

# COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

---

Cameron Musco

University of Massachusetts Amherst. Fall 2021.

Lecture 15

## Logistics:

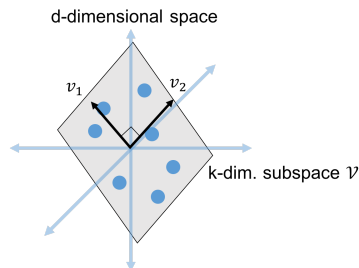
- We have almost finished grading the midterm. Will return grades tomorrow evening and tests in class on Thursday.

## Last Class:

- No-distortion embeddings for data lying in a  $k$ -dimensional subspace via an orthonormal basis  $\mathbf{V} \in \mathbb{R}^{d \times k}$  for that subspace.
- Using that  $\mathbf{V}^T \mathbf{V}$  is an identity matrix and  $\mathbf{V} \mathbf{V}^T$  is a projection matrix to argue this, and understand low-rank matrix approximation.
- ‘Dual view’ of low-rank approximation: data points that can be reconstructed from a few basis vectors vs. linearly dependent features.

## LAST CLASS: EMBEDDING WITH ASSUMPTIONS

**Set Up:** Assume that data points  $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$  lie in some  $k$ -dimensional subspace  $\mathcal{V}$  of  $\mathbb{R}^d$ .



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

$$\|\mathbf{V}^T \vec{x}_i - \mathbf{V}^T \vec{x}_j\|_2^2 = \|\vec{x}_i - \vec{x}_j\|_2^2.$$

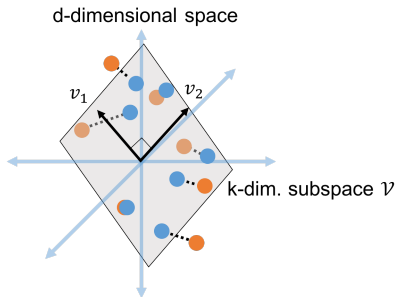
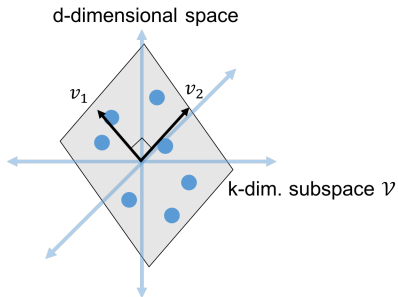
Letting  $\tilde{x}_i = \mathbf{V}^T \vec{x}_i$ , we have a perfect embedding from  $\mathcal{V}$  into  $\mathbb{R}^k$ .

# PROJECTION VIEW

**Claim:** If  $\vec{x}_1, \dots, \vec{x}_n$  lie in a  $k$ -dimensional subspace  $\mathcal{V}$  with orthonormal basis  $\mathbf{V} \in \mathbb{R}^{d \times k}$ , the data matrix can be written as

$$\mathbf{X} = \mathbf{XV}\mathbf{V}^T \text{ (Implies } \text{rank}(\mathbf{X}) \leq k \text{)}$$

- $\mathbf{V}\mathbf{V}^T$  is a **projection matrix**, which projects the rows of  $\mathbf{X}$  (the data points  $\vec{x}_1, \dots, \vec{x}_n$ ) onto the subspace  $\mathcal{V}$ .



**Quick Exercise 1:** Show that  $\mathbf{V}\mathbf{V}^T$  is idempotent. I.e.,  $(\mathbf{V}\mathbf{V}^T)(\mathbf{V}\mathbf{V}^T)\vec{y} = (\mathbf{V}\mathbf{V}^T)\vec{y}$  for any  $\vec{y} \in \mathbb{R}^d$ .

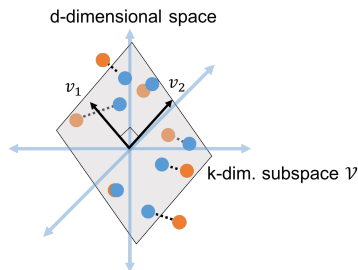
**Quick Exercise 2:** Show that  $\mathbf{V}\mathbf{V}^T(\mathbf{I} - \mathbf{V}\mathbf{V}^T) = \mathbf{0}$  ( the projection is orthogonal to its complement).

**Pythagorean Theorem:** For any orthonormal  $\mathbf{V} \in \mathbb{R}^{d \times k}$  and any  $\vec{y} \in \mathbb{R}^d$ ,

$$\|\vec{y}\|_2^2 = \|(\mathbf{V}\mathbf{V}^T)\vec{y}\|_2^2 + \|\vec{y} - (\mathbf{V}\mathbf{V}^T)\vec{y}\|_2^2.$$

## EMBEDDING WITH ASSUMPTIONS

**Main Focus of Today:** Assume that data points  $\vec{x}_1, \dots, \vec{x}_n$  lie close to any  $k$ -dimensional subspace  $\mathcal{V}$  of  $\mathbb{R}^d$ .



Letting  $\vec{v}_1, \dots, \vec{v}_k$  be an orthonormal basis for  $\mathcal{V}$  and  $\mathbf{V} \in \mathbb{R}^{d \times k}$  be the matrix with these vectors as its columns,  $\mathbf{V}^T \vec{x}_i \in \mathbb{R}^k$  is still a good embedding for  $x_i \in \mathbb{R}^d$ . The key idea behind low-rank approximation and principal component analysis (PCA).

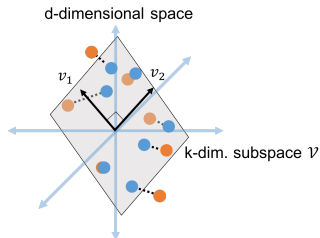
- How do we find  $\mathcal{V}$  and  $\mathbf{V}$ ?
- How good is the embedding?

# BEST FIT SUBSPACE

If  $\vec{x}_1, \dots, \vec{x}_n$  are close to a  $k$ -dimensional subspace  $\mathcal{V}$  with orthonormal basis  $\mathbf{V} \in \mathbb{R}^{d \times k}$ , the data matrix can be approximated as  $\mathbf{XV}^T$ .  $\mathbf{XV}$  gives optimal embedding of  $\mathbf{X}$  in  $\mathcal{V}$ .

How do we find  $\mathcal{V}$  (equivalently  $\mathbf{V}$ )?

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



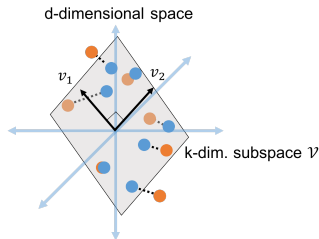
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$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



If  $\vec{x}_1, \dots, \vec{x}_n$  are close to a  $k$ -dimensional subspace  $\mathcal{V}$  with orthonormal basis  $\mathbf{V} \in \mathbb{R}^{d \times k}$ , the data matrix can be approximated as  $\mathbf{XV}\mathbf{V}^T$ .  $\mathbf{XV}$  gives optimal embedding of  $\mathbf{X}$  in  $\mathcal{V}$ .

How do we find  $\mathcal{V}$  (equivalently  $\mathbf{V}$ )?

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



# SOLUTION VIA EIGENDECOMPOSITION

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

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

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

$$\vec{v}_1 = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1} \|\mathbf{X}\vec{v}\|_2^2 \vec{v}^T \mathbf{X}^T \mathbf{X} \vec{v}.$$

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

...

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

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

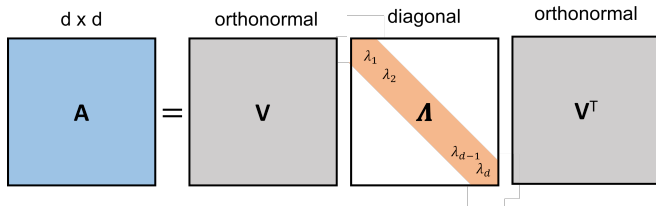
**Eigenvector:**  $\vec{x} \in \mathbb{R}^d$  is an eigenvector of a matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  if  $\mathbf{A}\vec{x} = \lambda\vec{x}$  for some scalar  $\lambda$  (the eigenvalue corresponding to  $\vec{x}$ ).

- That is,  $\mathbf{A}$  just ‘stretches’  $x$ .
- If  $\mathbf{A}$  is **symmetric**, can find  $d$  orthonormal eigenvectors  $\vec{v}_1, \dots, \vec{v}_d$ . Let  $\mathbf{V} \in \mathbb{R}^{d \times d}$  have these vectors as columns.

$$\mathbf{AV} = \begin{bmatrix} | & | & | & | \\ \mathbf{A}\vec{v}_1 & \mathbf{A}\vec{v}_2 & \cdots & \mathbf{A}\vec{v}_d \\ | & | & | & | \end{bmatrix} = \begin{bmatrix} | & | & | & | \\ \lambda_1\vec{v}_1 & \lambda_2\vec{v}_2 & \cdots & \lambda_d\vec{v}_d \\ | & | & | & | \end{bmatrix} = \mathbf{V}\mathbf{\Lambda}$$

Yields eigendecomposition:  $\mathbf{AVV}^T = \mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$ .

# REVIEW OF EIGENVECTORS AND EIGENDECOMPOSITION



Typically order the eigenvectors in decreasing order:

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d.$$

**Courant-Fischer Principal:** For symmetric  $\mathbf{A}$ , the eigenvectors are given via the greedy optimization:

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

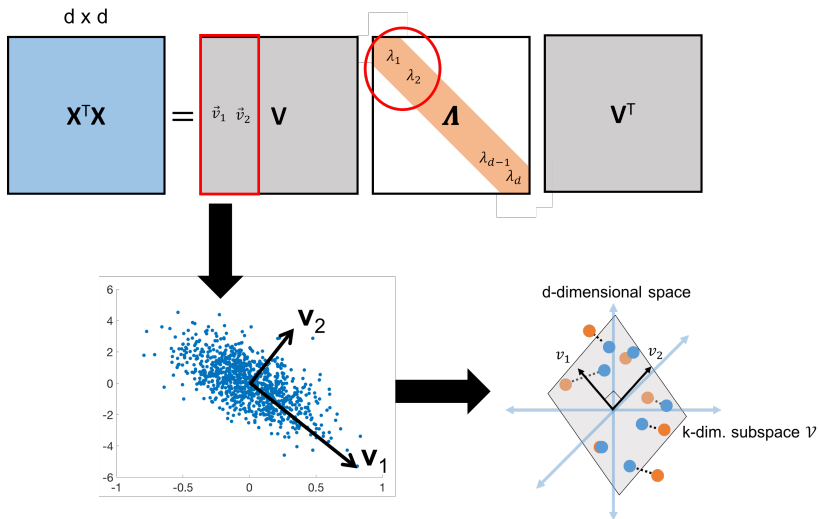
$$\vec{v}_2 = \arg \max_{\vec{v} \text{ with } \|\vec{v}\|_2=1, \langle \vec{v}, \vec{v}_1 \rangle = 0} \vec{v}^T \mathbf{A} \vec{v}.$$

...

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

- $\vec{v}_j^T \mathbf{A} \vec{v}_j = \lambda_j \cdot \vec{v}_j^T \vec{v}_j = \lambda_j$ , the  $j^{\text{th}}$  largest eigenvalue.
- The first  $k$  eigenvectors of  $\mathbf{X}^T \mathbf{X}$  (corresponding to the largest  $k$  eigenvalues) are exactly the directions of greatest variance in  $\mathbf{X}$  that we use for low-rank approximation.

# LOW-RANK APPROXIMATION VIA EIGENDECOMPOSITION



**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).

**How accurate is this low-rank approximation?** Can understand using eigenvalues 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  $\vec{v}_1, \dots, \vec{v}_k$ .