# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2023.

Lecture 18

## Logistics

- Problem Set 3 is due next **Friday at 11:59pm**.

- I made a small change to Problem 1.4: replacing $\sum_{i=1}^{n} \sigma_i(\mathbf{A})^2$ with $\sum_{i=1}^{\text{rank}(\mathbf{A})} \sigma_i(\mathbf{A})^2$. This don't change the solution to the problem, but as we will see will better match the conventions for SVD that I introduce today. L6RC A3)1

- Linear algebra review session **Monday 2-3pm**. Location TBD.

## Summary

### Last Class

- Finish up optimal low-rank approximation via eigendecomposition.
- 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 for non-linear dimensionality reduction.

3

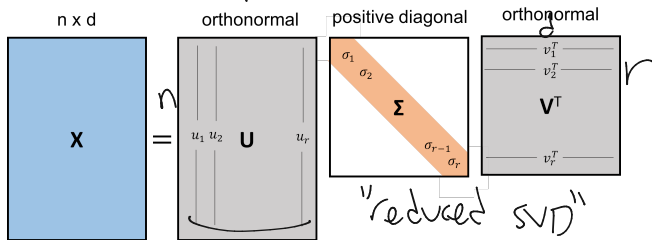# Singular Value Decomposition

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).

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

$X = U \Sigma V^T$



"reduced SVD"

## Connection of the SVD to Eigendecomposition

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

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

$X \in \mathbb{R}^{n \times d}$: data matrix, $U \in \mathbb{R}^{n \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{u}_1, \vec{u}_2, \dots$ (left singular vectors), $V \in \mathbb{R}^{d \times \text{rank}(X)}$: matrix with orthonormal columns $\vec{v}_1, \vec{v}_2, \dots$ (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 $X = U\Sigma V^T$:

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

$X \in \mathbb{R}^{n \times d}$: 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 $X = U\Sigma V^T$:

$$X^T X = V\Sigma U^T U\Sigma V^T = V\underline{\Sigma^2}V^T$$

$X \in \mathbb{R}^{n \times d}$: 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 $X = U\Sigma V^T$:

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

$X \in \mathbb{R}^{n \times d}$: 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: $XX^T = U\Sigma V^T V\Sigma U^T = U\Sigma^2 U^T$.

---

$X \in \mathbb{R}^{n \times d}$: 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 $X = U\Sigma V^T$:

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

※

$$\lambda_i(AB) = \lambda_i(BA)$$

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 $X^TX$ and the gram matrix $XX^T$ respectively.

$$\det(A - \lambda I) = 0$$

$$\sigma_i(X)^2 = \lambda_i(X^TX) = \lambda_i(XX^T)$$

$$(A - \lambda I)V$$

$$n\left[XX^T\right]_{ij} = \langle x_i, x_j \rangle$$

$$Av - \lambda v$$

$$\lambda v - \lambda v = 0$$

$X \in \mathbb{R}^{n \times d}$: 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.

$X^T X$

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).

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

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 $\underline{XV_k V_k^T}$ 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 $X = \underline{U\Sigma V^T}$:

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

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!

*(handwritten annotations: $V_k$ top $k$ eigvec? of $X^TX$; $XV_kV_k^T$ le st; LRA; $XV_kV_k^T = U_kU_k^TX$ ; $X$)*

> $X \in \mathbb{R}^{n \times d}$: 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 $X$.
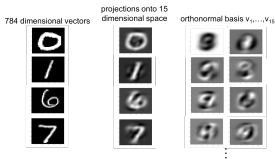
# The SVD and Optimal Low-Rank Approximation

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:

$$X_k = XV_kV_k^T = U_kU_k^TX$$

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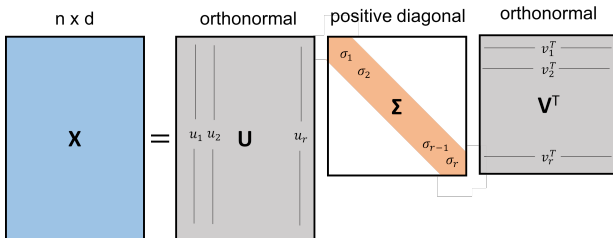


**Row (data point) compression**

**Column (feature) compression**

## 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:

$$X_k = XV_kV_k^T = U_kU_k^TX$$

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

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 $U_k$
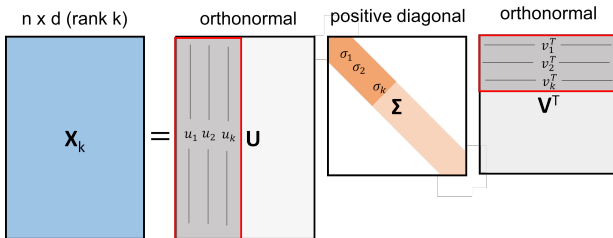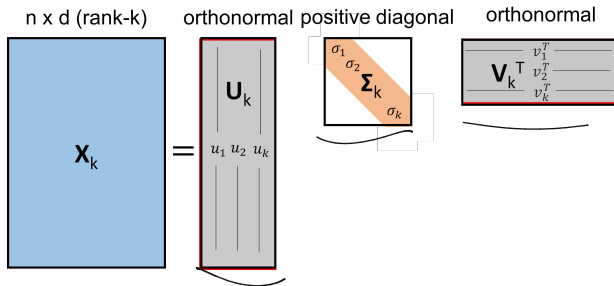
# The SVD and Optimal Low-Rank Approximation

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:

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

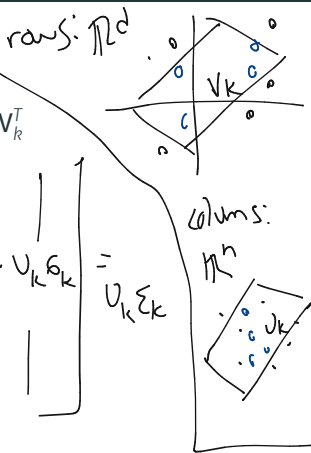Correspond to projecting the rows (data points) onto the span of $V_k$ or the columns (features) onto the span of $U_k$

# The SVD and Optimal Low-Rank Approximation

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

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



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

7

# The SVD and Optimal Low-Rank Approximation

The best low-rank approximation to $X$:
$$X_k = \arg\min_{\text{rank} - k \; B\in\mathbb{R}^{n\times d}}^{\text{rank}(B)=k} \|X - B\|_F \text{ is given by:}$$
$$X_k = XV_kV_k^T = U_kU_k^TX = U_k\Sigma_kV_k^T$$



$X \in \mathbb{R}^{n\times d}$: 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 $X$.

## SVD Review

$$SVD(X^TX) = EIG(X^TX) \quad X = U\Sigma V^T$$

$$V\Sigma^2 V^T$$

$$\left(X^TX\right)\left(X^TX\right)^T = \left(X^TX\right)^2 = V\Sigma^4 V^T$$

$$V\Lambda V^T$$

- 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 $X^TX$.
- $\Sigma \in \mathbb{R}^{r \times r}$ (diagonal) contains their eigenvalues.

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

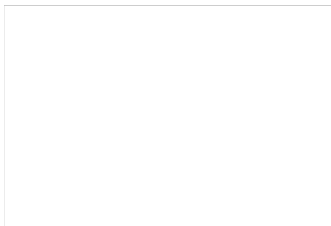$$n \begin{bmatrix} | & & | \\ v_1 & v. & v_n \\ | & & | \end{bmatrix}$$

9

# 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.

# 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.

Why might X be
close to low rank?
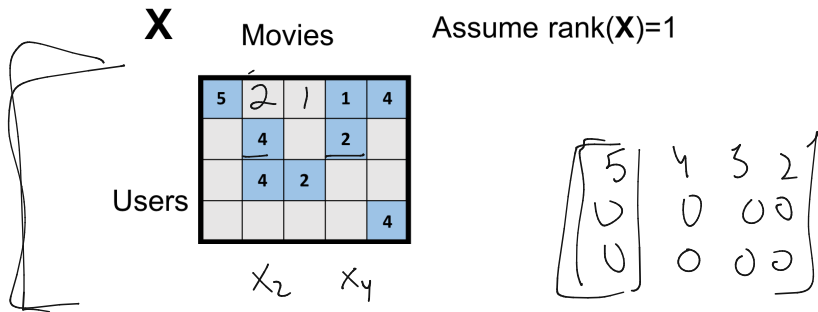
- users similar
- total ratings on movies
- genres

**X**  Movies

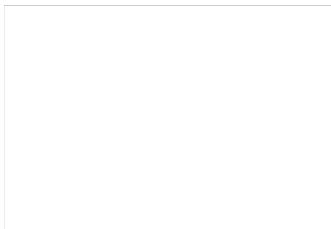| 5 |   |   | 1 | 4 |   |   |   |   |
|   |   | 3 |   |   |   |   | 5 |   |
|   |   |   |   |   |   |   |   |   |
|   |   |   |   | 4 |   |   |   |   |
|   |   |   |   |   |   |   |   |   |
|   |   | 5 |   |   |   |   |   | 5 |
| 1 |   |   | 2 |   |   |   |   |   |

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.
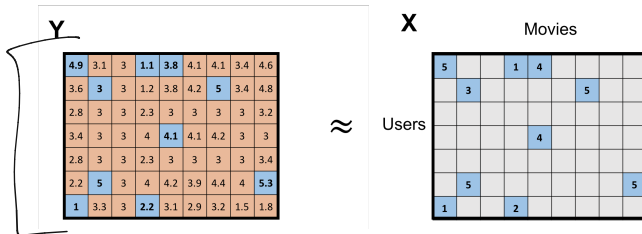


**X** Movies

Assume rank(**X**)=1

$$\begin{bmatrix} 5 & 2 & 1 & 1 & 4 \\ & 4 & & 2 & \\ & 4 & 2 & & \\ & & & & 4 \end{bmatrix}$$

Users

$X_2 \quad X_4$

$$\begin{bmatrix} 5 & 4 & 3 & 2 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

## 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.
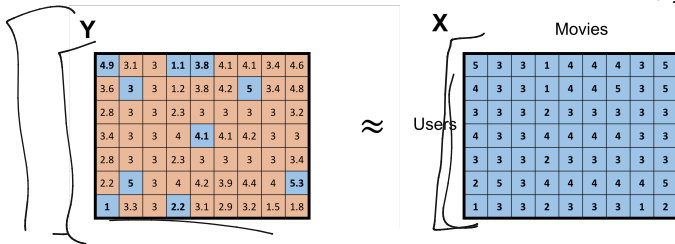


Solve: $Y = \underset{B \text{ s.t. } \mathrm{rank}(B) \leq k}{\arg\min} \sum_{\text{observed } (j,k)} \left[ X_{j,k} - B_{j,k} \right]^2$

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

## 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.

$\text{rank}(X) \leq \min(n, d)$



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

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.)
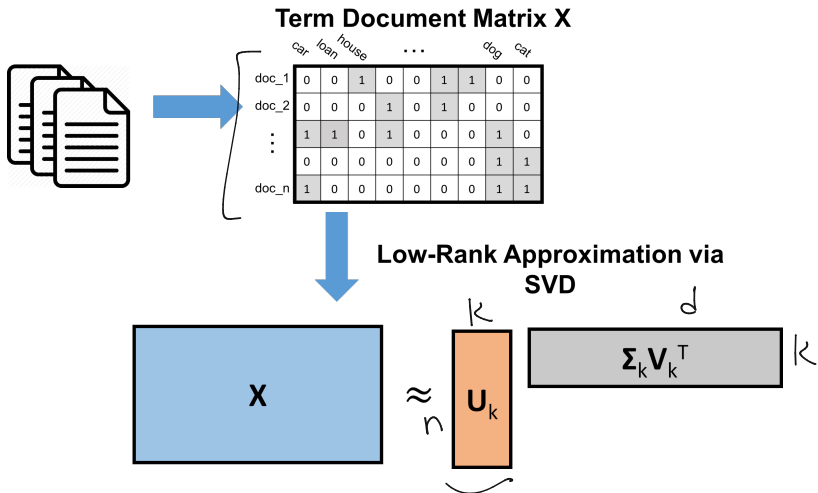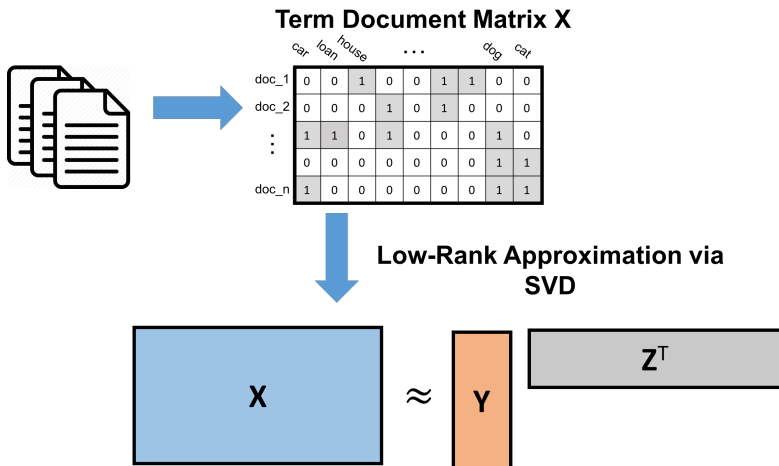- Nodes in a social network

## 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**

$$X \approx U_k \, \Sigma_k V_k^{\mathsf{T}}$$

**Term Document Matrix X**

Low-Rank Approximation via SVD

$X \approx Y \, Z^{\top}$

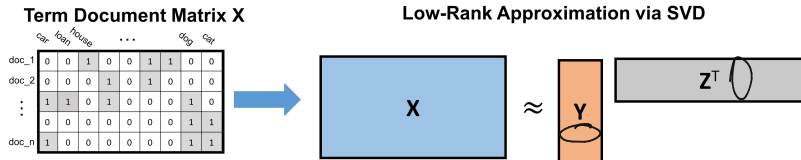**Term Document Matrix X**

**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**



**Low-Rank Approximation via SVD**



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

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

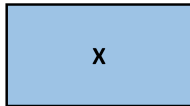- I.e., $\langle \vec{y}_i, \vec{z}_a \rangle \approx 1$ when $doc_i$ contains $word_a$.

# Example: Latent Semantic Analysis
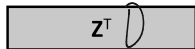
**Term Document Matrix X**



**Low-Rank Approximation via SVD**



each row of $y$ is k-dim representation of doc

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

$$X_{i,a} \approx (YZ^T)_{i,a} = \langle \vec{y}_i, \vec{z}_a \rangle.$$
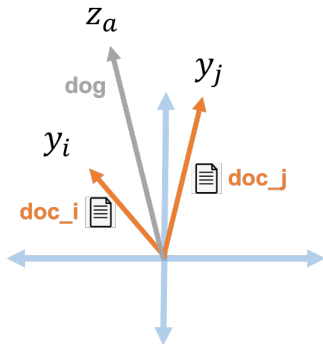
each column $Z$ is
" of word

- 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 $word_a$, $\underline{\langle \vec{y}_i, \vec{z}_a \rangle} \approx \underline{\langle \vec{y}_j, \vec{z}_a \rangle} \approx 1$.

$$y_i \approx y_j$$

If $doc_i$ and $doc_j$ both contain $word_a$, $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle \approx 1$
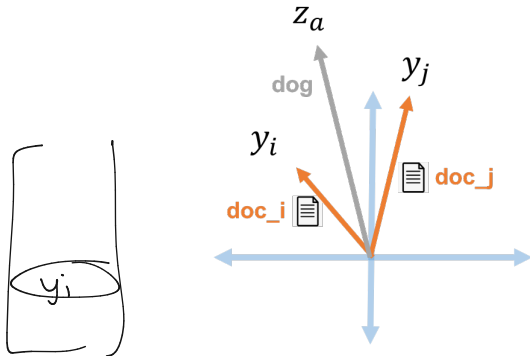
## Example: Latent Semantic Analysis

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.

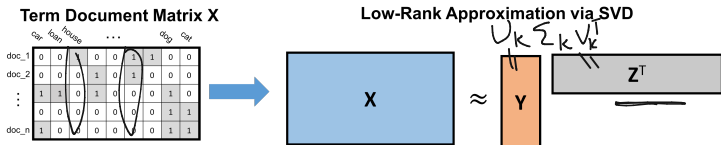# Example: Latent Semantic Analysis



**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.

# Example: Latent Semantic Analysis



**Term Document Matrix X** · · · · **Low-Rank Approximation via SVD**

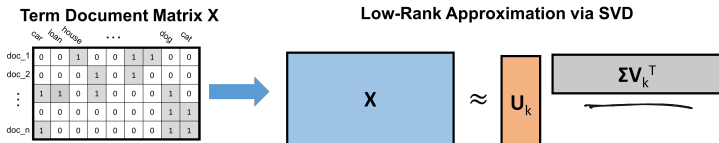$U_k \Sigma_k V_k^T$

$X \approx Y \quad Z^T$

- 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**

| | 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**

$$X \approx U_k \Sigma V_k^T$$

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

## Example: Word Embedding

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

- Embedding is via low-rank approximation of $X^T X$: where $(X^T X)_{a,b}$ 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 $X^T X$: 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 $word_b$.

## Example: Word Embedding

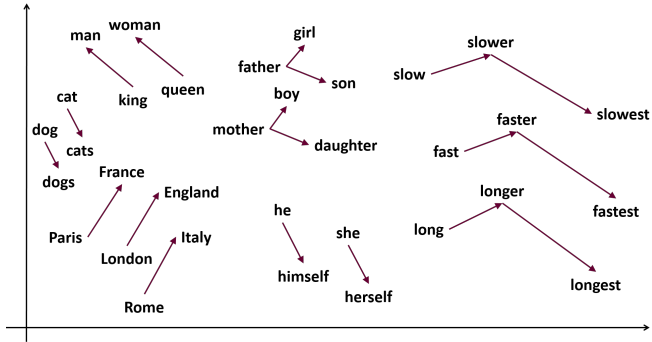LSA gives a way of embedding words into $k$-dimensional space.

- 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.
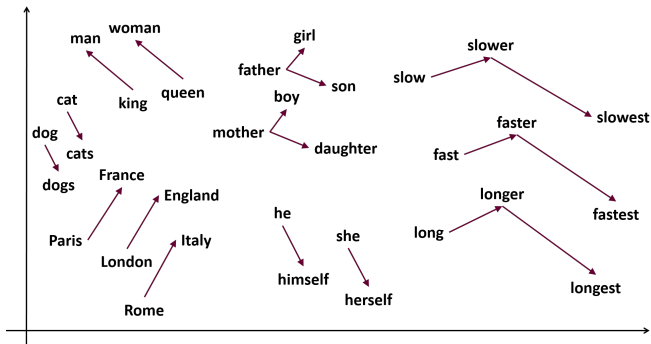
## Example: Word Embedding

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

- Embedding is via low-rank approximation of $X^TX$: 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 $word_b$.

- 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 $X^TX$ 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.