

# COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

---

Cameron Musco

University of Massachusetts Amherst. Spring 2020.

Lecture 16

## Last Class: Low-Rank Approximation, Eigendecomposition, and PCA

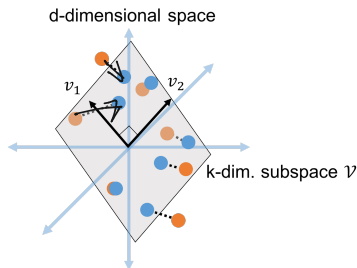
- Can approximate data lying close to in a  $k$ -dimensional subspace by projecting data points into that space.
- Finding the best  $k$ -dimensional subspace via eigendecomposition (PCA).  $[X^T X]_d \quad v_1 \dots v_k$
- Measuring error in terms of the eigenvalue spectrum.

## This Class: Finish Low-Rank Approximation and Connection to the singular value decomposition (SVD)

- Finish up PCA – runtime considerations and picking  $k$ .
- View of optimal low-rank approximation using the SVD.
- Applications of low-rank approximation beyond compression.

## BASIC SET UP

**Set Up:** Assume that data points  $\vec{x}_1, \dots, \vec{x}_n$  lie close to any  $k$ -dimensional subspace  $\mathcal{V}$  of  $\mathbb{R}^d$ . Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be the data matrix.



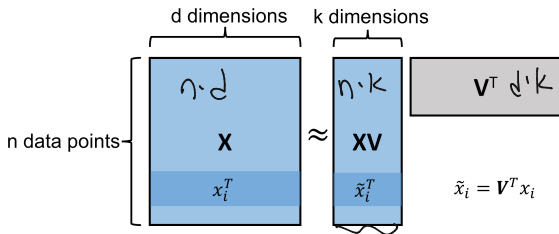
Let  $\vec{v}_1, \dots, \vec{v}_k$  be an orthonormal basis for  $\mathcal{V}$  and  $\mathbf{V} \in \mathbb{R}^{d \times k}$  be the matrix with these vectors as its columns.

- $\mathbf{W}\mathbf{W}^T \in \mathbb{R}^{d \times d}$  is the **projection matrix** onto  $\mathcal{V}$ .
- $\mathbf{X} \approx \mathbf{X}(\mathbf{W}\mathbf{W}^T)$ . Gives the closest approximation to  $\mathbf{X}$  with rows in  $\mathcal{V}$ .

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

## BASIC SET UP

**Set Up:** Assume that data points  $\vec{x}_1, \dots, \vec{x}_n$  lie close to any  $k$ -dimensional subspace  $\mathcal{V}$  of  $\mathbb{R}^d$ . Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be the data matrix.



Let  $\vec{v}_1, \dots, \vec{v}_k$  be an orthonormal basis for  $\mathcal{V}$  and  $\mathbf{V} \in \mathbb{R}^{d \times k}$  be the matrix with these vectors as its columns.

- $\mathbf{V}\mathbf{V}^T \in \mathbb{R}^{d \times d}$  is the **projection matrix** onto  $\mathcal{V}$ .
- $\mathbf{X} \approx \mathbf{X}(\mathbf{V}\mathbf{V}^T)$ . Gives the closest approximation to  $\mathbf{X}$  with rows in  $\mathcal{V}$ .

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# LOW-RANK APPROXIMATION VIA EIGENDECOMPOSITION

$\mathbf{V}$  minimizing  $\|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2$  is given by:

$$\arg \max_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \underbrace{\|\mathbf{X}\mathbf{V}\|_F^2}_{=} = \sum_{j=1}^k \underbrace{\|\mathbf{X}\vec{v}_j\|_2^2}_{=}$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# LOW-RANK APPROXIMATION VIA EIGENDECOMPOSITION

$\mathbf{V}$  minimizing  $\|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2$  is given by:

$$\arg \max_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X}\mathbf{V}\|_F^2 = \sum_{j=1}^k \|\mathbf{X}\vec{v}_j\|_2^2$$

**Solution via eigendecomposition:** Letting  $\mathbf{V}_k$  have columns  $\vec{v}_1, \dots, \vec{v}_k$  corresponding to the top  $k$  eigenvectors of the covariance matrix  $\mathbf{X}^T\mathbf{X}$ ,  $k$

$$\mathbf{V}_k = \arg \max_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X}\mathbf{V}\|_F^2$$

$$\begin{bmatrix} \mathbf{X}^T \mathbf{X} \end{bmatrix} \downarrow \begin{bmatrix} \vec{v}_1 \dots \vec{v}_k \end{bmatrix} \downarrow$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# LOW-RANK APPROXIMATION VIA EIGENDECOMPOSITION

$\mathbf{V}$  minimizing  $\|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2$  is given by:

$$\left[ \begin{array}{l} \arg \max \\ \text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k} \end{array} \right] \|\mathbf{X}\mathbf{V}\|_F^2 = \sum_{j=1}^k \|\mathbf{X}\vec{v}_j\|_2^2$$

**Solution via eigendecomposition:** Letting  $\mathbf{V}_k$  have columns  $\vec{v}_1, \dots, \vec{v}_k$  corresponding to the top  $k$  eigenvectors of the covariance matrix  $\mathbf{X}^T\mathbf{X}$ ,

$$\left[ \begin{array}{l} \vec{v}_1 \\ \vdots \\ \vec{v}_k \end{array} \right]^T \sum_i \mathbf{x}_i \mathbf{x}_i^T \mathbf{V}_k = \arg \max_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X}\mathbf{V}\|_F^2 \quad \sum \mathbf{x}_i = 0$$

$$\mathbb{E}(\mathbf{Y} \cdot \mathbb{E}(\mathbf{Y}))^2 = \mathbb{E} \mathbf{Y}^2$$

- Proof via Courant-Fischer and greedy maximization.
- Approximation error is  $\underbrace{\|\mathbf{X}\|_F^2}_{\text{total variance}} - \underbrace{\|\mathbf{X}\mathbf{V}_k\|_F^2}_{\text{variance captured}} = \sum_{i=k+1}^d \lambda_i(\mathbf{X}^T\mathbf{X})$ .

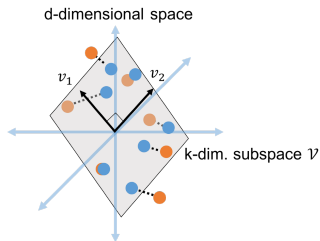
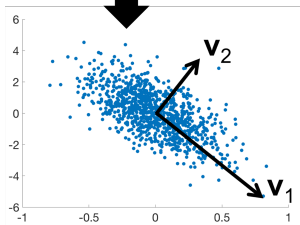
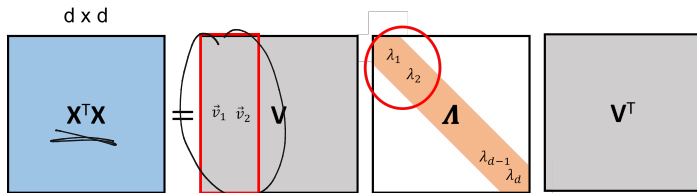
$$\text{tr}(\mathbf{X}^T\mathbf{X}) = \|\mathbf{X}\|_F^2$$

$$\|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# LOW-RANK APPROXIMATION VIA EIGENDECOMPOSITION

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq \lambda_d$$



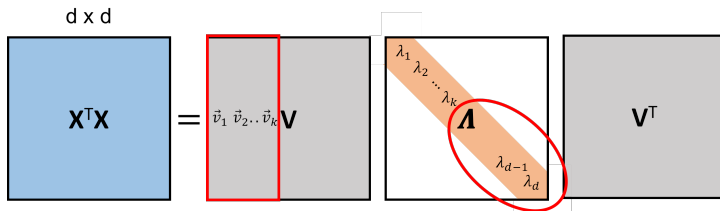


Plotting the **spectrum** of the covariance matrix  $\mathbf{X}^T\mathbf{X}$  (its eigenvalues) shows how compressible  $\mathbf{X}$  is using low-rank approximation (i.e., how close  $\vec{x}_1, \dots, \vec{x}_n$  are to a low-dimensional subspace).

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# SPECTRUM ANALYSIS

Plotting the **spectrum** of the covariance matrix  $\mathbf{X}^T\mathbf{X}$  (its eigenvalues) shows how compressible  $\mathbf{X}$  is using low-rank approximation (i.e., how close  $\vec{x}_1, \dots, \vec{x}_n$  are to a low-dimensional subspace).



error of optimal low rank approximation

$$\| \mathbf{X} - \mathbf{X} \mathbf{V}_k \mathbf{V}_k^T \|_F^2 = \sum_{i=k+1}^d \lambda_i (\mathbf{X}^T \mathbf{X})$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

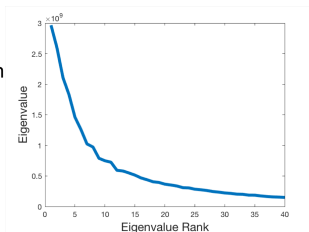
# SPECTRUM ANALYSIS

Plotting the **spectrum** of the covariance matrix  $X^T X$  (its eigenvalues) shows how compressible  $X$  is using low-rank approximation (i.e., how close  $\vec{x}_1, \dots, \vec{x}_n$  are to a low-dimensional subspace).

784 dimensional vectors



eigendecomposition



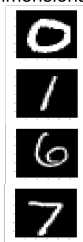
$\text{eig}(X^T X)$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $X \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $X^T X$ ,  $V_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

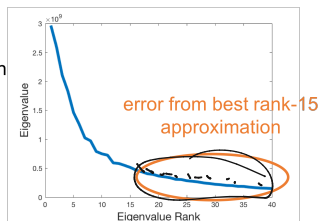
# SPECTRUM ANALYSIS

Plotting the **spectrum** of the covariance matrix  $\mathbf{X}^T\mathbf{X}$  (its eigenvalues) shows how compressible  $\mathbf{X}$  is using low-rank approximation (i.e., how close  $\vec{x}_1, \dots, \vec{x}_n$  are to a low-dimensional subspace).

784 dimensional vectors



eigendecomposition

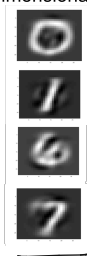


$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

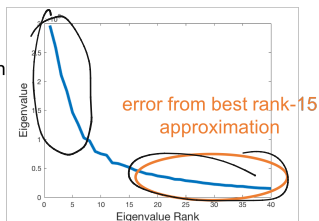
# SPECTRUM ANALYSIS

Plotting the **spectrum** of the covariance matrix  $\mathbf{X}^T\mathbf{X}$  (its eigenvalues) shows how compressible  $\mathbf{X}$  is using low-rank approximation (i.e., how close  $\vec{x}_1, \dots, \vec{x}_n$  are to a low-dimensional subspace).

784 dimensional vectors



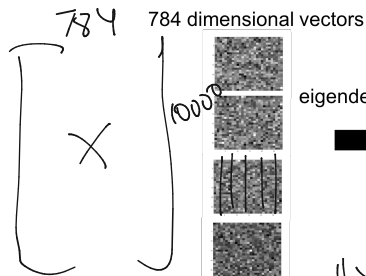
eigendecomposition



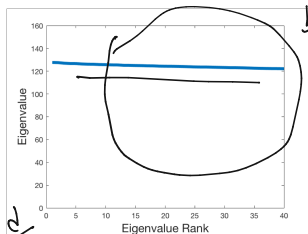
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# SPECTRUM ANALYSIS

Plotting the **spectrum** of the covariance matrix  $X^T X$  (its eigenvalues) shows how compressible  $X$  is using low-rank approximation (i.e., how close  $\vec{x}_1, \dots, \vec{x}_n$  are to a low-dimensional subspace).



$\lambda_i(X^T X)$  is the norm captured by  $v_i$



$\text{eig}(X^T X)$

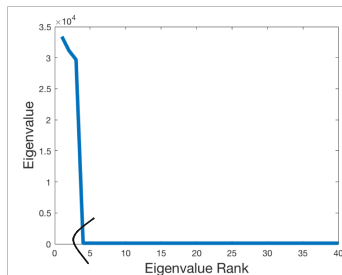
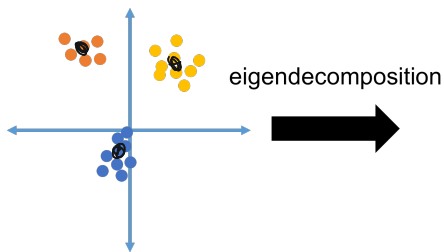
$$\|X\|_F^2 = \sum_{i=1}^n \lambda_i(X^T X)$$

$$\|X - X V_k V_k^T\|_F^2 = \sum_{i=k+1}^n \lambda_i(X^T X)$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $X \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $X^T X$ ,  $V_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# SPECTRUM ANALYSIS

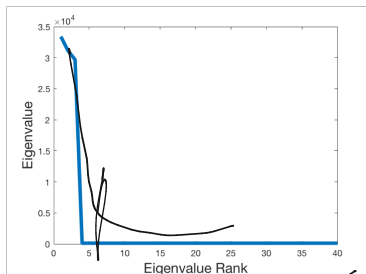
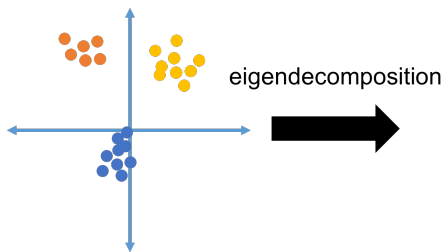
Plotting the **spectrum** of the covariance matrix  $\mathbf{X}^T\mathbf{X}$  (its eigenvalues) shows how compressible  $\mathbf{X}$  is using low-rank approximation (i.e., how close  $\vec{x}_1, \dots, \vec{x}_n$  are to a low-dimensional subspace).



$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# SPECTRUM ANALYSIS

Plotting the **spectrum** of the covariance matrix  $\mathbf{X}^T\mathbf{X}$  (its eigenvalues) shows how compressible  $\mathbf{X}$  is using low-rank approximation (i.e., how close  $\vec{x}_1, \dots, \vec{x}_n$  are to a low-dimensional subspace).



- Choose  $k$  to balance accuracy and compression.
- Often at an 'elbow'.

$$\text{tr}(\mathbf{X}^T\mathbf{X}) = \sum \lambda_i(\mathbf{X}^T\mathbf{X})$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .



# SPECTRUM ANALYSIS

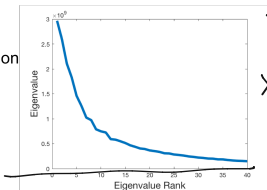
$$\lambda_1 > \lambda_2 \dots > \lambda_n$$

$$\sum_{i=k+1}^n x_i (x_i^T x)$$

784 dimensional vectors



eigendecomposition



$$X^T X v_j = \lambda_j v_j$$

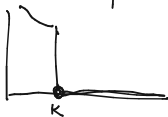
$$X^T X (-v_j) = -\lambda_j v_j = \lambda_j (-v_j)$$

$X^T X$   
non-negative.

**Exercise:** Show that the eigenvalues of  $X^T X$  are always positive.

Hint: Use that  $\lambda_j = \vec{v}_j^T X^T X \vec{v}_j$ .

positive semidefinite matrix



# INTERPRETATION IN TERMS OF CORRELATION

**Recall:** Low-rank approximation is possible when our data features are correlated.

10000\* bathrooms+ 10\* (sq. ft.)  $\approx$  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

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T \mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# INTERPRETATION IN TERMS OF CORRELATION

**Recall:** Low-rank approximation is possible when our data features are correlated.

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

$$X \approx X V_k V_k^T$$

↓  
-k  
[v<sub>1</sub> ... v<sub>k</sub>]  
V<sub>k</sub>

Our compressed dataset is  $C = \underbrace{X V_k}_{n \times k}$  where the columns of  $V_k$  are the top  $k$  eigenvectors of  $X^T X$ .

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $X \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $X^T X$ ,  $V_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# INTERPRETATION IN TERMS OF CORRELATION

**Recall:** Low-rank approximation is possible when our data features are correlated.

10000\* bathrooms + 10\* (sq. ft.)  $\approx$  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

Our compressed dataset is  $\mathbf{C} = \mathbf{X}\mathbf{V}_k$  where the columns of  $\mathbf{V}_k$  are the top  $k$  eigenvectors of  $\mathbf{X}^T\mathbf{X}$ .

What is the covariance of  $\mathbf{C}$ ?

$$\mathbf{C}^T\mathbf{C} = \mathbf{V}_k^T \mathbf{X}^T \mathbf{X} \mathbf{V}_k$$
$$\mathbf{V}_k^T \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \mathbf{V}_k$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# INTERPRETATION IN TERMS OF CORRELATION

**Recall:** Low-rank approximation is possible when our data features are correlated.

10000\* bathrooms+ 10\* (sq. ft.)  $\approx$  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

Our compressed dataset is  $\mathbf{C} = \mathbf{X}\mathbf{V}_k$  where the columns of  $\mathbf{V}_k$  are the top  $k$  eigenvectors of  $\mathbf{X}^T\mathbf{X}$ .

What is the covariance of  $\mathbf{C}$ ?  $\mathbf{C}^T\mathbf{C} = \mathbf{V}_k^T\mathbf{X}^T\mathbf{X}\mathbf{V}_k$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# INTERPRETATION IN TERMS OF CORRELATION

**Recall:** Low-rank approximation is possible when our data features are correlated.

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

$$V_k^T V = \begin{bmatrix} I & 0 \end{bmatrix}$$

$$V_k^T V \Lambda V^T V_k = C^T C$$

$$\begin{bmatrix} I & 0 \end{bmatrix}^T \begin{bmatrix} \Lambda \\ 0 \end{bmatrix} \begin{bmatrix} I \\ 0 \end{bmatrix} = \begin{bmatrix} \Lambda & \\ & 0 \end{bmatrix}$$

Our compressed dataset is  $C = X V_k$  where the columns of  $V_k$  are the top  $k$  eigenvectors of  $X^T X$ .

What is the covariance of  $C$ ?  $C^T C = V_k^T X^T X V_k = \underbrace{V_k^T V \Lambda V^T V_k}_{\text{circled}}$

$$(V_k^T V)_{ij} = \langle v_i, v_j \rangle$$

$$\begin{bmatrix} v_1^T \\ \vdots \\ v_k^T \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_d \end{bmatrix} = \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 & & \\ & & & 0 & \\ & & & & \ddots & \\ & & & & & 0 \end{bmatrix}$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $X \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $X^T X$ ,  $V_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# INTERPRETATION IN TERMS OF CORRELATION

**Recall:** Low-rank approximation is possible when our data features are correlated.

10000\* bathrooms+ 10\* (sq. ft.)  $\approx$  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

Our compressed dataset is  $\mathbf{C} = \mathbf{X}\mathbf{V}_k$  where the columns of  $\mathbf{V}_k$  are the top  $k$  eigenvectors of  $\mathbf{X}^T\mathbf{X}$ .

What is the covariance of  $\mathbf{C}$ ?  $\mathbf{C}^T\mathbf{C} = \mathbf{V}_k^T\mathbf{X}^T\mathbf{X}\mathbf{V}_k = \mathbf{V}_k^T\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T\mathbf{V}_k = \left[ \mathbf{\Lambda}_k \right]$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

# INTERPRETATION IN TERMS OF CORRELATION

**Recall:** Low-rank approximation is possible when our data features are correlated.

$10000 * \text{bathrooms} + 10 * (\text{sq. ft.}) \approx \text{list price}$

	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

Our compressed dataset is  $\mathbf{C} = \mathbf{X}\mathbf{V}_k$  where the columns of  $\mathbf{V}_k$  are the top  $k$  eigenvectors of  $\mathbf{X}^T\mathbf{X}$ .

**What is the covariance of  $\mathbf{C}$ ?**  $\mathbf{C}^T\mathbf{C} = \mathbf{V}_k^T\mathbf{X}^T\mathbf{X}\mathbf{V}_k = \mathbf{V}_k^T\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T\mathbf{V}_k = \mathbf{\Lambda}_k$

**Covariance becomes diagonal.** I.e., all correlations have been removed. Maximal compression.

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .



What is the runtime to compute an optimal low-rank approximation?

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T \mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

What is the runtime to compute an optimal low-rank approximation?

- Computing the covariance matrix  $\mathbf{X}^T\mathbf{X}$  requires  $O(nd^2)$  time.

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

## What is the runtime to compute an optimal low-rank approximation?

- Computing the covariance matrix  $\mathbf{X}^T\mathbf{X}$  requires  $O(nd^2)$  time.
- Computing its full eigendecomposition to obtain  $\vec{v}_1, \dots, \vec{v}_k$  requires  $O(d^3)$  time (similar to the inverse  $(\mathbf{X}^T\mathbf{X})^{-1}$ ).

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

$$d \ll n$$

What is the runtime to compute an optimal low-rank approximation?

- Computing the covariance matrix  $\mathbf{X}^T\mathbf{X}$  requires  $O(nd^2)$  time.
- Computing its full eigendecomposition to obtain  $\vec{v}_1, \dots, \vec{v}_k$  requires  $O(d^3)$  time (similar to the inverse  $(\mathbf{X}^T\mathbf{X})^{-1}$ ).

Many faster iterative and randomized methods. Runtime is roughly  $\tilde{O}(ndk)$  to output just the top  $k$  eigenvectors  $\vec{v}_1, \dots, \vec{v}_k$ .

- Will see in a few classes (power method, Krylov methods).

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : top eigenvectors of  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$ .

## 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).

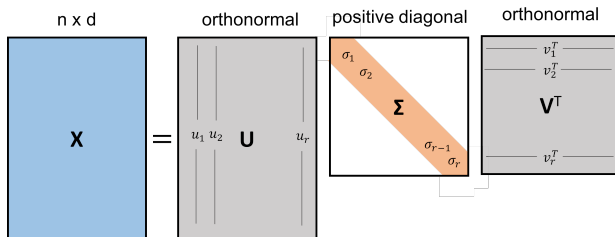
$$\begin{matrix} n & & d \\ \left[ \begin{array}{c} \mathbf{X} \end{array} \right] & = & \begin{matrix} n & & d \\ \left[ \begin{array}{c} \mathbf{U} \end{array} \right] \left[ \begin{array}{c} \sigma_1 \quad \sigma_2 \quad \dots \quad 0 \end{array} \right] \left[ \begin{array}{c} \mathbf{V}^T \end{array} \right] \end{matrix} \\ & & \left[ \begin{array}{c} \mathbf{U} \end{array} \right] \left[ \begin{array}{c} \sigma_1 \quad \sigma_2 \quad \dots \quad 0 \end{array} \right] \left[ \begin{array}{c} \mathbf{V}^T \end{array} \right] \\ & & \mathbf{X} = \sum_{i=1}^r \sigma_i \left[ \begin{array}{c} \mathbf{u}_i \mathbf{v}_i^T \end{array} \right]_n \\ & & \sigma_i = 0 \end{matrix}$$
$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$$

$\mathbf{U} \quad \mathbf{\Sigma} \quad \mathbf{V}$

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

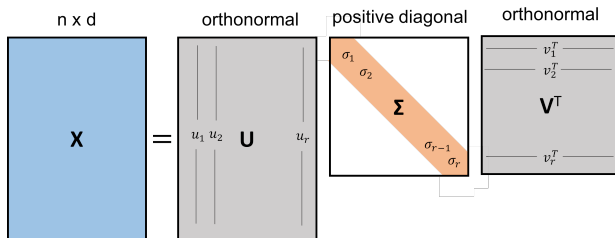
- $\mathbf{U}$  has orthonormal columns  $\vec{u}_1, \dots, \vec{u}_r \in \mathbb{R}^n$  (left singular vectors).
- $\mathbf{V}$  has orthonormal columns  $\vec{v}_1, \dots, \vec{v}_r \in \mathbb{R}^d$  (right singular vectors).
- $\mathbf{\Sigma}$  is diagonal with elements  $\sigma_1 \geq \sigma_2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

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

$\mathbf{I}$

$\mathbf{U}, \mathbf{V}$  orthonormal

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

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

$\downarrow$   
orthonormal      diagonal

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

ortho.      diag.

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : 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 left and right singular vectors are the eigenvectors of the covariance matrix  $\mathbf{X}^T\mathbf{X}$  and the gram matrix  $\mathbf{X}\mathbf{X}^T$  respectively.

*Singular values squared are the eigenvalues.*

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : 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 left and right 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  $\underline{\mathbf{X}\mathbf{V}_k\mathbf{V}_k^T}$  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 left and right 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 left and right 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  $\vec{u}_1, \dots, \vec{u}_k$ ?

Gives exactly the same approximation!

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

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

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

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:

$$X_k = \underline{X V_k V_k^T} = U_k \Sigma_k V_k^T = U_k U_k^T X$$

$$\underline{U_k U_k^T X} = U_k U_k^T U \Sigma V^T = U_k \begin{bmatrix} I & 0 \end{bmatrix} \Sigma V^T$$

$$= U_k \underline{\Sigma_k V_k^T}$$

$$\begin{bmatrix} u_1 \\ \vdots \\ u_k \end{bmatrix} \begin{bmatrix} \sigma_1 & \dots & \sigma_k \\ \vdots & & \vdots \\ v_1 & \dots & v_d \end{bmatrix}$$

# THE SVD AND OPTIMAL LOW-RANK APPROXIMATION

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

$V$  right sing. vect.  
 $U$  left sing. vect.

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:

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

$$X V_k V_k^T = U \Sigma \underbrace{V V_k V_k^T}_{V_k^T V_k = I_k} = U \Sigma \begin{bmatrix} I_k \\ 0 \end{bmatrix} V_k^T = U_k \Sigma_k V_k^T$$

rest of proof  
 ←

$$\begin{bmatrix} V_1^T \\ V_k^T \\ V_{k+1}^T \\ \vdots \\ V_d^T \end{bmatrix} \begin{bmatrix} | \\ | \\ | \\ \vdots \\ | \end{bmatrix} V_1 \dots V_k = \begin{bmatrix} | \\ | \\ | \\ \vdots \\ | \\ 0 \end{bmatrix} \begin{matrix} I_k \\ \vdots \\ 0 \end{matrix}$$

$$\begin{bmatrix} | \\ | \\ | \\ \vdots \\ | \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & \\ & \sigma_2 & & & \\ & & \sigma_k & & \\ & & & \sigma_{k+1} & \\ & & & & \sigma_d \end{bmatrix} \begin{bmatrix} | \\ | \\ | \\ \vdots \\ | \end{bmatrix} \begin{matrix} I_k \\ 0 \end{matrix} = \begin{bmatrix} | \\ | \\ | \\ \vdots \\ | \end{bmatrix} \begin{bmatrix} I_k \\ 0 \end{bmatrix} \begin{bmatrix} | \\ | \\ | \\ \vdots \\ | \end{bmatrix}$$

# THE SVD AND OPTIMAL LOW-RANK APPROXIMATION

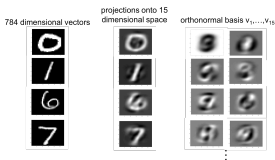
The best low-rank approximation to  $\mathbf{X}$ :

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

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

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

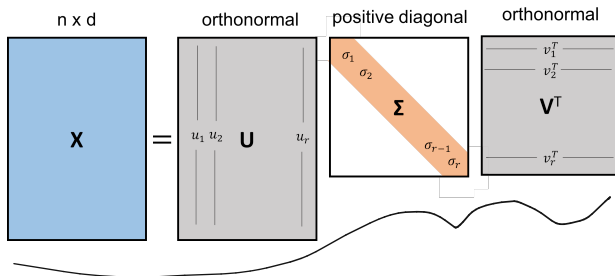
# THE SVD AND OPTIMAL LOW-RANK APPROXIMATION

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:

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

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

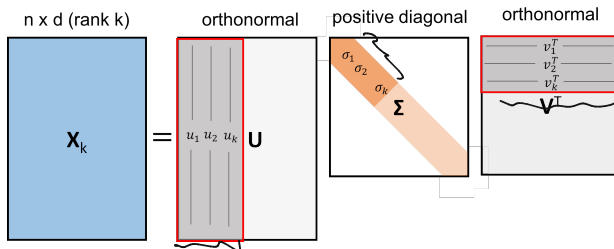


The best low-rank approximation to  $\mathbf{X}$ :

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

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

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



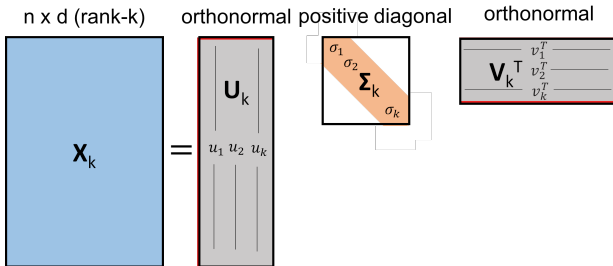
# THE SVD AND OPTIMAL LOW-RANK APPROXIMATION

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:

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

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 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:

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

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



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

**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.

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.

**X**                      Movies

Users

5			1	4				
	3					5		
				4				
	5							5
1			2					

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

## 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.

**Y**                      Movies

	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
Users	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

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

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  $d'$  dimensions. But what about when you want to embed objects other than vectors?



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

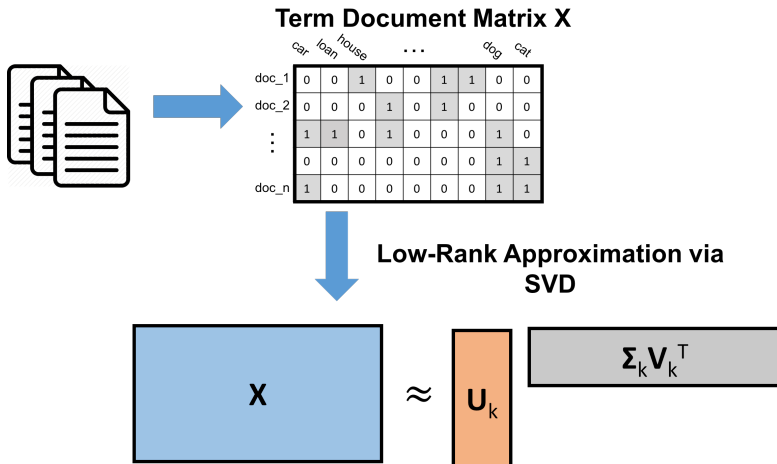
- Documents (for topic-based search and classification)
- Words (to identify synonyms, translations, etc.)
- Nodes in a social network

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

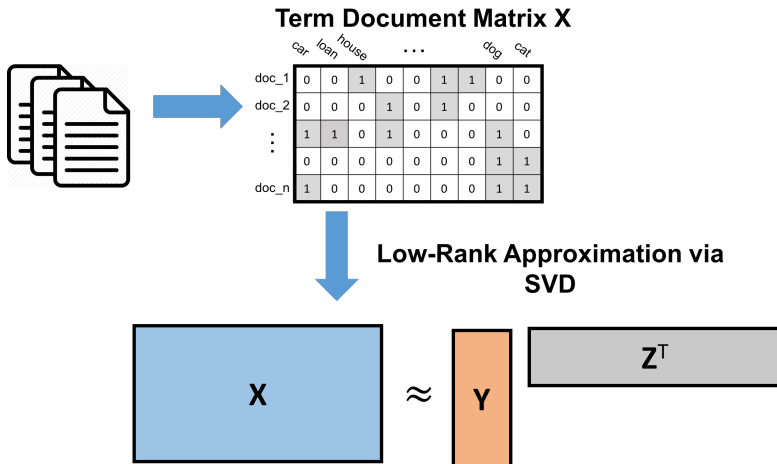
- Documents (for topic-based search and classification)
- Words (to identify synonyms, translations, etc.)
- Nodes in a social network

**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



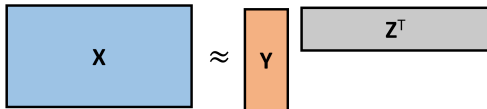
# EXAMPLE: LATENT SEMANTIC ANALYSIS

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



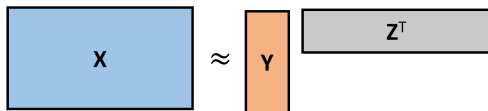
# EXAMPLE: LATENT SEMANTIC ANALYSIS

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
⋮	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



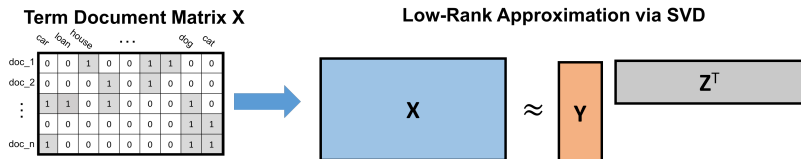
Low-Rank Approximation via SVD



- 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

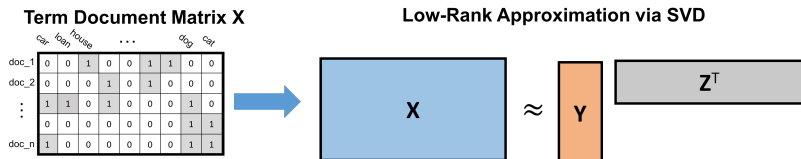


- If the error  $\|X - YZ^T\|_F$  is small, then on average,

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

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

## EXAMPLE: LATENT SEMANTIC ANALYSIS



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

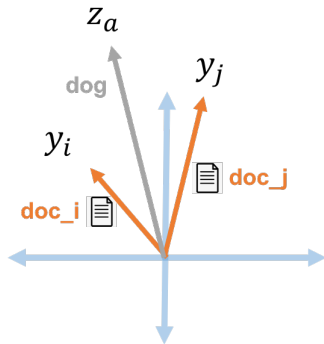
$$X_{i,a} \approx (YZ^T)_{i,a} = \langle \vec{y}_i, \vec{z}_a \rangle.$$

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



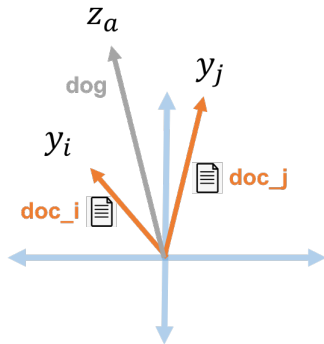
## EXAMPLE: LATENT SEMANTIC ANALYSIS

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



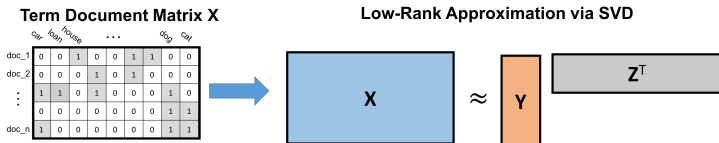
## EXAMPLE: LATENT SEMANTIC ANALYSIS

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



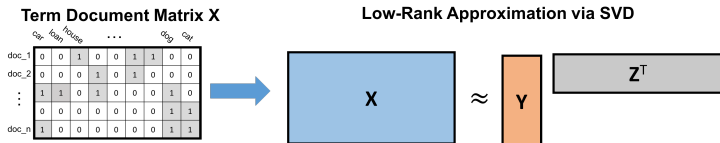
**Another View:** Each column of  $\mathbf{Y}$  represents a 'topic'.  $\vec{y}_i(j)$  indicates how much  $doc_i$  belongs to topic  $j$ .  $\vec{z}_a(j)$  indicates how much  $word_a$  associates with that topic.

## EXAMPLE: LATENT SEMANTIC ANALYSIS



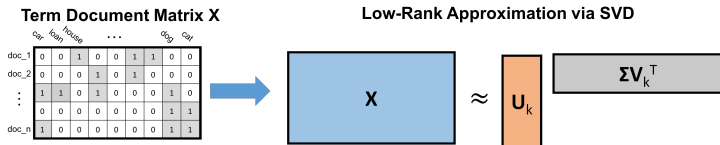
- Just like with documents,  $\vec{z}_a$  and  $\vec{z}_b$  will tend to have high dot product if  $word_i$  and  $word_j$  appear in many of the same documents.

## EXAMPLE: LATENT SEMANTIC ANALYSIS



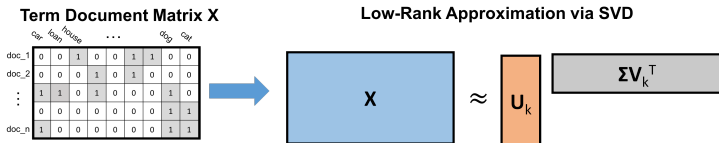
- Just like with documents,  $\vec{z}_a$  and  $\vec{z}_b$  will tend to have high dot product if  $word_i$  and  $word_j$  appear in many of the same documents.
- In an SVD decomposition we set  $Z = \Sigma_k V_K^T$ .

## EXAMPLE: LATENT SEMANTIC ANALYSIS



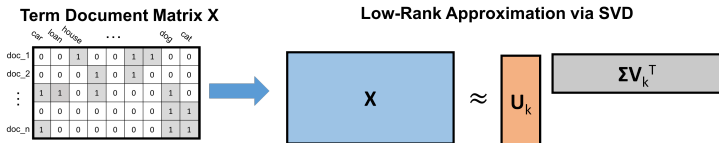
- Just like with documents,  $\vec{z}_a$  and  $\vec{z}_b$  will tend to have high dot product if  $word_i$  and  $word_j$  appear in many of the same documents.
- In an SVD decomposition we set  $Z = \mathbf{\Sigma}_k \mathbf{V}_k^T$ .
- The columns of  $\mathbf{V}_k$  are equivalently: the top  $k$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

## EXAMPLE: LATENT SEMANTIC ANALYSIS



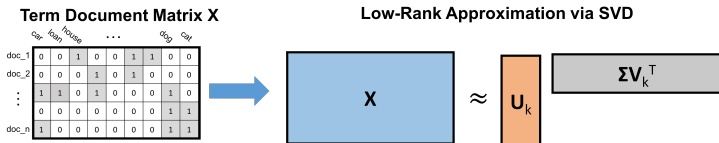
- Just like with documents,  $\vec{z}_a$  and  $\vec{z}_b$  will tend to have high dot product if  $word_i$  and  $word_j$  appear in many of the same documents.
- In an SVD decomposition we set  $Z = \mathbf{\Sigma}_k \mathbf{V}_k^T$ .
- The columns of  $\mathbf{V}_k$  are equivalently: the top  $k$  eigenvectors of  $\mathbf{X}^T \mathbf{X}$ . The eigendecomposition of  $\mathbf{X}^T \mathbf{X}$  is  $\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^T$ .

## EXAMPLE: LATENT SEMANTIC ANALYSIS



- Just like with documents,  $\vec{z}_a$  and  $\vec{z}_b$  will tend to have high dot product if  $word_i$  and  $word_j$  appear in many of the same documents.
- In an SVD decomposition we set  $Z = \sum_k \mathbf{v}_k^T$ .
- 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  $\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^T$ .
- **What is the best rank- $k$  approximation of  $\mathbf{X}^T \mathbf{X}$ ?** I.e.  
$$\arg \min_{\text{rank} = k \text{ B}} \|\mathbf{X}^T \mathbf{X} - \mathbf{B}\|_F$$

## EXAMPLE: LATENT SEMANTIC ANALYSIS



- Just like with documents,  $\vec{z}_a$  and  $\vec{z}_b$  will tend to have high dot product if  $word_i$  and  $word_j$  appear in many of the same documents.
- In an SVD decomposition we set  $Z = \mathbf{\Sigma}_k \mathbf{V}_k^T$ .
- 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  $\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^T$ .
- **What is the best rank- $k$  approximation of  $\mathbf{X}^T \mathbf{X}$ ?** I.e.  
$$\arg \min_{\text{rank} - k \mathbf{B}} \|\mathbf{X}^T \mathbf{X} - \mathbf{B}\|_F$$
- $\mathbf{X}^T \mathbf{X} = \mathbf{V}_k \mathbf{\Sigma}_k^2 \mathbf{V}_k^T = \mathbf{Z} \mathbf{Z}^T$ .



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.

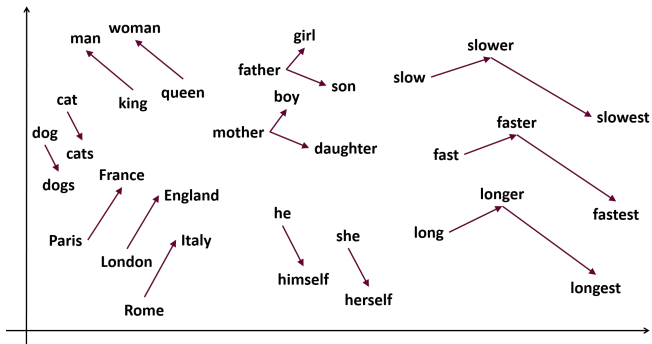
- Embedding is via low-rank approximation of  $\mathbf{X}^T\mathbf{X}$ : where  $(\mathbf{X}^T\mathbf{X})_{a,b}$  is the number of documents that both  $word_a$  and  $word_b$  appear in.
- Think about  $\mathbf{X}^T\mathbf{X}$  as a **similarity matrix** (gram matrix, kernel matrix) with entry  $(a,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

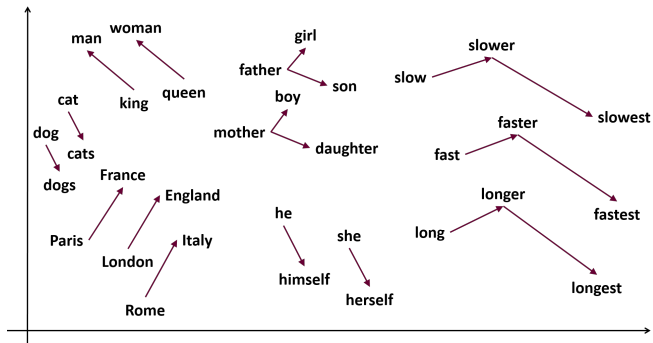
LSA gives a way of embedding words into  $k$ -dimensional space.

- 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.

# SUMMARY

## Summary:

$$\|X - XV_k V_k^T\|_F^2$$
~~$$\|X - XV_k\|$$~~

$$X \text{ - n.d. } \begin{bmatrix} \circ \end{bmatrix} \approx \begin{bmatrix} \circ \end{bmatrix}$$

$n \begin{bmatrix} k \\ \times V \end{bmatrix} \begin{bmatrix} d \\ V^T \end{bmatrix} k$   
 $(n+d)k$   
 $O(n \cdot d)$  time.  
 $Xy$   
 $XV_k V_k^T y = O((n+d)k)$

- Can use the SVD to understand optimal low-rank approximation in terms of the dual row/column projection view:  $XV_k V_k^T = U_k U_k^T X = U_k \Sigma_k V_k^T$ .
- A generalization of eigendecomposition: singular vectors are eigenvectors of  $XX^T$  and  $X^T X$ .
- Applications to low-rank approximation to matrix completion and entity embeddings.

$$\begin{bmatrix} X V \end{bmatrix} \begin{bmatrix} V^T \end{bmatrix} \begin{bmatrix} y \end{bmatrix}$$

$d \times k$

**Next Time:** Low-rank representations of graphs and networks.  $n \times k$   
 Beginning of spectral graph theory.  $(n+d)k$   
 $n \cdot d$