# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2024.
Lecture 14

- Midterms and Problem Set 2 are being graded now.
- Problem Set 3 will be released shortly, likely due 11/8.
- Quiz due monday.

## Summary

**Last Few Classes:** The Johnson-Lindenstrauss Lemma

- Reduce $n$ data points in any dimension $d$ to $O\left(\frac{\log n}{\epsilon^2}\right)$ dimensions and preserve all pairwise distances up to $1 \pm \epsilon$.

- Compression is linear via multiplication with a random, data oblivious, matrix (linear compression) $\quad [\ \Pi\ ][\vec{x}] \to \underbrace{(\tilde{x})}$

- Proved via the distributional JL-Lemma which shows that if $\mathbf{\Pi} \in \mathbb{R}^{m \times d}$ is a random matrix, $\mathbf{\Pi}\vec{y}_2 \approx \|\vec{y}\|$ for any $y$ with high probability. $\quad y = x_i - x_j$

- Proof of distributional JL via linearity of expectation, linearity of variance, stability of the Gaussian distribution, and an exponential concentration bound for Chi-Squared random variables.

## Summary

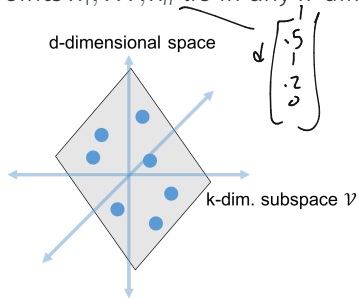Next Few Classes: Low-rank approximation, the SVD, and principal component analysis (PCA).

- Reduce $d$-dimesional 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).

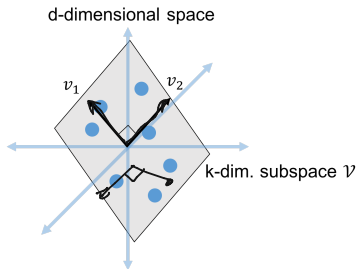Will be using a fair amount of linear algebra: orthogonal basis, column/row span, eigenvectors, etc.

4

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



d-dimensional space

k-dim. subspace $\mathcal{V}$

# Embedding with Assumptions

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



d-dimensional space

$v_1$  $v_2$

k-dim. subspace $\mathcal{V}$

$\mathbb{R}^d$

**Claim:** Let $\vec{v}_1, \ldots, \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and $\mathsf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns. For all $\vec{x}_i, \vec{x}_j$:
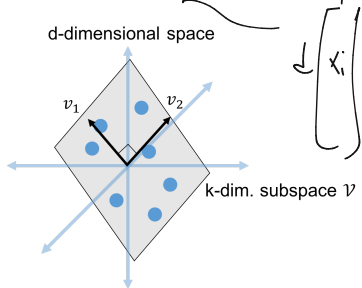
$$\|\mathsf{V}^T\vec{x}_i - \mathsf{V}^T\vec{x}_j\|_2 = \|\vec{x}_i - \vec{x}_j\|_2.$$

$k$

$\mathbb{R}^k$   $\mathbb{R}^d$

$$V = \begin{bmatrix} | & | & | \\ v_1 & v_2 & v_3 \cdots \\ | & | & | \end{bmatrix} \Bigg\} d$$

$k$

# Embedding with Assumptions

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



d-dimensional space

$v_1$   $v_2$

k-dim. subspace $\mathcal{V}$

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

$\|\Pi \vec{x}_i - \Pi \vec{x}_j\|$

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

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

5

# Dot Product Transformation

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

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

for all $i$, $\exists \; c_i \in \mathbb{R}^k$

$$\begin{bmatrix} x_i \end{bmatrix} = \begin{bmatrix} V & c_i \end{bmatrix}$$

$$x_i = V c_i$$

$$x_i = \vec{V}_1 \cdot c_i(1) + \vec{V}_2 \, c_i(2) \cdots + \vec{V}_k \, c_i(k)$$

$c_i$

$\|y\|_2 = \sqrt{\sum y_{(i)}^2}$

$\|y\|_2^2 = \sum y_{(i)}^2$

$= \langle y, y \rangle$

$= y^T y$

$$\|\mathbf{V}^T x_i - \mathbf{V}^T x_j\|_2^2$$

$$= \|\mathbf{V}^T V c_i - \mathbf{V}^T V c_j\|_2^2 = \|c_i - c_j\|_2^2$$

$\underbrace{\quad}_{I}$

$(\mathbf{V}^T V)_{ij} = \langle v_i, v_j \rangle = v_i^T v_j$

$= 1 \quad \text{if} \quad i = j$

$= 0 \quad \text{if} \quad i \neq j$

$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

$\|V c_i - V c_j\|_2$

$= \langle V(c_i - c_j), \; V(c_i - c_j) \rangle$

$= (c_i - c_j)^T V^T V (c_i - c_j)$

$\underset{I}{}$

$(c_i - c_j)^T (c_i - c_j)$

$\|c_i - c_j\|_2^2$

6

## Dot Product Transformation

**Claim:** Let $\vec{v}_1, \ldots, \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns. For all $\vec{x}_i, \vec{x}_j \in \mathcal{V}$:
$$\|\mathbf{V}^T \vec{x}_i - \mathbf{V}^T \vec{x}_j\|_2 = \|\vec{x}_i - \vec{x}_j\|_2.$$

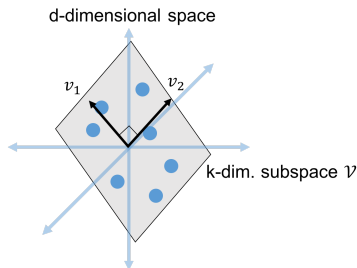$$\|y\|_2 = \sqrt{\sum_{i=1}^{d} y(i)^2}$$

$$\|y\|_2^2 = \sum_{i=1}^{d} y(i)^2 = \langle y, y \rangle$$

$$\left( \sum_{i=1}^{d} y(i) \cdot y(i) = \sum_{i=1}^{d} y(i)^2 \right.$$

$$\begin{bmatrix} \cdots & y^T & \cdots \end{bmatrix} \begin{bmatrix} | \\ y \\ | \end{bmatrix} = \sum_{i=1}^{|x|} y(i)^2$$
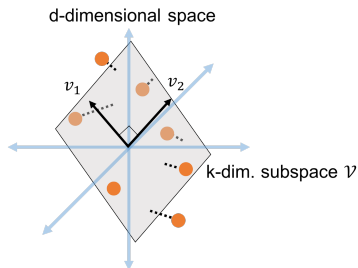
**Main Focus of Upcoming Classes:** Assume that data points $\vec{x}_1, \ldots, \vec{x}_n$ lie close to any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.



d-dimensional space

$v_1$   $v_2$

k-dim. subspace $\mathcal{V}$

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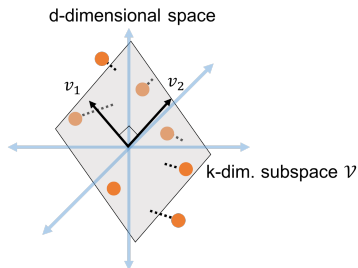


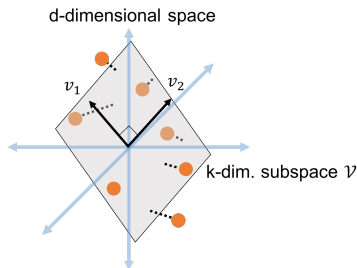d-dimensional space

$v_1$  $v_2$

k-dim. subspace $\mathcal{V}$

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



d-dimensional space

$v_1$ $v_2$

k-dim. subspace $\mathcal{V}$

Letting $\vec{v}_1, \ldots, \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and $\mathsf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns, $\underline{\mathsf{V}^T \vec{x}_i} \in \mathbb{R}^k$ is still a good embedding for $x_i \in \mathbb{R}^d$.
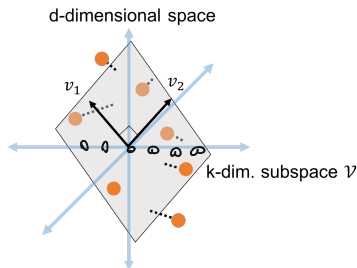
**Main Focus of Upcoming Classes:** Assume that data points $\vec{x}_1, \ldots, \vec{x}_n$ lie close to any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.



Letting $\vec{v}_1, \ldots, \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and $\mathsf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns, $\mathsf{V}^T \vec{x}_i \in \mathbb{R}^k$ 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 $\vec{x}_1, \ldots, \vec{x}_n$ lie close to any $k$-dimensional subspace $\mathcal{V}$ of $\mathbb{R}^d$.



d-dimensional space

$v_1$    $v_2$

k-dim. subspace $\mathcal{V}$

Letting $\vec{v}_1, \ldots, \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$ and $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns, $\mathbf{V}^T \vec{x}_i \in \mathbb{R}^k$ 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:** $\vec{x}_1, \ldots, \vec{x}_n$ lie in a $k$-dimensional subspace $\mathcal{V}$ $\Leftrightarrow$ the data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ has rank $\leq k$.
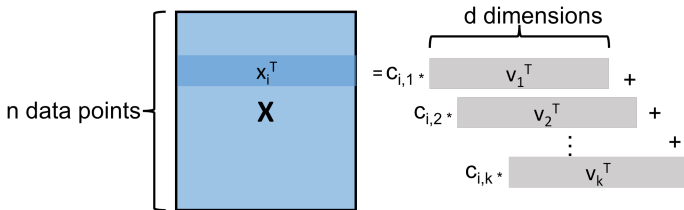


$\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \ldots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace $\mathcal{V}$. $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \ldots, \vec{v}_k$.

## Low-Rank Factorization

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

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

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



$\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$: data points, $X \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \ldots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace $\mathcal{V}$. $V \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \ldots, \vec{v}_k$.
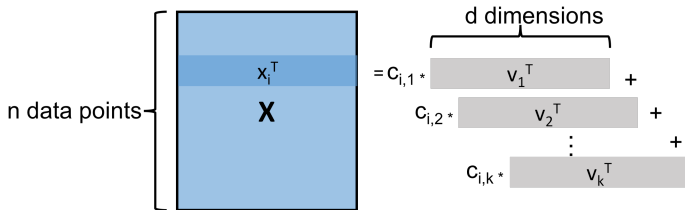
9

# Low-Rank Factorization

**Claim:** $\vec{x}_1, \ldots, \vec{x}_n$ 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 $\vec{v}_1, \ldots, \vec{v}_k$ be an orthonormal basis for $\mathcal{V}$, can write $\vec{x}_i$ as:

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

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



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

9

**Claim:** $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ 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
  $$\underbrace{\vec{x}_i = \mathbf{V}\vec{c}_i} = c_{i,1} \cdot \vec{v}_1 + \ldots + c_{i,k} \cdot \vec{v}_k.$$

---

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

**Claim:** $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ lie in a $k$-dimensional subspace $\mathcal{V}$ ⇔ 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

$$\vec{x}_i = \mathbf{V}\vec{c}_i = c_{i,1} \cdot \vec{v}_1 + \ldots + c_{i,k} \cdot \vec{v}_k.$$
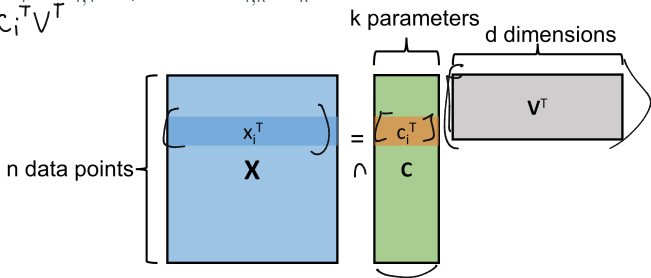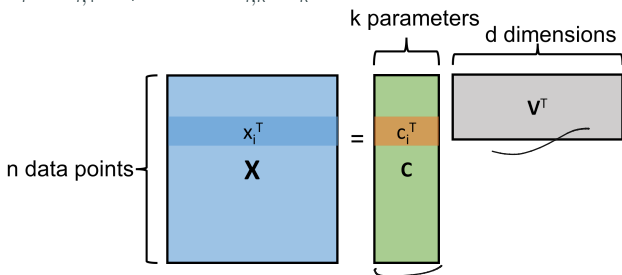
$$x_i^T = c_i^T \mathbf{V}^T$$



---

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

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

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

$\text{rank}(X) = 2$



k parameters     d dimensions
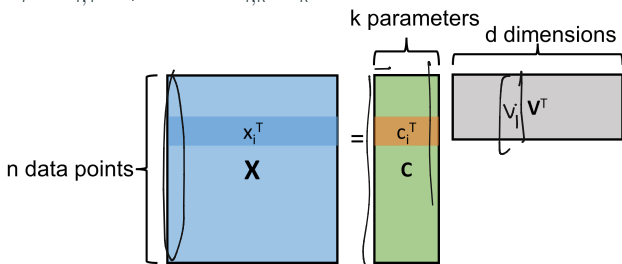
n data points

$x_i^T$

$X$

$=$

$c_i^T$

$C$

$V^T$

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

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

10

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

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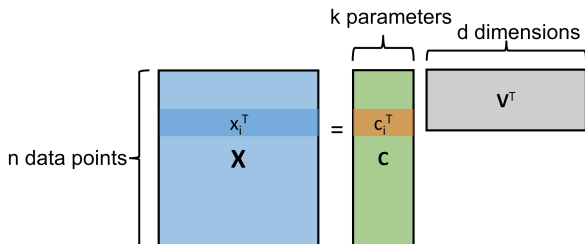  $\vec{x}_i = \mathsf{V}\vec{c}_i = c_{i,1} \cdot \vec{v}_1 + \ldots + c_{i,k} \cdot \vec{v}_k$.



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

- The rows of $\mathsf{X}$ are spanned by $k$ vectors: the columns of $\mathsf{V}$ $\implies$ the columns of $\mathsf{X}$ are spanned by $k$ vectors: the columns of $\mathsf{C}$.

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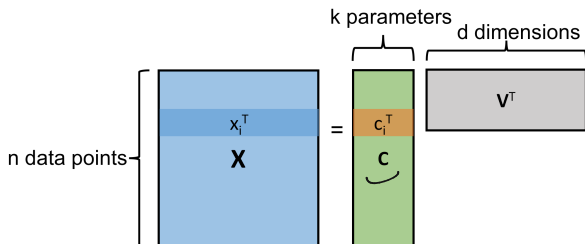
## Low-Rank Factorization

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



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

## Low-Rank Factorization

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



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

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

# Low-Rank Factorization

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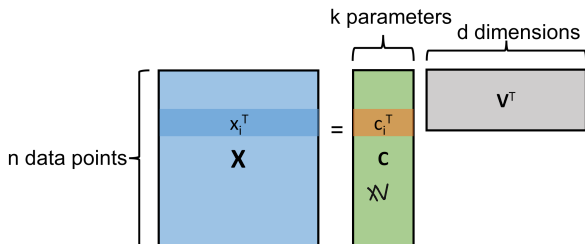
$V^T x_i = c_i$



k parameters

d dimensions

n data points

$x_i^T$

$X$

$=$

$c_i^T$

$C$

$XV$

$V^T$

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

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

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

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



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

- $X = CV^T \implies XV = CV^T V \implies XV = C$

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

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



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

- $X = CV^T \implies XV = CV^T V \implies XV = C$

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

## Projection View

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

$$X = CV^T.$$
$$\underset{XV}{\big\lfloor}$$

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

12

## Projection View

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

$$V^T V = I_k$$

$$X = XVV^T.$$

$$\left[ \begin{array}{c} V^T \end{array} \right] \left[ \begin{array}{c} V \end{array} \right] = \left[ I \right]$$

$$\overset{k}{\left[ \begin{array}{c} V \end{array} \right]} \left[ \begin{array}{c} V^T \end{array} \right] = \left[ VV^T \right]$$

---

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

## Projection View

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

$$X = XVV^T.$$

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

---

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

# Projection View

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

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

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



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

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

$$X = XVV^T.$$

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


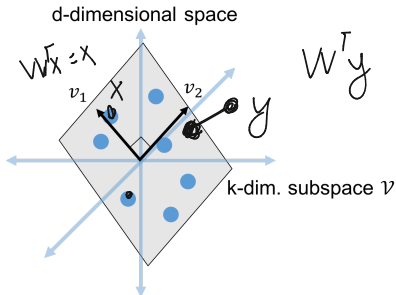
d-dimensional space

$v_1$  $v_2$

k-dim. subspace $\mathcal{V}$

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

# Projection View

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

$$X = XVV^T.$$

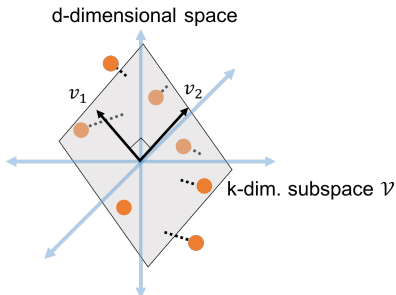- $VV^T$ is a projection matrix, which projects vectors onto the subspace $\mathcal{V}$.



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

## Low-Rank Approximation

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

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



d-dimensional space

$v_1$  $v_2$
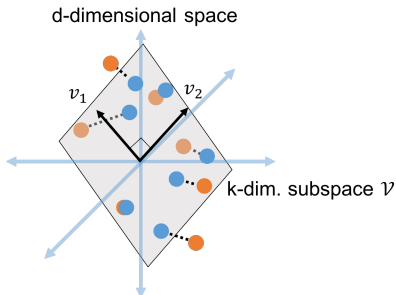
k-dim. subspace $\mathcal{V}$

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

# Low-Rank Approximation

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

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



d-dimensional space

$v_1$   $v_2$

k-dim. subspace $\mathcal{V}$

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

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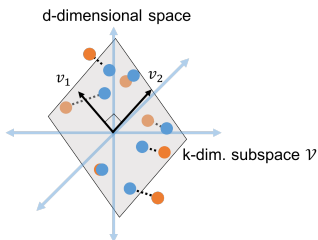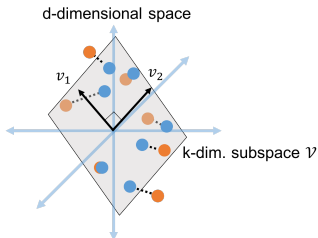
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathsf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace $\mathcal{V}$. $\mathsf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

# Low-Rank Approximation

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

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



d-dimensional space

$v_1$    $v_2$

k-dim. subspace $\mathcal{V}$

$$\|X - B\|_F^2 = \sum_{i\,j}(x_{ij} - B_{ij})^2$$
$$= \sum_{i=1}^{n}\|x_i - b_i\|_2^2$$

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

$$\mathbf{X}\mathbf{V}^\mathsf{T} = \underset{\mathbf{B} \text{ with rows in } \mathcal{V}}{\arg\min} \|\mathbf{X} - \mathbf{B}\|_F^2 = \sum_{i,j}(\mathbf{X}_{i,j} - \mathbf{B}_{i,j})^2 = \sum_{i}\|x_i - b_i\|_2^2$$

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

13

## Low-Rank Approximation

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

$$X \approx XVV^T.$$

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

---

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

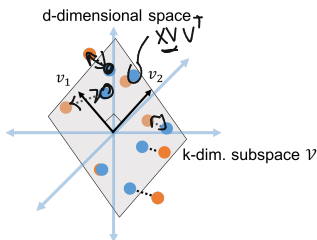## Low-Rank Approximation

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

$$\mathsf{X} \approx \mathsf{X}\mathsf{V}\mathsf{V}^T.$$

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

- Letting $(\mathsf{X}\mathsf{V}\mathsf{V}^T)_i, (\mathsf{X}\mathsf{V}\mathsf{V}^T)_j$ be the $i^{th}$ and $j^{th}$ projected data points,

$$\|\underbrace{(\mathsf{X}\mathsf{V}\mathsf{V}^T)_i}_{c_i} - \underbrace{(\mathsf{X}\mathsf{V}\mathsf{V}^T)_j}_{c_j}\|_2 = \|[(\mathsf{X}\mathsf{V})_i - (\mathsf{X}\mathsf{V})_j]\underbrace{\mathsf{V}^T}\|_2 = \|[(\mathsf{X}\mathsf{V})_i - (\mathsf{X}\mathsf{V})_j]\|_2.$$

$$\|c_i - c_j\|_2$$

---

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

## Low-Rank Approximation

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

$$X \approx XVV^T.$$

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

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

- Can use $XV \in \mathbb{R}^{n \times k}$ as a compressed approximate data set.

---

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

## Low-Rank Approximation

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

$$\mathsf{X} \approx \mathsf{X}\mathsf{V}\mathsf{V}^T.$$

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

- Letting $(\mathsf{X}\mathsf{V}\mathsf{V}^T)_i, (\mathsf{X}\mathsf{V}\mathsf{V}^T)_j$ be the $i^{th}$ and $j^{th}$ projected data points,

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

- Can use $\mathsf{X}\mathsf{V} \in \mathbb{R}^{n \times k}$ as a compressed approximate data set.

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

---

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

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

Why does this make sense intuitively?

Less Quick Exercise: (Pythagorean Theorem) Show that:

$V^TV = I$

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



$y$

$y - VV^Ty$

$V^Ty$

## A Step Back: Why Low-Rank Approximation?

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

## A Step Back: Why Low-Rank Approximation?

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

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

# A Step Back: Why Low-Rank Approximation?

**Question:** Why might we expect $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ 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 $v_1, \ldots, v_{15}$

## Dual View of Low-Rank Approximation

Question: Why might we expect $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ to lie close to a $k$-dimensional subspace?
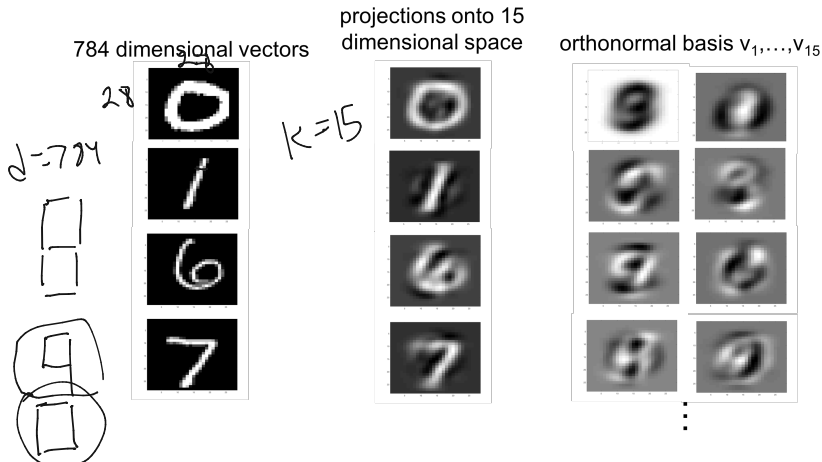
## Dual View of Low-Rank Approximation

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

- Equivalently, the columns of $X$ are approx. spanned by $k$ vectors.

## Dual View of Low-Rank Approximation

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

- Equivalently, the columns of $X$ 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 |

**Question:** Why might we expect $\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d$ 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 |

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

- Equivalently, the columns of $X$ 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 |

$X \approx \begin{Bmatrix} \\ \end{Bmatrix}$