# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2024.
Lecture 17

## Logistics

- Problem Set 3 is due tomorrow at 11:59pm.
- Due to Veteran's day and a short week this week, no quiz due Monday.

## Summary

### Last Class

- Finish up optimal low-rank approximation via eigendecomposition.

$X^T X$ · Eigenvalue spectrum as a way of measuring low-rank approximation error.

### This Class: The SVD and Application of Low-Rank Approximation Beyond Compression

- The Singular Value Decomposition (SVD) and its connection to eigendecomposition and low-rank approximation.

- Low-rank matrix completion (predicting missing measurements using low-rank structure).

- Entity embeddings (e.g., word embeddings, node embeddings).

# Low-Rank Approximation Review

**True or False?**

$V^k$

ⓐ

$$\min_{V \in \mathbb{R}^{d \times k} : V^T V = I} \|X - XVV^T\|_F^2 = \min_{B : \text{rank}(B) \leq k} \|X - B\|_F^2.$$

always equal

$B_2^*$  $B_1^*$

rows of X
we projected
subspace spanned
by V

ⓐ = ⓑ

ⓐ ≥ ⓑ   ($XVV^T$ has rank $\leq k$)

ⓑ ≥ ⓐ

ⓑ $\leq \|X - B\|_F^2$ for any $B : \text{rank}(B) \leq k$

sub. $XW^T = B$

ⓑ $\leq \|X - XVV^T\|_F^2$ for any $V$

Let $V^*$ be ortho. basis for rows of $B^*$

$B^* = \underset{B : \text{rank}(B) \leq k}{\arg\min} \|X - B\|_F^2$

$\|X - XV^*V^{*T}\|_F^2 \leq \|X - B^*\|_F^2$

ⓐ $= \min_V \|X - XVV^T\|_F^2 \leq \|X - XV^*V^{*T}\|_F^2 \leq \|X - B^*\|_F^2$ ⓑ

4

What is the value of

$$\min_{B:\text{rank}(B)\leq k} \|X - B\|_F^2? = \sum_{i=k+1}^{d} \lambda_i(X^TX)$$

$XV^xV^{xT}$

$\lambda_1(X^TX) \geq \lambda_2(X^TX) \geq \ldots \lambda_d(X^TX)$



$X_1$
$X_2$
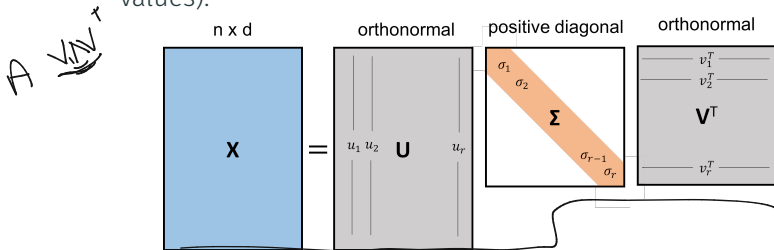$V^xV^{xT}X_1$ ⋯ $b_1^x$
$b_2^x$
$b_3^x$
$b_4^x$
$X_3$
$X_4$

# Singular Value Decomposition

The Singular Value Decomposition (SVD) generalizes the eigendecomposition to asymmetric (even rectangular) matrices.

## Singular Value Decomposition

The Singular Value Decomposition (SVD) generalizes the eigendecomposition to asymmetric (even rectangular) matrices. Any matrix $X \in \mathbb{R}^{n \times d}$ with $\text{rank}(X) = r$ can be written as $X = U\Sigma V^T$.

- $U$ has orthonormal columns $\vec{u}_1, \ldots, \vec{u}_r \in \mathbb{R}^n$ (left singular vectors).
- $V$ has orthonormal columns $\vec{v}_1, \ldots, \vec{v}_r \in \mathbb{R}^d$ (right singular vectors).
- $\Sigma$ is diagonal with elements $\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_r > 0$ (singular values).

## Connection of the SVD to Eigendecomposition

Writing $X \in \mathbb{R}^{n \times d}$ in its singular value decomposition $X = U\Sigma V^T$:

$$X^T X = \left( U \Sigma V^T \right)^T U \Sigma V^T = V \Sigma U^T U \Sigma V^T = V \Sigma^2 V^T$$

$X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\Sigma \in \mathbb{R}^{\text{rank}(X) \times \text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

## Connection of the SVD to Eigendecomposition

Writing $X \in \mathbb{R}^{n \times d}$ in its singular value decomposition $X = U\Sigma V^T$:

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

$X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\Sigma \in \mathbb{R}^{\text{rank}(X) \times \text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

## Connection of the SVD to Eigendecomposition

Writing $X \in \mathbb{R}^{n \times d}$ in its singular value decomposition $X = U\Sigma V^T$:

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

$X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\Sigma \in \mathbb{R}^{\text{rank}(X) \times \text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

## Connection of the SVD to Eigendecomposition

Writing $X \in \mathbb{R}^{n \times d}$ in its singular value decomposition $X = U\Sigma V^T$:

$$\Sigma^T = \Sigma \qquad X^T X = V\Sigma^T U^T U\Sigma V^T = V\Sigma^2 V^T \text{ (the eigendecomposition)}$$

right singular vectors = eigenvectors of $X^T X$

squared singular values = eigenvalues of $X^T X$

$X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\Sigma \in \mathbb{R}^{\text{rank}(X) \times \text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

## Connection of the SVD to Eigendecomposition

Writing $X \in \mathbb{R}^{n \times d}$ in its singular value decomposition $X = U\Sigma V^T$:

$$X^T X = V\Sigma U^T U\Sigma V^T = V\Sigma^2 V^T \text{ (the eigendecomposition)}$$

Similarly: $\underline{XX^T = U\Sigma V^T V\Sigma U^T} = \underline{U\Sigma^2 U^T}$.

---

$X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\Sigma \in \mathbb{R}^{\text{rank}(X) \times \text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

## Connection of the SVD to Eigendecomposition

Writing $X \in \mathbb{R}^{n \times d}$ in its singular value decomposition $X = U\Sigma V^T$:

$$X^T X = V\Sigma U^T U\Sigma V^T = V\Sigma^2 V^T \text{ (the eigendecomposition)}$$

Similarly: $XX^T = U\Sigma V^T V\Sigma U^T = U\Sigma^2 U^T$.

The left and right singular vectors are the eigenvectors of the covariance matrix $X^T X$ and the gram matrix $XX^T$ respectively.

---

$X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\Sigma \in \mathbb{R}^{\text{rank}(X) \times \text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

## Connection of the SVD to Eigendecomposition

Writing $X \in \mathbb{R}^{n \times d}$ in its singular value decomposition $X = U\Sigma V^T$:

$$X^T X = V\Sigma U^T U\Sigma V^T = V\Sigma^2 V^T \text{ (the eigendecomposition)}$$

Similarly: $XX^T = U\Sigma V^T V\Sigma U^T = U\Sigma^2 U^T$.

The left and right singular vectors are the eigenvectors of the covariance matrix $X^T X$ and the gram matrix $XX^T$ respectively.

So, letting $V_k \in \mathbb{R}^{d \times k}$ have columns equal to $\vec{v}_1, \ldots, \vec{v}_k$, we know that $XV_k V_k^T$ is the best rank-$k$ approximation to $X$ (given by PCA).

---

$X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\Sigma \in \mathbb{R}^{\text{rank}(X) \times \text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

## Connection of the SVD to Eigendecomposition

Writing $X \in \mathbb{R}^{n \times d}$ in its singular value decomposition $X = U\Sigma V^T$:

$$X^T X = V\Sigma U^T U\Sigma V^T = V\Sigma^2 V^T \text{ (the eigendecomposition)}$$

Similarly: $XX^T = U\Sigma V^T V\Sigma U^T = U\Sigma^2 U^T$.

The left and right singular vectors are the eigenvectors of the covariance matrix $X^T X$ and the gram matrix $XX^T$ respectively.

So, letting $V_k \in \mathbb{R}^{d \times k}$ have columns equal to $\vec{v}_1, \ldots, \vec{v}_k$, we know that $XV_k V_k^T$ is the best rank-$k$ approximation to $X$ (given by PCA).

What about $U_k U_k^T X$ where $U_k \in \mathbb{R}^{n \times k}$ has columns equal to $\vec{u}_1, \ldots, \vec{u}_k$?

---

$X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\Sigma \in \mathbb{R}^{\text{rank}(X) \times \text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

## Connection of the SVD to Eigendecomposition

Writing $X \in \mathbb{R}^{n \times d}$ in its singular value decomposition $X = U\Sigma V^T$:

$$X^T X = V\Sigma U^T U\Sigma V^T = V\Sigma^2 V^T \text{ (the eigendecomposition)}$$

Similarly: $XX^T = U\Sigma V^T V\Sigma U^T = U\Sigma^2 U^T$.

The left and right singular vectors are the eigenvectors of the covariance matrix $X^T X$ and the gram matrix $XX^T$ respectively.

So, letting $V_k \in \mathbb{R}^{d \times k}$ have columns equal to $\vec{v}_1, \ldots, \vec{v}_k$, we know that $XV_k V_k^T$ is the best rank-$k$ approximation to $X$ (given by PCA).

What about $U_k U_k^T X$ where $U_k \in \mathbb{R}^{n \times k}$ has columns equal to $\vec{u}_1, \ldots, \vec{u}_k$?
Gives exactly the same approximation! $\quad U_k U_k^T X = XV_k V_k^T$

> $X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\Sigma \in \mathbb{R}^{\text{rank}(X) \times \text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

# The SVD and Optimal Low-Rank Approximation
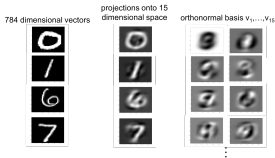
The best low-rank approximation to $X$:

$X_k = \arg\min_{\text{rank}-k \; B\in\mathbb{R}^{n\times d}} \|X - B\|_F$ is given by:

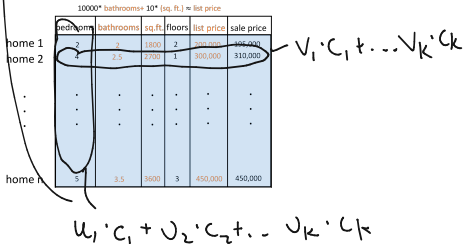$$X_k = XV_kV_k^T = U_kU_k^TX \quad \text{— projecting features optimally}$$

*projecting data points optimally*

Correspond to projecting the rows (data points) onto the span of $V_k$
or the columns (features) onto the span of $U_k$



**Row (data point) compression**

784 dimensional vectors — projections onto 15 dimensional space — orthonormal basis $v_1, \ldots, v_{16}$

**Column (feature) compression**

10000* bathrooms+ 10* (sq. ft.) ≈ list price

| | bedrooms | bathrooms | sq.ft. | floors | list price | sale price |
|---|---|---|---|---|---|---|
| home 1 | 2 | 1 | 1800 | 2 | 300,000 | 290,000 |
| home 2 | 4 | 2.5 | 2700 | 1 | 300,000 | 310,000 |
| . | . | . | . | . | . | . |
| . | . | . | . | . | . | . |
| . | . | . | . | . | . | . |
| home n | 5 | 3.5 | 3600 | 3 | 450,000 | 450,000 |

$\rightarrow v_1 \cdot c_1 + \ldots v_k \cdot c_k$

$u_1 \cdot c_1 + v_2 \cdot c_2 + \ldots v_k \cdot c_k$

$\sqrt[k]{\begin{bmatrix} X V_k \end{bmatrix}}$

$\overset{n}{\begin{bmatrix} X V_k V_k^T \end{bmatrix}} \approx \begin{bmatrix} V_k V_k^T X \end{bmatrix}$

8

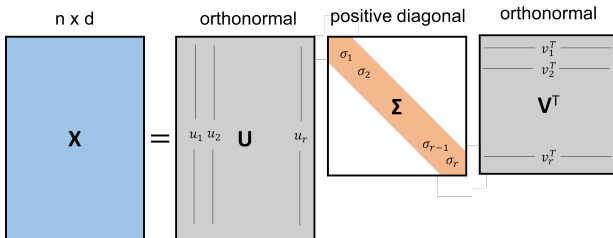## The SVD and Optimal Low-Rank Approximation

The best low-rank approximation to $X$:
$X_k = \arg\min_{\text{rank} - k \ B \in \mathbb{R}^{n \times d}} \|X - B\|_F$ is given by:

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

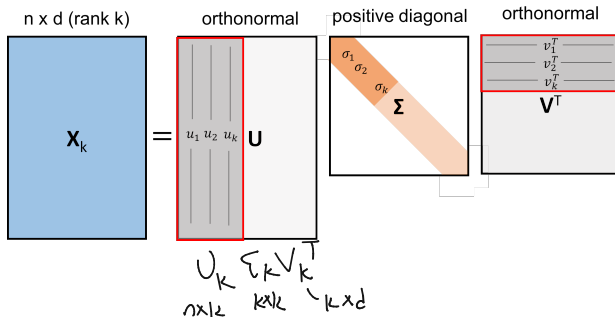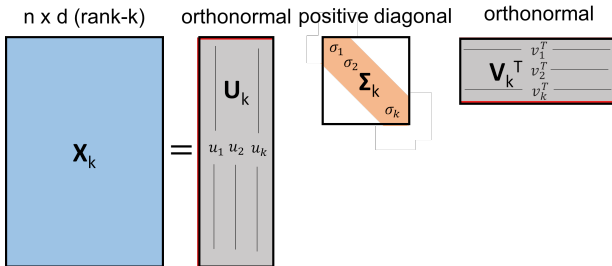Correspond to projecting the rows (data points) onto the span of $V_k$ or the columns (features) onto the span of $U_k$

# The SVD and Optimal Low-Rank Approximation

The best low-rank approximation to $X$:

$X_k = \arg\min_{\text{rank}-k\ B \in \mathbb{R}^{n \times d}} \|X - B\|_F$ is given by:

$$X_k = XV_kV_k^T = U_kU_k^TX \;=\; U_k\Sigma_kV_k^T$$

Correspond to projecting the rows (data points) onto the span of $V_k$ or the columns (features) onto the span of $U_k$



n x d (rank k)    orthonormal    positive diagonal    orthonormal

$X_k$ = $u_1\ u_2\ u_k$ $U$ | $\sigma_1$ $\sigma_2$ $\sigma_k$ $\Sigma$ | $v_1^T$ $v_2^T$ $v_k^T$ $V^T$

$$U_k\ \Sigma_k V_k^T$$
$$n \times k \qquad k \times k \qquad k \times d$$

The best low-rank approximation to $X$:
$X_k = \arg\min_{\text{rank}-k\ B \in \mathbb{R}^{n \times d}} \|X - B\|_F$ is given by:

$$X_k = XV_kV_k^T = U_kU_k^TX = U_k\Sigma_kV_k^T$$

Correspond to projecting the rows (data points) onto the span of $V_k$ or the columns (features) onto the span of $U_k$

# The SVD and Optimal Low-Rank Approximation

The best low-rank approximation to $X$:

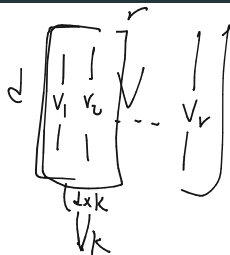$X_k = \arg\min_{\text{rank}-k\ B\in\mathbb{R}^{n\times d}} \|X - B\|_F$ is given by:

$$X_k = XV_kV_k^T = U_kU_k^TX = U_k\boldsymbol{\Sigma}_kV_k^T$$

$X = U\boldsymbol{\Sigma}V^T$

$XV_kV_k^T = U\boldsymbol{\Sigma}\underbrace{V^TV_k}V_k^T$

$$r\left[\quad \overset{d}{U^T}\quad\right] d\left[\overset{k}{V_k}\right] = r\begin{bmatrix} \overset{k}{1\ 0} \\ 0\ 1 \\ \hline 0 \end{bmatrix} \qquad i,j \text{ entry is } v_i^Tv_j$$

$$XV_kV_k^T = U\boldsymbol{\Sigma}\overset{n\times r}{\begin{bmatrix} I_k \\ 0 \end{bmatrix}}V_k^T = \overset{n\times k\ k\times k\ k\times d}{U\ r\begin{bmatrix} \boldsymbol{\Sigma}_k \\ 0 \end{bmatrix}}V_k^T = U_k\boldsymbol{\Sigma}_kV_k^T$$

$\overset{a}{\textcircled{a}} = \overset{c}{\textcircled{c}}$

$$\begin{bmatrix} \sigma_1 \cdots \sigma_k \\ \quad \sigma_{k+1} \\ \quad\quad \cdots \sigma_d \end{bmatrix}\begin{bmatrix} 1 \\ \cdots \\ 1 \\ \hline 0 \end{bmatrix} = \begin{bmatrix} \sigma_1\ 0 \\ 0\ \cdots\ \sigma_k \\ \hline 0 \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Sigma}_k \\ 0 \end{bmatrix}$$



> $X \in \mathbb{R}^{n\times d}$: data matrix, $U \in \mathbb{R}^{n\times\text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d\times\text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\boldsymbol{\Sigma} \in \mathbb{R}^{\text{rank}(X)\times\text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

9

# The SVD and Optimal Low-Rank Approximation

The best low-rank approximation to $X$:

$X_k = \arg\min_{\text{rank} - k\ B \in \mathbb{R}^{n \times d}} \|X - B\|_F$ is given by:

$$X_k = XV_kV_k^T = U_kU_k^TX = U_k\Sigma_kV_k^T$$

$U_kU_k^TX = U_kU_k^T \ U\Sigma V^T$

$\quad U_k[I_k \vdots 0]\Sigma V^T$

$\quad U_k[\Sigma_k \vdots 0]V^T$

$\quad U_k\Sigma_kV_k^T$

$$\begin{bmatrix} | & & | & | & & | \\ u_1 & \cdots & u_k & \cdots & u_y \\ | & & | & | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 \\ \ddots \\ & \sigma_k \\ \hline 0 \end{bmatrix} = v_1\sigma_1 + v_2\sigma_2 + \cdots v_k\sigma_k = U_k\Sigma_k$$

---

$X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \ldots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \ldots$ (right singular vectors), $\Sigma \in \mathbb{R}^{\text{rank}(X) \times \text{rank}(X)}$: positive diagonal matrix containing singular values of $X$.

## SVD Review

$$\begin{bmatrix} X \end{bmatrix}\begin{bmatrix} V_k \end{bmatrix} = \begin{bmatrix} XV_k \end{bmatrix}\begin{bmatrix} V_k^T \end{bmatrix} \quad \text{typically dense}$$

- Every $X \in \mathbb{R}^{n \times d}$ can be written in its SVD as $U\Sigma V^T$.

- $U \in \mathbb{R}^{n \times r}$ (orthonormal) contains the eigenvectors of $XX^T$.
  $V \in \mathbb{R}^{d \times r}$ (orthonormal) contains the eigenvectors of $X^T X$.
  $\Sigma \in \mathbb{R}^{r \times r}$ (diagonal) contains their eigenvalues. *square roots of*

- $U_k U_k^T X = XV_k V_k^T = U_k \Sigma_k V_k^T = \underset{B \text{ s.t. } \mathrm{rank}(B) \leq k}{\arg\min} \|X - B\|_F.$

$$U_k \Sigma_k V_k^T = SVdS(X, k)$$

11

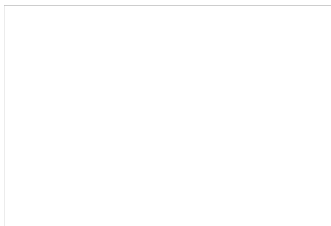# Applications of Low-Rank Approximation Beyond Compression

## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix).

## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.
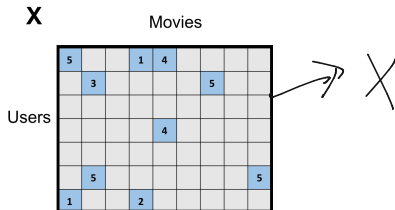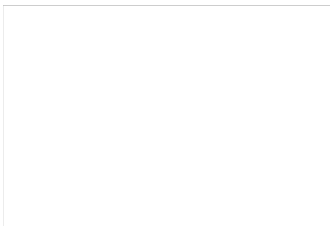
**X**      Movies

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 5 | 3 | 3 | 1 | 4 | 4 | 4 | 3 | 5 |
| 4 | 3 | 3 | 1 | 4 | 4 | 5 | 3 | 5 |
| 3 | 3 | 3 | 2 | 3 | 3 | 3 | 3 | 3 |
| 4 | 3 | 3 | 4 | 4 | 4 | 4 | 3 | 3 |
| 3 | 3 | 3 | 2 | 3 | 3 | 3 | 3 | 3 |
| 2 | 5 | 3 | 4 | 4 | 4 | 4 | 4 | 5 |
| 1 | 3 | 3 | 2 | 3 | 3 | 3 | 1 | 2 |

Users

## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.
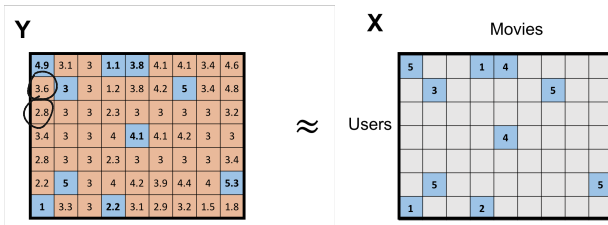


**X** Movies

Assume rank(**X**)=1

$r_2 = 2 \cdot r_1$

$r_1$

$r_2 \rightarrow$

Users

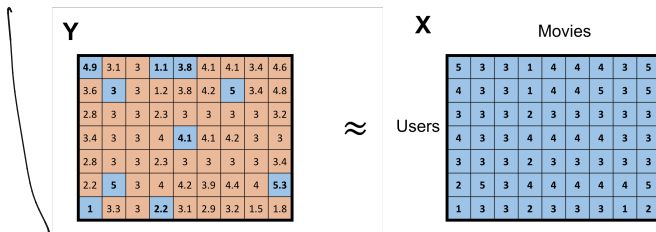| 5 | 2 | 1 | 1 | 4 |
| | 4 | | 2 | |
| | 4 | 2 | | |
| | | | | 4 |

## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.



**Solve:** $Y = \underset{B \text{ s.t. } \mathrm{rank}(B) \leq k}{\arg\min} \sum_{\text{observed } (j,k)} \left[ X_{j,k} - B_{j,k} \right]^2$

## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.



**Solve:** $Y = \underset{B \text{ s.t. } \operatorname{rank}(B) \leq k}{\arg \min} \sum_{\text{observed } (j,k)} \left[ X_{j,k} - B_{j,k} \right]^2$

## Matrix Completion

Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem.



**Solve:** $Y = \underset{B \text{ s.t. } \mathrm{rank}(B) \leq k}{\arg\min} \sum_{\text{observed } (j,k)} \left[ X_{j,k} - B_{j,k} \right]^2$

Under certain assumptions, can show that $Y$ well approximates $X$ on both the observed and (most importantly) unobserved entries.

## Entity Embeddings

Dimensionality reduction embeds $d$-dimensional vectors into $k$ dimensions. But what about when you want to embed objects other than vectors?

- Documents (for topic-based search and classification)
- Words (to identify synonyms, translations, etc.)
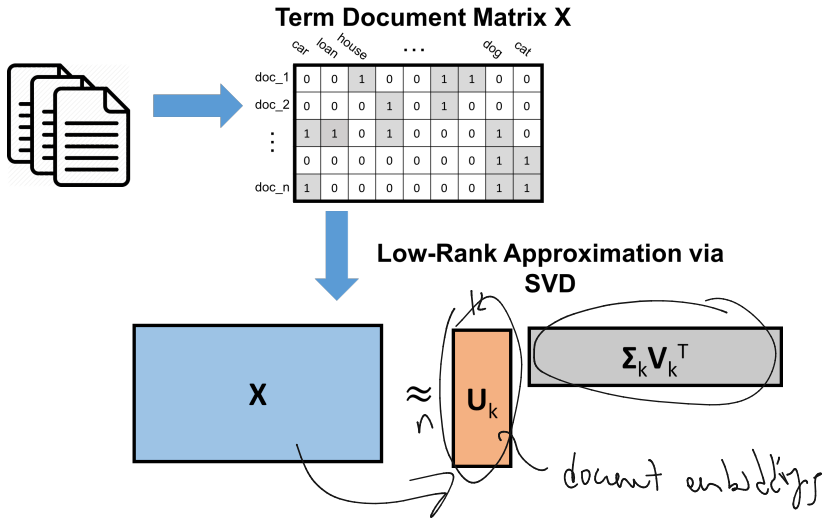- Nodes in a social network

## Entity Embeddings

Dimensionality reduction embeds $d$-dimensional vectors into $k$ dimensions. But what about when you want to embed objects other than vectors?
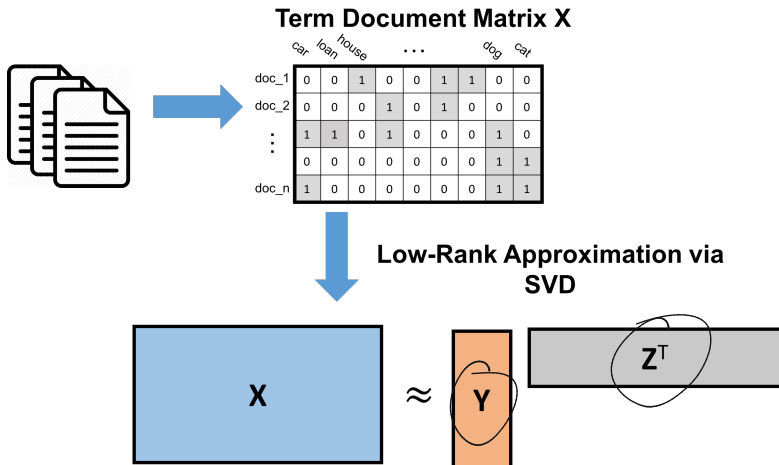
- Documents (for topic-based search and classification)
- Words (to identify synonyms, translations, etc.)
- Nodes in a social network

**Classic Approach:** Convert each item into a (very) high-dimensional feature vector and then apply low-rank approximation.

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

document embeddings

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

$$X \approx Y \, Z^T$$

# Example: Latent Semantic Analysis

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

- If the error $\|X - YZ^T\|_F$ is small, then on average,

$$X_{i,a} \approx (YZ^T)_{i,a} = \langle \vec{y}_i, \vec{z}_a \rangle.$$

# Example: Latent Semantic Analysis



**Term Document Matrix X**

| | car | loan | house | ... | | | dog | cat | |
|---|---|---|---|---|---|---|---|---|---|
| doc_1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 |
| doc_2 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| ⋮ | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| doc_n | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |

**Low-Rank Approximation via SVD**

$$X \approx Y \ Z^T$$

- If the error $\|X - YZ^T\|_F$ is small, then on average,

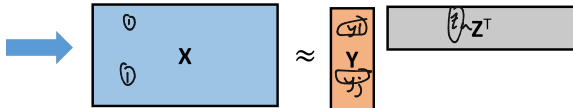$$X_{i,a} \approx (YZ^T)_{i,a} = \langle \vec{y}_i, \vec{z}_a \rangle.$$

- I.e., $\langle \vec{y}_i, \vec{z}_a \rangle \approx 1$ when $doc_i$ contains $word_a$.

# Example: Latent Semantic Analysis



**Term Document Matrix X**      **Low-Rank Approximation via SVD**
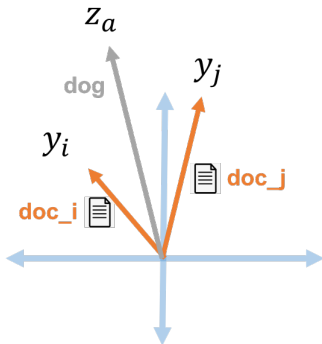
- If the error $\|X - YZ^T\|_F$ is small, then on average,

$$X_{i,a} \approx (YZ^T)_{i,a} = \langle \vec{y}_i, \vec{z}_a \rangle.$$

- I.e., $\langle \vec{y}_i, \vec{z}_a \rangle \approx 1$ when $doc_i$ contains $word_a$.

- If $doc_i$ and $doc_j$ both contain $word_a$, $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle \approx 1$.

# Example: Latent Semantic Analysis

If $doc_i$ and $doc_j$ both contain $word_a$, $\langle \vec{y_i}, \vec{z_a} \rangle \approx \langle \vec{y_j}, \vec{z_a} \rangle \approx 1$
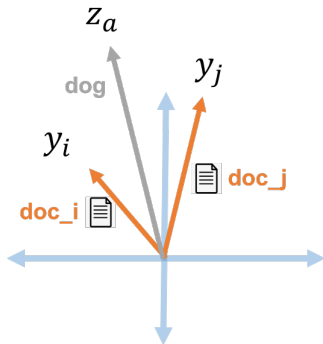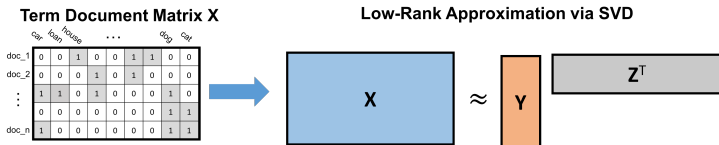
## Example: Latent Semantic Analysis

If $doc_i$ and $doc_j$ both contain $word_a$, $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle \approx 1$
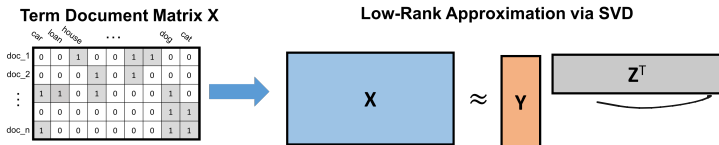


**Another View:** Each column of **Y** represents a 'topic'. $\vec{y}_i(j)$ indicates how much $doc_i$ belongs to topic $j$. $\vec{z}_a(j)$ indicates how much $word_a$ associates with that topic.

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

- Just like with documents, $\vec{z}_a$ and $\vec{z}_b$ will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.

# Example: Latent Semantic Analysis



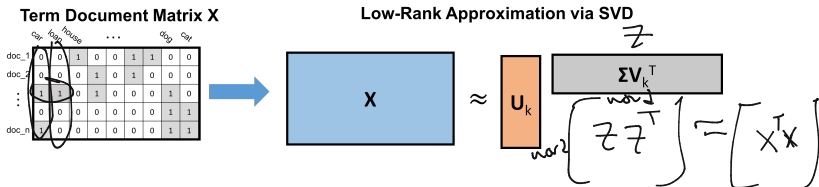**Term Document Matrix X**

**Low-Rank Approximation via SVD**

- Just like with documents, $\vec{z}_a$ and $\vec{z}_b$ will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.

- In an SVD decomposition we set $Z^T = \Sigma_k V_K^T$.

- The columns of $V_k$ are equivalently: the top $k$ eigenvectors of $X^T X$.

# Example: Latent Semantic Analysis



**Term Document Matrix X**

**Low-Rank Approximation via SVD**

- Just like with documents, $\vec{z}_a$ and $\vec{z}_b$ will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.

- In an SVD decomposition we set $Z^T = \Sigma_k V_K^T$.

- The columns of $V_k$ are equivalently: the top $k$ eigenvectors of $X^T X$.

- ~~Claim:~~ Exercise: $ZZ^T$ is the best rank-$k$ approximation of $X^T X$. I.e.,
  $\arg\min_{\text{rank}-k \text{ B}} \|X^T X - B\|_F$

## Example: Word Embedding

LSA gives a way of embedding words into $k$-dimensional space.

- Embedding is via low-rank approximation of $X^TX$: where $(X^TX)_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.

## Example: Word Embedding

LSA gives a way of embedding words into $k$-dimensional space.

- Embedding is via low-rank approximation of $X^TX$: where $(X^TX)_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.
- Think about $X^TX$ as a similarity matrix (gram matrix, kernel matrix) with entry $(a, b)$ being the similarity between $word_a$ and $word_b$.

## Example: Word Embedding

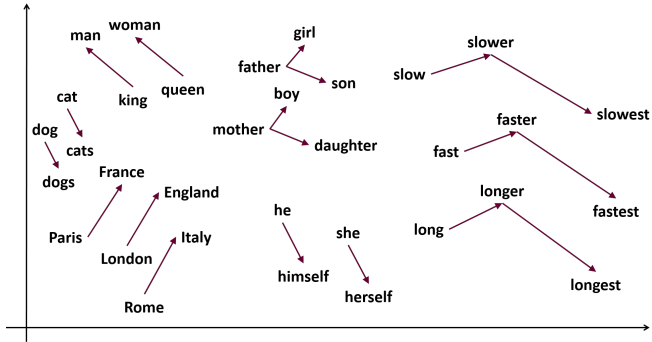LSA gives a way of embedding words into $k$-dimensional space.

- Embedding is via low-rank approximation of $X^T X$: where $(X^T X)_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.
- Think about $X^T X$ as a similarity matrix (gram matrix, kernel matrix) with entry $(a, b)$ being the similarity between $word_a$ and $word_b$.
- Many ways to measure similarity: number of sentences both occur in, number of times both appear in the same window of $w$ words, in similar positions of documents in different languages, etc.
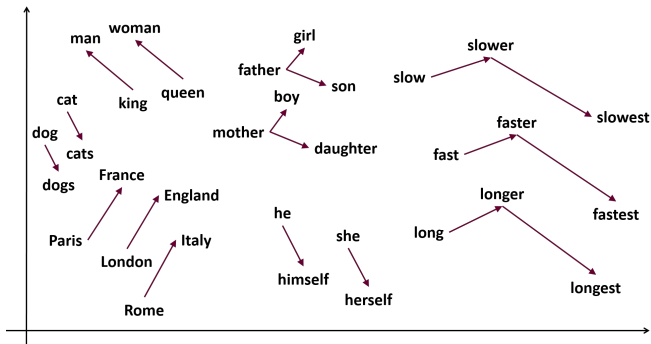
## Example: Word Embedding

LSA gives a way of embedding words into $k$-dimensional space.

- Embedding is via low-rank approximation of $X^T X$: where $(X^T X)_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.

- Think about $X^T X$ as a similarity matrix (gram matrix, kernel matrix) with entry $(a, b)$ being the similarity between $word_a$ and $word_b$.

- Many ways to measure similarity: number of sentences both occur in, number of times both appear in the same window of $w$ words, in similar positions of documents in different languages, etc.

- Replacing $X^T X$ with these different metrics (sometimes appropriately transformed) leads to popular word embedding algorithms: word2vec, GloVe, fastText, etc.

## Example: Word Embedding



**Note:** word2vec is typically described as a neural-network method, but can be viewed as just a low-rank approximation of a specific similarity matrix. *Neural word embedding as implicit matrix factorization,* Levy and Goldberg.

Questions?