

$$\begin{array}{l} \boxed{P_r(\text{some event}) > \underline{1-\delta}} \\ \downarrow \\ \underline{P_r(\text{not some event}) < \underline{\delta}} \end{array}$$

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Fall 2020.

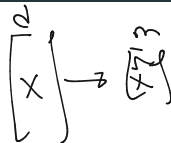
Lecture 13

- **Midterm is this Thursday - Friday.**
- **Office Hours:** I'll hold extra office hours, tomorrow from 2-3pm. The TAs will also hold their regular hours (see course page). Recordings of my office hours will be posted on Piazza.
- **Logistics:** Sometime on Thursday/Friday, you will download the exam in Gradescope, and should upload a pdf either of typed or handwritten answers 2 hours later. There will be a 15 minute buffer to upload in. Must submit by 11:59pm on Friday.
- **Questions:** Via private Piazza message. We'll try to answer frequently between 8am-10pm. If you don't get an answer, state any assumptions/interpretations you make clearly and move forward.
- **Academic Honestly:** You may not discuss the exam with any other students. Any cheating on the exam will result in failing the class. Please don't do this! It is much easier to catch than you might think, and the consequences seriously outweigh the benefits.

- **Midterm is this Thursday - Friday.**
- You must show your work/derive any answers to get full credit. Even on multiple choice questions.
- The exam is open notes. If you use outside resources (this should not be necessary) make sure to cite them.
- Very important to do some practice problems and to try them first with no resources, to simulate the exam.
- Make sure you can recognize when to apply the fundamentals: union bound, linearity of expectation and variance, Markov's inequality, Chebyshev's inequality, indicator random variables.
- Understand the goal of each algorithm/data structure. I.e., what problem it solves with what guarantees. No need to memorize proofs.

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.

$$\lfloor \uparrow \rfloor [x] \rightarrow [\tilde{x}]$$

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.

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.

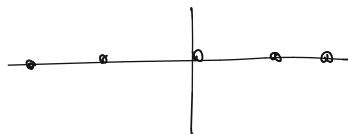
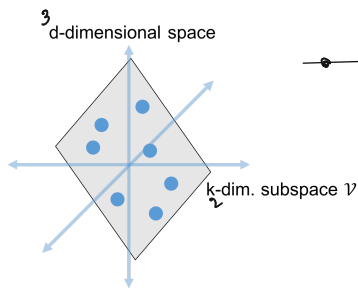
Will be using a fair amount of linear algebra: orthogonal basis, column/row span, eigenvectors, etc,

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 .



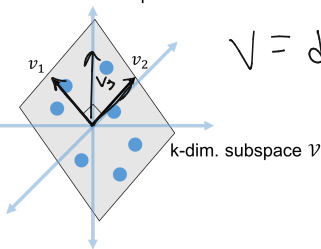
EMBEDDING WITH ASSUMPTIONS

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

$$\forall i, \|\vec{v}_i\|_2 = 1$$

$$\forall i, j \quad \langle \vec{v}_i, \vec{v}_j \rangle = 0$$

d-dimensional space



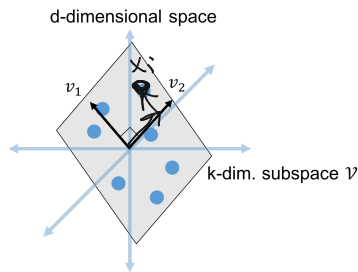
$$V = \begin{bmatrix} | & | & | & | \\ v_1 & v_2 & \dots & v_k \\ | & | & | & | \end{bmatrix}$$

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 :

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

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 .



$$V = \begin{bmatrix} | & & | \\ v_1 & \dots & v_k \\ | & & | \end{bmatrix}^k$$
$$V^T = \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix}^d$$

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 :

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

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

DOT PRODUCT TRANSFORMATION

$$\|y\|_2^2 = \sum_{i=1}^n y(i)^2 \quad \text{[0 0 y 0]} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} = y(1)y(1) + y(2)y(2) + \dots = \|y\|_2^2$$

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

① $\|V^T \vec{x}_i - V^T \vec{x}_j\|_2^2 = \|\vec{x}_i - \vec{x}_j\|_2^2$

$$d \begin{bmatrix} \vec{v}_1 & \dots & \vec{v}_k \end{bmatrix} \begin{bmatrix} c_1 \\ \vdots \\ c_k \end{bmatrix} = \begin{bmatrix} x_i \\ \vdots \\ x_j \end{bmatrix}$$



② x_i can be written as $v_1 c_1 + v_2 c_2 + \dots + v_k c_k$ for some values c_1, \dots, c_k

$$\boxed{x_i = V \vec{c}_i \text{ for } c_i \in \mathbb{R}^k}$$

③ $\|V^T V c_i - V^T V c_j\|_2^2 = \|V c_i - V c_j\|_2^2$

$V^T V = k \times k$ identity matrix

$$[V^T V]_{ij} = \langle v_i, v_j \rangle = 0 \text{ if } i \neq j \text{ (ortho)}$$

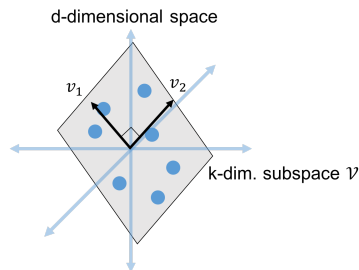
$$= 1 \text{ if } i = j \text{ (norm)}$$

④ $\|c_i - c_j\|_2^2 = \|V c_i - V c_j\|_2^2$

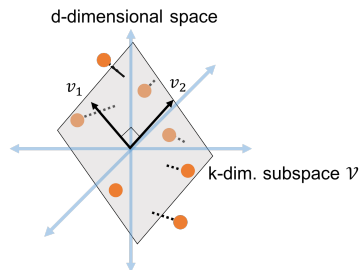
for any y , $\|y\|_2^2 = \langle y, y \rangle = y^T y$

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

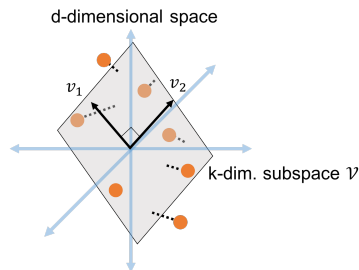
Main Focus of Upcoming Classes: Assume that data points $\vec{x}_1, \dots, \vec{x}_n$ lie **close to** any k -dimensional subspace \mathcal{V} of \mathbb{R}^d .



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 .

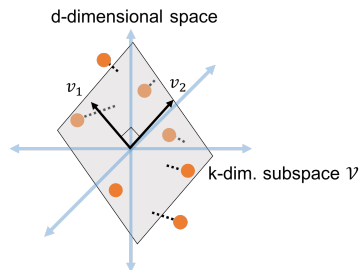


Main Focus of Upcoming Classes: Assume that data points $\vec{x}_1, \dots, \vec{x}_n$ lie **close to** any k -dimensional subspace \mathcal{V} of \mathbb{R}^d .



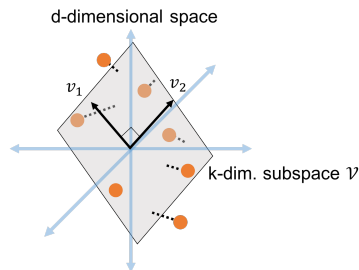
Letting $\vec{v}_1, \dots, \vec{v}_k$ be an orthonormal basis for \mathcal{V} and $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns, $\mathbf{V}^T \vec{x}_i \in \mathbb{R}^k$ is **still a good embedding** for $x_i \in \mathbb{R}^d$.

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 $\leq k$.

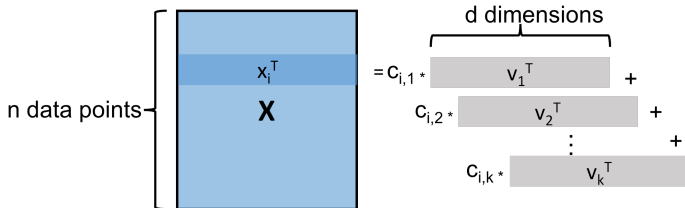
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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

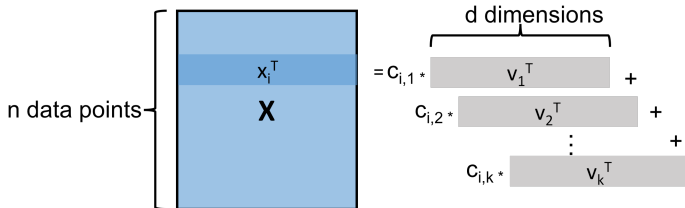
LOW-RANK FACTORIZATION

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

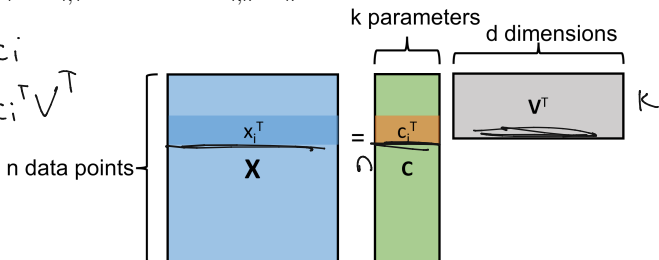
$\vec{x}_1, \dots, \vec{x}_n$: data points (in \mathbb{R}^d), \mathcal{V} : k -dimensional subspace of \mathbb{R}^d , $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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

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

$$\begin{aligned} x_i &= \mathbf{V}c_i \\ x_i^T &= c_i^T \mathbf{V}^T \end{aligned}$$

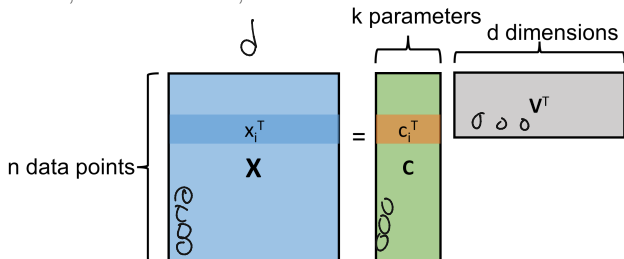


$\vec{x}_1, \dots, \vec{x}_n$: data points (in \mathbb{R}^d), \mathcal{V} : k -dimensional subspace of \mathbb{R}^d , $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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

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



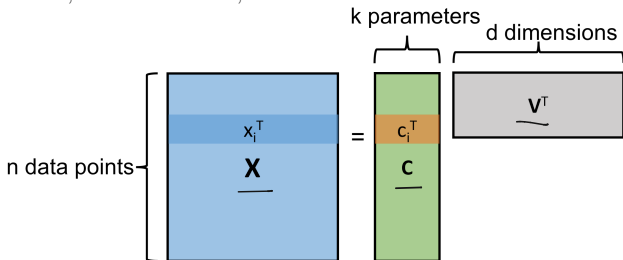
- \mathbf{X} can be represented by $(n + d) \cdot k$ parameters vs. $n \cdot d$.

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

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

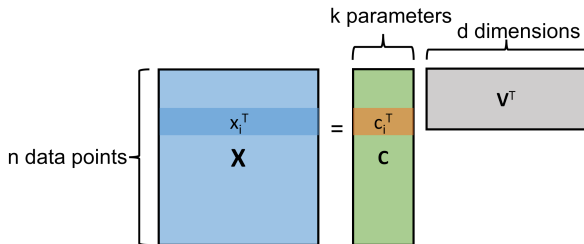


- \mathbf{X} can be represented by $(n + d) \cdot k$ parameters vs. $n \cdot d$.
- The rows of \mathbf{X} are spanned by k vectors: the columns of $\mathbf{V} \implies$ the columns of \mathbf{X} are spanned by k vectors: the columns of \mathbf{C} .

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

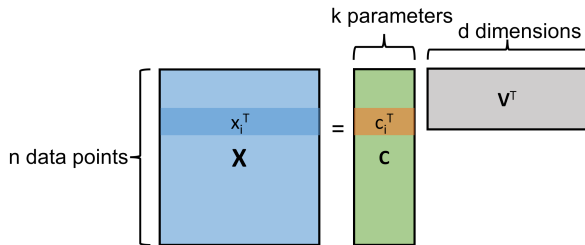
Claim: If $\vec{x}_1, \dots, \vec{x}_n$ lie in a k -dimensional subspace with orthonormal basis $\mathbf{V} \in \mathbb{R}^{d \times k}$, the data matrix can be written as $\mathbf{X} = \mathbf{C}\mathbf{V}^T$.



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

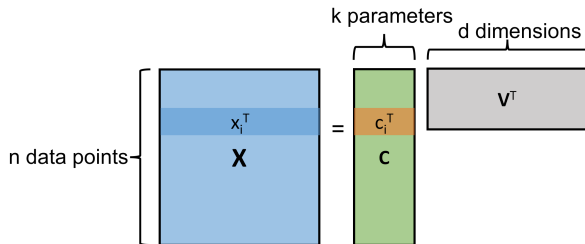


Exercise: What is this coefficient matrix \mathbf{C} ? **Hint:** Use that $\mathbf{V}^T\mathbf{V} = \mathbf{I}$.

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



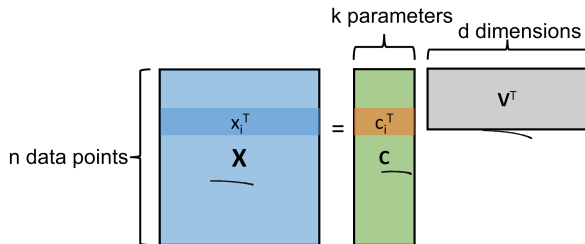
Exercise: What is this coefficient matrix \mathbf{C} ? **Hint:** Use that $\mathbf{V}^T\mathbf{V} = \mathbf{I}$.

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

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



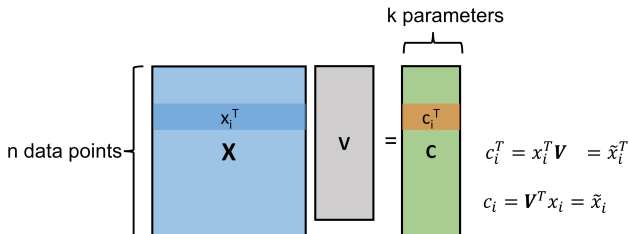
Exercise: What is this coefficient matrix \mathbf{C} ? **Hint:** Use that $\mathbf{V}^T\mathbf{V} = \mathbf{I}$.

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

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

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

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

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

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

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

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

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

$$\begin{matrix} d & k \\ \left[\begin{array}{c} \vdots \\ \mathbf{X} \\ \vdots \end{array} \right] & \left[\begin{array}{c} \vdots \\ \mathbf{V}^T \\ \vdots \end{array} \right]_k = \left[\begin{array}{c} \vdots \\ \mathbf{W} \\ \vdots \end{array} \right]_d \end{matrix} \quad \mathbf{X} = \mathbf{X}\mathbf{V}\mathbf{V}^T.$$

- $\mathbf{V}\mathbf{V}^T$ 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} .

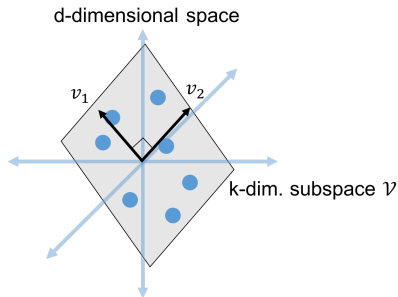
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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

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

- $\mathbf{V}\mathbf{V}^T$ 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} . $\vec{x}_i \in \mathcal{V}$



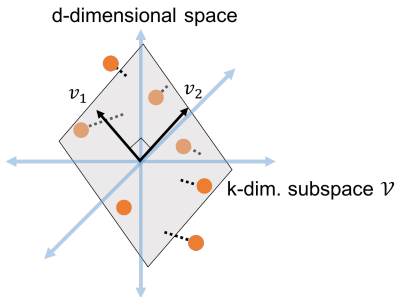
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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

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

- $\mathbf{V}\mathbf{V}^T$ 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} .



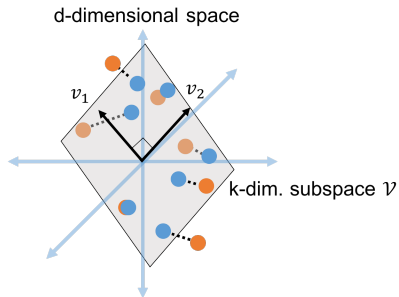
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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} .

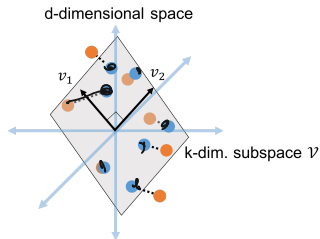


$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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:

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

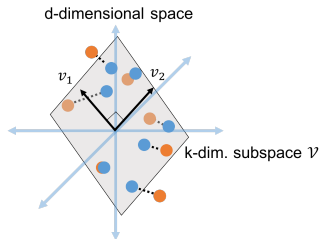


$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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:**

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



$$\begin{bmatrix} \vec{x}_1 \\ \vdots \\ \vec{x}_n \end{bmatrix} \begin{bmatrix} \vec{v}_1^T \\ \vdots \\ \vec{v}_k^T \end{bmatrix}$$

every column of $\mathbf{XV}\mathbf{V}^T$ is a linear combination of columns of \mathbf{XV}

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

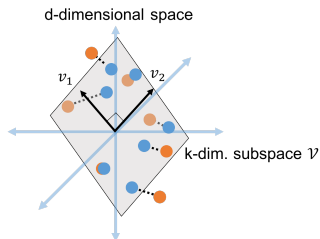
$$\begin{bmatrix} \vec{x}_1 \\ \vdots \\ \vec{x}_n \end{bmatrix} \begin{bmatrix} \vec{v}_1^T \\ \vdots \\ \vec{v}_k^T \end{bmatrix} = \begin{bmatrix} \vec{x}_1 \vec{v}_1^T \\ \vdots \\ \vec{x}_n \vec{v}_k^T \end{bmatrix} \quad \text{rank} \leq k$$

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthonormal 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**:

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



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

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

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

$$\mathbf{X} \approx \mathbf{XV}\mathbf{V}^T. \quad \vec{x}_1, \dots$$

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

1. If $\vec{x}_1, \dots, \vec{x}_n$ lie in a k -dim subspace, no distortion embedding into \mathbb{R}^k
2. If $\vec{x}_1, \dots, \vec{x}_n$ lie close to a k -dim subspace, $\underbrace{\mathbf{XV}\mathbf{V}^T}_{n \times k} \approx \mathbf{X}$
3. How do we find \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 $\vec{v}_1, \dots, \vec{v}_k$.

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^{th} and j^{th} projected data points,

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

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

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

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^{th} and j^{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 $\vec{v}_1, \dots, \vec{v}_k$.

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?

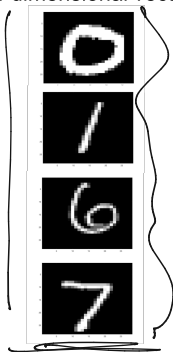
- The rows of \mathbf{X} can be approximately reconstructed from a basis of k vectors.

A STEP BACK: WHY LOW-RANK APPROXIMATION?

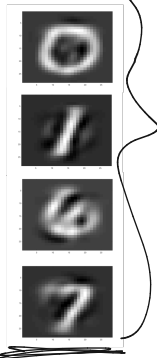
Question: Why might we expect $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ to lie close to a k -dimensional subspace?

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

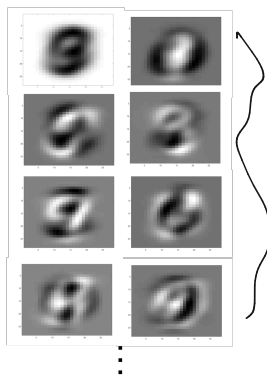
784 dimensional vectors



projections onto 15 dimensional space



orthonormal basis v_1, \dots, v_{15}



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?

- 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* 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