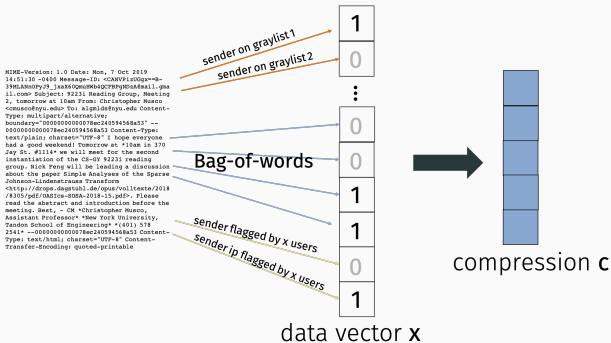


CS-GY 9223 I: Lecture 9

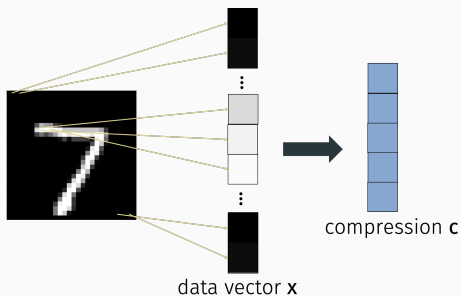
Low-rank approximation and singular value decomposition

NYU Tandon School of Engineering, Prof. Christopher Musco

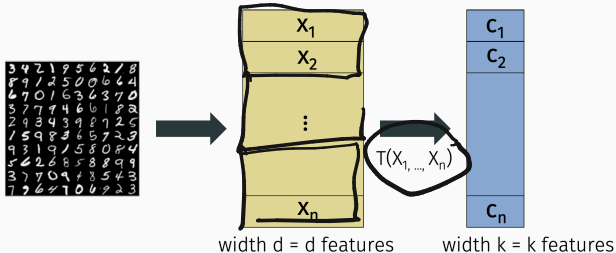
Return to data compression:



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.

SPECTRAL METHODS

Advantages of data **independent** methods: → faster

- data dependent methods often only work for specific tasks or data sets
- give very general theoretical guarantees

Advantages of data **dependent** methods:

- better compression for some data sets
→ take advantage of data structure

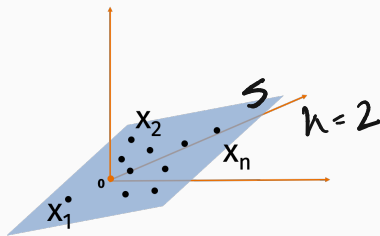
Feature
sparsity.

- preserve more complex structure.

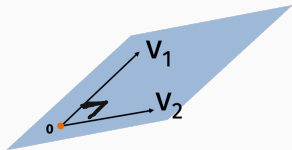
LOW-RANK DATA

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 $k < d$.

↓
dim. k



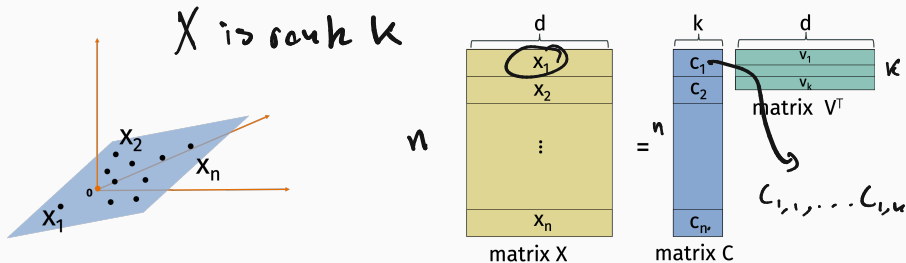
Let $\mathbf{v}_1, \dots, \mathbf{v}_k$ be orthogonal unit vectors spanning S .



For all i , we can write:

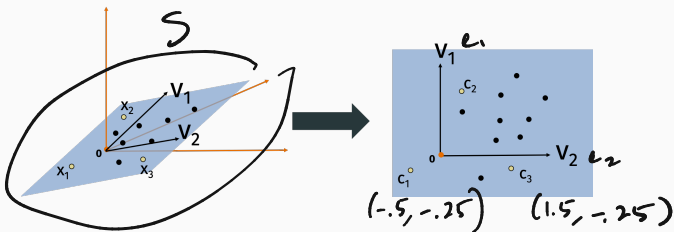
$$\underline{\mathbf{x}_i} = \underline{c_{i,1}} \underline{\mathbf{v}_1} + \dots + \underline{c_{i,k}} \underline{\mathbf{v}_k}.$$

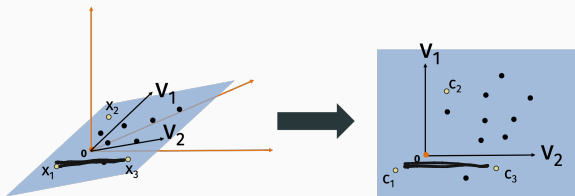
LOW-RANK DATA



What are c_1, \dots, c_n ?

rank k factorization

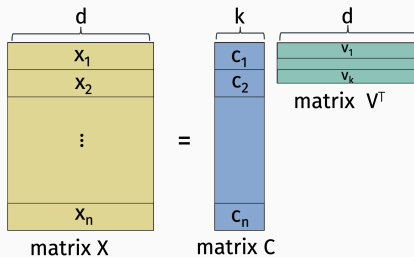




Lots of information preserved:

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

$$\underline{X = CV^T} \Rightarrow XV = \cancel{CV^T V}$$

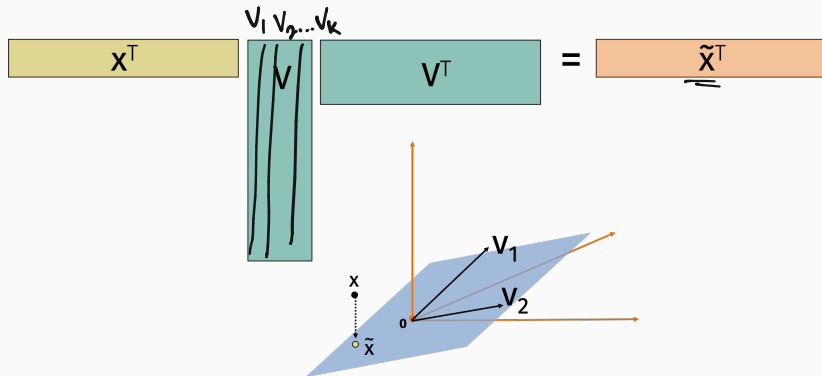
Since V 's columns are an orthonormal basis, $V^T V = I$.

So $X = \cancel{XV} V^T$

Identity matrix

PROJECTION MATRICES

VV^T is a symmetric projection matrix. $V^TV = I$



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

LOW-RANK APPROXIMATION

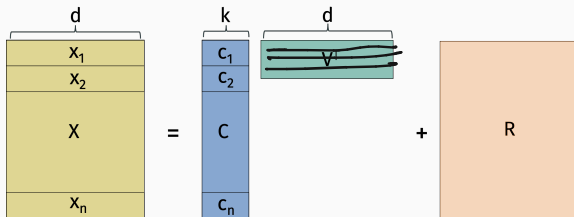
When X 's rows lie close to a k dimensional subspace, we can still approximate

$$X \approx XW^T.$$

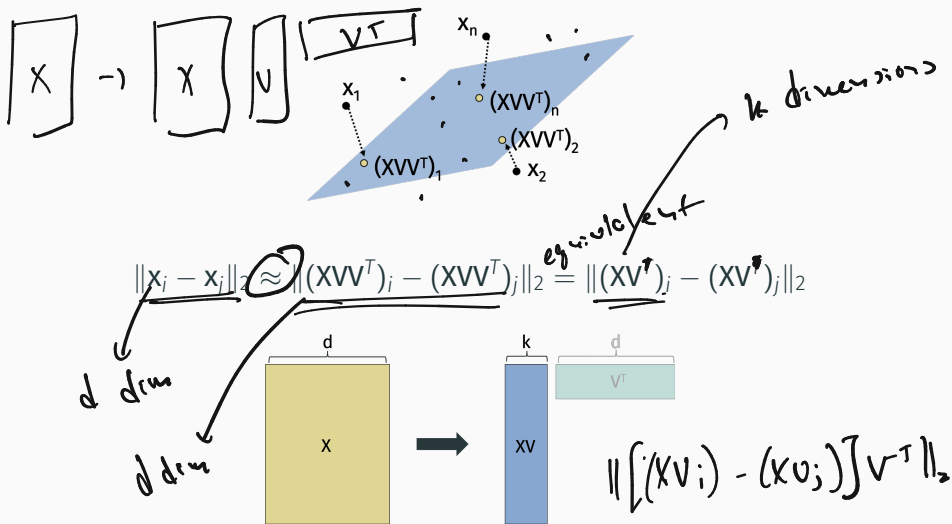
XW^T is a low-rank approximation for X .

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

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



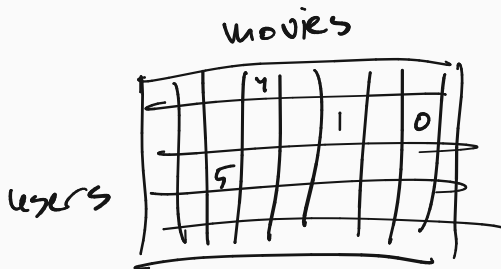
LOW-RANK APPROXIMATION



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

WHY IS DATA APPROXIMATELY LOW-RANK?

- images: have "patches" of similar colors \rightarrow low-rank structure.



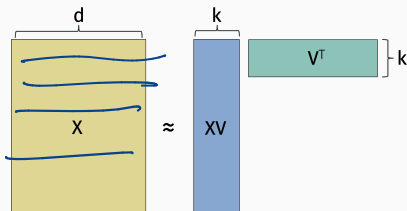
row con



rank reduction
movies

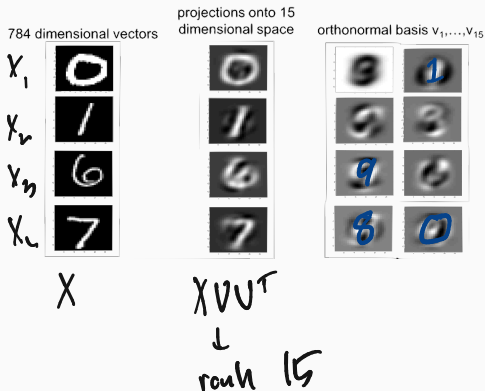
DUAL VIEW

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 .



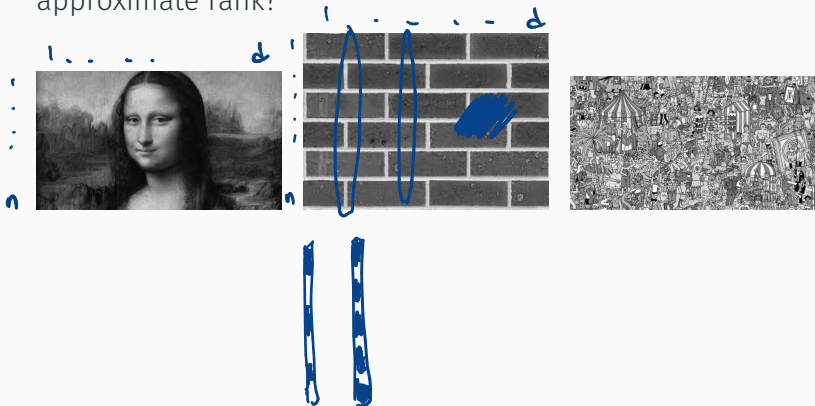
COLUMN REDUNDANCY

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

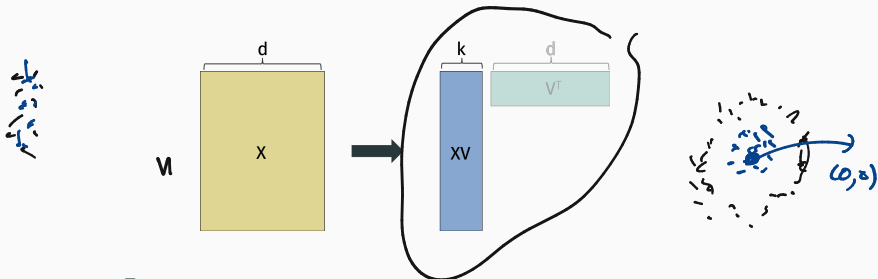
	bedrooms	bathrooms	sq.ft	floors	list price	sale price
home 1	2	2	1800	2	200,000	195,000
home 2	4	2.5	2700	1	300,000	310,000
.
.
.
home n	5	3.5	3600	3	450,000	450,000

OTHER REASONS FOR LOW-RANK STRUCTURE

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



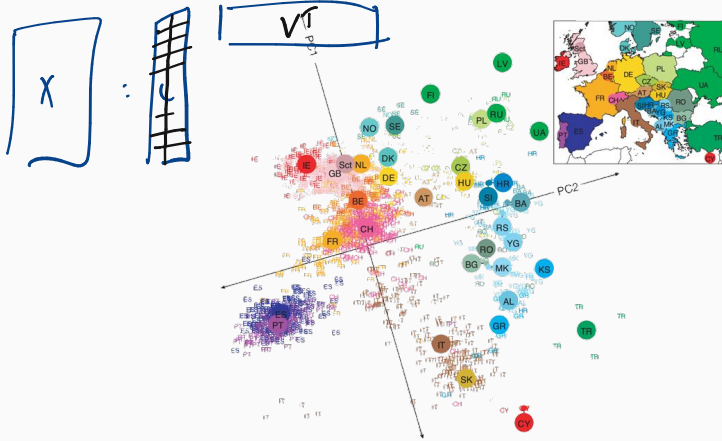
APPLICATIONS OF LOW-RANK APPROXIMATION



- $XV \cdot V^T$ 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 x_i contains genetic information for one person in Europe. Set $k = 2$ and plot $(XV)_i$ for each i on a 2-d plane. Color points by what country they are from.

COMPUTATIONAL QUESTION

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

$$\begin{aligned} \underline{\underline{\|X - XVV^T\|_F^2}} &= \min_C \|X - CV^T\|_F^2 \rightarrow \sum_{ij} (X_{ij} - (CV^T)_{ij})^2 \\ &= \sum_{i=1}^n \sum_j (X_{ij} - (CV^T)_{ij})^2 \\ &= \sum_{i=1}^n \|X_i^T - X_i^T VV^T\|_2^2 \quad (1) \end{aligned}$$

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

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

Pythagorean theorem

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

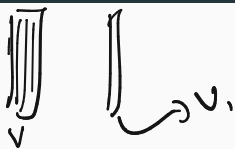
$$\max_{\text{orthonormal } V \in \mathbb{R}^{d \times k}} \|XVV^T\|_F^2 = \|XV\|_F^2 \quad (2)$$

X_i^T
row

$X_i^T VV^T$
project row

$$\begin{aligned} (X_i^T - X_i^T VV^T) \perp X_i^T VV^T \\ \|X_i^T - X_i^T VV^T\|_2^2 = \|X_i^T\|_2^2 - \|X_i^T VV^T\|_2^2 \end{aligned}$$

RANK 1 CASE



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

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

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

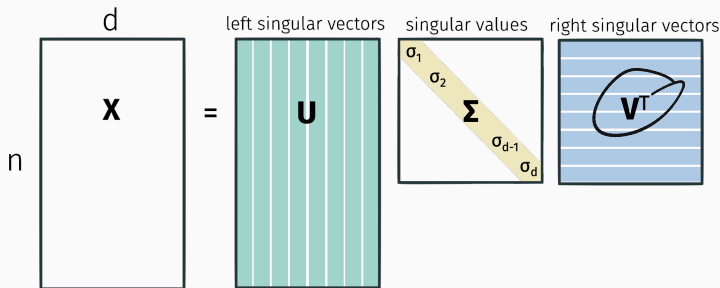
What about higher k ?

SINGULAR VALUE DECOMPOSITION

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

Any matrix X can be written:

SVD



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

- $\rightarrow u_i$
• U contains the orthonormal eigenvectors of XX^T .
- V contains the orthonormal eigenvectors of $X^T X$. $= V^T \Sigma^2 V$
- σ_i^2 = $\lambda_i(XX^T)$ = $\lambda_i(X^T X)$
↳ it's eigenvalue

This can be checked directly:

$$XX^T = U \underbrace{\Sigma V^T V \Sigma}_{I} U^T = U \Sigma^2 U^T$$

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


$$XX^T u_i = \lambda u_i$$

↓
scalar

$$XX^T u_i = U \Sigma^2 \underbrace{u_i^T u_i}_{e_i} = \underline{\sigma_i^2} u_i$$

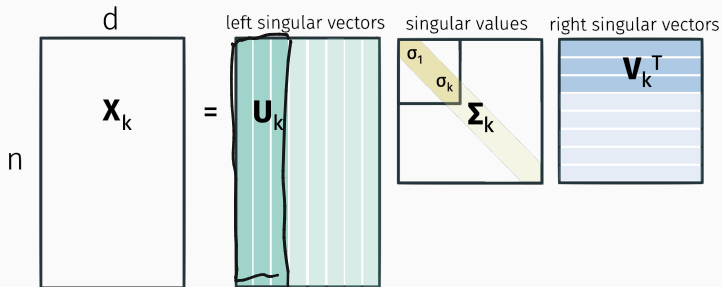
$e_i \rightarrow$ basis vector

← σ_i^2



SINGULAR VALUE DECOMPOSITION

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



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

low rank approx.

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

$$U_k \Sigma_k V_k^T$$

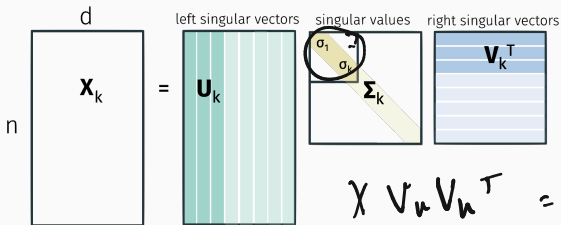
SINGULAR VALUE DECOMPOSITION

- \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



$$X V_k V_k^T = \underbrace{U_k \Sigma_k}_{C} V_k^T$$

Observation 1: The optimal compression XV_k has orthogonal columns.

$$X = CV^T$$

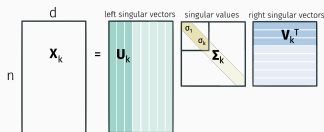
↓

also has orthogonal columns

USEFUL OBSERVATIONS

$$V^T V_k \neq I$$

$$= \begin{bmatrix} 1 & & \\ & \dots & \\ & & 0 \dots 0 \end{bmatrix}$$



Observation 2: The optimal low-rank approximation error

$$E_k = \|X - U_k U_k^T X\|_F^2 = \|X - \underline{X V_k V_k^T}\|_F^2 \text{ can be written:}$$



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

$$\|U \Sigma V^T - U \Sigma V^T V_k V_k^T\|_F^2$$

$$= \|U \Sigma V^T - U \bar{\Sigma}_k V^T\|_F^2$$

$$\hookrightarrow \bar{\Sigma}_k = \begin{bmatrix} \sigma_1 & & \\ & \sigma_k & \\ & & 0 \dots 0 \end{bmatrix}$$

$$\|U(\Sigma - \bar{\Sigma}_k)U^T\|_F^2$$

$$= \|\Sigma - \bar{\Sigma}_k\|_F^2$$

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

SPECTRAL PLOTS

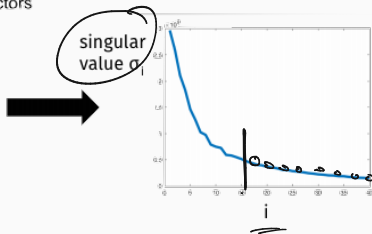
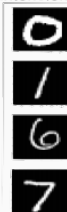
Observation 2: The optimal low-rank approximation error

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



SPECTRAL PLOTS

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:

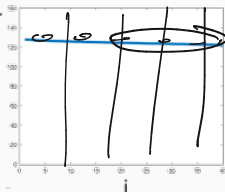
$$E_k = \sum_{i=k+1}^d \sigma_i^2.$$

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

784 dimensional vectors



singular
value σ_i

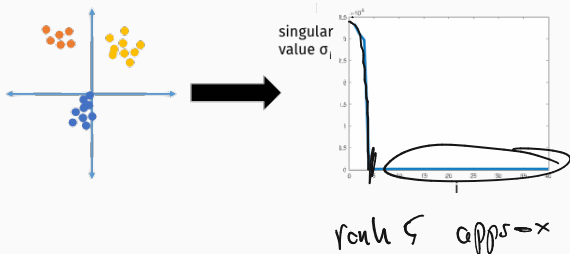


Observation 2: The optimal low-rank approximation error

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



COMPUTING THE SVD

Suffices to compute V . Then $U\Sigma = XV$.

- Compute $X^T X \rightarrow O(nd^2)$
- Find eigendecomposition $V\Lambda V^T = X^T X$.
- Compute $L = XV$. Set $\sigma_j = \|L_j\|_2$ and $U_j = L_j / \|L_j\|_2$.

$$X = \begin{matrix} n \\ \downarrow \\ d \end{matrix}$$

$$X^T X = \begin{matrix} d \\ \downarrow \\ d \end{matrix}$$

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

$$\begin{aligned} XV &= U\Sigma V^T V \\ &= U\Sigma \\ &= L \end{aligned}$$

$$U^T U = I \quad \beta_i =$$

Total runtime $\approx O(nd^2)$

Eigendecomposition of $d \times d$ matrix

QR Algorithm: $\rightarrow \approx O(d^3)$
 $O(d^3 + d^2 \log \log(1/\epsilon))$

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).

POWER METHOD

Today: What about when $k=1$? \rightarrow d dimensions

Goal: Find some $\underline{z} \approx \underline{v}_1$ \rightarrow u_1, σ_1
unit vector

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

$$\begin{matrix} n \\ \downarrow \\ X \end{matrix}$$

$$\begin{matrix} \boxed{X^T} \\ \downarrow \\ X \end{matrix} \beta_z$$

$$z = \begin{bmatrix} .1 \\ .2 \\ \vdots \\ 0.5 \end{bmatrix}$$

Power method:

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

• For $i = 1, \dots, T$

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

\rightarrow time: $O(nd + nd) = O(2nd)$

• $n_i = \|\underline{z}^{(i)}\|_2$

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

Return \underline{z}_T

$$\dots X^T X, X^T X \cdot X^T X \underline{z}^{(0)}$$

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

wey faster than $O(4d^2)$

POWER METHOD INTUITION

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

X is four right singular vectors

$$\underline{z}^{(0)} = \underbrace{c_1 v_1 + c_2 v_2 + \dots + c_d v_d}_{\rightarrow = 0}$$

Update step: $z^{(1)} = X^T \cdot (X z^{(0)}) = \underline{V \Sigma^2 V^T z^{(0)}}$ (then normalize)

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

$\begin{bmatrix} \sigma_1^2 \\ \vdots \\ \sigma_d^2 \end{bmatrix}$

$$z^{(1)} = \frac{1}{n_1} [c_1 \cdot \sigma_1^2 v_1 + c_2 \cdot \sigma_2^2 v_2 + \dots + c_d \cdot \sigma_d^2 v_d]$$

$$V \Sigma^2 V^T z^{(0)}$$

$$\begin{bmatrix} \sigma_1^2 c_1 \\ \vdots \\ \sigma_d^2 c_d \end{bmatrix}$$

$$V \Sigma^2 V^T z^{(0)}$$

$$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}^T \begin{bmatrix} z_0 \\ \vdots \\ z_n \end{bmatrix} = \begin{matrix} \rightarrow l_1 \\ \rightarrow l_2 \\ \vdots \\ \rightarrow l_n \end{matrix}$$

$$\begin{aligned} v_1 \cdot z^{(0)} &= c_1 \cdot 1 = c_1 \\ &= v_1 \cdot (l_1 v_1 + \dots + l_n v_n) \\ &= \cancel{l_1 v_1 \cdot v_1} + \dots + \cancel{l_n v_n \cdot v_n} \end{aligned}$$

POWER METHOD INTUITION

Claim: $z^{(0)} = c_1 v_1 + \dots + c_d v_d$

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

POWER METHOD FORMAL CONVERGENCE

Theorem (Basic Power Method Convergence) $\gamma > 0$

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:

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

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

faster than $O(nd^2)$

POWER METHOD FORMAL CONVERGENCE

First observation: For all i If $\text{Moose } z^{(0)} \sim \mathcal{N}(0, I)$

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

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

Corollary:

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

POWER METHOD FORMAL CONVERGENCE

$$(1-\gamma)^{1/\gamma} \approx 1/e \quad (1-\gamma)^{\log(s)/\gamma} \approx (1/e)^{\log(s)} = \frac{1}{s}$$

$$z^{(T)} = \frac{1}{\prod_{i=1}^T n_i} c_1 \cdot \sigma_1^{2T} v_1 + \frac{1}{\prod_{i=1}^T n_i} c_2 \cdot \sigma_2^{2T} v_2 + \dots + \frac{1}{\prod_{i=1}^T n_i} c_d \cdot \sigma_d^{2T} v_d$$

$\underbrace{\hspace{10em}}_{\alpha_1} \quad \underbrace{\hspace{10em}}_{\alpha_2} \quad \dots \quad \underbrace{\hspace{10em}}_{\alpha_d}$

Want to show: α_1 is \gg than α_j for all $j \geq 2$.

$$\frac{\alpha_j}{\alpha_1} = \frac{c_j}{c_1} \cdot \frac{b_j^{2T}}{b_1^{2T}} \leq \frac{c_j}{c_1} \cdot \frac{b_2^{2T}}{b_1^{2T}} \leq O(d^3 \left(\frac{b_2}{b_1}\right)^{2T})$$

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

$$\frac{b_2}{b_1} = \gamma: \frac{b_1 - b_2}{b_1} = 1 - \frac{b_2}{b_1} = 1 - \gamma \leq \frac{e}{2d} \text{ if } T = O\left(\frac{\log d/b_1}{\gamma}\right)$$

POWER METHOD FORMAL CONVERGENCE

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

$$\alpha_j \leq \frac{\epsilon}{2d}$$

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

$$\begin{aligned} \|v_1 - z^{(T)}\|_2 &= \|v_1 - \alpha_1 v_1 - \alpha_2 v_2 - \dots - \alpha_d v_d\|_2 \\ &\leq \|v_1 - \alpha_1 v_1\|_2 + \|\alpha_2 v_2 + \dots + \alpha_d v_d\|_2 \\ &\leq \epsilon/2 + \epsilon/2 \\ &\leq \|v_1 - z^{(T)}\|_2 \leq \epsilon. \end{aligned}$$

↓

$$\alpha_1 v_1 + \dots + \alpha_d v_d$$

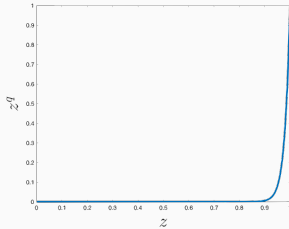
POWER METHOD – NO GAP DEPENDENCE

Theorem (Gapless Power Method Convergence)

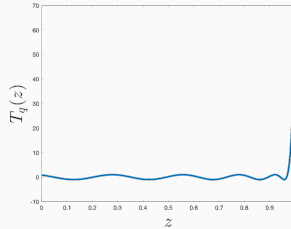
If Power Method is initialized with a random Gaussian vector then, with high probability, after $T = O\left(\frac{\log d/\epsilon}{\epsilon}\right)$ steps, we obtain a \mathbf{z} satisfying:

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

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

$\ll O(nd^2)$

to obtain a nearly optimal low-rank approximation.