

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2019.

Lecture 13

- Pass/Fail Deadline is 10/29 for undergraduates and 10/31 for graduates. We will have your Problem Set 2 and midterm grades back before then.
- Will release Problem Set 3 next week due ~ 11/11.

- Pass/Fail Deadline is 10/29 for undergraduates and 10/31 for graduates. We will have your Problem Set 2 and midterm grades back before then.
- Will release Problem Set 3 next week due \sim 11/11.
- MAP Feedback:
 - Going to adjust a bit how I take questions in class.
 - Will try to more clearly identify important information (what will appear on exams or problem sets) v.s. motivating examples.
 - Will try to use iPad more to write out proofs in class.

Last Few Classes: Low-Rank Approximation and PCA

Last Few Classes: Low-Rank Approximation and PCA

- Discussed how to compress a dataset that lies close to a k -dimensional subspace.
- Optimal compression by projecting onto the top k eigenvectors of the covariance matrix $\mathbf{X}^T\mathbf{X}$ (PCA).
- Saw how to calculate the error of the approximation – interpret the **spectrum** of $\mathbf{X}^T\mathbf{X}$.

Last Few Classes: Low-Rank Approximation and PCA

- Discussed how to compress a dataset that lies close to a k -dimensional subspace.
- Optimal compression by projecting onto the top k eigenvectors of the covariance matrix $\mathbf{X}^T\mathbf{X}$ (PCA).
- Saw how to calculate the error of the approximation – interpret the **spectrum** of $\mathbf{X}^T\mathbf{X}$.

This Class: Low-rank approximation and connection to singular value decomposition.

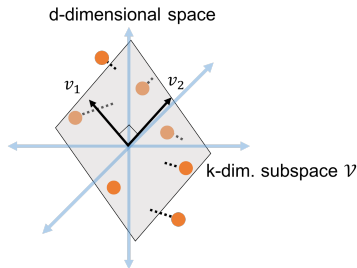
Last Few Classes: Low-Rank Approximation and PCA

- Discussed how to compress a dataset that lies close to a k -dimensional subspace.
- Optimal compression by projecting onto the top k eigenvectors of the covariance matrix $\mathbf{X}^T\mathbf{X}$ (PCA).
- Saw how to calculate the error of the approximation – interpret the **spectrum** of $\mathbf{X}^T\mathbf{X}$.

This Class: Low-rank approximation and connection to singular value decomposition.

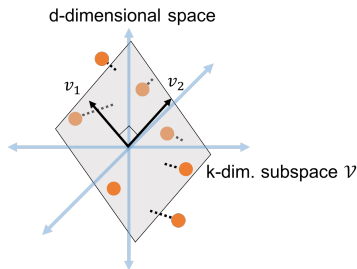
- Show how PCA can be interpreted in terms of the singular value decomposition (SVD) of \mathbf{X} .
- Applications to word embeddings, graph embeddings, document classification, recommendation systems.

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.



$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

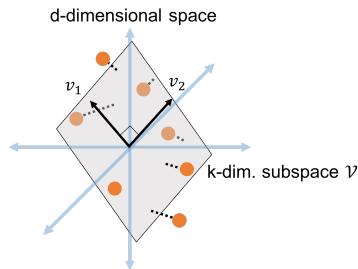
Set Up: Assume that data points $\vec{x}_1, \dots, \vec{x}_n$ lie close to any k -dimensional subspace \mathcal{V} of \mathbb{R}^d . Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the data matrix.



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.

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

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.

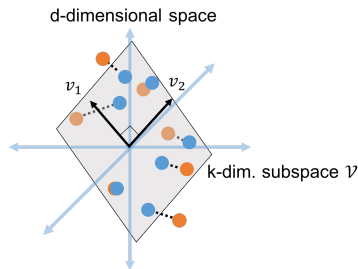


Let $\vec{v}_1, \dots, \vec{v}_k$ be an orthonormal basis for \mathcal{V} and $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- $\mathbf{W}^T \in \mathbb{R}^{d \times d}$ is the **projection matrix** onto \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 $\vec{v}_1, \dots, \vec{v}_k$.

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.

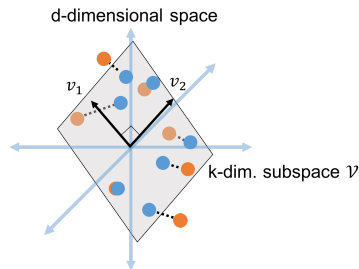


Let $\vec{v}_1, \dots, \vec{v}_k$ be an orthonormal basis for \mathcal{V} and $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- $\mathbf{W}^T \in \mathbb{R}^{d \times d}$ is the **projection matrix** onto \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 $\vec{v}_1, \dots, \vec{v}_k$.

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.



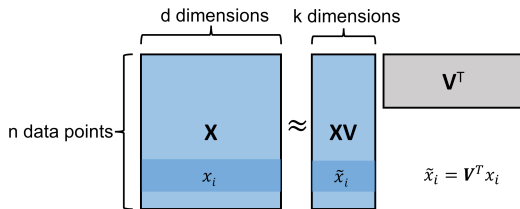
Let $\vec{v}_1, \dots, \vec{v}_k$ be an orthonormal basis for \mathcal{V} and $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

- $\mathbf{W}^T \in \mathbb{R}^{d \times d}$ is the **projection matrix** onto \mathcal{V} .
- $\mathbf{X} \approx \mathbf{X}(\mathbf{W}^T)$. 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 $\vec{v}_1, \dots, \vec{v}_k$.

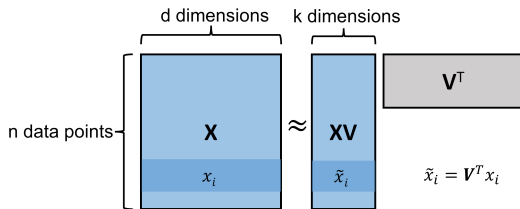
REVIEW OF LAST TIME

Low-Rank Approximation: Approximate $X \approx XV^T$.



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

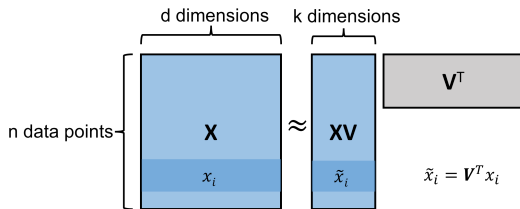
Low-Rank Approximation: Approximate $X \approx XVV^T$.



- XV^T is a **rank- k matrix** – all its rows fall in \mathcal{V} .

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

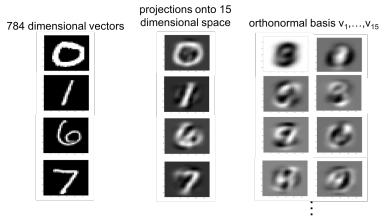
Low-Rank Approximation: Approximate $X \approx XVV^T$.



- XV^T is a **rank- k matrix** – all its rows fall in \mathcal{V} .
- X 's rows are approximately spanned by the columns of V .
- X 's columns are approximately spanned by the columns of XV .

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

DUAL VIEW OF LOW-RANK APPROXIMATION



Row (data point) compression

Column (feature) compression

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

	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

OPTIMAL LOW-RANK APPROXIMATION

Given $\vec{x}_1, \dots, \vec{x}_n$ (the rows of \mathbf{X}) we want to find an orthonormal span $\mathbf{V} \in \mathbb{R}^{d \times k}$ (spanning a k -dimensional subspace \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 $\vec{v}_1, \dots, \vec{v}_k$.

OPTIMAL LOW-RANK APPROXIMATION

Given $\vec{x}_1, \dots, \vec{x}_n$ (the rows of \mathbf{X}) we want to find an orthonormal span $\mathbf{V} \in \mathbb{R}^{d \times k}$ (spanning a k -dimensional subspace \mathcal{V}).

$$\arg \min_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

OPTIMAL LOW-RANK APPROXIMATION

Given $\vec{x}_1, \dots, \vec{x}_n$ (the rows of \mathbf{X}) we want to find an orthonormal span $\mathbf{V} \in \mathbb{R}^{d \times k}$ (spanning a k -dimensional subspace \mathcal{V}).

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

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

OPTIMAL LOW-RANK APPROXIMATION

Given $\vec{x}_1, \dots, \vec{x}_n$ (the rows of \mathbf{X}) we want to find an orthonormal span $\mathbf{V} \in \mathbb{R}^{d \times k}$ (spanning a k -dimensional subspace \mathcal{V}).

$$\arg \min_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2 = \arg \max_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2 = \sum_{i=1}^n \|\mathbf{V}\mathbf{V}^T \vec{x}_i\|_2^2$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

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

SOLUTION VIA EIGENDECOMPOSITION

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

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

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

SOLUTION VIA EIGENDECOMPOSITION

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

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

Surprisingly, can find the columns of \mathbf{V} , $\vec{v}_1, \dots, \vec{v}_k$ greedily.

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

SOLUTION VIA EIGENDECOMPOSITION

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

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

Surprisingly, can find the columns of \mathbf{V} , $\vec{v}_1, \dots, \vec{v}_k$ greedily.

$$\vec{v}_1 = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1} \vec{v}^T \mathbf{X}^T \mathbf{X} \vec{v}.$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

SOLUTION VIA EIGENDECOMPOSITION

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

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

Surprisingly, can find the columns of \mathbf{V} , $\vec{v}_1, \dots, \vec{v}_k$ greedily.

$$\vec{v}_1 = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1} \vec{v}^T \mathbf{X}^T \mathbf{X} \vec{v}.$$

$$\vec{v}_2 = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1, \langle \vec{v}, \vec{v}_1 \rangle = 0} \vec{v}^T \mathbf{X}^T \mathbf{X} \vec{v}.$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

SOLUTION VIA EIGENDECOMPOSITION

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

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

Surprisingly, can find the columns of \mathbf{V} , $\vec{v}_1, \dots, \vec{v}_k$ greedily.

$$\vec{v}_1 = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1} \vec{v}^T \mathbf{X}^T \mathbf{X} \vec{v}.$$

$$\vec{v}_2 = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1, \langle \vec{v}, \vec{v}_1 \rangle = 0} \vec{v}^T \mathbf{X}^T \mathbf{X} \vec{v}.$$

...

$$\vec{v}_k = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1, \langle \vec{v}, \vec{v}_j \rangle = 0 \ \forall j < k} \vec{v}^T \mathbf{X}^T \mathbf{X} \vec{v}.$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

SOLUTION VIA EIGENDECOMPOSITION

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

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

Surprisingly, can find the columns of \mathbf{V} , $\vec{v}_1, \dots, \vec{v}_k$ greedily.

$$\vec{v}_1 = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1} \vec{v}^T \mathbf{X}^T \mathbf{X} \vec{v}.$$

$$\vec{v}_2 = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1, \langle \vec{v}, \vec{v}_1 \rangle = 0} \vec{v}^T \mathbf{X}^T \mathbf{X} \vec{v}.$$

...

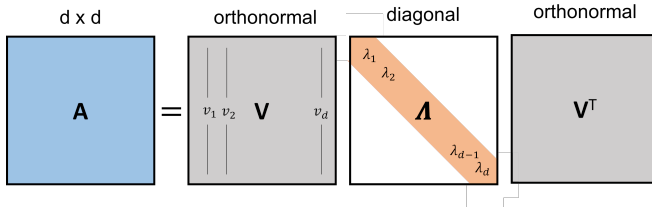
$$\vec{v}_k = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1, \langle \vec{v}, \vec{v}_j \rangle = 0 \ \forall j < k} \vec{v}^T \mathbf{X}^T \mathbf{X} \vec{v}.$$

The top k eigenvectors of $\mathbf{X}^T \mathbf{X}$ by the **Courant-Fischer Principal**.

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

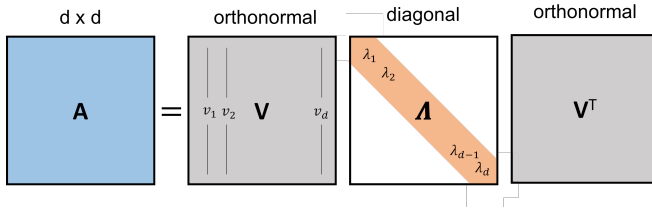
EIGENDECOMPOSITION

Any symmetric matrix \mathbf{A} can be decomposed as $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$, where the columns \mathbf{V} are d orthonormal eigenvectors $\vec{v}_1, \dots, \vec{v}_d$.



EIGENDECOMPOSITION

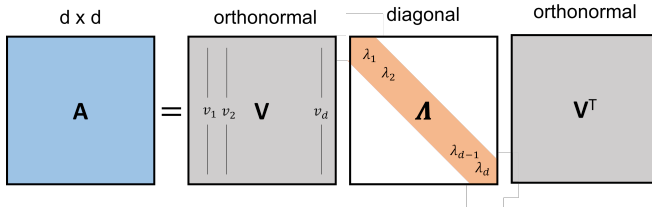
Any symmetric matrix \mathbf{A} can be decomposed as $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$, where the columns \mathbf{V} are d orthonormal eigenvectors $\vec{v}_1, \dots, \vec{v}_d$.



Typically order the eigenvalues in decreasing order: $\lambda_1 \geq \lambda_2 \geq \dots \lambda_d$.

EIGENDECOMPOSITION

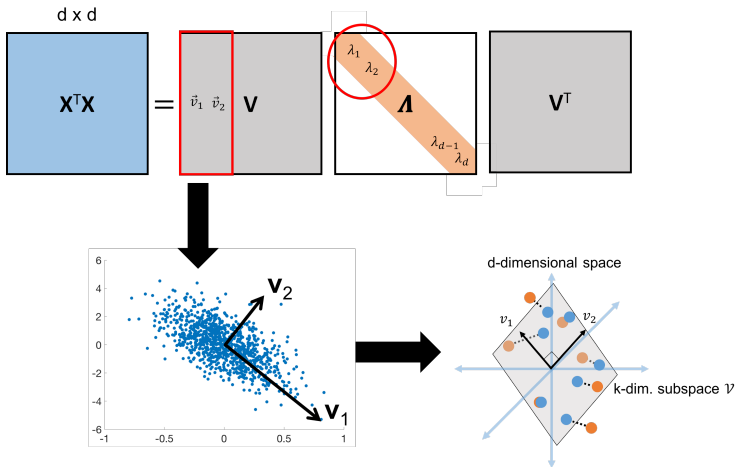
Any symmetric matrix \mathbf{A} can be decomposed as $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$, where the columns \mathbf{V} are d orthonormal eigenvectors $\vec{v}_1, \dots, \vec{v}_d$.



Typically order the eigenvalues in decreasing order: $\lambda_1 \geq \lambda_2 \geq \dots \lambda_d$.

When $\mathbf{A} = \mathbf{X}^T\mathbf{X}$ all eigenvalues are ≥ 0 . Why?

LOW-RANK APPROXIMATION VIA EIGENDECOMPOSITION



$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

Upshot: Letting \mathbf{V}_k have columns $\vec{v}_1, \dots, \vec{v}_k$ corresponding to the top k eigenvectors of the covariance matrix $\mathbf{X}^T\mathbf{X}$, \mathbf{V}_k is the orthogonal basis minimizing

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

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

Upshot: Letting \mathbf{V}_k have columns $\vec{v}_1, \dots, \vec{v}_k$ corresponding to the top k eigenvectors of the covariance matrix $\mathbf{X}^T\mathbf{X}$, \mathbf{V}_k is the orthogonal basis minimizing

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

This is principal component analysis (PCA).

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

Upshot: Letting \mathbf{V}_k have columns $\vec{v}_1, \dots, \vec{v}_k$ corresponding to the top k eigenvectors of the covariance matrix $\mathbf{X}^T\mathbf{X}$, \mathbf{V}_k is the orthogonal basis minimizing

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

This is principal component analysis (PCA).

Last Time: Saw how to determine accuracy by looking at the eigenvalues (the ‘spectrum’) of $\mathbf{X}^T\mathbf{X}$.

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \dots, \vec{v}_k$.

SINGULAR VALUE DECOMPOSITION

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

SINGULAR VALUE DECOMPOSITION

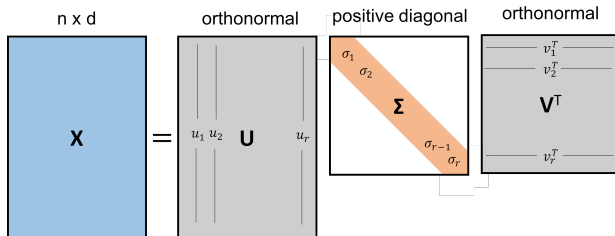
The Singular Value Decomposition (SVD) generalizes the eigendecomposition to asymmetric (even rectangular) matrices. Any matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ with $\text{rank}(\mathbf{X}) = r$ can be written as $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$.

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

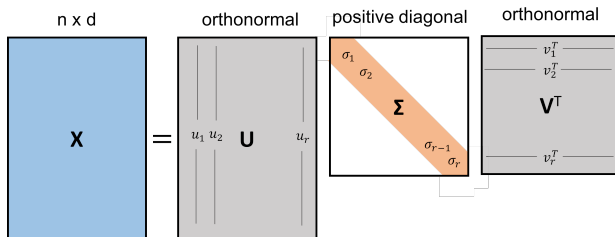
- \mathbf{U} has orthonormal columns $\vec{u}_1, \dots, \vec{u}_r \in \mathbb{R}^n$ (left singular vectors).
- \mathbf{V} has orthonormal columns $\vec{v}_1, \dots, \vec{v}_r \in \mathbb{R}^d$ (right singular vectors).
- $\mathbf{\Sigma}$ is diagonal with elements $\sigma_1 \geq \sigma_2 \geq \dots \geq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

$\mathbf{X} \in \mathbb{R}^{n \times d}$: 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 have that $\mathbf{X}\mathbf{V}_k\mathbf{V}_k^T$ is the best rank- k approximation to \mathbf{X} (given by PCA approximation).

$\mathbf{X} \in \mathbb{R}^{n \times d}$: 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 have that $\mathbf{X}\mathbf{V}_k\mathbf{V}_k^T$ is the best rank- k approximation to \mathbf{X} (given by PCA approximation).

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 have that $\mathbf{X}\mathbf{V}_k\mathbf{V}_k^T$ is the best rank- k approximation to \mathbf{X} (given by PCA approximation).

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!

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

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

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:

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

THE SVD AND OPTIMAL LOW-RANK APPROXIMATION

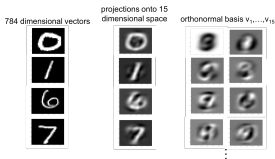
The best low-rank approximation to \mathbf{X} :

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

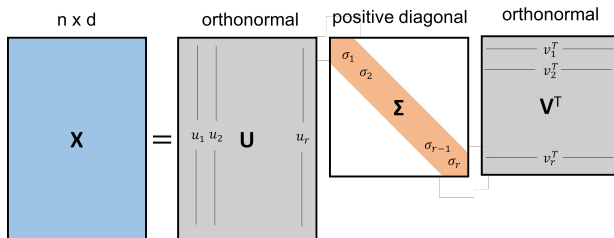
	bedrooms	bathrooms	sq.ft.	floors	list price	sale price
home 1	2	2	1800	2	200,000	195,000
home 2	4	2.5	2700	1	300,000	310,000
.
.
.
home n	5	3.5	3600	3	450,000	450,000

The 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

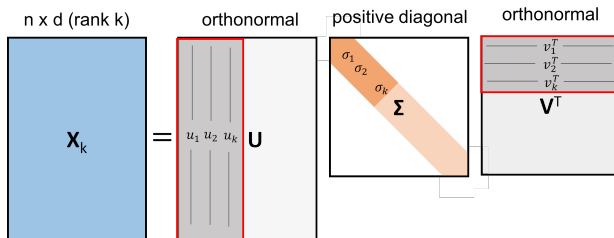


The best low-rank approximation to \mathbf{X} :

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



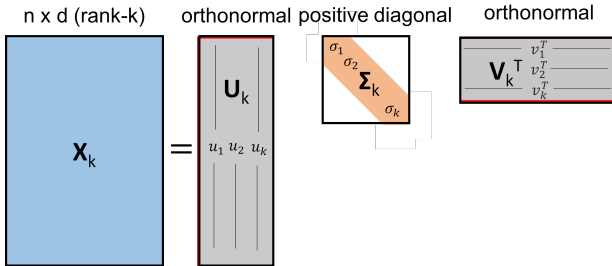
THE SVD AND OPTIMAL LOW-RANK APPROXIMATION

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:

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

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 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} .

SVD is a 'swiss army knife'.

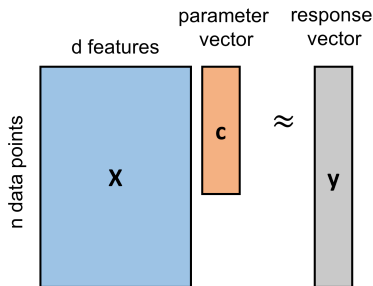
SVD is a 'swiss army knife'.

Classic Linear Regression: Given $\mathbf{X} \in \mathbb{R}^{n \times d}$ where $n > d$ (we have more data points than parameters), and response vector $\vec{y} \in \mathbb{R}^d$, want to find $\vec{c} \in \mathbb{R}^d$ minimizing $\|\mathbf{X}\vec{c} - \vec{y}\|_2$.

THE SVD AND LINEAR REGRESSION

SVD is a 'swiss army knife'.

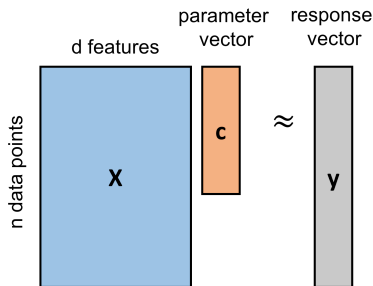
Classic Linear Regression: Given $X \in \mathbb{R}^{n \times d}$ where $n > d$ (we have more data points than parameters), and response vector $\vec{y} \in \mathbb{R}^d$, want to find $\vec{c} \in \mathbb{R}^d$ minimizing $\|X\vec{c} - \vec{y}\|_2$.



THE SVD AND LINEAR REGRESSION

SVD is a 'swiss army knife'.

Classic Linear Regression: Given $X \in \mathbb{R}^{n \times d}$ where $n > d$ (we have more data points than parameters), and response vector $\vec{y} \in \mathbb{R}^d$, want to find $\vec{c} \in \mathbb{R}^d$ minimizing $\|X\vec{c} - \vec{y}\|_2$.



E.g., $c_1 \cdot (\# \text{ baths}) + c_2 \cdot (\text{sq.ft.}) + c_3 \cdot (\# \text{ floors}) + \dots \approx \text{home price}$

Classic Linear Regression: Given $\mathbf{X} \in \mathbb{R}^{n \times d}$ where $n > d$ (we have more data points than parameters) and response vector $\vec{y} \in \mathbb{R}^d$, want to find $\vec{c} \in \mathbb{R}^d$ minimizing $\|\mathbf{X}\vec{c} - \vec{y}\|_2$.

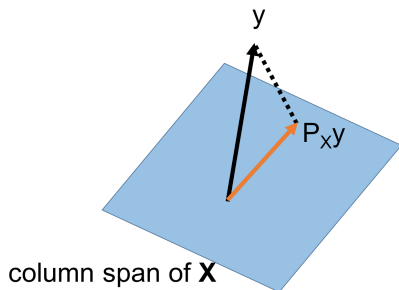
Classic Linear Regression: Given $\mathbf{X} \in \mathbb{R}^{n \times d}$ where $n > d$ (we have more data points than parameters) and response vector $\vec{y} \in \mathbb{R}^d$, want to find $\vec{c} \in \mathbb{R}^d$ minimizing $\|\mathbf{X}\vec{c} - \vec{y}\|_2$.

Optimal solution is to choose \vec{c} so that $\mathbf{X}\vec{c} = \mathbf{P}_X\vec{y}$ – the projection of \vec{y} onto the column span of \mathbf{X} .

THE SVD AND LINEAR REGRESSION

Classic Linear Regression: Given $\mathbf{X} \in \mathbb{R}^{n \times d}$ where $n > d$ (we have more data points than parameters) and response vector $\vec{y} \in \mathbb{R}^d$, want to find $\vec{c} \in \mathbb{R}^d$ minimizing $\|\mathbf{X}\vec{c} - \vec{y}\|_2$.

Optimal solution is to choose \vec{c} so that $\mathbf{X}\vec{c} = \mathbf{P}_X \vec{y}$ – the projection of \vec{y} onto the column span of \mathbf{X} .



Classic Linear Regression: Given $\mathbf{X} \in \mathbb{R}^{n \times d}$ where $n > d$ (we have more data points than parameters) and response vector $\vec{y} \in \mathbb{R}^d$, want to find $\vec{c} \in \mathbb{R}^d$ minimizing $\|\mathbf{X}\vec{c} - \vec{y}\|_2$.

Optimal solution is to choose \vec{c} so that $\mathbf{X}\vec{c} = \mathbf{P}_X\vec{y}$ – the projection of \vec{y} onto the column span of \mathbf{X} .

Writing the SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ we have:

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

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

X Movies

Users

5		1	4				
	3				5		
			4				
	5						5
1		2					

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

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

Classical approach is to convert each item into a high-dimensional feature vector and then apply low-rank approximation

EXAMPLE: LATENT SEMANTIC ANALYSIS

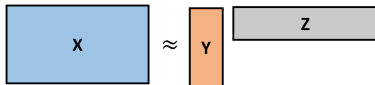
Corpus of Documents



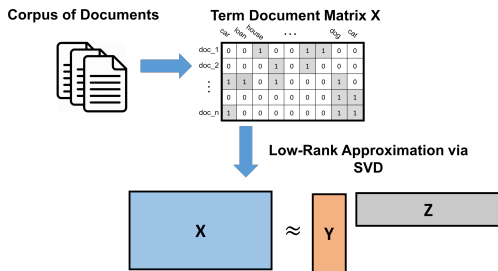
Term Document Matrix X

	cat	fish	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

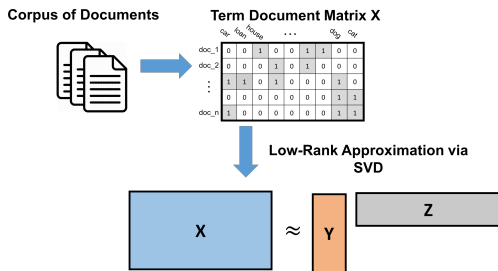


EXAMPLE: LATENT SEMANTIC ANALYSIS

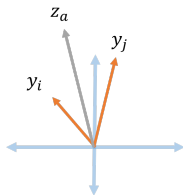


- $\langle \vec{y}_i, \vec{z}_a \rangle \approx 1$ when doc_i contains $word_a$.
- If doc_i and doc_j both contain $word_a$, $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle = 1$.

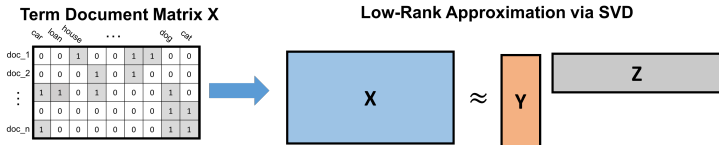
EXAMPLE: LATENT SEMANTIC ANALYSIS



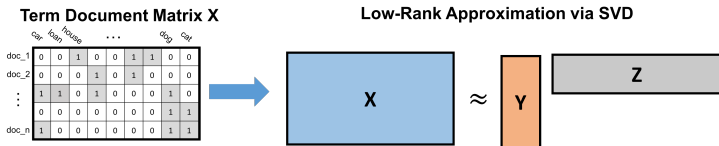
- $\langle \vec{y}_i, \vec{z}_a \rangle \approx 1$ when doc_i contains $word_a$.
- If doc_i and doc_j both contain $word_a$, $\langle \vec{y}_i, \vec{z}_a \rangle \approx \langle \vec{y}_j, \vec{z}_a \rangle = 1$.



EXAMPLE: LATENT SEMANTIC ANALYSIS

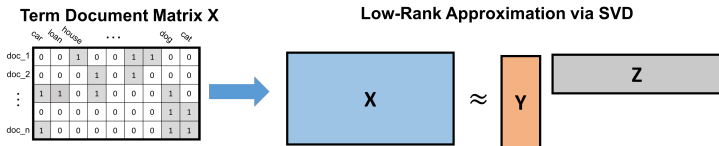


EXAMPLE: LATENT SEMANTIC ANALYSIS



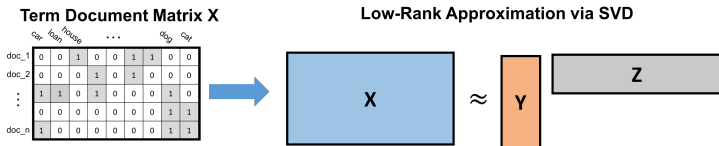
- The columns $\vec{z}_1, \vec{z}_2, \dots$ give representations of words, with \vec{z}_i and \vec{z}_j tending to have high dot product if $word_i$ and $word_j$ appear in many of the same documents.
- Z corresponds to the top k right singular vectors: the eigenvectors of XX^T .

EXAMPLE: LATENT SEMANTIC ANALYSIS



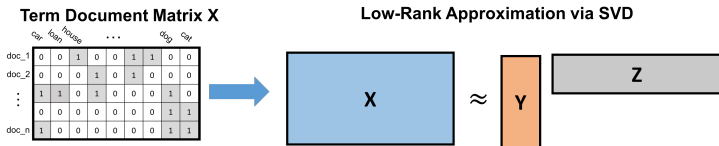
- The columns $\vec{z}_1, \vec{z}_2, \dots$ give representations of words, with \vec{z}_i and \vec{z}_j tending to have high dot product if $word_i$ and $word_j$ appear in many of the same documents.
- Z corresponds to the top k right singular vectors: the eigenvectors of XX^T . *Intuitively, what is XX^T ?*

EXAMPLE: LATENT SEMANTIC ANALYSIS



- The columns $\vec{z}_1, \vec{z}_2, \dots$ give representations of words, with \vec{z}_i and \vec{z}_j tending to have high dot product if $word_i$ and $word_j$ appear in many of the same documents.
- Z corresponds to the top k right singular vectors: the eigenvectors of XX^T . *Intuitively, what is XX^T ?*
- $(XX^T)_{i,j} = \#$ documents that $word_i$ and $word_j$ co-occur in.

EXAMPLE: LATENT SEMANTIC ANALYSIS



- The columns $\vec{z}_1, \vec{z}_2, \dots$ give representations of words, with \vec{z}_i and \vec{z}_j tending to have high dot product if $word_i$ and $word_j$ appear in many of the same documents.
- Z corresponds to the top k right singular vectors: the eigenvectors of XX^T . *Intuitively, what is XX^T ?*
- $(XX^T)_{i,j} = \#$ documents that $word_i$ and $word_j$ co-occur in.
- A document based similarity matrix.

Not obvious how to convert a word into a feature vector that captures the meaning of that word.

- In LSA, feature vector is the set of documents that word appears in.
- SVD of term-document matrix \mathbf{X} corresponds to eigendecomposition of document based similarity matrix \mathbf{XX}^T .

EXAMPLE: WORD EMBEDDING

Not obvious how to convert a word into a feature vector that captures the meaning of that word.

- In LSA, feature vector is the set of documents that word appears in.
- SVD of term-document matrix \mathbf{X} corresponds to eigendecomposition of document based similarity matrix $\mathbf{X}\mathbf{X}^T$.
- Many alternative similarities: how often do $word_i, word_j$ appear in the same sentence, in the same window of w words, in similar positions of documents in different languages, etc.

EXAMPLE: WORD EMBEDDING

Not obvious how to convert a word into a feature vector that captures the meaning of that word.

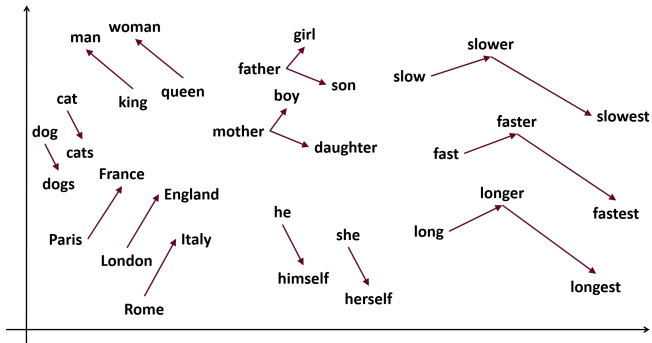
- In LSA, feature vector is the set of documents that word appears in.
- SVD of term-document matrix \mathbf{X} corresponds to eigendecomposition of document based similarity matrix \mathbf{XX}^T .
- Many alternative similarities: how often do $word_i, word_j$ appear in the same sentence, in the same window of w words, in similar positions of documents in different languages, etc.
- Replacing \mathbf{XX}^T with these different metrics (sometimes appropriately transformed) leads to popular word embedding algorithms: word2vec, GloVe, fastText, etc.

EXAMPLE: WORD EMBEDDING

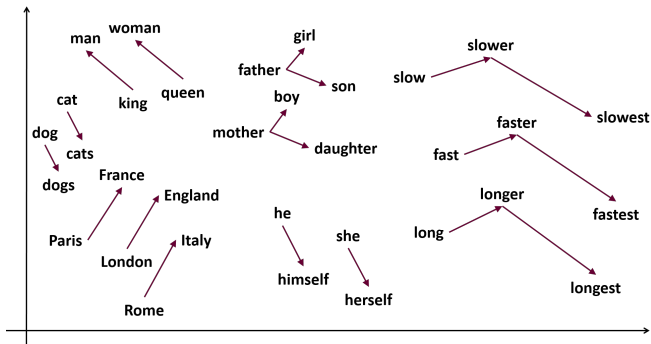
Not obvious how to convert a word into a feature vector that captures the meaning of that word.

- In LSA, feature vector is the set of documents that word appears in.
- SVD of term-document matrix \mathbf{X} corresponds to eigendecomposition of document based similarity matrix \mathbf{XX}^T .
- Many alternative similarities: how often do $word_i, word_j$ appear in the same sentence, in the same window of w words, in similar positions of documents in different languages, etc.
- Replacing \mathbf{XX}^T with these different metrics (sometimes appropriately transformed) leads to popular word embedding algorithms: word2vec, GloVe, fastText, etc.
- Perform low-rank approximation of similarity matrix directly.

EXAMPLE: WORD EMBEDDING



EXAMPLE: WORD EMBEDDING



word2vec was originally described as a neural-network method, but Levy and Goldberg show that it is simply low-rank approximation of a specific similarity matrix. *Neural word embedding as implicit matrix factorization.*

Next Time: Build on the idea of low-rank approximation of similarity matrix low-rank approximation to perform **non-linear** dimensionality reduction for data that is not close to a low-dimensional linear subspace.

Questions?