# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2022.
Lecture 17

- Problem Set 3 is posted. Due Monday 11/14, 11:59pm.
- Quiz this week due Monday at 8pm.

## Summary

### Last Class: Optimal Low-Rank Approximation

- When data lies close to $\mathcal{V}$, the optimal embedding in that space is given by projecting onto that space.

$$\mathbf{X}\mathbf{V}\mathbf{V}^T = \underset{\mathbf{B} \text{ with rows in } \mathcal{V}}{\arg\min} \|\mathbf{X} - \mathbf{B}\|_F^2.$$
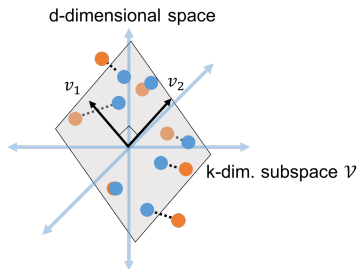
- Optimal $\mathbf{V}$ maximizes $\|\mathbf{X}\mathbf{V}\mathbf{V}^T\|_F$ and can be found greedily. Equivalently by computing the top $k$ eigenvectors of $\mathbf{X}^T\mathbf{X}$.

### This Class:

- How do we assess the error of this optimal $\mathbf{V}$.
- Connection to the singular value decomposition.

# Basic Set Up

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



d-dimensional space

$v_1$    $v_2$

k-dim. subspace $\mathcal{V}$

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

- $\mathsf{V}\mathsf{V}^T \in \mathbb{R}^{d \times d}$ is the projection matrix onto $\mathcal{V}$.

- $\mathsf{X} \approx \mathsf{X}(\mathsf{V}\mathsf{V}^T)$. Gives the closest approximation to $\mathsf{X}$ with rows in $\mathcal{V}$.

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

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

$$\underset{\text{orthonormal } \mathsf{V} \in \mathbb{R}^{d \times k}}{\arg\min} \|\mathsf{X} - \mathsf{X}\mathsf{V}\mathsf{V}^T\|_F^2 = \underset{\text{orthonormal } \mathsf{V} \in \mathbb{R}^{d \times k}}{\arg\max} \|\mathsf{X}\mathsf{V}\|_F^2 = \sum_{j=1}^{k} \|\mathsf{X}\vec{v}_j\|_2^2$$

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

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

- Proof via Courant-Fischer and greedy maximization.

- How accurate is this low-rank approximation? Can understand using eigenvalues of $\mathsf{X}^T\mathsf{X}$.

---

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

# Spectrum Analysis

Let $\vec{v}_1, \ldots, \vec{v}_k$ be the top $k$ eigenvectors of $X^T X$ (the top $k$ principal components). Approximation error is:

$$\|X - XV_k V_k^T\|_F^2 = \|X\|_F^2 \operatorname{tr}(X^T X) - \|XV_k V_k^T\|_F^2 \operatorname{tr}(V_k^T X^T X V_k)$$

$$= \sum_{i=1}^{d} \lambda_i(X^T X) - \sum_{i=1}^{k} \vec{v}_i^T X^T X \vec{v}_i$$

$$= \sum_{i=1}^{d} \lambda_i(X^T X) - \sum_{i=1}^{k} \lambda_i(X^T X) = \sum_{i=k+1}^{d} \lambda_i(X^T X)$$

- **Exercise:** For any matrix $A$, $\|A\|_F^2 = \sum_{i=1}^{d} \|\vec{a}_i\|_2^2 = \operatorname{tr}(A^T A)$ (sum of diagonal entries = sum eigenvalues).

---

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

# Spectrum Analysis

**Claim:** The error in approximating $X$ with the best rank $k$ approximation (projecting onto the top $k$ eigenvectors of $X^TX$ is:

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



d x d

$X^TX$ = $\vec{v}_1\ \vec{v}_2..\vec{v}_k$ $V$ $\Lambda$ $V^T$

$\lambda_1$
$\lambda_2$
$\lambda_k$
$\lambda_{d-1}$
$\lambda_d$

error of optimal low rank approximation

784 dimensional vect

eige

- Choose $k$ to balance accuracy/compression – often at an 'elbow'.

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

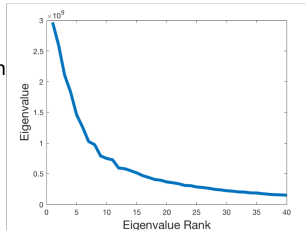Plotting the spectrum of $X^TX$ (its eigenvalues) shows how compressible $X$ is using low-rank approximation (i.e., how close $\vec{x}_1, \ldots, \vec{x}_n$ are to a low-dimensional subspace).

784 dimensional vectors



eigendecomposition

784 dimensional vectors



eigendec

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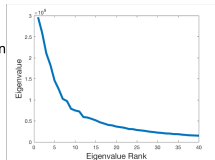
# Spectrum Analysis

784 dimensional vectors



eigendecomposition

Exercises:

1. Show that the eigenvalues of $X^T X$ are always positive. **Hint:** Use that $\lambda_j = \vec{v}_j^T X^T X \vec{v}_j$.

2. Show that for symmetric $A$, the trace is the sum of eigenvalues: $\text{tr}(A) = \sum_{i=1}^{n} \lambda_i(A)$. **Hint:** First prove the cyclic property of trace, that for any $MN$, $\text{tr}(MN) = \text{tr}(NM)$ and then apply this to $A$'s eigendecomposition

## Summary

- Many (most) datasets can be approximated via projection onto a low-dimensional subspace.

- Find this subspace via a maximization problem:

$$\max_{\text{orthonormal } V} \|XV\|_F^2.$$

- Greedy solution via eigendecomposition of $X^TX$.

- Columns of $V$ are the top eigenvectors of $X^TX$.

- Error of best low-rank approximation (compressibility of data) is determined by the tail of $X^TX$'s eigenvalue spectrum.

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

| | bedrooms | bathrooms | sq.ft. | floors | list price | sale price |
|---|---|---|---|---|---|---|
| | | 10000* bathrooms+ 10* (sq. ft.) ≈ list 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 $\mathsf{C} = \mathsf{X}\mathsf{V}_k$ where the columns of $\mathsf{V}_k$ are the top $k$ eigenvectors of $\mathsf{X}^T\mathsf{X}$.

Observe that $\mathsf{C}^T\mathsf{C} = \mathbf{\Lambda}_k$

$C^T C$ **is diagonal.** I.e., all columns are orthogonal to each other, and correlations have been removed. Maximal compression.

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

## Algorithmic Considerations

Runtime to compute an optimal low-rank approximation:

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

Many faster iterative and randomized methods. Runtime is roughly $\tilde{O}(ndk)$ to output just to top $k$ eigenvectors $\vec{v}_1, \ldots, \vec{v}_k$.
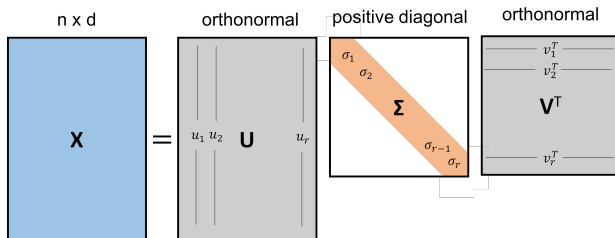
- Will see in a few classes (power method, Krylov methods).
- One of the most intensively studied problems in numerical computation.

---

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

## Singular Value Decomposition

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

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



13

## Connection of the SVD to Eigendecomposition

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

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

Similarly: $XX^T = U\Sigma V^T V\Sigma U^T = U\Sigma^2 U^T$.

The left and right singular vectors are the eigenvectors of the covariance matrix $X^T X$ and the gram matrix $XX^T$ respectively.

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

What about $U_k U_k^T X$ where $U_k \in \mathbb{R}^{n \times k}$ has columns equal to $\vec{u}_1, \ldots, \vec{u}_k$?
Gives exactly the same approximation!

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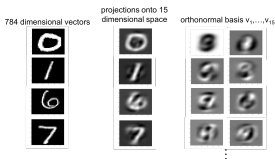
The best low-rank approximation to $\mathbf{X}$:

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

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

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

**Row (data point) compression**          **Column (feature) compression**

n x d          orthonormal   positive diagonal   orthonormal   n x d (rank k)   orthonorma

## The SVD and Optimal Low-Rank Approximation

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:

$$X_k = XV_kV_k^T = U_kU_k^TX = U_k\boldsymbol{\Sigma}_kV_k^T$$

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

## The SVD and Optimal Low-Rank Approximation

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:

$$X_k = XV_kV_k^T = U_kU_k^TX = U_k\Sigma_kV_k^T$$

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