

# COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

---

Cameron Musco

University of Massachusetts Amherst. Fall 2021.

Lecture 14

- We will be grading midterms soon, and plan to return before the add/drop deadline.
- No quiz this week.
- No office hours today.

## 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.

$$[V] \begin{bmatrix} 1 \\ x \end{bmatrix} = \begin{bmatrix} x \end{bmatrix}$$

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.

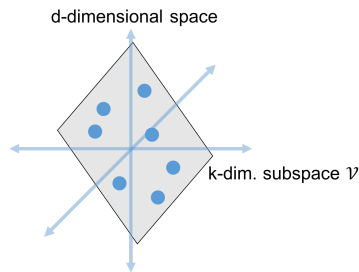
Will be using a fair amount of linear algebra: orthogonal basis, column/row span, eigenvectors, etc.

## RANDOMIZED ALGORITHMS UNIT TAKEAWAYS

- Randomization is an important tool in working with large datasets.
- Lets us solve 'easy' problems that get really difficult on massive datasets. Fast/space efficient look up (hash tables and bloom filters), distinct items counting, frequent items counting, near neighbor search (LSH), etc.
- The analysis of randomized algorithms leads to complex output distributions, which we can't compute exactly.
- We've covered many of the key ideas used through a small number of example applications/algorithms.
- We use concentration inequalities to bound these distributions and behaviors like accuracy, space usage, and runtime.
- Concentration inequalities and probability tools used in randomized algorithms are also fundamental in statistics, machine learning theory, probabilistic modeling of complex systems, etc.

## EMBEDDING WITH ASSUMPTIONS

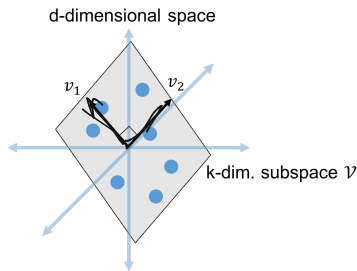
Assume that data points  $\vec{x}_1, \dots, \vec{x}_n$  lie in any  $k$ -dimensional subspace  $\mathcal{V}$  of  $\mathbb{R}^d$ .



# EMBEDDING WITH ASSUMPTIONS

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

$$V = \begin{bmatrix} | & | \\ \vec{v}_1 & \vec{v}_2 \\ | & | \end{bmatrix}$$



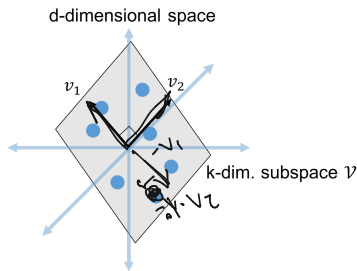
$$V^T = \begin{bmatrix} \vec{v}_1^T \\ \vec{v}_2^T \end{bmatrix}$$

**Claim:** Let  $\vec{v}_1, \dots, \vec{v}_k$  be an orthonormal basis for  $\mathcal{V}$  and  $V \in \mathbb{R}^{d \times k}$  be the matrix with these vectors as its columns. For all  $\vec{x}_i, \vec{x}_j$ :

$$\| \underbrace{V^T}_{2 \times d} \vec{x}_i - \underbrace{V^T}_{2 \times d} \vec{x}_j \|_2 = \| \vec{x}_i - \vec{x}_j \|_2.$$

## 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  $V \in \mathbb{R}^{d \times k}$  be the matrix with these vectors as its columns. For all  $\vec{x}_i, \vec{x}_j$ :

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

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



# DOT PRODUCT TRANSFORMATION

$$V = \begin{bmatrix} v_1 & v_2 & \dots & v_k \end{bmatrix} \in \mathbb{R}^{d \times k}$$

$$\begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_k^T \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_k \end{bmatrix} = \begin{bmatrix} v_1^T v_1 & v_1^T v_2 & \dots & v_1^T v_k \\ v_2^T v_1 & v_2^T v_2 & \dots & v_2^T v_k \\ \vdots & \vdots & \ddots & \vdots \\ v_k^T v_1 & v_k^T v_2 & \dots & v_k^T v_k \end{bmatrix}$$

**Claim:** Let  $\vec{v}_1, \dots, \vec{v}_k$  be an orthonormal basis for  $\mathcal{V}$  and

$V \in \mathbb{R}^{d \times k}$  be the matrix with these vectors as its columns. For

all  $\vec{x}_i, \vec{x}_j \in \mathcal{V}$ :

(A)  $\|V^T \vec{x}_i - V^T \vec{x}_j\|_2 = \|\vec{x}_i - \vec{x}_j\|_2$

$$\begin{aligned} (V(c_i - c_j))^T &= (c_i - c_j)^T V^T \end{aligned}$$

for any  $\vec{x}_i$ ,  $\vec{x}_i = V c_i$  for some  $c_i \in \mathbb{R}^k$

$$\vec{x}_i = v_1 \cdot c_i(1) + v_2 \cdot c_i(2) + \dots + v_k \cdot c_i(k)$$

$$\|V^T V c_i - V^T V c_j\|_2 = \|V c_i - V c_j\|_2$$

$$V^T V = I$$

$$(V^T V)_{ij} = v_i^T v_j \quad i \neq j, = 0$$

$$i=j = 1 = \|v_i\|_2^2$$

(For any  $y$ ,  $\|y\|_2^2 = y^T y$ )

$$y = V(c_i - c_j)$$

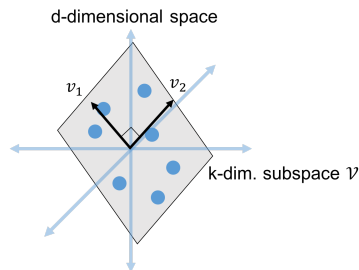
$$\|V(c_i - c_j)\|_2^2 = (c_i - c_j)^T V^T V (c_i - c_j)$$

$$= (c_i - c_j)^T I (c_i - c_j) = \|c_i - c_j\|_2^2$$

(B)  $\|c_i - c_j\|_2 = \|V(c_i - c_j)\|_2$

# 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  $\mathbb{R}^d$ .



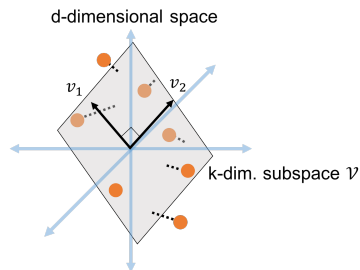
$$\|V^T x_i - V^T x_j\|_2^2 = \|x_i - x_j\|_2^2$$

$\downarrow$   
 $V^T(x_i - x_j)$

$$(x_i - x_j)^T V V^T (x_i - x_j)$$

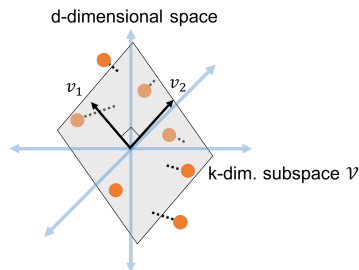
$(V^T V)$

**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  $\mathbb{R}^d$ .



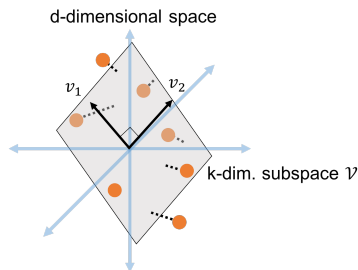
## EMBEDDING WITH ASSUMPTIONS

**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  $\mathbb{R}^d$ .



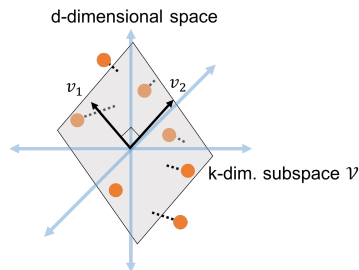
Letting  $\vec{v}_1, \dots, \vec{v}_k$  be an orthonormal basis for  $\mathcal{V}$  and  $\mathbf{V} \in \mathbb{R}^{d \times k}$  be the matrix with these vectors as its columns,  $\mathbf{V}^T \vec{x}_i \in \mathbb{R}^k$  is **still a good embedding** for  $x_i \in \mathbb{R}^d$ .

**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).

**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  $\leq k$ .

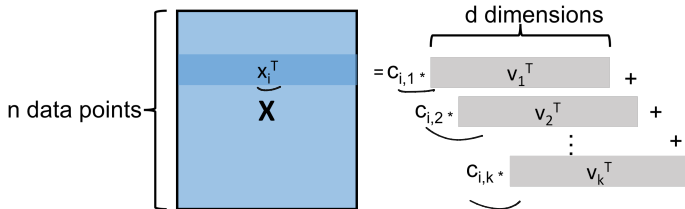
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$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

# LOW-RANK 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 any  $\vec{x}_i$  as:

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



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



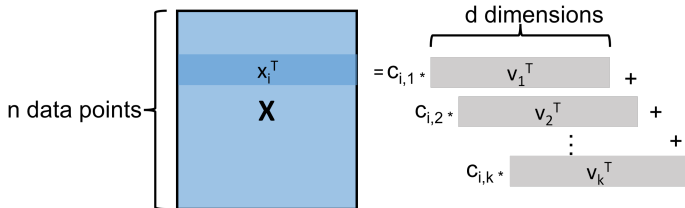
# LOW-RANK FACTORIZATION

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

$\vec{x}_1, \dots, \vec{x}_n$ : 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  $\vec{v}_1, \dots, \vec{v}_k$ .

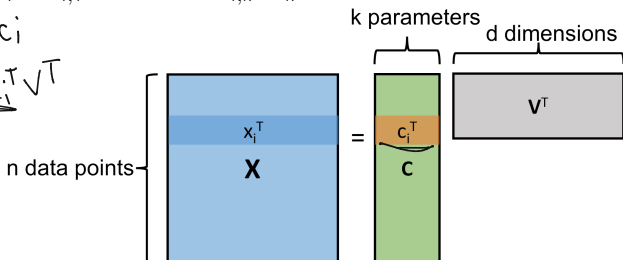
**Claim:**  $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$  lie in a  $k$ -dimensional subspace  $\mathcal{V} \Leftrightarrow$  the data matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  has rank  $\leq k$ .

- Every data point  $\vec{x}_i$  (row of  $\mathbf{X}$ ) can be written as

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

$$x_i = \mathbf{V}c_i$$

$$\underline{x_i^T} = \underline{c_i^T} \mathbf{V}^T$$

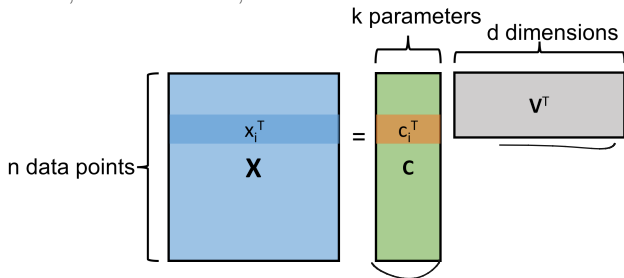


$\vec{x}_1, \dots, \vec{x}_n$ : data points (in  $\mathbb{R}^d$ ),  $\mathcal{V}$ :  $k$ -dimensional subspace of  $\mathbb{R}^d$ ,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots,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

**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

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

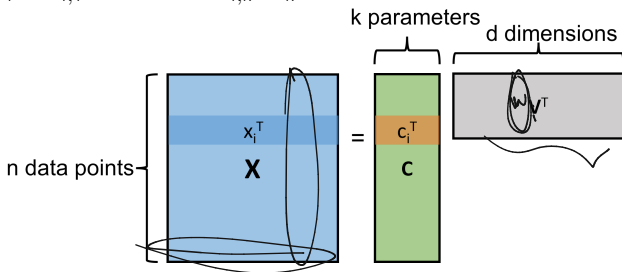


- $\mathbf{X}$  can be represented by  $(n + d) \cdot k$  parameters vs.  $n \cdot d$ .

$\vec{x}_1, \dots, \vec{x}_n$ : 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  $\vec{v}_1, \dots, \vec{v}_k$ .

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

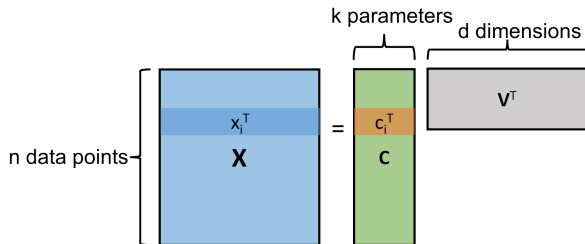


- $\mathbf{X}$  can be represented by  $(n + d) \cdot k$  parameters vs.  $n \cdot d$ .
- The rows of  $\mathbf{X}$  are spanned by  $k$  vectors: the columns of  $\mathbf{V} \implies$  the columns of  $\mathbf{X}$  are spanned by  $k$  vectors: the columns of  $\mathbf{C}$ .

$\vec{x}_1, \dots, \vec{x}_n$ : 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

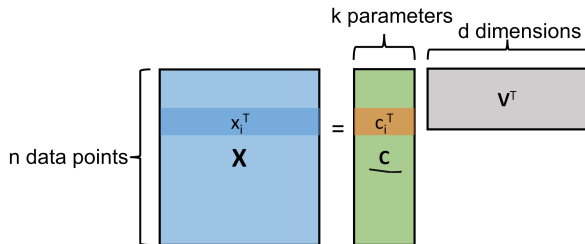
**Claim:** If  $\vec{x}_1, \dots, \vec{x}_n$  lie in a  $k$ -dimensional subspace with orthonormal basis  $\mathbf{V} \in \mathbb{R}^{d \times k}$ , the data matrix can be written as  $\mathbf{X} = \mathbf{C}\mathbf{V}^T$ .



$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : 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 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  $\mathbf{X} = \mathbf{C}\mathbf{V}^T$ .



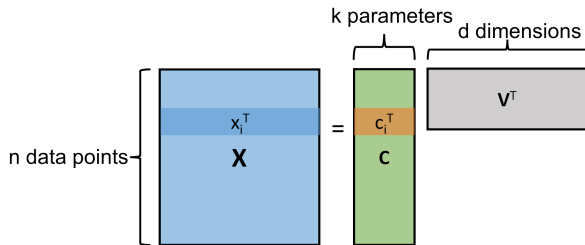
Exercise: What is this coefficient matrix  $\mathbf{C}$ ? Hint: Use that  $\mathbf{V}^T\mathbf{V} = \mathbf{I}$ .

$$\mathbf{C} = \mathbf{X}\mathbf{V}$$

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

# 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  $\mathbf{V}^T\mathbf{V} = \mathbf{I}$ .

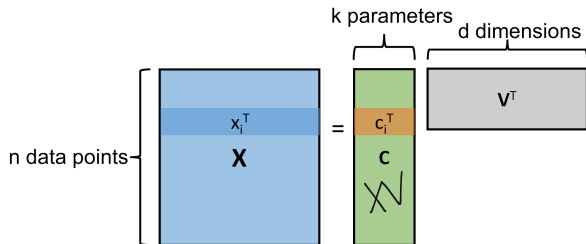
$$\bullet \mathbf{X} = \mathbf{C}\mathbf{V}^T \implies \mathbf{X}\mathbf{V} = \mathbf{C}\mathbf{V}^T\mathbf{V}$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : 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 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  $\mathbf{X} = \mathbf{C}\mathbf{V}^T$ .



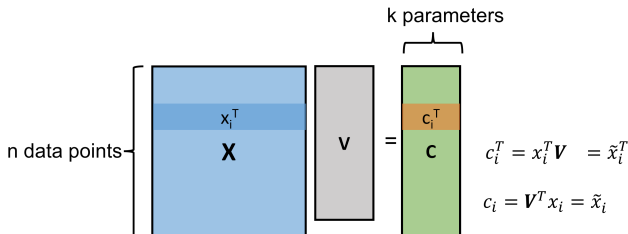
**Exercise:** What is this coefficient matrix  $\mathbf{C}$ ? **Hint:** Use that  $\mathbf{V}^T\mathbf{V} = \mathbf{I}$ .

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

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

# 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  $\vec{v}_1, \dots, \vec{v}_k$ .

**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

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

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

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

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

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

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

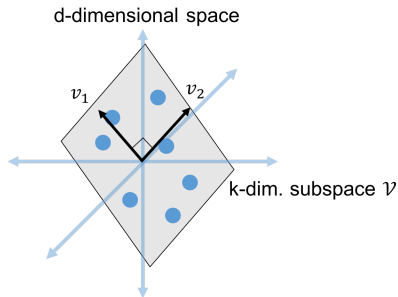
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$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

# 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

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

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



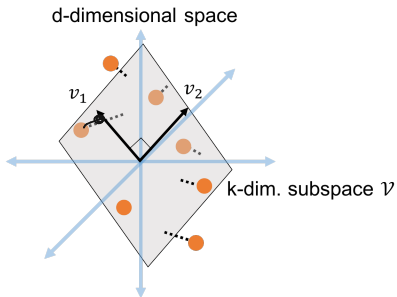
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$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

# 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

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

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



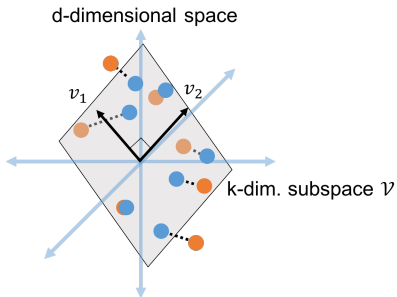
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$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

# 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

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

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



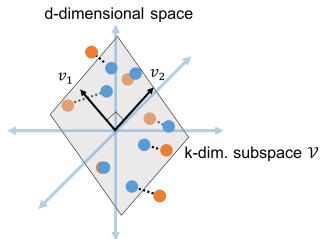
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$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



# LOW-RANK APPROXIMATION

**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:

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



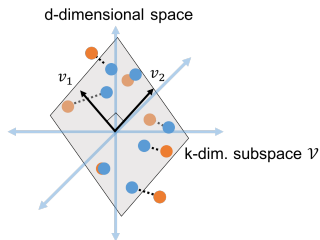
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$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

# LOW-RANK APPROXIMATION

**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:

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

*(Handwritten note:  $n \times d$  above  $\mathbf{X}$ )*



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

*(Handwritten note:  $n \times k$  under  $\mathbf{X}$ ,  $k \times d$  under  $\mathbf{V}$ )*

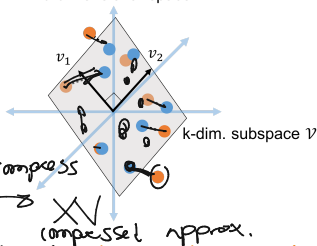
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ : data points,  $\mathbf{X} \in \mathbb{R}^{n \times d}$ : data matrix,  $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$ : orthogonal basis for subspace  $\mathcal{V}$ .  $\mathbf{V} \in \mathbb{R}^{d \times k}$ : matrix with columns  $\vec{v}_1, \dots, \vec{v}_k$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

# LOW-RANK APPROXIMATION

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

$$X \approx XVV^T$$

d-dimensional space



$$\|A\|_F^2 = \sum_{ij} A_{ij}^2$$

$$V, W$$

$$W^T = WW^T$$

rows lie in  $\mathcal{V}$   
they can be represented  
by  $k$  dimension vectors

original  $X$   $\rightarrow$  project to  $k$   
 $XV$  approx.

Note:  $XV$  has rank  $k$ . It is a low-rank approximation of  $X$ .

$$XV = \left( \arg \min_{B \text{ with rows in } \mathcal{V}} \|X - B\|_F^2 \right) = \sum_{ij} (X_{i,j} - B_{i,j})^2 = \sum_{i=1}^n \|x_i - b_i\|_2^2$$

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

**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:

$$\underline{\mathbf{X} \approx \mathbf{X}\mathbf{V}\mathbf{V}^T}.$$

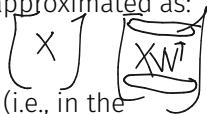
This is the closest approximation to  $\mathbf{X}$  with rows in  $\mathcal{V}$  (i.e., in the column span of  $\mathbf{V}$ ).

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

# 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,

$$\|(\mathbf{XV}^T)_i - (\mathbf{XV}^T)_j\|_2 = \underbrace{\|[(\mathbf{XV})_i - (\mathbf{XV})_j] \mathbf{V}^T\|_2}_{\|y^T \mathbf{V}^T\|_2 = y^T \mathbf{V}^T \mathbf{V} = y^T y = \|y\|_2^2} = \|[(\mathbf{XV})_i - (\mathbf{XV})_j]\|_2.$$

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

**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.

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

**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:

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

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

- Letting  $(\mathbf{XVV}^T)_i, (\mathbf{XVV}^T)_j$  be the  $i^{\text{th}}$  and  $j^{\text{th}}$  projected data points,
 
$$\|(\mathbf{XVV}^T)_i - (\mathbf{XVV}^T)_j\|_2 = \|[(\mathbf{XV})_i - (\mathbf{XV})_j]\mathbf{V}^T\|_2 = \|[(\mathbf{XV})_i - (\mathbf{XV})_j]\|_2.$$
- 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

$$d \begin{bmatrix} \overset{d}{\mathbf{V}\mathbf{V}^T} \\ \mathbf{y} \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{V} \\ \mathbf{X} \end{bmatrix} \mathbf{V}^T$$

$$\begin{bmatrix} (\mathbf{V}\mathbf{V}^T)\mathbf{y} \\ \mathbf{y}^T \mathbf{V}^T \end{bmatrix}$$

**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:

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





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

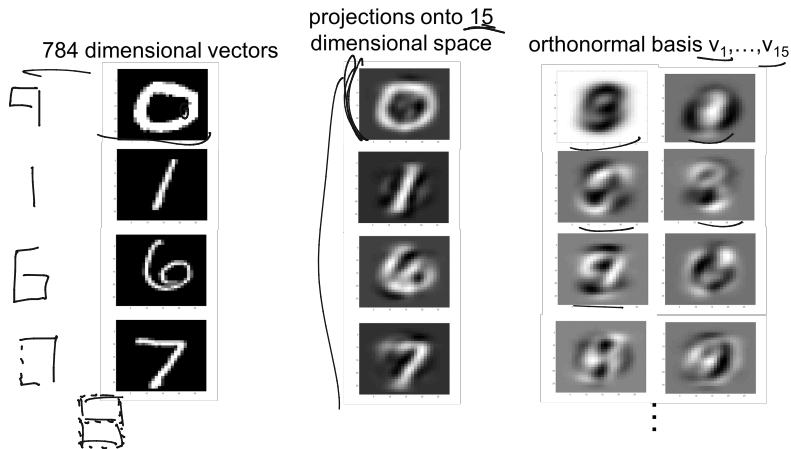
**Question:** Why might we expect  $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$  to lie close to a  $k$ -dimensional subspace?

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

## 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.



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

## 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?

$$X = (XV)V^T$$

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

# 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

# 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

# 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



# 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:

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

	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