

COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2023.

Lecture 15

- Midterm grades and solutions are posted on Moodle.
- We'll hand out the midterms at the end of class.
- The class average was $\approx \underline{30/39 = 77\%}$.

See Piazza post for more details. If you aren't happy with your grade, I'm happy to chat about strategies moving forward.

Quiz Question

Question 5

Not complete

Points out of 1.00

Flag question

Edit question

Suppose $x=(1,2,3,4)$ and let $y=(y_1,y_2,y_3,y_4)$ be a random vector where each y_i is independent and is distributed according to a Normal distribution with mean 0 and variance 1. What is the expected value of $\langle x, y \rangle^2$?

Answer:

Check

$$\|x\|_2^2 = 30$$

$$1 + 2^2 + 3^2 + 4^2 = \underline{\underline{30}}$$

$$x = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$$

$$y = \begin{bmatrix} N(0,1) \\ N(0,1) \\ \vdots \\ N(0,1) \end{bmatrix}$$

$$E[\langle x, y \rangle^2]$$

$$\text{Var}(\langle x, y \rangle) = E[\langle x, y \rangle^2] - E[\langle x, y \rangle]^2$$

$$E[y_1] + E[2y_2] + E[3y_3] + E[4y_4]$$

$$\text{Var}(\langle x, y \rangle) = \text{Var}(y_1) + \text{Var}(2y_2) + \text{Var}(3y_3) + \text{Var}(4y_4)$$

Summary

Last Few Classes: The Johnson-Lindenstrauss Lemma

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

- **Compression is linear** via multiplication with a random, **data oblivious**, matrix (linear compression)

$$\begin{bmatrix} \pi \\ \downarrow \\ X \end{bmatrix} = \begin{bmatrix} \downarrow \\ X \end{bmatrix}$$

- 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.
- 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

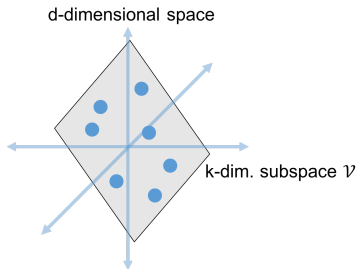
Next Few Classes: Low-rank approximation, the SVD, and principal component analysis (PCA).

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

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

Embedding with Assumptions

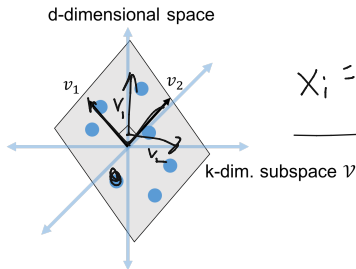
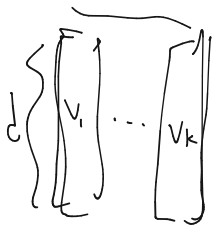
Assume that data points $\vec{x}_1, \dots, \vec{x}_n$ lie in any k -dimensional subspace \mathcal{V} of \mathbb{R}^d .



$$x_i = \begin{pmatrix} | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \end{pmatrix} \in \mathbb{R}^d$$

Embedding with Assumptions

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



$$\vec{x}_i = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_k \vec{v}_k$$

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

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

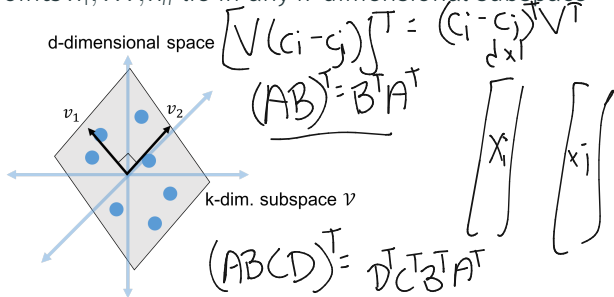
$$k \begin{bmatrix} \vec{v}_1 \\ \vdots \\ \vec{v}_k \end{bmatrix}^T \begin{bmatrix} \vec{x}_i \\ \vdots \\ \vec{x}_j \end{bmatrix} = \begin{bmatrix} \vec{v}_1^T \vec{x}_i \\ \vdots \\ \vec{v}_k^T \vec{x}_i \end{bmatrix}$$

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 .

$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}^T \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}^T \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$k \times 1$ $k \times n$



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**.

$\langle z, y \rangle = z^T y = \sum_{k=1}^d z(k) y(k)$
 $\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}^T \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}^T \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$

Dot Product Transformation

$y = x_i - x_j$ $\|V^T y\| \leq \|V^T\| \|y\| = \|y\|$
 $\xrightarrow{X \rightarrow V}$ 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 $\vec{x}_i, \vec{x}_j \in \mathcal{V}$: $\langle y, y \rangle = y^T y$

$V^T = 2 \times V$

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

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

$\exists c_i, c_j \in \mathbb{R}^k$ s.t. $x_i = V c_i = v_1 c_i(1) + v_2 c_i(2) \dots$

$x_j = V c_j$

$\Rightarrow \|V^T V c_i - V^T V c_j\|_2 = \|c_i - c_j\|_2 = \|x_i - x_j\|_2$

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

$\begin{bmatrix} V_1^T \\ \vdots \\ V_k^T \end{bmatrix} \begin{bmatrix} v_1 \dots v_k \end{bmatrix} = (V^T V)_{ij} = \langle v_i, v_j \rangle$

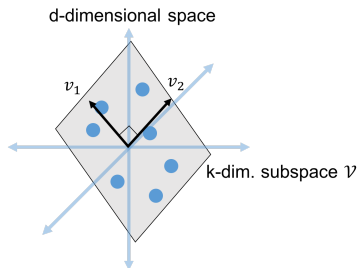
$i \neq j \Rightarrow \langle v_i, v_j \rangle = 0$
 $i = j \Rightarrow \langle v_i, v_i \rangle = \|v_i\|_2^2 = 1$

$\sum_{k=1}^d v_i(k) \cdot v_i(k) = \sum_{k=1}^d v_i(k)^2 = \|v_i\|_2^2 = 1$

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

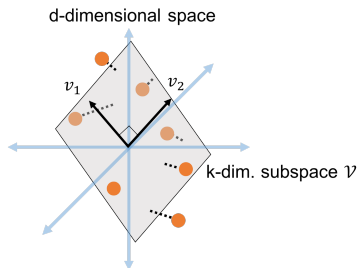
Embedding with Assumptions

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 .



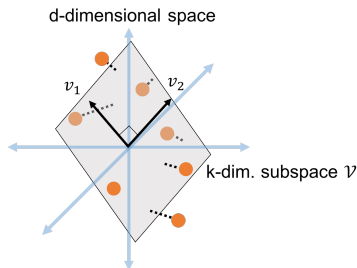
Embedding with Assumptions

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 .



Embedding with Assumptions

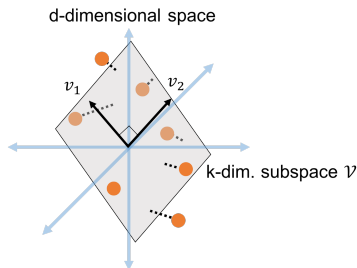
Main Focus of Upcoming Classes: Assume that data points $\vec{x}_1, \dots, \vec{x}_n$ lie **close to** any k -dimensional subspace \mathcal{V} of \mathbb{R}^d .



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 $x_i \in \mathbb{R}^d$.

Embedding with Assumptions

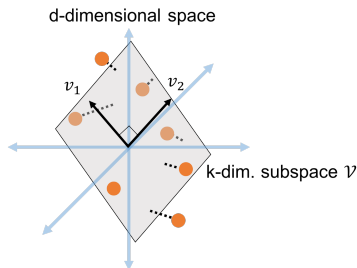
Main Focus of Upcoming Classes: Assume that data points $\vec{x}_1, \dots, \vec{x}_n$ lie **close to** any k -dimensional subspace \mathcal{V} of \mathbb{R}^d .



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

Embedding with Assumptions

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



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

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

Low-Rank Factorization

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

A hand-drawn diagram illustrating the data matrix \mathbf{X} . The matrix is enclosed in large square brackets. On the left side, the number n is written vertically. On the top side, the number d is written horizontally. Inside the brackets, the matrix is represented as a list of row vectors: x_1^T , x_2^T , followed by a vertical ellipsis, and x_n^T . To the right of the matrix, the equation $\text{rank}(\mathbf{X}) \leq k$ is written in a cursive hand-drawn style.

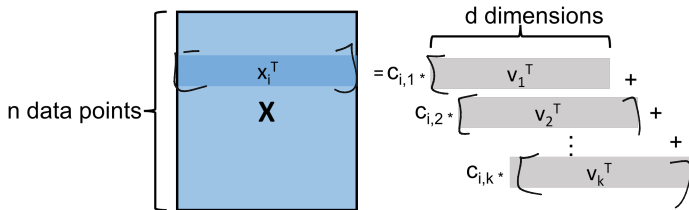
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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

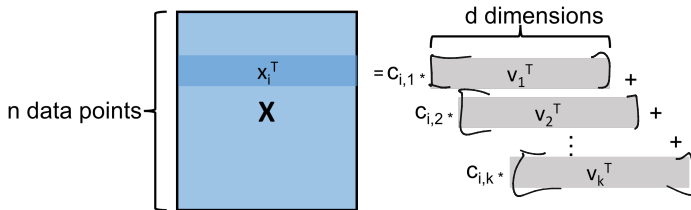
Low-Rank Factorization

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

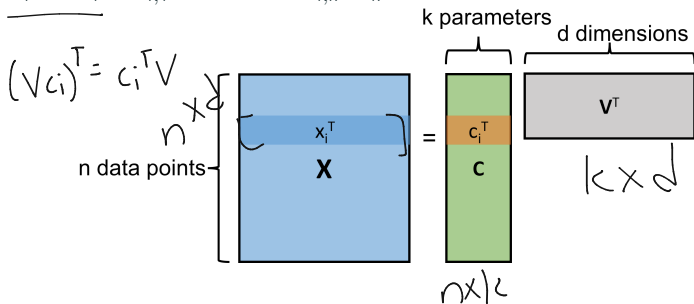
Claim: $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ lie in a k -dimensional subspace $\mathcal{V} \Leftrightarrow$ the data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ has rank $\leq k$.

- 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 $\leq k$.

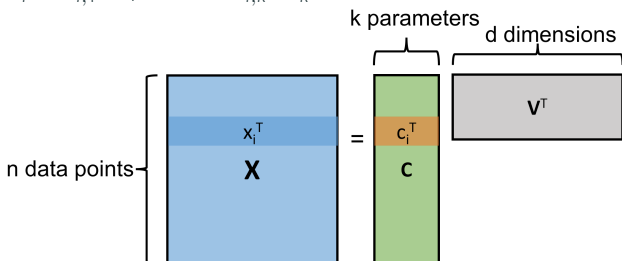
- Every data point \vec{x}_i (row of \mathbf{X}) can be written as $\vec{x}_i = \mathbf{V} \vec{c}_i = c_{i,1} \cdot \vec{v}_1 + \dots + c_{i,k} \cdot \vec{v}_k$.



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



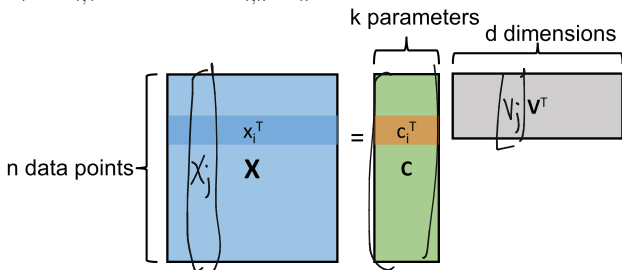
- \mathbf{X} can be represented by $(n + d) \cdot k$ parameters vs. $n \cdot d$.

$$\underbrace{n \cdot k}_{\text{C's}} \quad \underbrace{d \cdot k}_{\text{V}}$$

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

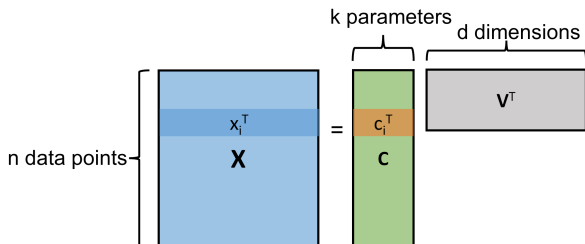


- \mathbf{X} can be represented by $(n + d) \cdot k$ parameters vs. $n \cdot d$.
- The rows of \mathbf{X} are spanned by k vectors: the columns of $\mathbf{V} \Rightarrow$ the columns of \mathbf{X} are spanned by k vectors: the columns of \mathbf{C} .

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

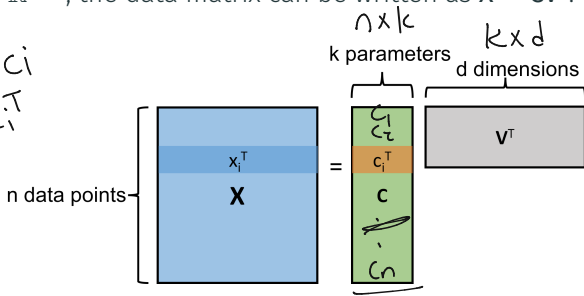


$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthonormal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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

$$\begin{aligned} \mathbf{V}^T \mathbf{x}_i &= \mathbf{c}_i \\ \mathbf{x}_i^T \mathbf{V} &= \mathbf{c}_i^T \\ \mathbf{X} \mathbf{V} &= \mathbf{C} \end{aligned}$$



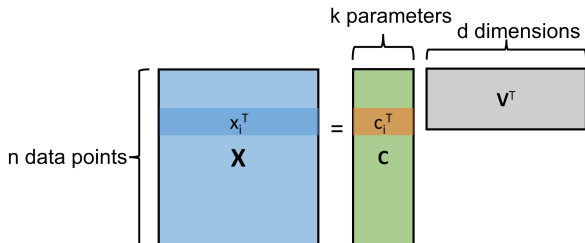
Exercise: What is this coefficient matrix \mathbf{C} ? Hint: Use that $\mathbf{V}^T \mathbf{V} = \mathbf{I}$.

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

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



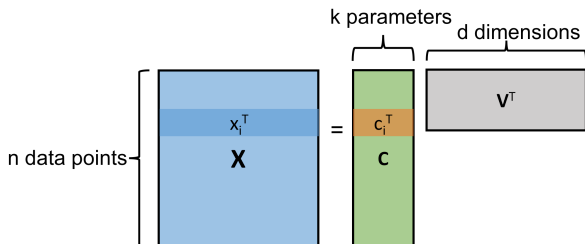
Exercise: What is this coefficient matrix \mathbf{C} ? **Hint:** Use that $\mathbf{V}^T\mathbf{V} = \mathbf{I}$.

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

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



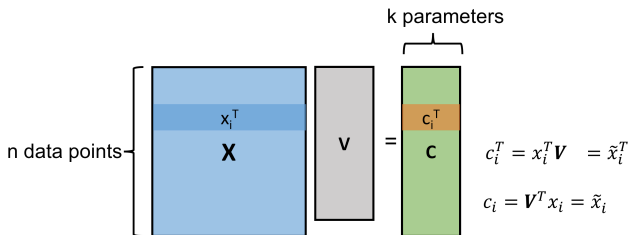
Exercise: What is this coefficient matrix \mathbf{C} ? Hint: Use that $\mathbf{V}^T\mathbf{V} = \mathbf{I}$.

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

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

Projection View

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

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

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

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

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

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

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 vectors onto the subspace \mathcal{V} .

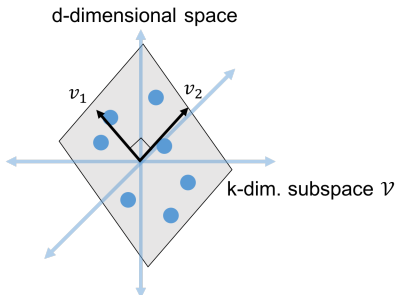
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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

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

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



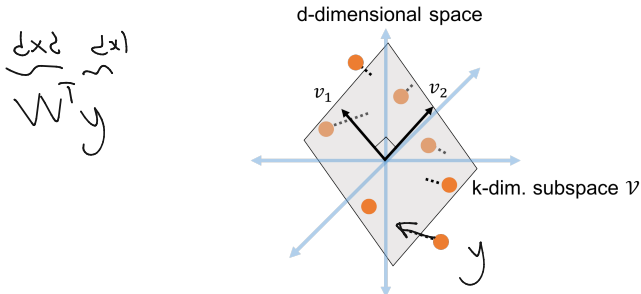
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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 vectors onto the subspace \mathcal{V} .



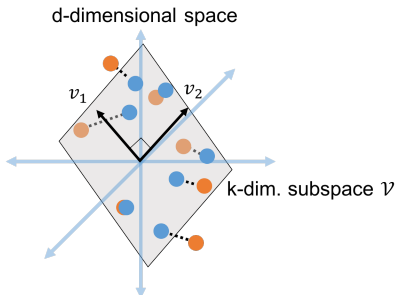
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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 vectors onto the subspace \mathcal{V} .

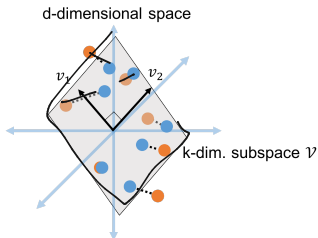


$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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:**

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

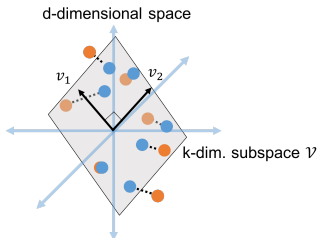


$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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:

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



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

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