

CS 331, Fall 2024  
Lecture 17 (10/28)

Today: - Low-rank approx  
- SVD/PCA  
- Power method

## Low-rank approximation (Part VI, Section 4.1)

Quiz: What are missing entries?

$$\begin{pmatrix} 7 & ? & ? & 14 & 21 \\ ? & 8 & 12 & ? & 6 \\ ? & ? & 6 & 2 & ? \end{pmatrix}$$

Obviously, an unfair question.

Nonetheless, we believe data "in the wild" is structured... maybe guessable?

Ans: 
$$\begin{pmatrix} 7 & 28 & 42 & 14 & 21 \\ 2 & 8 & 12 & 4 & 6 \\ 1 & 4 & 6 & 2 & 3 \end{pmatrix}$$

Highly-structured. Every col = multiple of  $\begin{pmatrix} 7 \\ 2 \\ 1 \end{pmatrix}$

Rank-1 matrix:  $A = \underbrace{U}_{n \times 1} \underbrace{V^T}_{d \times 1} \in \mathbb{R}^{n \times d}$

$$A = \begin{pmatrix} A_{:1} & A_{:2} & \dots & A_{:d} \end{pmatrix} = \begin{pmatrix} v_1 u & v_2 u & \dots & v_d u \end{pmatrix}$$

More generally, rank- $r$

$$A = \begin{pmatrix} u_1 v_1^T \end{pmatrix} + \begin{pmatrix} u_2 v_2^T \end{pmatrix} + \dots + \begin{pmatrix} u_r v_r^T \end{pmatrix}$$

Compactly, rank- $r$  decomposition

$$A = UV^T = \sum_{k \in \mathcal{K}(r)} u_k v_k^T \in \mathbb{R}^{n \times d}$$

$$U = (u_1 \ u_2 \ \dots \ u_r) \in \mathbb{R}^{n \times r}$$

$$V = (v_1 \ v_2 \ \dots \ v_r) \in \mathbb{R}^{d \times r}$$

Most interestingly when  $r \ll \min(n, d)$

We can always achieve  $r = \min(n, d)$

by choosing  $U = I$  or  $V = I$

$$I = \text{"identity matrix"} = \begin{pmatrix} 1 & & & 0 \\ & 1 & & \\ & & \ddots & \\ 0 & & & 1 \end{pmatrix}$$

Who cares? One motivation:

Netflix prize (matrix completion)

$$A = \begin{matrix} n \\ \text{users} \end{matrix} \left( \begin{array}{cccc} ? & ? & ? & \vdots & ? \\ ? & ? & ? & \vdots & ? \\ ? & ? & ? & \vdots & ? \\ ? & ? & ? & \vdots & ? \end{array} \right)$$

d movies

guess unknown  
movie ratings

Simplified model: everyone agrees

$$V^T = \begin{pmatrix} 9.3 & 9.2 & \dots & 1.9 & 1.5 \end{pmatrix}$$

e.g.      Showshank    Godfather                  Disaster    Superheroes  
                Redemption                                  movie

$$A_{i:} \approx V \cdot \underbrace{u_i}_w \qquad A \approx UV^T$$

how much does  
user  $i$  appreciate movies?

More sophisticated model:

$$A \approx UV^T = \begin{pmatrix} u_1 v_1^T \\ \vdots \\ u_r v_r^T \end{pmatrix} + \dots + \begin{pmatrix} u_r v_r^T \end{pmatrix}$$

$$v_1^T = (\dots 9.5 \dots 9.3 \dots)$$

Silence of the lambs      Get out

"horror enjoyer"

$$v_r^T = (\dots \dots 8.7 \dots 9.9 \dots)$$

John Wick      Kill Bill

"action enjoyer"

$$A_{i:} = \sum_{k \in [r]} [u_k]_i v_k \in \mathbb{R}^d$$

ratings of user  $i$       explained by topic  $k$

In real life, not exactly low-rank

$$A = UV^T + N = n \begin{matrix} \boxed{r} \\ \boxed{r} \end{matrix} + \begin{matrix} \boxed{d} \\ \boxed{r} \end{matrix} + \begin{matrix} \boxed{\text{small?}} \end{matrix}$$

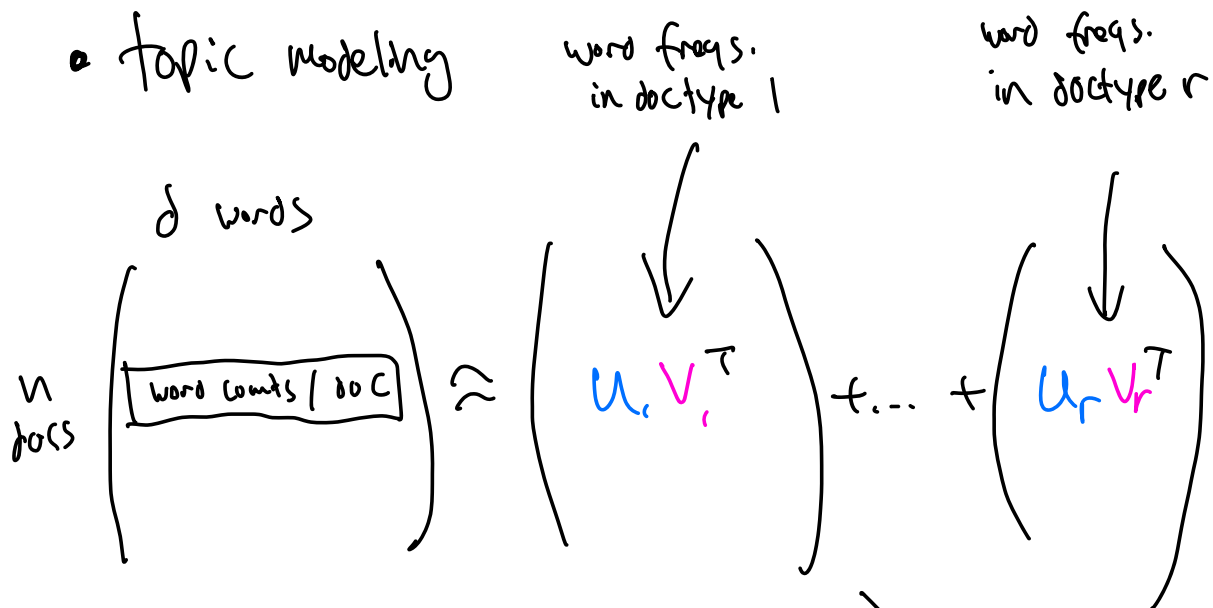
"explanatory factors"      noise      "simple" explanation      sampling error, rounds, etc.

Goal in low-rank approximation (LRA):

Make "unexplained"  $A - UV^T$  as small as possible

Other motivations:

- topic modeling



- word embeddings (see HW!)

- Save on computation, e.g. if  $r = O(1)$

$$\underbrace{Ax}_{\text{takes } O(nd) \text{ time}} \approx \underbrace{U}_{\text{takes } O(d) \text{ time}} \underbrace{V^T x}_{\text{takes } O(n) \text{ time}} = \underbrace{U}_{\text{takes } O(n) \text{ time}} (V^T x)$$

# Singular value decomposition (Part VI, Section 4.2)

We have 2 complete characterizations of

2-norm

$$\|A - UV^T\|$$

$U \in \mathbb{R}^{n \times n}$   
 $V \in \mathbb{R}^{d \times d}$

"unitarily invariant"  
norm. incl. basically all  
common size norms...  
Frobenius, operator, nuclear norms

Detour to SVD needed.

Say that  $U \in \mathbb{R}^{d \times d}$  is orthonormal if

$$U^T U = I = \begin{pmatrix} 1 & & & 0 \\ & 1 & & \\ & & \ddots & \\ 0 & & & 1 \end{pmatrix}$$

Interpretation: let  $U = \begin{pmatrix} u_1 & u_2 & \dots & u_n \end{pmatrix}$

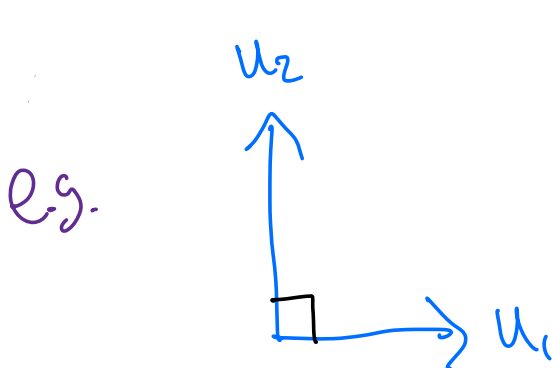
Then,  $\forall (i,j) \in [n] \times [n]$

$$\left[ U^T U \right]_{ij} = \underbrace{u_i^T u_j}_{\cos(\theta_{ij})} = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

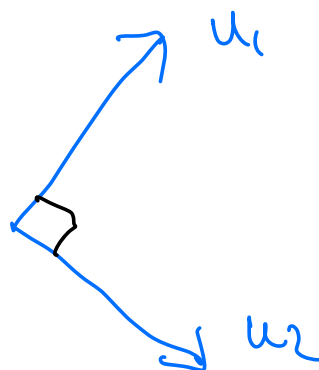
Recall  $u^T u = \sum u_i^2 = \underbrace{\|u\|_2^2}_{\text{length of } u}$  (Pythagoras)

Hence columns of  $U$

- unit length
- pairwise perpendicular



"Standard basis"  $U = I$



Orthonormal basis



Also possible for rectangular "full" = orthonormal

$$U \in \mathbb{R}^{n \times d}, \quad n \geq d \quad U^T U = I$$

$$U = \begin{pmatrix} u_1 & u_2 & \dots & u_d \end{pmatrix}$$

Unit length,  
pairwise perp

subset of full  
basis  $\{u_i\}_{i \in [d]}$

e.g.

$$U = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}$$

$e_1 \quad e_2 \quad e_4$

is orthonormal,

$\text{Span}(U)$  is "subspace"  
of  $\mathbb{R}^5$  with dimension 3

$$Ux = x_1 e_1 + x_2 e_2 + x_4 e_4$$

(no mass on  $e_3, e_5$  allowed)

SVD: All matrices  $A \in \mathbb{R}^{n \times d}$ ,  $n \geq d$   
 can be decomposed as

$$A = U \Sigma V^T = \sum_{i \in (d)} \sigma_i u_i v_i^T$$

$$U = \begin{pmatrix} u_1 & \dots & u_d \end{pmatrix} \in \mathbb{R}^{n \times d}$$

$$V = \begin{pmatrix} v_1 & \dots & v_d \end{pmatrix} \in \mathbb{R}^{d \times d}$$

} orthogonal matrices

$$\Sigma = \begin{pmatrix} \sigma_1 & & & 0 \\ & \sigma_2 & & \\ & & \ddots & \\ 0 & & & \sigma_d \\ & & & & 0 \end{pmatrix} \in \mathbb{R}_{\geq 0}^{d \times d}$$

nonneg. diagonal matrix

$$\# \text{ nonzero} = \text{rank}(A) = \dim(\text{span}(A))$$

Interpretation: SVD sends  $x \rightarrow Ax$

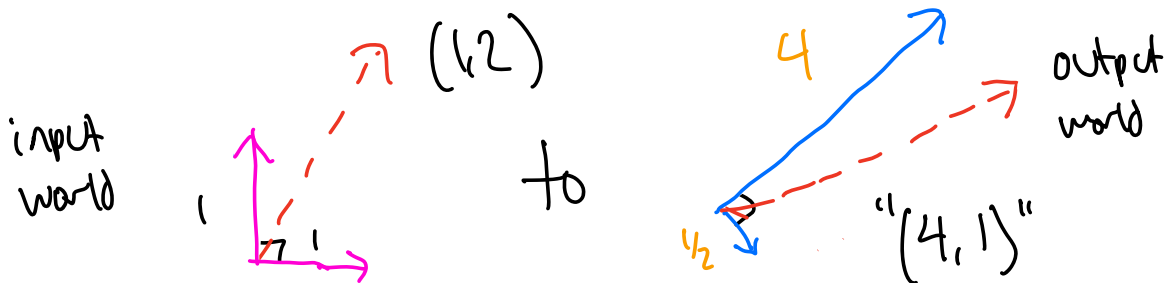
$V_i \in \mathbb{R}^d$  = "input world" to  $\sigma_i U_i \in \mathbb{R}^n$  = "output world"

$$X = \sum_{i \in \mathcal{I}} c_i v_i = V C \quad (\text{coeffs. } c_i \in \mathbb{R}^d)$$

$$\rightarrow Ax = U \underbrace{\sum_{i \in \mathcal{I}} v_i^T v_i}_{= I} C = \sum_{i \in \mathcal{I}} \sigma_i c_i u_i$$

e.g.

$$A = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} 4 & 0 \\ 0 & \frac{1}{2} \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



Eckart-Young-Mirsky:

For any  $A \in \mathbb{R}^{n \times d}$ ,  $n \geq d$ , unitarily-invariant  $\|\cdot\|$

$$\text{argmin}_{\text{rank-}r \text{ } M} \|A - M\|$$

$$\text{achieved by } M = \sum_{k \in \mathcal{K}(r)} \sigma_k u_k v_k^T$$

$$\text{where } A = U \Sigma V^T \text{ (SVD)}$$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d$$

AKA, optimal low-rank approx.

= do SVD, keep top  $r$  components,

all else  $\Rightarrow 0$ .

Special case: Symmetric matrix  $M \in \mathbb{R}^{d \times d}$

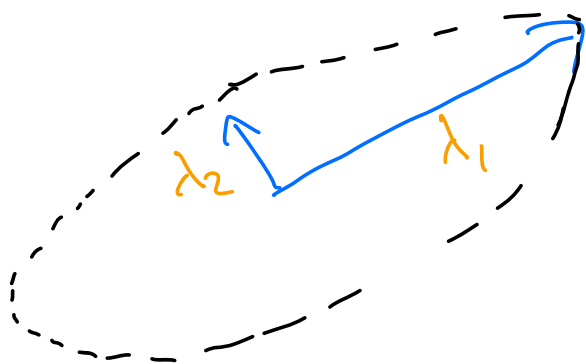
$$M = \underbrace{U \Lambda U^T}_{\text{"eigendecomposition"}}, \text{ but } \Lambda \in \mathbb{R}^d$$

diagonal,  
could be neg.

Even more special: if  $\Lambda \in \mathbb{R}_{\geq 0}^{d \times d}$  already

$M$  is "positive semidefinite" (PSD)

Input world = output world!



$$M u_i = \lambda_i u_i$$

geometric interpretation  
of eigvecs/eigvals

"every PSD matrix is an ellipse"

If  $A = U \Sigma V^T$ , we have

$$A^T A = V \Sigma U^T U \Sigma V^T = V \Sigma^2 V^T$$

$$A A^T = U \Sigma V^T V \Sigma U^T = U \Sigma^2 U^T$$

Both PSD, eigenvectors give SVD of  $A$ .

Enough to give algos to recover:

top eigvec of PSD matrix  $M$  (PCA)

Power method (Part VI, Section 4.3)

Let PSD  $M \in \mathbb{R}^{d \times d}$

$$M = U \underbrace{\Delta}_{\text{diag}(\lambda)} U^T = \sum_{i \in [d]} \lambda_i u_i u_i^T$$

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$

Goal: recover  $u_1$  to high accuracy.

Idea: suppose  $U = I$ ,  $\lambda_1 \geq 2000 d^2 \lambda_2$

Then  $Mg \approx$  multiple of  $e_1$

random  
normal vector

In fact,  $Ug \sim g$  (still works for  $U \neq I$ )

But  $\lambda_1 \geq 2000 d^2 \lambda_2$  really strong...

What if only  $\lambda_1 \geq 1.1 \lambda_2$ ?

No problem: for  $p = O(\log d)$

$$M^p g = \cancel{U \Delta U^T U \Delta U^T U \Delta U^T \dots}$$

$$= \underbrace{U \Delta^p U^T}_g g$$

gap now  $1.1^p \geq 2000 d^2$

normalize  $M^p g$

Runtime:  $O(d^2 \log(d))$

gives  $\cos \theta(\hat{u}, u_1) \geq 0.99$   
w.p.  $\geq 0.99$