

COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Spring 2026.

Lecture 15


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

- office hours after class

↳ SVD, other applications of SVD.

Summary

Last Class

- Review of optimal low-rank approximation via eigendecomposition. 
- Linear algebra exercises and practice proofs.
- The Singular Value Decomposition (SVD) and its connection to eigendecomposition and low-rank approximation

This Class:

- Applications of SVD beyond low-rank approximation.
- Applications of low-rank approximation beyond compression
- 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? $X \in \mathbb{R}^{n \times d}$

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

TRUE

Proof by contradiction:

$$\min \|X - XV^T\|_F < \min \|X - B\|_F$$

X

this has
rank $\leq k$

$$\min \|X - XV^T\|_F > \min \|X - B\|_F$$

X

+ \bar{V} be a basis for B 's rows

$$\bar{V} \in \mathbb{R}^{d \times k}$$

$$\min \|X - XV^T\|_F < \|X - X\bar{V}\bar{V}^T\|_F \leq \|X - B\|_F$$

Low-Rank Approximation Review

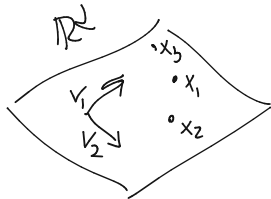
What is the value of

$$\min_{B: \text{rank}(B) \leq k} \|X - B\|_F^2?$$

$$\begin{aligned} & \sum_{i=k+1}^d \lambda_i (X^T X) = \sum_{i=k+1}^r \sigma_i(X)^2 \\ & \sum_{i=k+1}^d \lambda_i (X X^T) \end{aligned}$$

Column and Row Spans

Let $X \in \mathbb{R}^{n \times d}$ have its ~~columns~~ ^{rows} spanned by the ~~rows~~ ^{columns} of an orthonormal matrix $V \in \mathbb{R}^{d \times k}$. Show that the columns of X are spanned by the columns of $XV \in \mathbb{R}^{n \times k}$.



$$X =$$

columns

$$X_i = V c_i$$

$$X_i^T = c_i^T V^T$$

$$X = \begin{bmatrix} c_1^T V^T \\ c_2^T V^T \\ \vdots \end{bmatrix} = \begin{bmatrix} c \\ n \end{bmatrix} \begin{bmatrix} V^T \end{bmatrix}_k$$

XV^T :

X but with rows projected to subspace spanned by V 's columns $\text{col}(V)$

$$C = XV$$

$$C V^T = X V V^T = X$$

$XV v_i$

Matrix Completion

Assume that A is a rank-1 matrix, and that you have access to a subset of its entries, shown below:

$$c_1 = \frac{2}{3} \cdot c_2$$

$$A = \begin{bmatrix} \textcircled{6} & \textcircled{9} & 4 \\ x & 2 & y \\ 3 & z & 2 \end{bmatrix}$$

4.5

$$r_3 = r_1 \cdot \frac{1}{2}$$

What are the values of x , y , and z ?

$$x = \frac{6}{4.5} = \frac{4}{3}$$

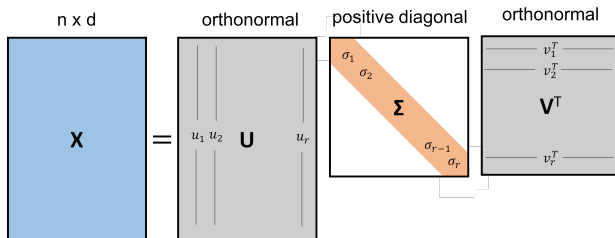
$$x = \frac{2}{3} \cdot 2 = \frac{4}{3}$$

Singular Value Decomposition

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

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

$$A = V\Lambda V^T$$



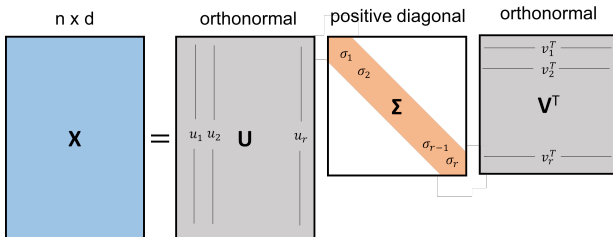
The SVD and Optimal Low-Rank Approximation

The best low-rank approximation to \mathbf{X} :

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

$$\mathbf{X}_k = \underline{\mathbf{X}\mathbf{V}_k\mathbf{V}_k^T} = \underline{\mathbf{U}_k\mathbf{U}_k^T\mathbf{X}}$$

Correspond to projecting the rows (data points) onto the span of \mathbf{V}_k or the columns (features) onto the span of \mathbf{U}_k



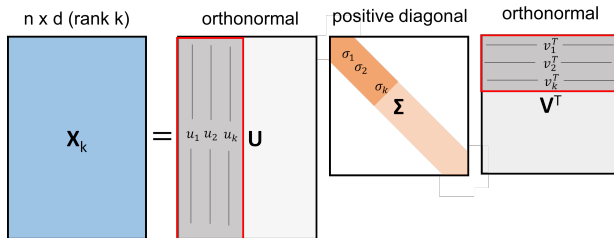
The SVD and Optimal Low-Rank Approximation

The best low-rank approximation to X :

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



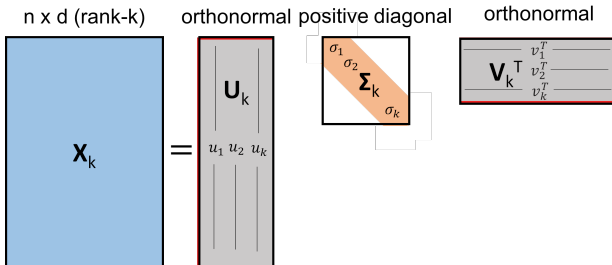
The SVD and Optimal Low-Rank Approximation

The best low-rank approximation to X :

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



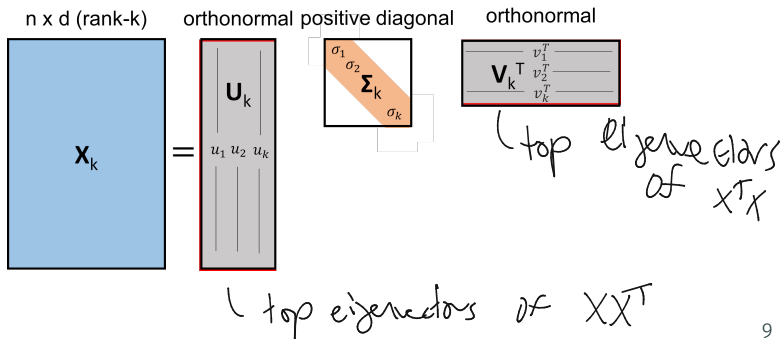
The SVD and Optimal Low-Rank Approximation

The best low-rank approximation to X :

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

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

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



Another Application of SVD

SVD is the 'Swiss army knife' of linear algebra. Lots of applications beyond low-rank approximation.

Another Application of SVD

SVD is the 'Swiss army knife' of linear algebra. Lots of applications beyond low-rank approximation.

Consider the following optimization problem: given $\mathbf{A} \in \mathbb{R}^{n \times d}$ and $\mathbf{b} \in \mathbb{R}^n$, compute $\arg \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2$.

What is this problem known as?

↳ least squares
linear regression

$$n \begin{bmatrix} \mathbf{A} \\ \hline \mathbf{a}_i^T \end{bmatrix} \begin{bmatrix} \mathbf{x} \end{bmatrix} \approx n \begin{bmatrix} \mathbf{b} \end{bmatrix}$$

$$\sum_{i=1}^n ([\mathbf{Ax}]_i - b_i)^2 = \sum_{i=1}^n (\underbrace{\mathbf{a}_i^T}_{\text{data point}} \mathbf{x} - b_i)^2$$

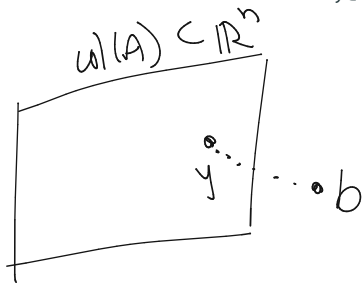
\downarrow linear function \downarrow scalar function

Another Application of SVD

Equivalent Formulation: Letting $y = Ax$, computing $\arg \min_x \|Ax - b\|_2^2$

is equivalent to finding: $\text{colspan}(A) : \{y \mid \overbrace{y = Ax \text{ for some } x}^{\text{for some } x}\}$

$$\arg \min_{y \in \text{colspan}(A)} \|y - b\|_2^2.$$



Another Application of SVD

Equivalent Formulation: Letting $\mathbf{y} = \mathbf{Ax}$, computing $\arg \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2$ is equivalent to finding:

$$\arg \min_{\mathbf{y} \in \text{colspan}(\mathbf{A})} \|\mathbf{y} - \mathbf{b}\|_2^2.$$

What is the solution to this problem?

Another Application of SVD

Equivalent Formulation: Letting $y = Ax$, computing $\arg \min_x \|Ax - b\|_2^2$ is equivalent to finding:

$$\arg \min_{y \in \text{colspan}(A)} \|y - b\|_2^2.$$

What is the solution to this problem? $y = \mathbf{U}\mathbf{U}^T\mathbf{b}$ where \mathbf{U} is an orthonormal basis for the columns of \mathbf{A}

projection onto the column span

Taking \mathbf{U} to be the left singular vectors of $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ and letting

$\mathbf{x}^* = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{b}$ we have:

$$y = Ax^* = \underbrace{U\Sigma V^T}_A \underbrace{V\Sigma^{-1}U^T b}_{x^*} = \underline{\underline{UU^T b}}$$

Another Application of SVD

Equivalent Formulation: Letting $\mathbf{y} = \mathbf{Ax}$, computing $\arg \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2$ is equivalent to finding:

$$\arg \min_{\mathbf{y} \in \text{colspan}(\mathbf{A})} \|\mathbf{y} - \mathbf{b}\|_2^2.$$

What is the solution to this problem? $\mathbf{y} = \mathbf{UU}^T\mathbf{b}$ where \mathbf{U} is an orthonormal basis for the columns of \mathbf{A}

Taking \mathbf{U} to be the left singular vectors of $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ and letting

$\mathbf{x}^* = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{b}$ we have:

$$\mathbf{y} = \mathbf{Ax} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{b} = \mathbf{UU}^T\mathbf{b}.$$

So $\mathbf{x}^* = \arg \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2$.

Another Application of SVD

Equivalent Formulation: Letting $\mathbf{y} = \mathbf{Ax}$, computing $\arg \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2$ is equivalent to finding:

$$\arg \min_{\mathbf{y} \in \text{colspan}(\mathbf{A})} \|\mathbf{y} - \mathbf{b}\|_2^2.$$

What is the solution to this problem? $\mathbf{y} = \mathbf{UU}^T\mathbf{b}$ where \mathbf{U} is an orthonormal basis for the columns of \mathbf{A}

Taking \mathbf{U} to be the left singular vectors of $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ and letting $\mathbf{x}^* = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{b}$ we have:

$$\mathbf{y} = \mathbf{Ax} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{b} = \mathbf{UU}^T\mathbf{b}.$$

So $\mathbf{x}^* = \arg \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2$.

$\mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T$ is the **Moore–Penrose pseudoinverse** of \mathbf{A} . Can be read off from the SVD.

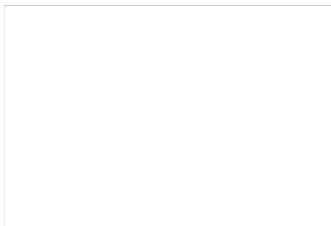
Applications of Low-Rank Approximation Beyond Compression

Matrix Completion

Consider a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ 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 $\mathbf{X} \in \mathbb{R}^{n \times d}$ 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

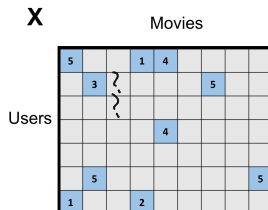
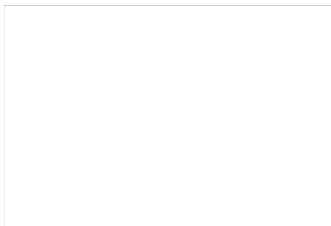
Users

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

Matrix Completion

Consider a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ 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

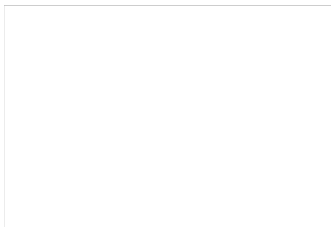
Consider a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank- k (i.e., well approximated by a rank k matrix).
Classic example: the Netflix prize problem.

X Movies Assume $\text{rank}(\mathbf{X})=1$

	5	2	1	1	4
		4		2	
		4	2		
Users					4

Matrix Completion

Consider a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ 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

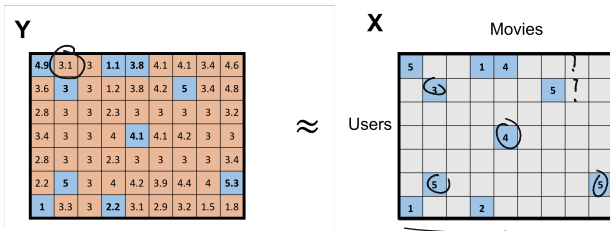
Users

5		1	4				
	3				5		
			4				
	5						5
1		2					

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

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.

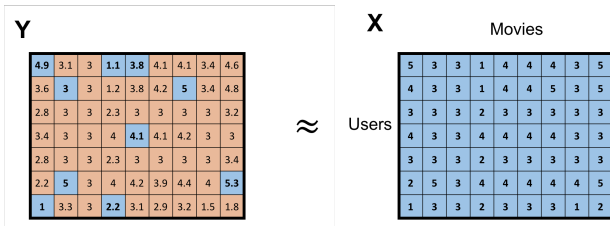


Solve: $Y = \arg \min_{B \text{ s.t. } \text{rank}(B) \leq k} \sum_{\text{observed } (j,k)} [X_{j,k} - B_{j,k}]^2$

$\|X - B\|_F$

Matrix Completion

Consider a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ 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.



$$\text{Solve: } \mathbf{Y} = \arg \min_{\mathbf{B} \text{ s.t. } \text{rank}(\mathbf{B}) \leq k} \sum_{\text{observed } (j,k)} [X_{j,k} - \mathbf{B}_{j,k}]^2$$

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

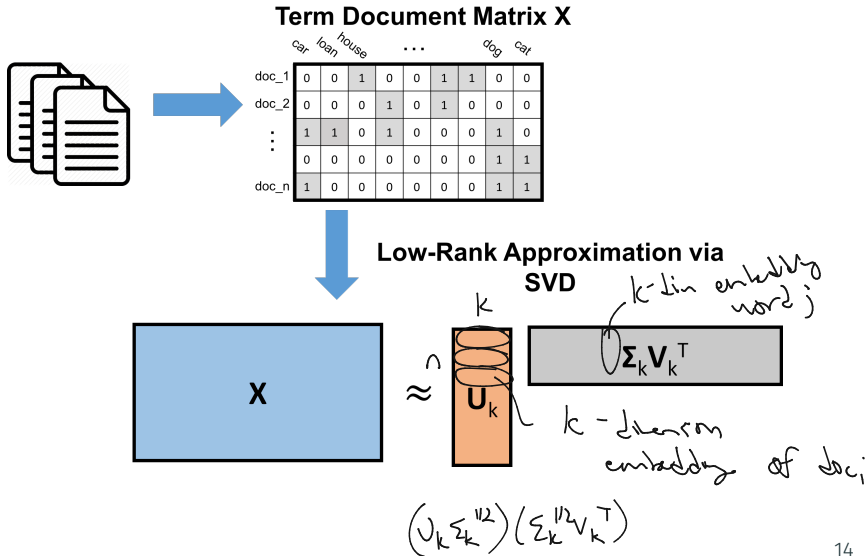
Entity Embeddings

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

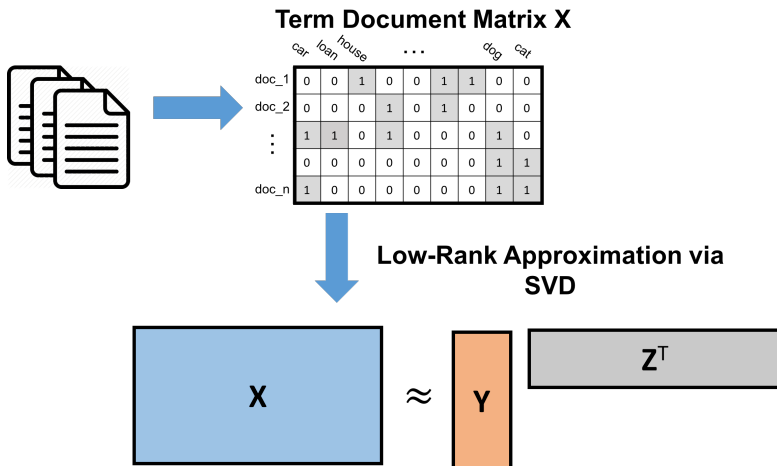
- Documents (for topic-based search and classification)
- Words (to identify synonyms, translations, etc.)
- Nodes in a social network

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

Example: Latent Semantic Analysis



Example: Latent Semantic Analysis



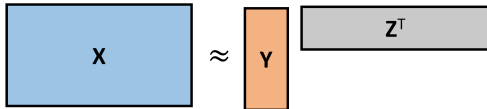
Example: Latent Semantic Analysis

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



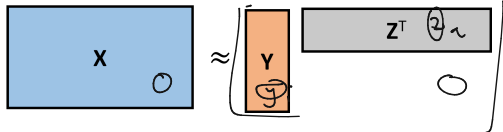
Example: Latent Semantic Analysis

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
⋮									
doc_n	1	0	0	0	0	0	0	1	1



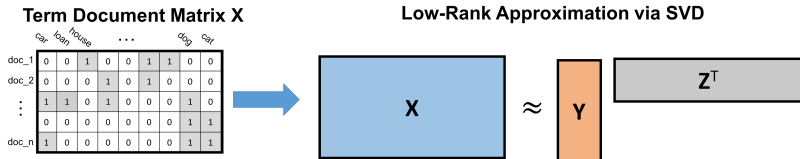
Low-Rank Approximation via SVD



- 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



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

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

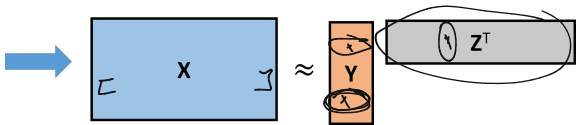
- I.e., $\langle \vec{y}_i, \vec{z}_a \rangle \approx 1$ when doc_i contains $word_a$.

Example: Latent Semantic Analysis

Term Document Matrix X

	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
doc_n	1	0	0	0	0	0	0	1	1

Low-Rank Approximation via SVD



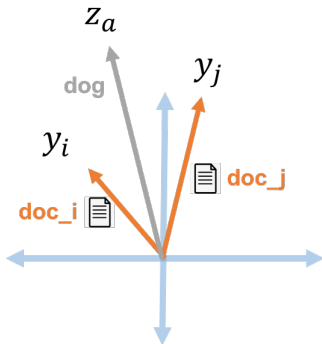
- If the error $\|X - YZ^T\|_F$ is small, then on average,

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

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

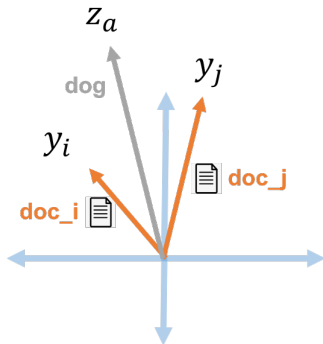
Example: Latent Semantic Analysis

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



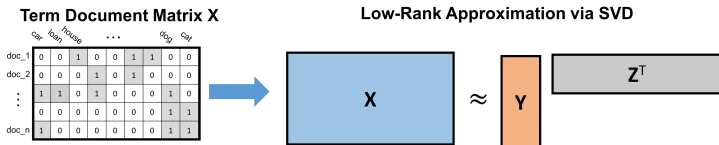
Example: Latent Semantic Analysis

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



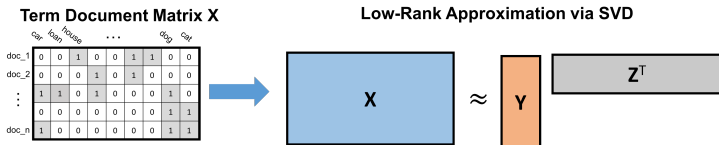
row τ
Another View: Each column of V represents a 'topic'. $\vec{y}_i(j)$ indicates how much doc_i belongs to topic j . $\vec{z}_a(j)$ indicates how much $word_a$ associates with that topic.

Example: Latent Semantic Analysis



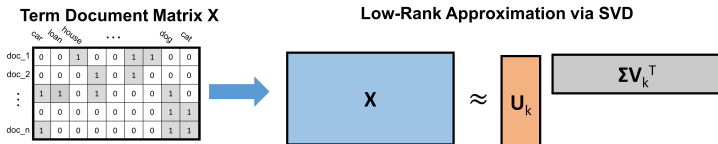
- Just like with documents, \vec{z}_a and \vec{z}_b will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.

Example: Latent Semantic Analysis



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

Example: Latent Semantic Analysis



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

Example: Word Embedding

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

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

Hand-drawn diagram illustrating the matrix $X^T X$. The matrix is shown as a square with rows labeled "word" and columns labeled "docs". The first row is circled and contains the values 0, 1, 0, 0. The matrix is labeled X^T .

Example: Word Embedding

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

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

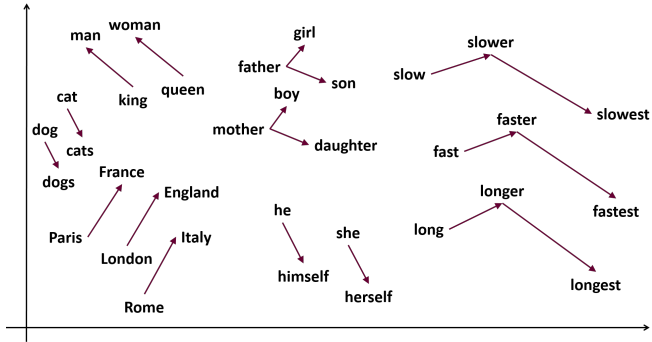
- Embedding is via low-rank approximation of $\mathbf{X}^T\mathbf{X}$: where $(\mathbf{X}^T\mathbf{X})_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.
- Think about $\mathbf{X}^T\mathbf{X}$ as a **similarity matrix** (gram matrix, kernel matrix) with entry (a, b) being the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

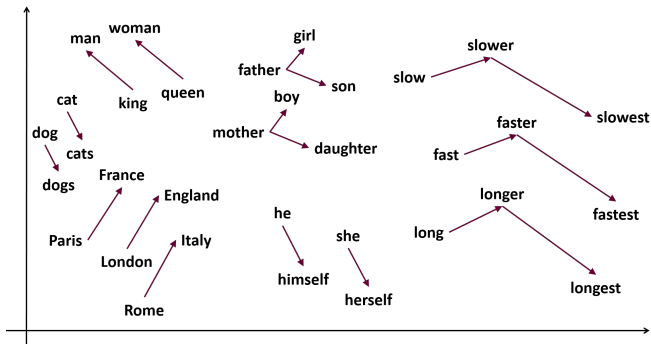
LSA gives a way of embedding words into k -dimensional space.

- Embedding is via low-rank approximation of $\mathbf{X}^T\mathbf{X}$: where $(\mathbf{X}^T\mathbf{X})_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.
- Think about $\mathbf{X}^T\mathbf{X}$ 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 $\mathbf{X}^T\mathbf{X}$ with these different metrics (sometimes appropriately transformed) leads to popular word embedding algorithms: word2vec, GloVe, fastText, etc.

Example: Word Embedding



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?