

# COMPSCI 514: Algorithms for Data Science

---

Cameron Musco

University of Massachusetts Amherst. Fall 2022.

Lecture 15

- Midterms should be graded by end of the week.
- Will release grades once we are done and hand them back in class next week.
- Quiz due Monday 8pm as usual this week.

# Summary

## Last Few Classes:

### The Johnson-Lindenstrauss Lemma

- Reduce  $n$  data points in **any dimension  $d$**  to  $O\left(\frac{\log n/\delta}{\epsilon^2}\right)$  dimensions and preserve (with probability  $\geq 1 - \delta$ ) **all pairwise distances** up to  $1 \pm \epsilon$ .
- **Compression is linear** via multiplication with a random, **data oblivious**, matrix (linear compression)

### High-Dimensional Geometry

- Why high-dimensional space is so different than low-dimensional space.
- How the JL Lemma can still work, and why it is optimal.

# Summary

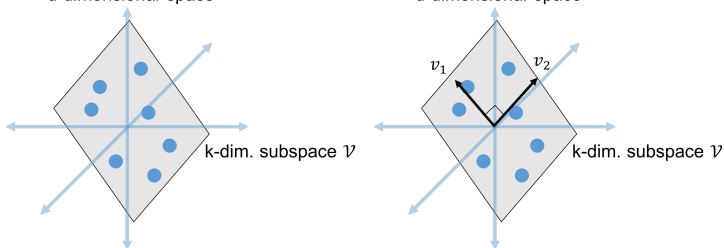
Next Few Classes: Low-rank approximation, the SVD, and principal component analysis (PCA).

- Reduce  $d$ -dimensional data points to a smaller dimension  $m$ .
- Like JL, **compression is linear** – by applying a matrix.
- Chose this matrix carefully, taking into account **structure of the dataset**.
- Can give better compression than random projection (although not directly comparable).

Will be using a fair amount of linear algebra: orthogonal basis, column/row span, eigenvectors, etc.

# Embedding with Assumptions

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



**Claim:** 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. For all  $\vec{x}_i, \vec{x}_j$ :

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

- $\mathbf{V}^T \in \mathbb{R}^{k \times d}$  is a linear embedding of  $\vec{x}_1, \dots, \vec{x}_n$  into  $k$  dimensions with **no distortion**.

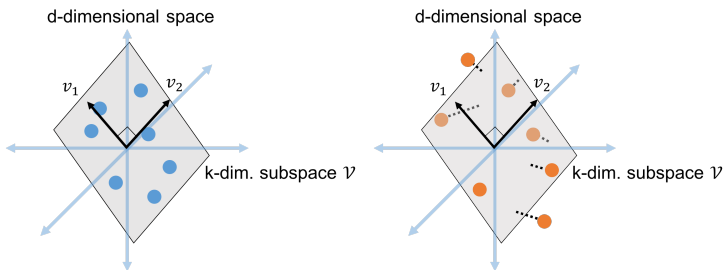
## Dot Product Transformation

**Claim:** 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. For all  $\vec{x}_i, \vec{x}_j \in \mathcal{V}$ :

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

# Embedding with Assumptions

**Main Focus of Upcoming Classes:** 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?

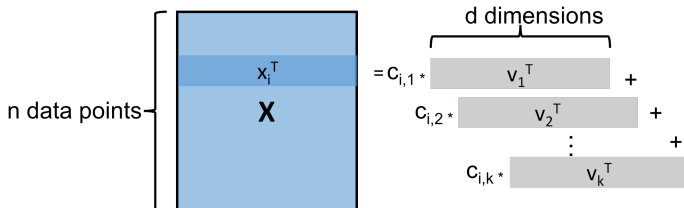
# Low-Rank Factorization

**Claim:**  $\vec{x}_1, \dots, \vec{x}_n$  lie in a  $k$ -dimensional subspace  $\mathcal{V} \Leftrightarrow$  the data matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  has rank  $\leq k$ .

- Letting  $\vec{v}_1, \dots, \vec{v}_k$  be an orthonormal basis for  $\mathcal{V}$ , can write  $\vec{x}_i$  as:

$$\vec{x}_i = \mathbf{V}\vec{c}_i = c_{i,1} \cdot \vec{v}_1 + c_{i,2} \cdot \vec{v}_2 + \dots + c_{i,k} \cdot \vec{v}_k.$$

- So  $\vec{v}_1, \dots, \vec{v}_k$  span the rows of  $\mathbf{X}$  and thus  $\text{rank}(\mathbf{X}) \leq k$ .

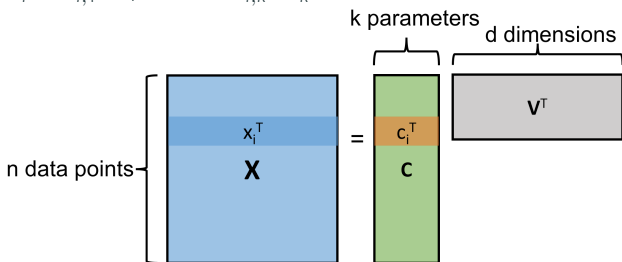


$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$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



**Claim:**  $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$  lie in a  $k$ -dimensional subspace  $\mathcal{V} \Leftrightarrow$  the data matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  has rank  $\leq k$ .

- Every data point  $\vec{x}_i$  (row of  $\mathbf{X}$ ) can be written as  $\vec{x}_i = \mathbf{V}\vec{c}_i = c_{i,1} \cdot \vec{v}_1 + \dots + c_{i,k} \cdot \vec{v}_k$ .

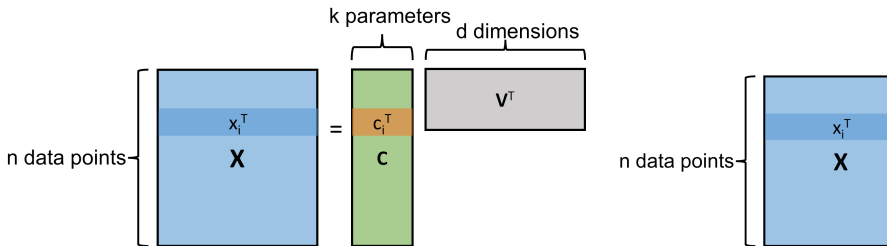


- $\mathbf{X}$  can be represented by  $(n + d) \cdot k$  parameters vs.  $n \cdot d$ .
- The rows of  $\mathbf{X}$  are spanned by  $k$  vectors: the columns of  $\mathbf{V} \implies$  the columns of  $\mathbf{X}$  are spanned by  $k$  vectors: the columns of  $\mathbf{C}$ .

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

# Low-Rank Factorization

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



**Exercise:** What is this coefficient matrix  $\mathbf{C}$ ? **Hint:** Use that  $\mathbf{V}^T\mathbf{V} = \mathbf{I}$ .

$$\bullet \mathbf{X} = \mathbf{C}\mathbf{V}^T \implies \mathbf{X}\mathbf{V} = \mathbf{C}\mathbf{V}^T\mathbf{V} \implies \mathbf{X}\mathbf{V} = \mathbf{C}$$

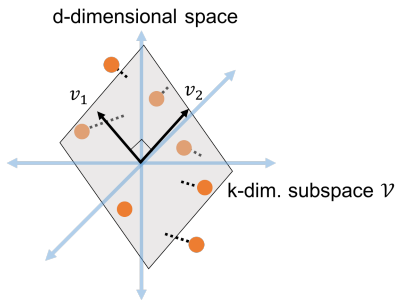
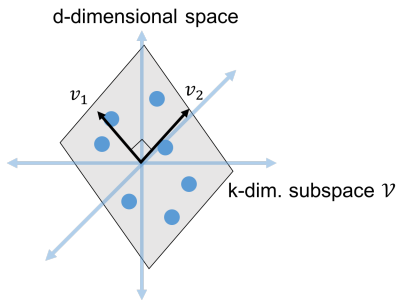
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$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

# 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

$$\mathbf{X} = \mathbf{C}\mathbf{V}^T\mathbf{X}\mathbf{V}^T.$$

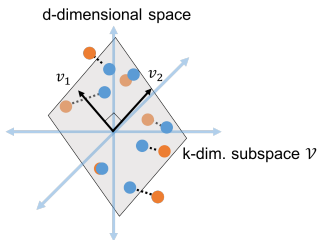
- $\mathbf{V}^T$  is a **projection matrix**, which projects vectors onto the subspace  $\mathcal{V}$ .



# Low-Rank Approximation

**Claim:** If  $\vec{x}_1, \dots, \vec{x}_n$  lie **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:**

$$\mathbf{X} \approx \mathbf{XV}^T$$



**Note:**  $\mathbf{XV}^T$  has rank  $k$ . It is a **low-rank approximation** of  $\mathbf{X}$ .

$$\mathbf{XV}^T = \underset{\mathbf{B} \text{ with rows in } \mathcal{V}}{\arg \min} \|\mathbf{X} - \mathbf{B}\|_F^2 = \sum_{i,j} (X_{i,j} - B_{i,j})^2.$$

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

# Low-Rank Approximation

**So Far:** If  $\vec{x}_1, \dots, \vec{x}_n$  lie 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:

$$\mathbf{X} \approx \mathbf{XV}^T.$$

This is the closest approximation to  $\mathbf{X}$  with rows in  $\mathcal{V}$  (i.e., in the column span of  $\mathbf{V}$ ).

- Letting  $(\mathbf{XV}^T)_i, (\mathbf{XV}^T)_j$  be the  $i^{\text{th}}$  and  $j^{\text{th}}$  projected data points,  
$$\|(\mathbf{XV}^T)_i - (\mathbf{XV}^T)_j\|_2 = \|[(\mathbf{XV})_i - (\mathbf{XV})_j]\mathbf{V}^T\|_2 = \|[(\mathbf{XV})_i - (\mathbf{XV})_j]\|_2.$$
- Can use  $\mathbf{XV} \in \mathbb{R}^{n \times k}$  as a compressed approximate data set.

Key question is how to find the subspace  $\mathcal{V}$  and correspondingly  $\mathbf{V}$ .

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

# Properties of Projection Matrices

**Quick Exercise:** 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  $\vec{y} \in \mathbb{R}^d$ .

Why does this make sense intuitively?

**Less Quick Exercise: (Pythagorean Theorem)** Show that:

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

# A Step Back: Why Low-Rank Approximation?

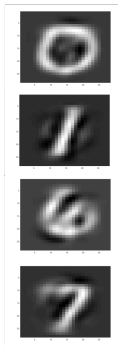
**Question:** Why might we expect  $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$  to lie close to a  $k$ -dimensional subspace?

- The rows of  $X$  can be approximately reconstructed from a basis of  $k$  vectors.

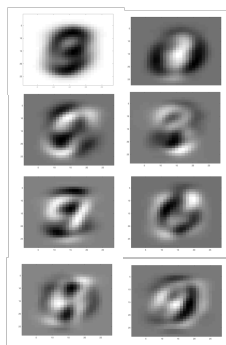
784 dimensional vectors



projections onto 15 dimensional space



orthonormal basis  $v_1, \dots, v_{15}$



# Dual View of Low-Rank Approximation

**Question:** Why might we expect  $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$  to lie close to a  $k$ -dimensional subspace?

- Equivalently, the columns of  $\mathbf{X}$  are approx. spanned by  $k$  vectors.

**Linearly Dependent Variables:**

	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

	bedrooms
home 1	2
home 2	4
•	•
•	•
•	•
home n	5 <sup>16</sup>