

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Spring 2020.

Lecture 14

Move Online:

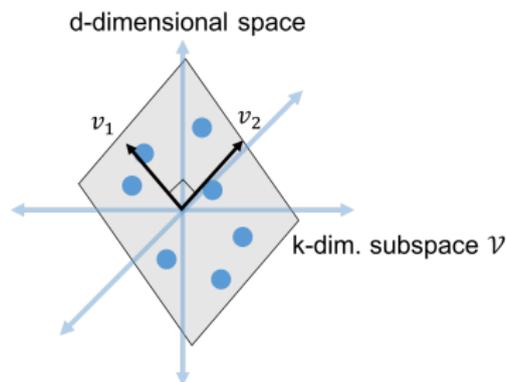
- Check out Piazza post for details about moving online.
- Lectures will be streamed and recorded. Feel free to ask questions using audio or by typing into chat. **Mute when not talking.**
- Feel free to turn on video, although it will be automatically off at the beginning of each lecture.
- Office hours will be over Zoom, after class on Tuesdays. **Different Zoom link.**
- Message me if you want to attend office hours but can't.
- Problem set rules will remain the same: you can submit in groups of up to three, but do not have to.

Midterm:

- Midterm grades are posted in Moodle. Average was a 30/37.
- Email me if you'd like to see your graded midterm.
- I won't release an answer key, but you can ask about midterm solutions in office hours or on Piazza.
- If you were not happy with your performance I'm happy to talk about it, and see if there are any adjustments we can make to get things on track.

LAST CLASS: EMBEDDING WITH ASSUMPTIONS

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



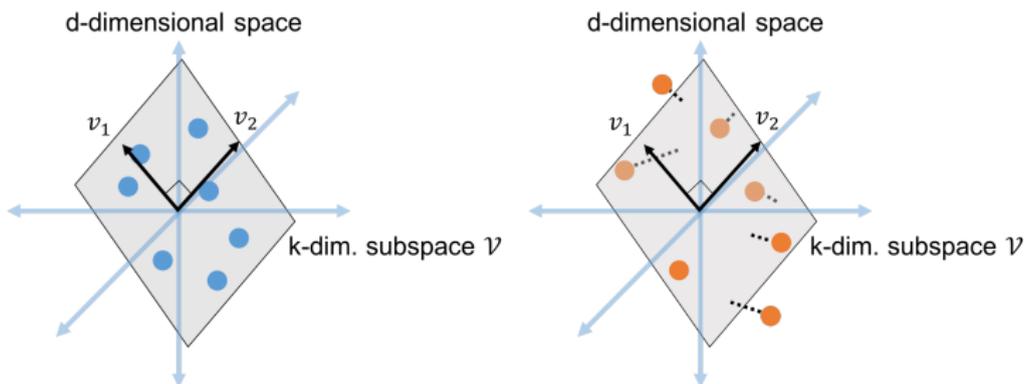
Let $\vec{v}_1, \dots, \vec{v}_k$ be an orthonormal basis for \mathcal{V} and $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns.

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

Letting $\tilde{x}_i = \mathbf{V}^T \vec{x}_i$, we have a perfect embedding from \mathcal{V} into \mathbb{R}^k .

EMBEDDING WITH ASSUMPTIONS

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



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

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

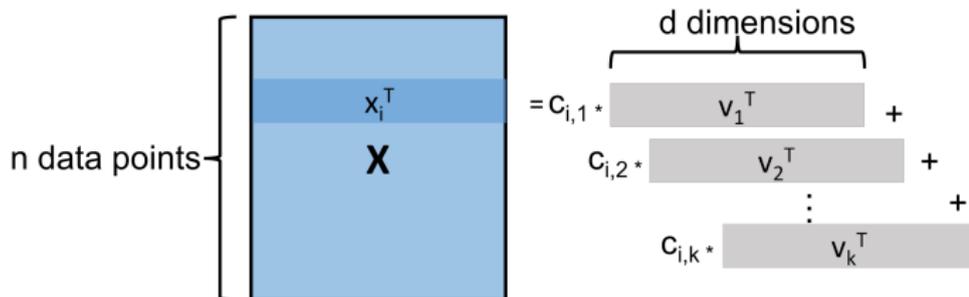
LOW-RANK FACTORIZATION

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

- Letting $\vec{v}_1, \dots, \vec{v}_k$ be an orthonormal basis for \mathcal{V} , can write 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 $\text{rank}(\mathbf{X}) \leq k$.

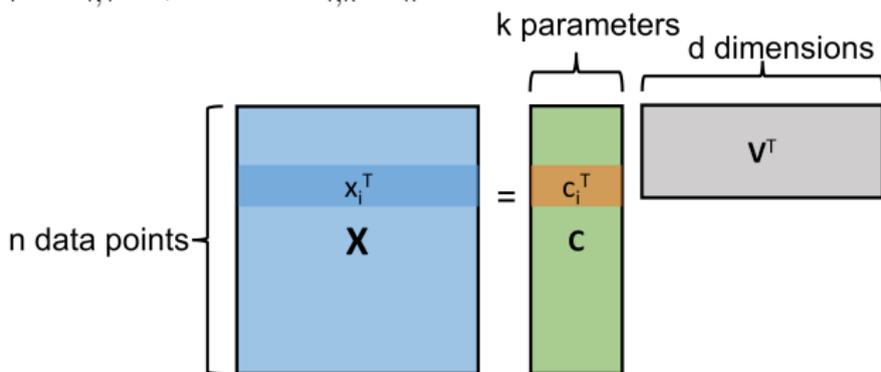


$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

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

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

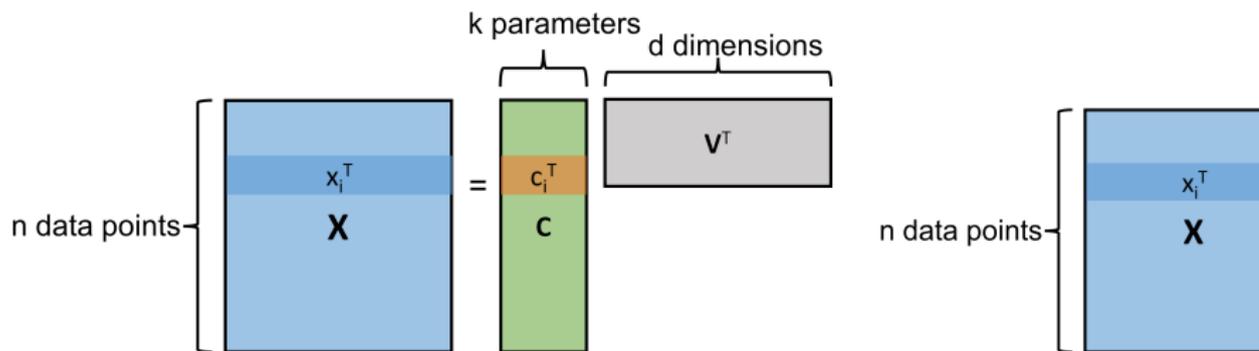


- \mathbf{X} can be represented by $(n + d) \cdot k$ parameters vs. $n \cdot d$.
- The rows of \mathbf{X} are spanned by k vectors: the columns of $\mathbf{V} \implies$ the columns of \mathbf{X} are spanned by k vectors: the columns of \mathbf{C} .

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

LOW-RANK FACTORIZATION

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



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

- $\mathbf{X} = \mathbf{C}\mathbf{V}^T \implies \mathbf{X}\mathbf{V} = \mathbf{C}\mathbf{V}^T\mathbf{V}$
- $\mathbf{V}^T\mathbf{V} = \mathbf{I}$, the identity (since \mathbf{V} is orthonormal) $\implies \mathbf{X}\mathbf{V} = \mathbf{C}$.

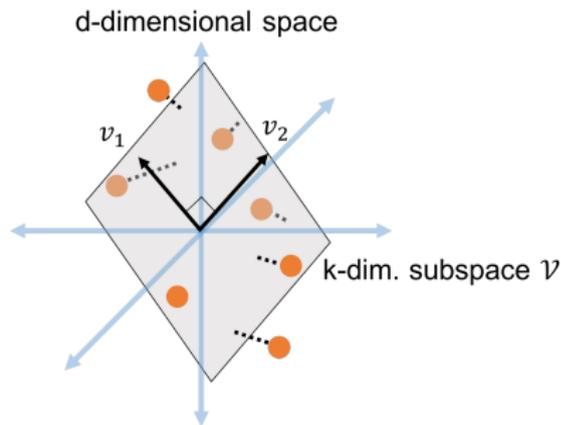
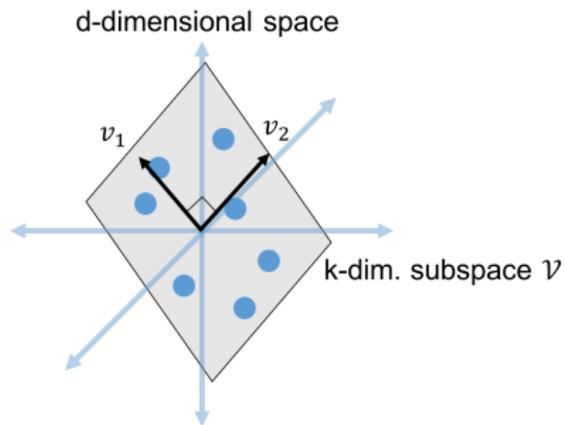
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} , $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

PROJECTION VIEW

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

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

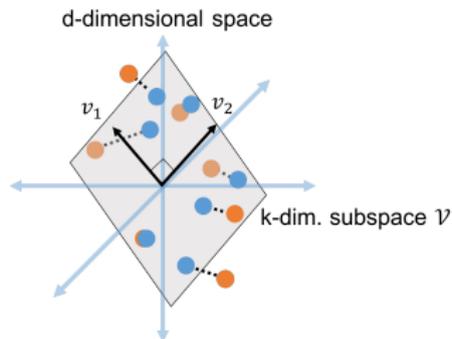
- $\mathbf{V}\mathbf{V}^T$ is a **projection matrix**, which projects the rows of \mathbf{X} (the data points $\vec{x}_1, \dots, \vec{x}_n$) onto the subspace \mathcal{V} .



LOW-RANK APPROXIMATION

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

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



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

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

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

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:

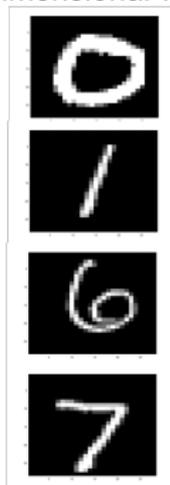
$$\|\vec{y}\|_2^2 = \|(\mathbf{V}\mathbf{V}^T)\vec{y}\|_2^2 + \|\vec{y} - (\mathbf{V}\mathbf{V}^T)\vec{y}\|_2^2.$$

A STEP BACK: WHY LOW-RANK APPROXIMATION?

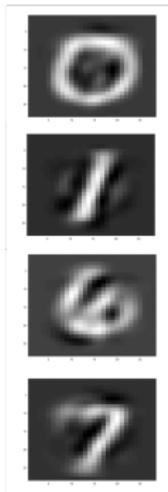
Question: Why might we expect $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ to lie close to a k -dimensional subspace?

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

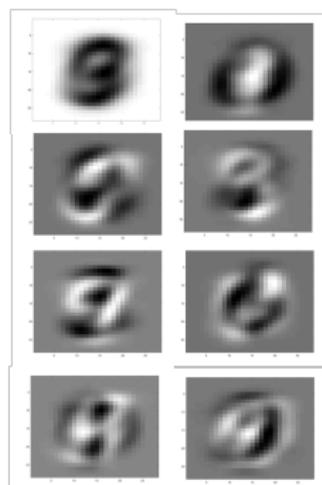
784 dimensional vectors



projections onto 15 dimensional space



orthonormal basis v_1, \dots, v_{15}



DUAL VIEW OF LOW-RANK APPROXIMATION

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

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

Linearly Dependent Variables:

	bedrooms	bathrooms	sq.ft.	floors	list price	sale price
home 1	2	2	1800	2	200,000	195,000
home 2	4	2.5	2700	1	300,000	310,000
.
.
.
home n	5	3.5	3600	3	450,000	450,000

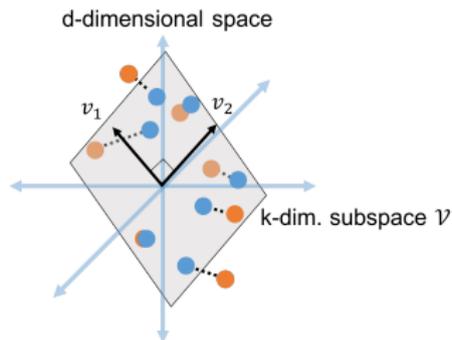
	bedrooms
home 1	2
home 2	4
.	.
.	.
.	.
home n	5 ¹³

BEST FIT SUBSPACE

If $\vec{x}_1, \dots, \vec{x}_n$ are 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{XV}^T . \mathbf{XV} gives optimal embedding of \mathbf{X} in \mathcal{V} .

How do we find \mathcal{V} (equivalently \mathbf{V})?

$$\arg \min_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X} - \mathbf{XV}^T\|_F^2 = \sum_{i,j} (\mathbf{X}_{i,j} - (\mathbf{XV}^T)_{i,j})^2 = \sum_{i=1}^n \|\vec{x}_i - \mathbf{V}^T \vec{x}_i\|_2^2$$



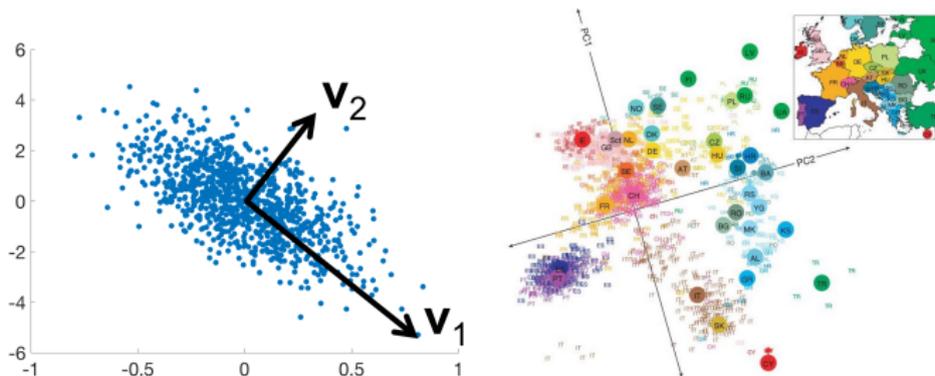
$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$: data points, $\mathbf{X} \in \mathbb{R}^{n \times d}$: data matrix, $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$: orthogonal basis for subspace \mathcal{V} . $\mathbf{V} \in \mathbb{R}^{d \times k}$: matrix with columns $\vec{v}_1, \dots, \vec{v}_k$.

BEST FIT SUBSPACE

\mathbf{V} minimizing $\|\mathbf{X} - \mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2$ is given by:

$$\arg \max_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X}\mathbf{V}\mathbf{V}^T\|_F^2 = \sum_{i=1}^n \|\mathbf{V}\mathbf{V}^T \vec{x}_i\|_2^2 \quad \arg \max_{\text{orthonormal } \mathbf{V} \in \mathbb{R}^{d \times k}} \|\mathbf{X}\mathbf{V}\|_F^2 = \sum_{i=1}^n \|\mathbf{V}^T \vec{x}_i\|_2^2 =$$

Columns of \mathbf{V} are 'directions of greatest variance' in the data.



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

- Many datasets lie close to a k -dimensional subspace.
- Can take advantage of this to do data-dependent linear dimensionality reduction (low-rank approximation).
- Dual view: both rows (data points) and columns (features) are approximated spanned by a small number of vectors.

- **Step 1:** Find this subspace by finding the directions of greatest variance in the data.
- **Step 2:** Get best approximation to the data points in this subspace via **projection** matrix $\mathbf{V}\mathbf{V}^T$. $\mathbf{V} \in \mathbb{R}^{d \times k}$ used as linear mapping from d -dimensional to k -dimensional space.