## COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2021.
Lecture 16

$$A/A^- \qquad 85\% \qquad\qquad C - 50 - 5\%$$

$$B \qquad 65 - 75\%$$

- Problem Set 3 is posted. Due Monday 11/8, 11:59pm.
- I strongly encourage you to work together on the problems, rather than split them up.
- Midterms can be collected after class today. Solutions were posted in Moodle. The class average was a 34/40.
- Quiz this week due Monday at 8pm.

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.

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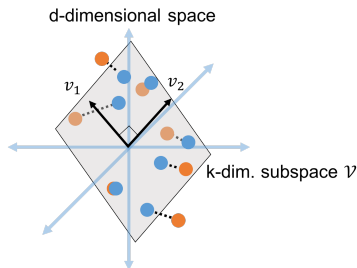
· Optimal $\mathbf{V}$ maximizes $\|\mathbf{X}\mathbf{V}\mathbf{V}^T\|_F$ and can be found greedily. Equivilantly by computing the top $k$ eigenvectors of $\mathbf{X}^T\mathbf{X}$.

### 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.
  Equivilantly 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.

2

**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 $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the data matrix.



Let $\vec{v}_1, \ldots, \vec{v}_k$ 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 $\mathcal{V}$.

---

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

3

$V$ minimizing $\|X - XVV^T\|_F^2$ is given by:

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

---

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

$V$ minimizing $\|X - XVV^T\|_F^2$ is given by:

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

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

$$\underline{V_k} = \underset{\text{orthonormal } V \in \mathbb{R}^{d \times k}}{\arg\max} \|XV\|_F^2$$

---

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

$V$ minimizing $\|X - XVV^T\|_F^2$ is given by:

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

**Solution via eigendecomposition:** Letting $V_k$ have columns $\vec{v}_1, \ldots, \vec{v}_k$ corresponding to the top $k$ eigenvectors of $X^T X$,
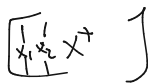
$$V_k = \underset{\text{orthonormal } V \in \mathbb{R}^{d \times k}}{\arg\max} \|XV\|_F^2$$

· Proof via Courant-Fischer and greedy maximization.

---

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

V minimizing $\|X - XVV^T\|_F^2$ is given by:

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

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

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

- Proof via Courant-Fischer and greedy maximization.
- How accurate is this low-rank approximation? Can understand using eigenvalues of $X^T X$.

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

4

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

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

$$V_k = \begin{bmatrix} & & & & \\ v_1 & v_2 & \cdots & v_k \\ & & & & \end{bmatrix}$$

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

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:

$$\|X - XV_k V_k^T\|_F^2 = \underbrace{\|X\|_F^2} - \underbrace{\|XV_k V_k^T\|_F^2}$$

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

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:

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

$$\|XV_k V_k^T\|_F^2 = \|XV_k\|_F^2$$

$$n \begin{bmatrix} XV_k V_k^T \\ \subseteq x_i^T V_k V_k^T \end{bmatrix}$$

$$\|x_i^T V_k V_k^T\|_2^2 = \|x_i^T V_k\|_2^2$$

$$\overset{\frown}{A}$$

$$V_k^T V_k = I$$

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

5

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:

$$\|X - XV_k V_k^T\|_F^2 = \underline{\|X\|_F^2} - \underline{\|XV_k\|_F^2}$$

$$(A^TA)_{ii} = \|a_i\|_2^2$$

$n \times d$

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

$$= \text{tr}(AA^T)$$

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

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:

$$\|X\|_F^2 - \|XVk\|_F^2$$

$$\|X - XV_kV_k^T\|_F^2 = \text{tr}(X^TX) - \text{tr}(V_k^TX^TXV_k)$$

- **Exercise:** For any matrix $A$, $\|A\|_F^2 = \sum_{i=1}^{d} \|\vec{a}_i\|_2^2 = \text{tr}(A^TA)$ (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 $X^TX$, $V_k \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \ldots, \vec{v}_k$.

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:

$$\|X - XV_k V_k^T\|_F^2 = \text{tr}(X^T X) - \text{tr}(V_k^T X^T X V_k)$$
$$= \sum_{i=1}^{d} \lambda_i(X^T X) - \sum_{i=1}^{k} \vec{v}_i^T X^T X \vec{v}_i$$

*Handwritten annotations:*

$$k \left[ \begin{array}{c} v_1^T \\ v_k \\ v_k^T \end{array} \right] \prod \left[ X^T X \right] \prod \left[ \begin{array}{c} v_k \\ v_1 .. v_k \end{array} \right]$$

$1 \times d \quad d \times d \quad d \times 1$

$$V_i^T (X^T X) V_i = V_i^T (\lambda_i \cdot V_i)$$
$$= \lambda_i \cdot v_i^T v_i = \lambda_i$$

- **Exercise:** For any matrix $A$, $\|A\|_F^2 = \sum_{i=1}^{d} \|\vec{a}_i\|_2^2 = \text{tr}(A^T A)$ (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

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:

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

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

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

(annotations above equation: $\|X\|_F^2$, $\|XV_k\|_F^2$)

- **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 $\vec{v}_1, \ldots, \vec{v}_k$.

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:

$$\|X - XV_k V_k^T\|_F^2 = \text{tr}(X^T X) - \text{tr}(V_k^T X^T X V_k)$$
$$= \sum_{i=1}^{d} \lambda_i(X^T X) - \sum_{i=1}^{k} \vec{v}_i^T X^T X \vec{v}_i$$
$$= \sum_{i=1}^{d} \lambda_i(X^T X) - \sum_{i=1}^{k} \lambda_i(X^T X) = \sum_{i=k+1}^{d} \lambda_i(X^T X)$$

$$\lambda_1(X^T X) \geq \lambda_2(X^T X) \geq \ldots \geq \lambda_d(X^T X)$$

- **Exercise:** For any matrix $A$, $\|A\|_F^2 = \sum_{i=1}^{d} \|\vec{a}_i\|_2^2 = \text{tr}(A^T A)$ (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 $\vec{v}_1, \ldots, \vec{v}_k$.

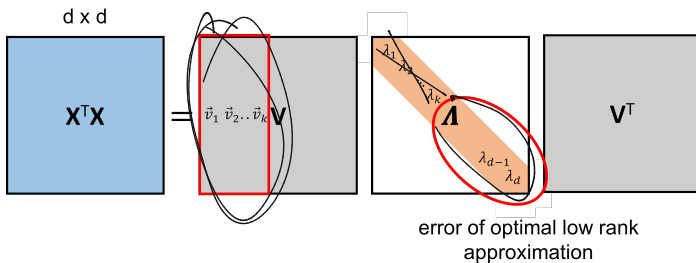Claim: The error in approximating X with the best rank $k$ approximation (projecting onto the top $k$ eigenvectors of $X^T X$ is:

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

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

**Claim:** The error in approximating $\mathsf{X}$ with the best rank $k$ approximation (projecting onto the top $k$ eigenvectors of $\mathsf{X}^T\mathsf{X}$ is:

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



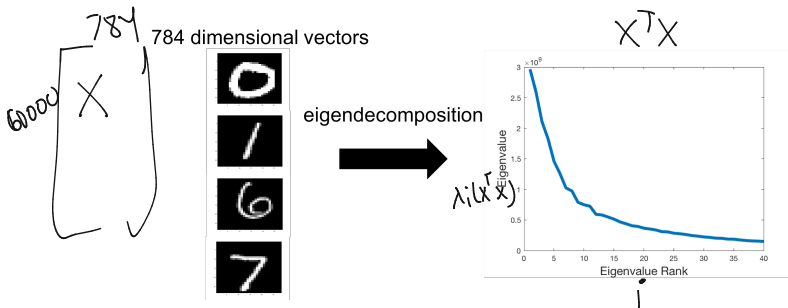error of optimal low rank approximation

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

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

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



784 dimensional vectors

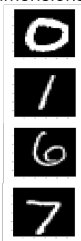eigendecomposition

$X^T X$

$\lambda_i(X^T X)$

Eigenvalue Rank

---

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

6

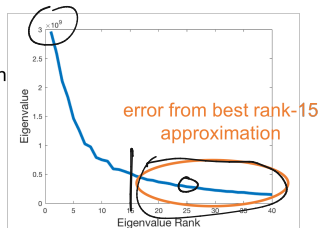**Claim:** The error in approximating X with the best rank $k$ approximation (projecting onto the top $k$ eigenvectors of $X^TX$ is:

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

784 dimensional vectors
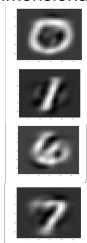


eigendecomposition

error from best rank-15 approximation

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

6

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

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

784 dimensional vectors



eigendecomposition

error from best rank-15 approximation

$\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$: 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 $\vec{v}_1, \ldots, \vec{v}_k$.
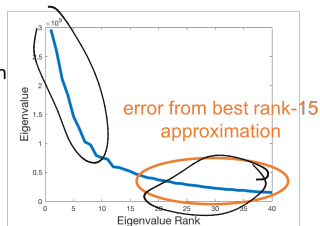
6

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

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

784 dimensional vectors



eigendecomposition

error from best rank-15 approximation

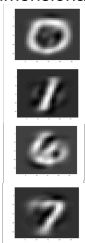· Choose $k$ to balance accuracy/compression – often at an 'elbow'.

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

6

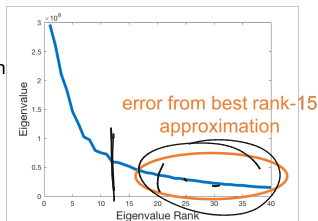Plotting the spectrum of $X^T X$ (its eigenvalues) shows how compressible $X$ is using low-rank approximation (i.e., how close $\vec{x}_1, \ldots, \vec{x}_n$ are to a low-dimensional subspace).

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

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

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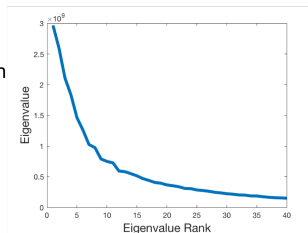
Plotting the spectrum of $X^TX$ (its eigenvalues) shows how compressible $X$ is using low-rank approximation (i.e., how close $\vec{x}_1, \ldots, \vec{x}_n$ are to a low-dimensional subspace).

$$\underbrace{\|X - XW^T\|}_{\text{large}}$$

$$\|X - XVV^T\|_F$$

$$\sum \|x_i - VV^Tx_i\|_2$$

$$\frac{\log n}{\varepsilon^2}$$

784 dimensional vectors

eigendecomposition



$X^TX$

---

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

7

Plotting the spectrum of $X^T X$ (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).
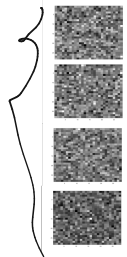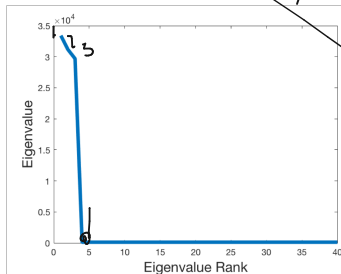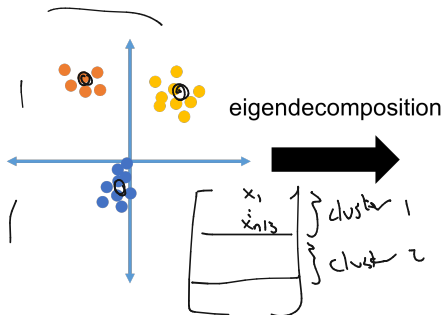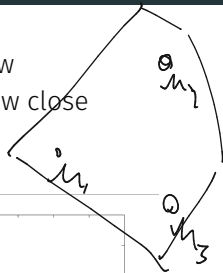


eigendecomposition

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

784 dimensional vectors



eigendecomposition

### Exercises:

*positive semidefinite*

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

784 dimensional vectors



eigendecomposition

### Exercises:

1. Show that the eigenvalues of $X^TX$ 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.

- 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 $X^T X$.
- Columns of $V$ are the top eigenvectors of $X^T X$.
- Error of best low-rank approximation (compressibility of data) is determined by the tail of $X^T X$'s eigenvalue spectrum.

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



| | 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 |

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

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

**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 $\mathsf{C} = \mathsf{X}\mathsf{V}_k$ where the columns of $\mathsf{V}_k$ are the top $k$ eigenvectors of $\mathsf{X}^T\mathsf{X}$.

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

10

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

10000* bedrooms+ 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 |

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

Observe that $C^TC = V_k^T X^T X V_k = V_k^T V \Lambda V^T V_k$

$$X^TX = V \Lambda V^T$$

$$V_k^T V \Lambda V^T V_k = \begin{bmatrix} I : 0 \end{bmatrix} \begin{bmatrix} \Lambda \end{bmatrix} \begin{bmatrix} I \\ 0 \end{bmatrix} = \begin{bmatrix} \lambda_k \end{bmatrix}_k^k$$

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

10

**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 $C = XV_k$ where the columns of $V_k$ are the top $k$ eigenvectors of $X^T X$.

Observe that $C^T C = \Lambda_k$

$C^T C$ 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, $X \in \mathbb{R}^{n \times d}$: 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 $\vec{v}_1, \ldots, \vec{v}_k$.

10

Runtime to compute an optimal low-rank approximation:

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

$$\sqrt{\boxed{\phantom{X}X^T\phantom{X}}}\,\boxed{\phantom{X}}\binom{X}{n} = \boxed{X^TX}\,d$$

Runtime to compute an optimal low-rank approximation:

- Computing $X^TX$ requires $O(nd^2)$ time.

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

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 $(X^T X)^{-1}$).

---

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

Runtime to compute an optimal low-rank approximation:

- Computing $X^TX$ 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 $(X^TX)^{-1}$).

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 $X^TX$, $V_k \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \ldots, \vec{v}_k$.