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

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

## COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2020.

Lecture 17

① eigenvalues of $X^T X$ = squared sing. values of $X$

② = singular values of $X^T X$

③ any sym. A

singular values of $A$ = absolute values of eigenvals of $A$.

- Problem Set 3 deadline extended until Monday 10/26, 8pm.
- Week 9 Quiz will now be due Tuesday 10/27, 8pm. ?

symmetric
$A = W \Lambda W^T$

What are eigenvalues of $A^2$?
$= \Lambda^2$?

r symmetric

$$A = W \Lambda W^T$$

where $W$ is orthogonal

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

$U = W \cdot S$     $S = \text{sign}(\Lambda)$

$\Sigma = |\Lambda|$

$V = W$

$A^2 = W \Lambda W^T W \Lambda W^T$
$= W \Lambda^2 W^T$

1

**Last Few Classes: Low-Rank Approximation and PCA**   $\rightarrow A = V \Lambda V^T$
$= V S \Sigma V^T$
where $\Sigma = |\Lambda|$
$S = sign(\Lambda)$

- Compress data that lies close to a $k$-dimensional subspace.
- Equivalent to finding a low-rank approximation of the data matrix $X$: $X \approx XVV^T$ for orthonormal $V \in \mathbb{R}^{d \times k}$.

$\Lambda \begin{bmatrix} -6 \\ 2 \\ 3 \end{bmatrix} \begin{bmatrix} 6 \\ 2 \\ 3 \end{bmatrix}$

- Optimal solution via PCA (eigendecomposition of $X^T X$ or equivalently, SVD of $X$).
- Singular vectors of $X$ are the eigenvectors of $XX^T$ and $X^T X$. Singular values squared are the eigenvalues.

$U, V$ are orthonormal
$\Sigma$ is positive diagonal

$$X = U \Sigma V^T$$
$$XX^T = U \Sigma^2 U^T \longrightarrow \text{eigendecomp.}$$
$$X^T X = V \Sigma^2 V^T$$

SVD of $X^T X$?

SVD of $X^T X$
positive semidefinite matrix
all eigens are +

if I have any symmetric $A$ with eigendecomp
$V \Lambda V^T =$
SVD of $A$?

2

### Last Few Classes: Low-Rank Approximation and PCA

- Compress data that lies close to a $k$-dimensional subspace.
- Equivalent to finding a low-rank approximation of the data matrix $X$: $X \approx XVV^T$ for orthonormal $V \in \mathbb{R}^{d \times k}$.
- Optimal solution via PCA (eigendecomposition of $X^T X$ or equivalently, SVD of $X$).
- Singular vectors of $X$ are the eigenvectors of $XX^T$ and $X^T X$. Singular values squared are the eigenvalues.

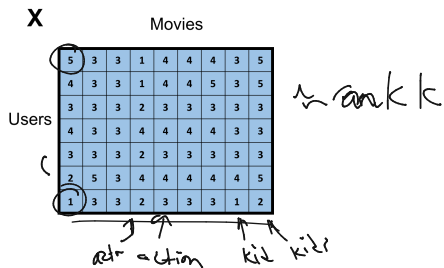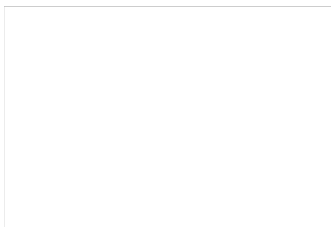### This Class: Applications of low-rank approx. beyond compression.

- Matrix completion and collaborative filtering
- Entity embeddings (word embeddings, node embeddings, etc.)
- Low-rank approximation fornon-linear dimensionality reduction.
- Spectral graph theory, spectral clustering.

2

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).
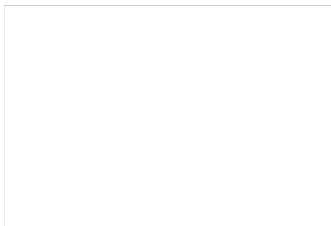
Consider a matrix $X \in \mathbb{R}^{n \times d}$ which we cannot fully observe but believe is close to rank-$k$ (i.e., well approximated by a rank $k$ matrix). Classic example: the Netflix prize problem. $1, 2, 3, 4, 5$

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

LRA : solve L by PCA

$$Y = \underset{\text{rank-}B}{\arg\min} \ \|X - B\|_F^2 = \sum_{i,j} (x_{ij} - B_{ij})^2$$

**X** Movies

Users

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

**Solve:** $Y = \underset{\text{rank} -k \ \mathbf{B}}{\arg\min} \underbrace{\sum_{\text{observed } (j,k)}}_{} \left[ \mathbf{X}_{j,k} - \mathbf{B}_{j,k} \right]^2$

3

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

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.

$$B = \left[ U W^T \right] k$$



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

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

$$y = \underset{z}{\arg\min} \; f(z)$$

3

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.



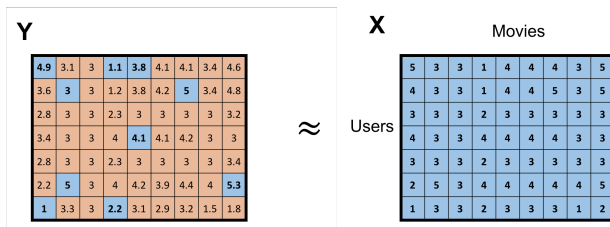**Solve:** $Y = \underset{\text{rank} -k \ B}{\arg\min} \sum_{\text{observed } (j, \ell)} \left[ X_{j, \ell} - B_{j, \ell} \right]^2$

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

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

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

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.)
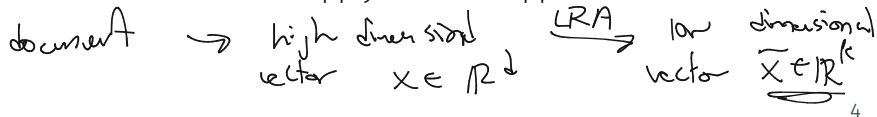- Nodes in a social network

*with a neural net*

*tune*

$$\begin{bmatrix} \ \end{bmatrix} \underset{\mathbb{R}^d}{} \overset{NN}{\longrightarrow} \begin{bmatrix} \ \end{bmatrix} \underset{\mathbb{R}^k}{}$$

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

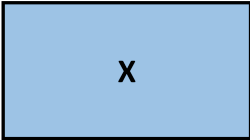document → high dimensional vector $x \in \mathbb{R}^d$ $\xrightarrow{LRA}$ low dimensional vector $\widetilde{x} \in \mathbb{R}^k$

**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 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |

bag of words

the

**Low-Rank Approximation via SVD**

$$X \approx U_k \; \Sigma_k V_k^\top$$

Term Document Matrix X

Low-Rank Approximation via SVD

$X \approx Y Z^T$

**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,

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

**Term Document Matrix X**

**Low-Rank Approximation via SVD**
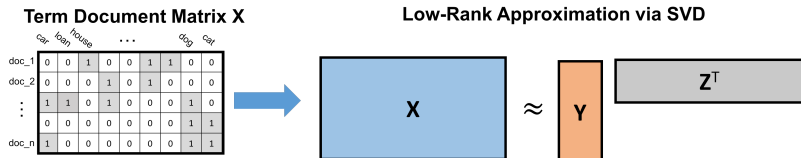
- If the error $\|X - YZ^T\|_F$ is small, then on average,
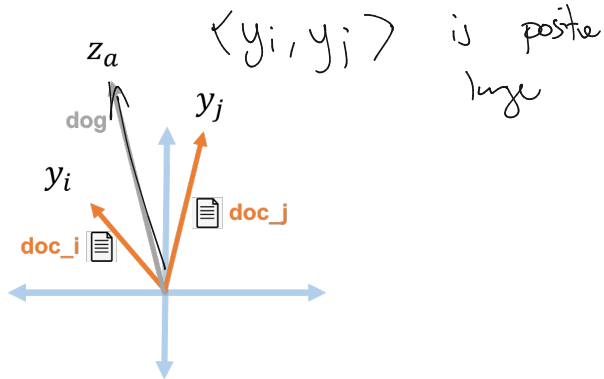
$$X_{i,a} \approx (YZ^T)_{i,a} = \langle \vec{y}_i, \vec{z}_a \rangle.$$

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

**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**

$$\mathbf{X} \approx \mathbf{Y}\ \mathbf{Z}^T$$

- If the error $\|\mathbf{X} - \mathbf{Y}\mathbf{Z}^T\|_F$ is small, then on average,

$$\mathbf{X}_{i,a} \approx (\mathbf{Y}\mathbf{Z}^T)_{i,a} = \langle \vec{y}_i, \vec{z}_a \rangle.$$

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

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



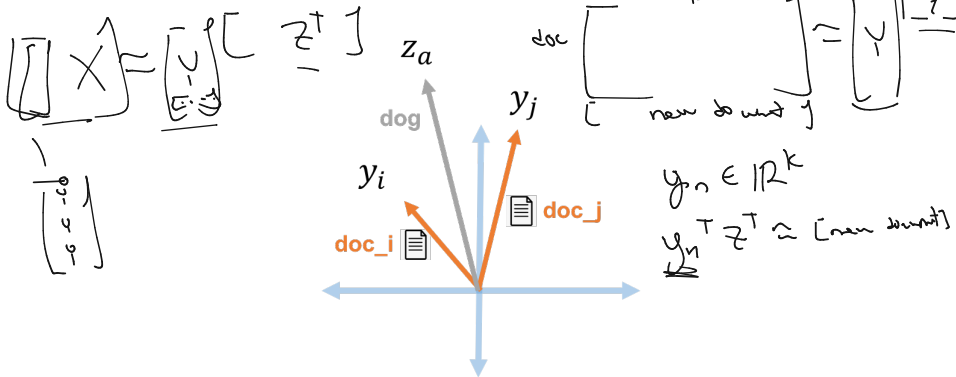$\langle y_i, y_j \rangle$ is postie large

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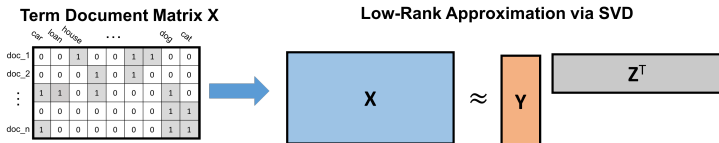


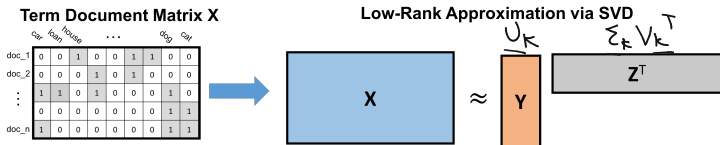**Another View:** Each column of $Y$ represents a 'topic'. $\vec{y}_i(j)$ indicates how much $doc_i$ belongs to topic $j$. $\vec{z}_a(j)$ indicates how much $word_a$ associates with that topic.

**Term Document Matrix X**

**Low-Rank Approximation via SVD**

$X \approx Y \, Z^T$

- 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.

**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 $Z^T = \mathbf{\Sigma}_k V_K^T$.
- The columns of $V_k$ are equivalently: the top $k$ eigenvectors of $X^T X$.

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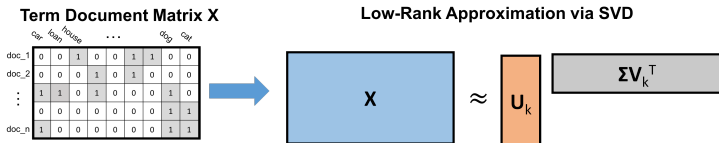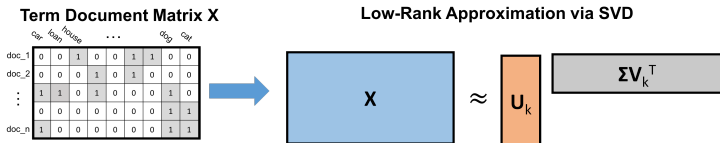
**Low-Rank Approximation via SVD**

- 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 $\mathbf{Z}^T = \boldsymbol{\Sigma}_k \mathbf{V}_k^T$.

- The columns of $\mathbf{V}_k$ are equivalently: the top $k$ eigenvectors of $\mathbf{X}^T\mathbf{X}$. The eigendecomposition of $\mathbf{X}^T\mathbf{X}$ is $\mathbf{X}^T\mathbf{X} = \mathbf{V}\boldsymbol{\Sigma}^2\mathbf{V}^T$.

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

**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 $\mathbf{Z}^T = \mathbf{\Sigma}_k \mathbf{V}_K^T$.

- The columns of $\mathbf{V}_k$ are equivalently: the top $k$ eigenvectors of $\mathbf{X}^T\mathbf{X}$. The eigendecomposition of $\mathbf{X}^T\mathbf{X}$ is $\mathbf{X}^T\mathbf{X} = \mathbf{V}\mathbf{\Sigma}^2\mathbf{V}^T$. SVD $X^TX$

- What is the best rank-$k$ approximation of $\mathbf{X}^T\mathbf{X}$? I.e. $\arg\min_{\text{rank}-k\ \mathbf{B}} \|\mathbf{X}^T\mathbf{X} - \mathbf{B}\|_F = V_k \Sigma_k^2 V_k^T = Z Z^T$
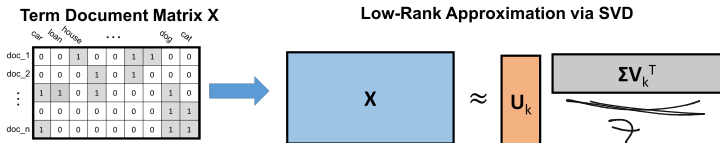
8

**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 $Z^T = \boldsymbol{\Sigma}_k V_K^T$.

- The columns of $V_k$ are equivalently: the top $k$ eigenvectors of $X^T X$. The eigendecomposition of $X^T X$ is $X^T X = V \boldsymbol{\Sigma}^2 V^T$.

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

- $X^T X = V_k \boldsymbol{\Sigma}_k^2 V_k^T = ZZ^T$.

$Z \in w \begin{bmatrix} k \\ \end{bmatrix}$   each row is an embedding of a work

$ZZ^T \approx X^TX$   $w \times w$ matrix

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.

docs

words $\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$ words $\begin{bmatrix} & & \\ \end{bmatrix}$ doc

$0 \cdot 0 + 1 \cdot 1 + 1 \cdot 1 = 2$

$(X^TX)_{ab} =$ # docs that $word_a$ + $word_b$ both appear in

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 $word_b$.

$$Z Z^T \text{ is best LRA to this similarity matrix.}$$

## EXAMPLE: WORD EMBEDDING
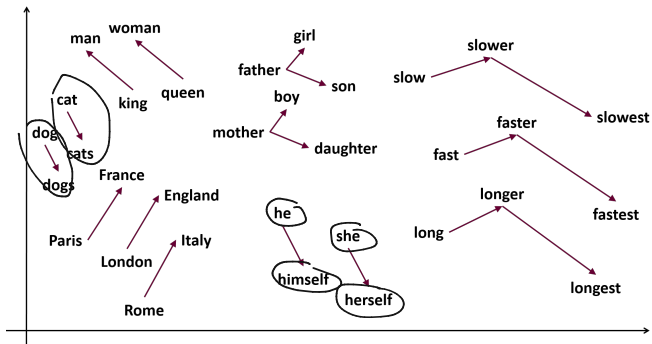
LSA gives a way of embedding words into $k$-dimensional space.

- 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.
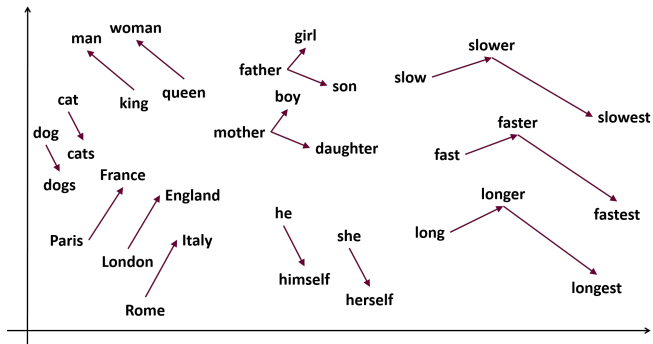
## EXAMPLE: WORD EMBEDDING

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.

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