

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Andrew McGregor

Lecture 17

Last Class: Low-Rank Approximation, Eigendecomposition, PCA

- For any symmetric square matrix A , we can write $A = V\Lambda V^T$ where columns of V are orthonormal eigenvectors.
- Can approximate data lying close to in a k -dimensional subspace by projecting data points into that space.
- Can find the best k -dimensional subspace via eigendecomposition applied to $\mathbf{X}^T\mathbf{X}$ (PCA).
- Measuring error in terms of the eigenvalue spectrum.

This Class: SVD and Applications

- SVD and connection to eigenvalue value decomposition.
- Applications of low-rank approximation beyond compression.

SINGULAR VALUE DECOMPOSITION

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

SINGULAR VALUE DECOMPOSITION

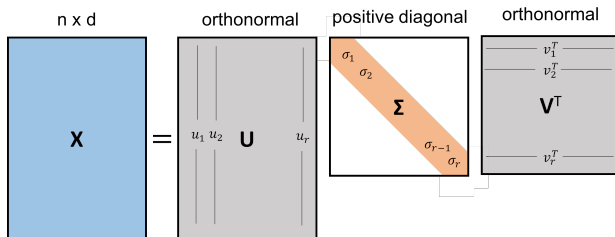
The Singular Value Decomposition (SVD) generalizes the eigendecomposition to asymmetric (even rectangular) matrices. Any matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ with $\text{rank}(\mathbf{X}) = r$ can be written as $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$.

- \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 $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ (singular values).

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 $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$.

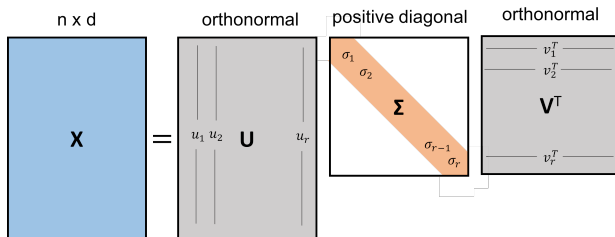
- \mathbf{U} has orthonormal columns $\vec{u}_1, \dots, \vec{u}_r \in \mathbb{R}^n$ (left singular vectors).
- \mathbf{V} has orthonormal columns $\vec{v}_1, \dots, \vec{v}_r \in \mathbb{R}^d$ (right singular vectors).
- $\mathbf{\Sigma}$ is diagonal with elements $\sigma_1 \geq \sigma_2 \geq \dots \geq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



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 $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ (singular values).



The 'swiss army knife' of modern linear algebra.

CONNECTION OF THE SVD TO EIGENDECOMPOSITION

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

$$\mathbf{X}^T \mathbf{X} =$$

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

CONNECTION OF THE SVD TO EIGENDECOMPOSITION

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

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

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

CONNECTION OF THE SVD TO EIGENDECOMPOSITION

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

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

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

CONNECTION OF THE SVD TO EIGENDECOMPOSITION

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

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

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

CONNECTION OF THE SVD TO EIGENDECOMPOSITION

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

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

Similarly: $\mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^T$.

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

CONNECTION OF THE SVD TO EIGENDECOMPOSITION

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

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

Similarly: $\mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^T$.

The right and left singular vectors are the eigenvectors of the covariance matrix $\mathbf{X}^T \mathbf{X}$ and the gram matrix $\mathbf{X} \mathbf{X}^T$ respectively.

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

CONNECTION OF THE SVD TO EIGENDECOMPOSITION

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

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

Similarly: $\mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^T$.

The right and left singular vectors are the eigenvectors of the covariance matrix $\mathbf{X}^T \mathbf{X}$ and the gram matrix $\mathbf{X} \mathbf{X}^T$ respectively.

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

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

CONNECTION OF THE SVD TO EIGENDECOMPOSITION

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

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

Similarly: $\mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^T$.

The right and left singular vectors are the eigenvectors of the covariance matrix $\mathbf{X}^T \mathbf{X}$ and the gram matrix $\mathbf{X} \mathbf{X}^T$ respectively.

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

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

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

CONNECTION OF THE SVD TO EIGENDECOMPOSITION

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

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

Similarly: $\mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^T$.

The right and left singular vectors are the eigenvectors of the covariance matrix $\mathbf{X}^T \mathbf{X}$ and the gram matrix $\mathbf{X} \mathbf{X}^T$ respectively.

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

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

Exercise: $\mathbf{U}_k \mathbf{U}_k^T \mathbf{X} = \mathbf{X} \mathbf{V}_k \mathbf{V}_k^T = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^T$

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

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:

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

THE SVD AND OPTIMAL LOW-RANK APPROXIMATION

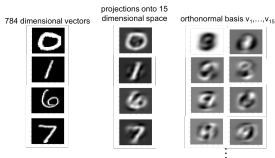
The best low-rank approximation to \mathbf{X} :

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

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

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

Row (data point) compression



Column (feature) compression

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

	bedrooms	bathrooms	sq.ft.	floors	list price	sale price
home 1	2	2	1800	2	200,000	195,000
home 2	4	2.5	2700	1	300,000	310,000
.
.
.
home n	5	3.5	3600	3	450,000	450,000

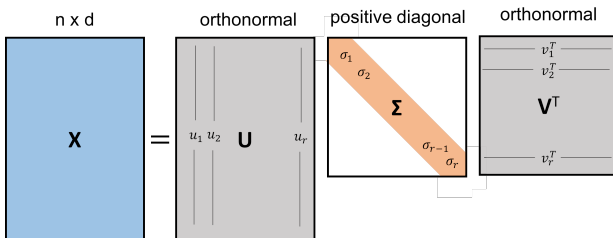
THE SVD AND OPTIMAL LOW-RANK APPROXIMATION

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:

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

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



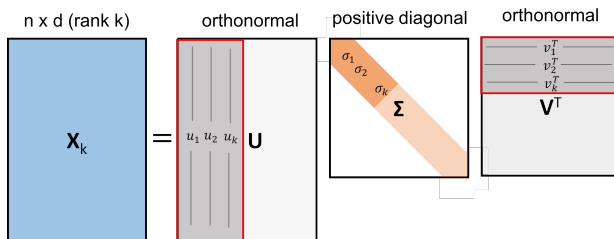
THE SVD AND OPTIMAL LOW-RANK APPROXIMATION

The best low-rank approximation to \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:

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

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



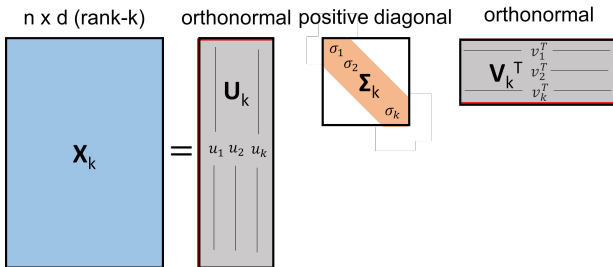
THE SVD AND OPTIMAL LOW-RANK APPROXIMATION

The best low-rank approximation to \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:

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

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



- Let $\vec{v}_1, \vec{v}_2, \dots, \in \mathbb{R}^d$ be orthonormal eigenvectors of $\mathbf{X}^T \mathbf{X}$.

BASIC IDEA TO PROVE EXISTENCE OF SVD

- Let $\vec{v}_1, \vec{v}_2, \dots, \in \mathbb{R}^d$ be orthonormal eigenvectors of $\mathbf{X}^T \mathbf{X}$.
- Let $\sigma_i = \|\mathbf{X} \vec{v}_i\|_2$ and define unit vector $\vec{u}_i = \frac{\mathbf{X} \vec{v}_i}{\sigma_i}$.

BASIC IDEA TO PROVE EXISTENCE OF SVD

- Let $\vec{v}_1, \vec{v}_2, \dots, \in \mathbb{R}^d$ be orthonormal eigenvectors of $\mathbf{X}^T \mathbf{X}$.
- Let $\sigma_i = \|\mathbf{X} \vec{v}_i\|_2$ and define unit vector $\vec{u}_i = \frac{\mathbf{X} \vec{v}_i}{\sigma_i}$.
- **Exercise:** Show $\vec{u}_1, \vec{u}_2, \dots$ are orthonormal.

BASIC IDEA TO PROVE EXISTENCE OF SVD

- Let $\vec{v}_1, \vec{v}_2, \dots, \in \mathbb{R}^d$ be orthonormal eigenvectors of $\mathbf{X}^T \mathbf{X}$.
- Let $\sigma_i = \|\mathbf{X} \vec{v}_i\|_2$ and define unit vector $\vec{u}_i = \frac{\mathbf{X} \vec{v}_i}{\sigma_i}$.
- **Exercise:** Show $\vec{u}_1, \vec{u}_2, \dots$ are orthonormal.
- This establishes that $\mathbf{X} \mathbf{V} = \mathbf{U} \mathbf{\Sigma}$ and that \mathbf{V} and \mathbf{U} have the required properties to show $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$.

BASIC IDEA TO PROVE EXISTENCE OF SVD

- Let $\vec{v}_1, \vec{v}_2, \dots, \in \mathbb{R}^d$ be orthonormal eigenvectors of $\mathbf{X}^T \mathbf{X}$.
- Let $\sigma_i = \|\mathbf{X} \vec{v}_i\|_2$ and define unit vector $\vec{u}_i = \frac{\mathbf{X} \vec{v}_i}{\sigma_i}$.
- **Exercise:** Show $\vec{u}_1, \vec{u}_2, \dots$ are orthonormal.
- This establishes that $\mathbf{XV} = \mathbf{U}\Sigma$ and that \mathbf{V} and \mathbf{U} have the required properties to show $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^T$.
- To see rest of the details, see https://math.mit.edu/classes/18.095/2016IAP/lec2/SVD_Notes.pdf

Rest of Class: Examples of how low-rank approximation is applied in a variety of data science applications.

Rest of Class: Examples of how low-rank approximation is applied in a variety of data science applications.

- Used for many reasons other than dimensionality reduction/data 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

Movies

Users

5			1	4				
	3					5		
				4				
	5							5
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.

\mathbf{X} Movies

Users

5			1	4				
	3					5		
				4				
	5							5
1			2					

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

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.

		Movies								
Users	Y	4.9	3.1	3	1.1	3.8	4.1	4.1	3.4	4.6
	3.6	3	3	1.2	3.8	4.2	5	3.4	4.8	
	2.8	3	3	2.3	3	3	3	3	3.2	
	3.4	3	3	4	4.1	4.1	4.2	3	3	
	2.8	3	3	2.3	3	3	3	3	3.4	
	2.2	5	3	4	4.2	3.9	4.4	4	5.3	
	1	3.3	3	2.2	3.1	2.9	3.2	1.5	1.8	

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

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

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

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

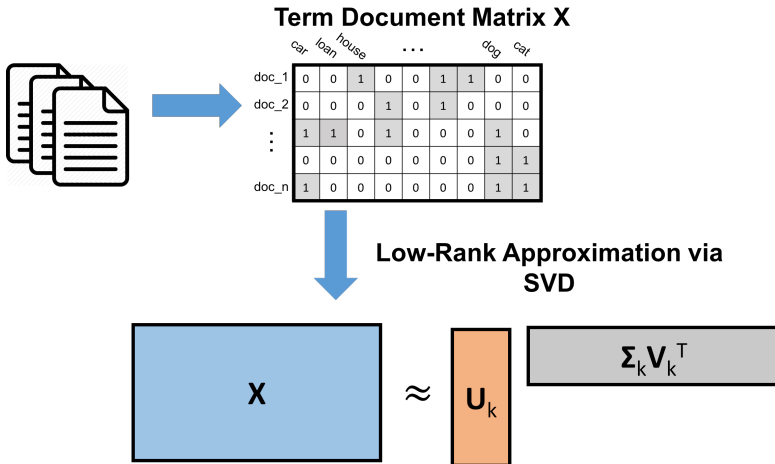
- Documents (for topic-based search and classification)
- Words (to identify synonyms, translations, etc.)
- Nodes in a social network

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

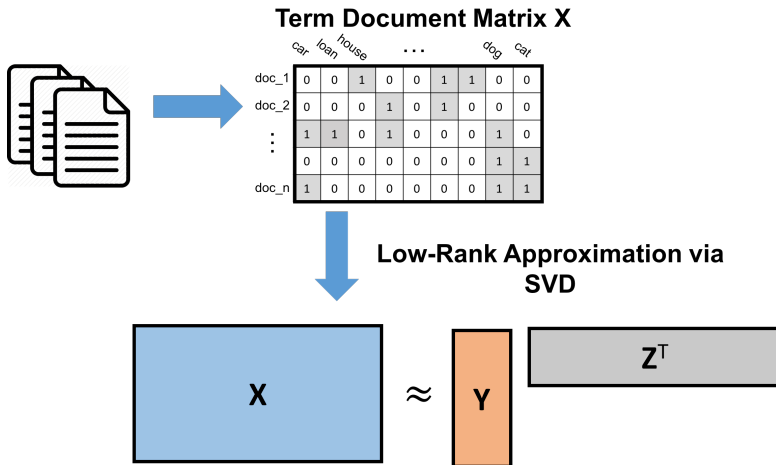
- Documents (for topic-based search and classification)
- Words (to identify synonyms, translations, etc.)
- Nodes in a social network

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

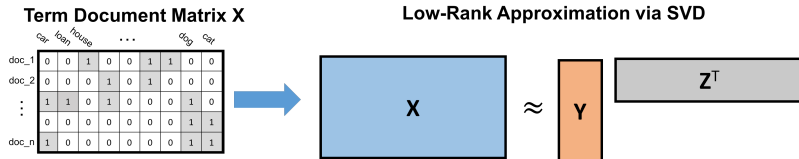
EXAMPLE: LATENT SEMANTIC ANALYSIS



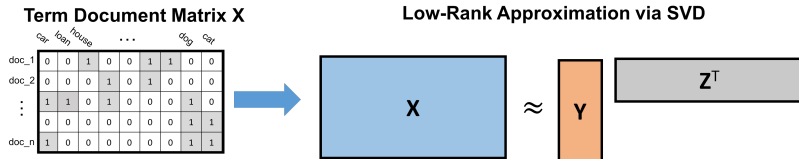
EXAMPLE: LATENT SEMANTIC ANALYSIS



EXAMPLE: LATENT SEMANTIC ANALYSIS



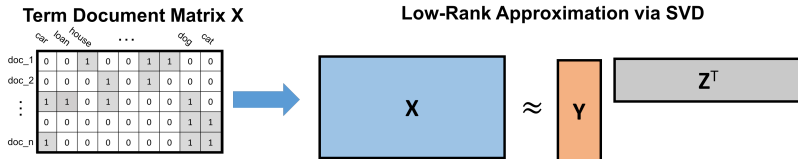
EXAMPLE: LATENT SEMANTIC ANALYSIS



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

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

EXAMPLE: LATENT SEMANTIC ANALYSIS

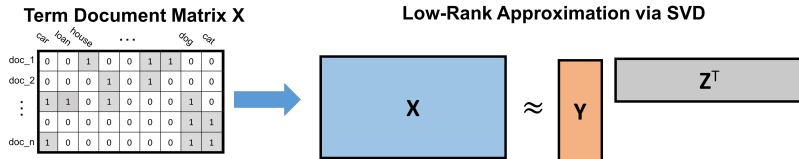


- If the error $\|\mathbf{X} - \mathbf{Y}\mathbf{Z}^T\|_F$ is small, then on average,

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

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

EXAMPLE: LATENT SEMANTIC ANALYSIS



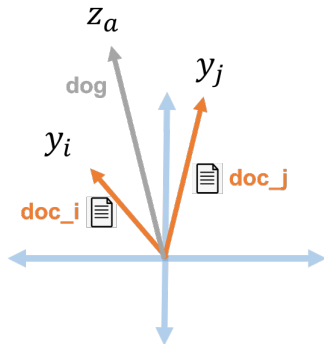
- If the error $\|\mathbf{X} - \mathbf{Y}\mathbf{Z}^T\|_F$ is small, then on average,

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

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

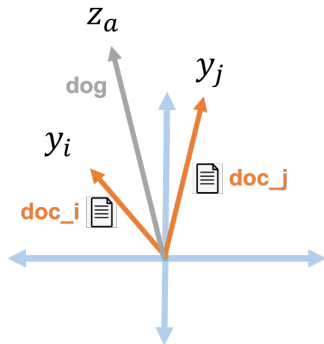
EXAMPLE: LATENT SEMANTIC ANALYSIS

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



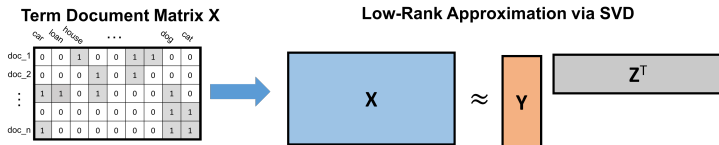
EXAMPLE: LATENT SEMANTIC ANALYSIS

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



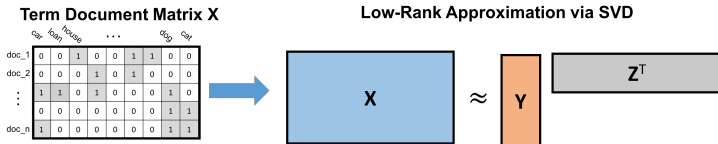
Another View: Each column of \mathbf{Y} represents a 'topic'. $\vec{y}_i(j)$ indicates how much doc_i belongs to topic j . $\vec{z}_a(j)$ indicates how much $word_a$ associates with that topic.

EXAMPLE: LATENT SEMANTIC ANALYSIS



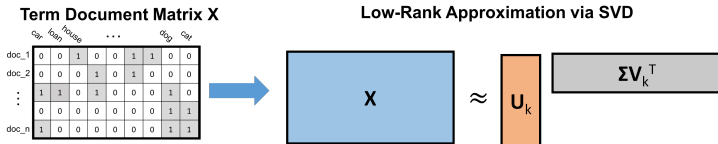
- Just like with documents, \vec{z}_a and \vec{z}_b will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.

EXAMPLE: LATENT SEMANTIC ANALYSIS



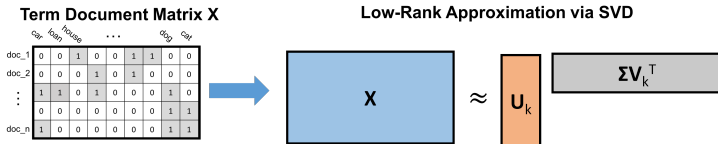
- Just like with documents, \vec{z}_a and \vec{z}_b will tend to have high dot product if $word_a$ and $word_b$ appear in many of the same documents.
- In an SVD decomposition we set $\mathbf{Z}^T = \mathbf{\Sigma}_k \mathbf{V}_K^T$.
- The columns of \mathbf{V}_k are equivalently: the top k eigenvectors of $\mathbf{X}^T \mathbf{X}$.

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 $\mathbf{Z}^T = \Sigma_k \mathbf{V}_k^T$.
- 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 $\mathbf{X}^T \mathbf{X} = \mathbf{V} \Sigma^2 \mathbf{V}^T$.

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 $\mathbf{Z}^T = \mathbf{\Sigma}_k \mathbf{V}_k^T$.
- 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 $\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^T$.
- The best rank- k approximation of $\mathbf{X}^T \mathbf{X}$ is

$$\arg \min_{\text{rank}-k \text{ } \mathbf{B}} \|\mathbf{X}^T \mathbf{X} - \mathbf{B}\|_F = \mathbf{V}_k \mathbf{\Sigma}_k^2 \mathbf{V}_k^T = \mathbf{Z} \mathbf{Z}^T$$

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.

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.

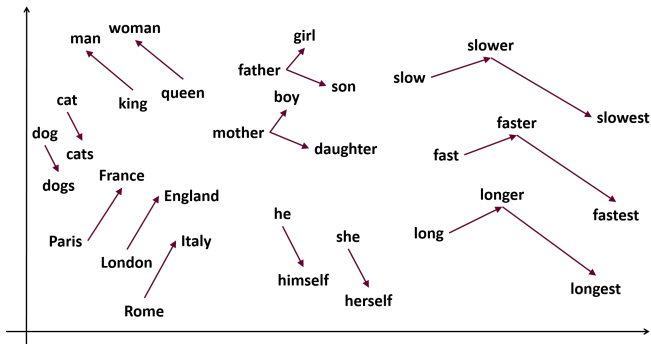
- Embedding is via low-rank approximation of $\mathbf{X}^T \mathbf{X}$: where $(\mathbf{X}^T \mathbf{X})_{a,b}$ is the number of documents that both $word_a$ and $word_b$ appear in.
- Think about $\mathbf{X}^T \mathbf{X}$ as a **similarity matrix** (gram matrix, kernel matrix) with entry (a, b) being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

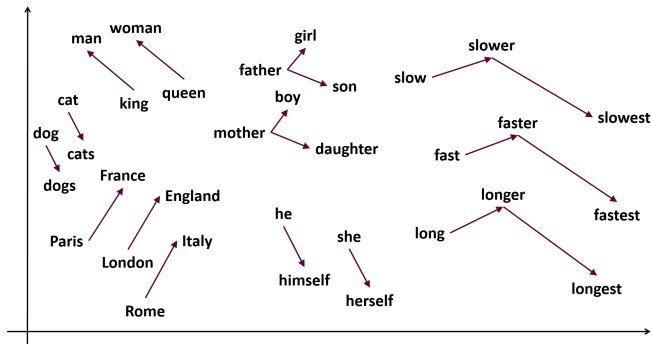
LSA gives a way of embedding words into k -dimensional space.

- 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 it is really just low-rank approximation of a specific similarity matrix. *Neural word embedding as implicit matrix factorization*, Levy and Goldberg.