basic Power Method Convergence)

$\sigma_1 = \sigma_2 : \gamma = 0$

$\sigma_2 \ll \sigma_1 : \gamma$ larger

*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:*

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

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

$\leq d$

43

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

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

$c^{(0)}$

$$z^{(0)} = c_1^{(0)}v_1 + c_2^{(0)}v_2 + \ldots + c_d^{(0)}v_d$$

$V^T c^{(1)}$

$$z^{(1)} = c_1^{(1)}v_1 + c_2^{(1)}v_2 + \ldots + c_d^{(1)}v_d$$

$V$ $\bigg]$

$$\vdots$$

$$z^{(i)} = c_1^{(i)}v_1 + c_2^{(i)}v_2 + \ldots + c_d^{(i)}v_d$$

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

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

44

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

$$c_1^{(i-1)} \dots c_d^{(i-1)}$$

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

$$\mathbf{z}^{(i-1)}$$

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

$$c^{(i)} = \Sigma^2 c^{(i-1)}$$

$$c^{(i-1)} = V^T z^{(i-1)}$$

$$c^{(i)} = V^T z^{(i)} = V^T X^T X z^{(i-1)} \cdot \frac{1}{n_i}$$

$$V\Sigma U^T U \Sigma V^T = V\Sigma^2 V^T \qquad \underbrace{X^T X \Sigma^2 \underbrace{V^T z^{(i-1)}}_{c^{(i-1)}}}$$

45

$$\sigma_1^2 > \sigma_2^2 > \sigma_3^2 > \ldots$$

**Claim:** After $T$ updates:

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

$\approx 0$

Let $\alpha_j = \left( \frac{1}{\prod_{i=1}^{T} n_i} \right) c_j^{(0)} \sigma_j^{2T}$. **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 $\left| \frac{\alpha_j}{\alpha_1} \right| \leq \sqrt{\frac{\epsilon}{d}}$ then:

$$\| \mathbf{z}^{(T)} - \mathbf{U}_1 \|_2^2 \leq 2\epsilon$$

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

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

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

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

$$|\alpha_1| \geq 1 - \epsilon$$

$$\geq (1 - \epsilon)$$

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

Lets proves that $\left|\frac{\alpha_j}{\alpha_1}\right| \leq \sqrt{\frac{\epsilon}{d}}$ where $\alpha_j = \frac{1}{\prod_{l=1}^{T} n_l} c_j^{(0)} \sigma_j^{2T}$

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

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

**Assumption:** Starting coefficients are all roughly equal.

For all $j$ $\qquad O(1/d^{1.5}) \leq \left|c_j^{(0)}\right| \leq 1.$

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.

$$\frac{|\alpha_j|}{|\alpha_1|} = \frac{\sigma_j^{2T}}{\sigma_1^{2T}} \cdot \frac{|c_j^{(0)}|}{|c_1^{(0)}|} \leq \left(\frac{\sigma_j}{\sigma_1}\right)^{2T} \cdot d^{1.5} \leq \left(\frac{\sigma_2}{\sigma_1}\right)^{2T} \cdot d^{1.5}$$

$\nearrow \leq 1$

$\geq \frac{1}{d^{1.5}}$

$\epsilon \ (1-y)^{2T} \cdot d^{1.5} \leq \sqrt{\frac{\epsilon}{d}}$

Need $T = \frac{\log(d/\epsilon)}{y}$

$\left((1-y)^{1/y}\right)^{\log(d/\epsilon)} \rightarrow \left(\frac{1}{e}\right)^{\log(d/\epsilon)} \leq \sqrt{\frac{\epsilon}{d}}$

48

Need to prove: Starting coefficients are all <u>roughly</u> equal.

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

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

$$\frac{z^0}{\|z^0\|_2}$$

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:

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

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

49

Need to show that with high probability, every entry of
$\frac{\mathbf{g}}{\|\mathbf{g}\|_2} \geq c \cdot \frac{1}{d^{1.5}}$.

Part 1: With probablility 999/100,

$$\|\mathbf{g}\|_2 \leq \sigma(\sqrt{d})$$

$$\|\mathbf{g}\|_2^2 \leq \underline{2d}$$

$$= \sum g_i^2$$

$$\mathbb{E}[\|\mathbf{g}\|_2^2] = \sum \mathbb{E}[g_i^2] = \underline{d}$$

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

$1 - c$

**Part 2:** With probablility $1 - c/d$,

$$\text{for any } i, \quad |\underline{g_i}| \geq O\left(\frac{c}{d}\right).$$



Standard normal distribution

$\frac{2c}{d} \cdot , \gamma = O(c/d)$

$\geq 1 - O(c/d)$

$\cdot \gamma$

$-\frac{c}{d} \quad 0 \quad \frac{c}{d}$

Applying union bound completes the result.

51

## 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:* $O(\gamma 1^\perp)$

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

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

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

### 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:*

$$\|X - Xzz^T\|_F^2 \le (1+\epsilon)\|X - Xv_1v_1^T\|_F^2$$

**Intuition:** For a good low-rank approximation, we don't actually need to converge to $v_1$ if $\sigma_1$ and $\sigma_2$ are the same or very close. Would suffice to return either $v_1$ or $v_2$, or some linear combination of the two.

53

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

**Power method:**

- Choose $G \in \mathbb{R}^{d \times k}$ be a random Gaussian matrix.
- $Z_0 = \text{orth}(G)$.
- For $i = 1, \ldots, T$
    - $Z^{(i)} = X^T \cdot (XZ^{(i-1)})$
    - $Z^{(i)} = \text{orth}(Z^{(i)})$

  Return $Z^{(T)}$

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

$$\|X - XZZ^T\|_F^2 \leq (1 + \epsilon)\|X - XV_kV_k^T\|_F^2.$$

Possible to "accelerate" these methods.

$$\frac{\log d/\epsilon}{\epsilon}$$

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

$$\|X - XZZ^T\|_F^2 \leq (1 + \epsilon)\|X - XV_kV_k^T\|_F^2.$$

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

$$\frac{\log d/(\epsilon)}{y} \qquad \frac{\log(d/\epsilon)}{\sqrt{y}}$$

$X^{\tau}X \; X^{\tau}X \ldots X^{\tau}X$

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



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

56

$$\left(X^{T}X\right)^{t-1} g$$

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

Along the way we computed:

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

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

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

57

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

$$\|X - X_{vv^T}\|_F^2$$

**Lanczos method:**

- Let $\underline{\mathbf{Q}} \in \mathbb{R}^{d \times \mathfrak{q}}$ be an orthonormal span for the vectors in $\mathcal{K}$.
- Solve $\min_{\mathbf{v}=\mathbf{Q}\mathbf{w}} \|\mathbf{X} - \mathbf{X}\mathbf{v}\mathbf{v}^T\|_F^2$.     return $Q\omega$
  - Find <u>best</u> vector in the Krylov subspace, instead of just using last vector.
  - Can be done in $O\left(\underline{nd\mathfrak{q}} + \underline{d\mathfrak{q}^2}\right)$ time.

$$\underline{2nd \; \mathfrak{q}}$$

58

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

→ degree 8

Power method returns:

$$\left( \mathbf{X}^T\mathbf{X} - 2(\mathbf{X}^T\mathbf{X})^2 + 3(\mathbf{X}^T\mathbf{X})^3 \cdots \right)\mathbf{g}$$

$\mathbf{v}_{\mathbf{X}\mathbf{g}}$.

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

$$p^* = \operatorname*{arg\,min}_{\text{degree } q\, p} \|\mathbf{X} - \mathbf{X}\mathbf{v}_p\mathbf{v}_p^T\|_F^2.$$

$C_0\, \mathbf{g}^n \qquad C_1\, \mathbf{X}^T\mathbf{X}\, \mathbf{g} + \quad \cdots \quad + C_q\, (\mathbf{X}^T\mathbf{X})^8\, \mathbf{g}$

$p(\mathbf{X}^T\mathbf{X})\, \mathbf{g} \quad \text{for} \quad \text{deg} \quad 8 \quad \text{poly} \quad P.$

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

$$\Delta = \frac{\text{poly}(\epsilon)}{\text{poly}(d)}$$



$$q = \sqrt{\frac{\log(d/\epsilon)}{\gamma}}$$

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

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

$$nd \qquad\qquad\qquad nd$$

60

Again 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:

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

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 $G \in \mathbb{R}^{d \times k}$ be a random Gaussian matrix.
- $\mathcal{K}_q = \left[ G, (X^TX) \cdot G, (X^TX)^2 \cdot G, \ldots, (X^TX)^q \cdot \underline{\underline{G}} \right]$

$orth\left( p\left( X^TX \right) \cdot G \right)$

$\downarrow$

approx to $V_k$.

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

$d(k+p)$

$\sqrt{\gamma}$

$k+p$ vectors

$$\boxed{\frac{\lambda_k - \lambda_{k+1}}{\lambda_k}}$$

$\frac{\lambda_k - \lambda_{k+p}}{\lambda_k}$