

# PCA

CS 446

2020-12-27 (a50ee91)

# Plan for today

- ▶ **Unsupervised learning?**
- ▶ PCA.
- ▶ PCA example.
- ▶ PCA proofs.

# Supervised learning

So far, we've done **supervised learning**:

Given  $((\mathbf{x}_i, \mathbf{y}_i))_{i=1}^n$ , find predictor  $f$  with  $f(\mathbf{x}_i) \approx \mathbf{y}_i$ .

*Linear regression, deep networks, k-nn, decision trees, ...*

# Supervised learning

So far, we've done **supervised learning**:

Given  $((\mathbf{x}_i, \mathbf{y}_i))_{i=1}^n$ , find predictor  $f$  with  $f(\mathbf{x}_i) \approx \mathbf{y}_i$ .

*Linear regression, deep networks, k-nn, decision trees, ...*

Most methods used **(regularized) ERM**:

minimize  $\hat{\mathcal{R}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), \mathbf{y}_i)$ , hope  $\mathcal{R}$  is small.

*least squares, logistic regression, deep networks, SVM, perceptron, ...*

# Unsupervised learning

Now we only receive  $(\mathbf{x}_i)_{i=1}^n$ , and the goal is...?

# Unsupervised learning

Now we only receive  $(\mathbf{x}_i)_{i=1}^n$ , and the goal is...?

- ▶ Encoding data in some compact representation (and decoding this).
- ▶ Data analysis; recovering “hidden structure” in data (e.g., recovering cliques or clusters).
- ▶ Features for supervised learning.
- ▶ ...?

# Unsupervised learning

Now we only receive  $(\mathbf{x}_i)_{i=1}^n$ , and the goal is...?

- ▶ Encoding data in some compact representation (and decoding this).
- ▶ Data analysis; recovering “hidden structure” in data (e.g., recovering cliques or clusters).
- ▶ Features for supervised learning.
- ▶ ...?

The task is less clear-cut, and lacks a single accepted formalization.

(Side note: can *still* use the “pytorch meta-algorithm”.)

# Principal Component Analysis (PCA) motivation



# Principal Component Analysis (PCA) motivation

Let's formulate a *simplistic linear unsupervised method*.

- ▶ Encoding (and decoding) data in some compact representation.  
Let's linearly map data in  $\mathbb{R}^d$  to  $\mathbb{R}^k$  and back.
- ▶ Data analysis;  
recovering "hidden structure" in data.  
Let's find if data mostly lies on a low-dimensional subspace.
- ▶ Features for supervised learning.  
Let's feed the  $\mathbb{R}^k$ -dimensional encoding to supervised methods.

Let  $M \in \mathbb{R}^{n \times d}$  be given.  $((s_i, \mathbf{u}_i, \mathbf{v}_i))_{i=1}^r$  is an **SVD of  $M$**  if:

- ▶  $M$  has rank  $r$ ;
- ▶  $s_1 \geq s_2 \geq \dots \geq s_r > 0$ ;
- ▶  $(\mathbf{u}_i)_{i=1}^r$  are orthonormal (orthogonal and unit length), and span the column space of  $M$ ;
- ▶  $(\mathbf{v}_i)_{i=1}^r$  are orthonormal, and span the row space of  $M$ .
- ▶  $M = \sum_i s_i \mathbf{u}_i \mathbf{v}_i^\top$ .

Let's also collect these into matrices:  $S := \text{diag}(s_1, \dots, s_r) \in \mathbb{R}^{r \times r}$ , and

$$U := \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{u}_1 & \cdots & \mathbf{u}_r \\ \downarrow & & \downarrow \end{bmatrix} \in \mathbb{R}^{n \times r}, \quad V := \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{v}_1 & \cdots & \mathbf{v}_r \\ \downarrow & & \downarrow \end{bmatrix} \in \mathbb{R}^{d \times r},$$

whereby  $M = USV^\top$ . Additionally define **best rank- $k$  approximations**  
 $S_k := \text{diag}(s_1, \dots, s_k) \in \mathbb{R}^{k \times k}$ ,

$$U_k := \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{u}_1 & \cdots & \mathbf{u}_k \\ \downarrow & & \downarrow \end{bmatrix} \in \mathbb{R}^{n \times k}, \quad V_k := \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{v}_1 & \cdots & \mathbf{v}_k \\ \downarrow & & \downarrow \end{bmatrix} \in \mathbb{R}^{d \times k},$$

and  $M_k := U_k S_k V_k^\top$ . (**Side note:** not unique!)

# PCA (Principal component analysis)

**Input:** Data as rows of  $\mathbb{R}^{n \times d} \ni \mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ , integer  $k$ .

**Output:** Encoder  $\mathbf{V}_k$ , decoder  $\mathbf{V}_k^T$ , encoded data  $\mathbf{X}\mathbf{V}_k = \mathbf{U}_k\mathbf{S}_k$ .

# PCA (Principal component analysis)

**Input:** Data as rows of  $\mathbb{R}^{n \times d} \ni \mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ , integer  $k$ .

**Output:** Encoder  $\mathbf{V}_k$ , decoder  $\mathbf{V}_k^\top$ , encoded data  $\mathbf{X}\mathbf{V}_k = \mathbf{U}_k\mathbf{S}_k$ .

The goal in unsupervised learning is unclear.

We'll try to define this as "best encoding/decoding in Frobenius sense":

$$\min_{\substack{\mathbf{D} \in \mathbb{R}^{k \times d} \\ \mathbf{E} \in \mathbb{R}^{d \times k}}} \|\mathbf{X} - \mathbf{X}\mathbf{E}\mathbf{D}\|_F^2 = \left\| \mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^\top \right\|_F^2 = \left\| \mathbf{X} - \mathbf{U}_k\mathbf{S}_k\mathbf{V}_k^\top \right\|_F^2 = \|\mathbf{X} - \mathbf{X}_k\|_F^2.$$

# PCA (Principal component analysis)

**Input:** Data as rows of  $\mathbb{R}^{n \times d} \ni \mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ , integer  $k$ .

**Output:** Encoder  $\mathbf{V}_k$ , decoder  $\mathbf{V}_k^\top$ , encoded data  $\mathbf{X}\mathbf{V}_k = \mathbf{U}_k\mathbf{S}_k$ .

The goal in unsupervised learning is unclear.

We'll try to define this as “best encoding/decoding in Frobenius sense”:

$$\min_{\substack{\mathbf{D} \in \mathbb{R}^{k \times d} \\ \mathbf{E} \in \mathbb{R}^{d \times k}} \|\mathbf{X} - \mathbf{X}\mathbf{E}\mathbf{D}\|_F^2 = \left\| \mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^\top \right\|_F^2 = \left\| \mathbf{X} - \mathbf{U}_k\mathbf{S}_k\mathbf{V}_k^\top \right\|_F^2 = \|\mathbf{X} - \mathbf{X}_k\|_F^2.$$

Note  $\mathbf{V}_k\mathbf{V}_k^\top$  performs orthogonal projection onto subspace spanned by  $\mathbf{V}_k$ ; thus we are finding “best  $k$ -dimensional projection of the data”.

# PCA properties

Let's try to capture “best low-rank approximation” and “best linear encoder/decoder”.

# PCA properties

Let's try to capture “best low-rank approximation” and “best linear encoder/decoder”.

**Theorem.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  with SVD  $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$  and integer  $k \leq r$  be given.

$$\begin{aligned} \min_{\substack{D \in \mathbb{R}^{k \times d} \\ E \in \mathbb{R}^{d \times k}}} \|\mathbf{X} - \mathbf{XED}\|_F^2 &= \min_{\substack{D \in \mathbb{R}^{d \times k} \\ D^\top D = I}} \|\mathbf{X} - \mathbf{XDD}^\top\|_F^2 \\ &= \|\mathbf{X} - \mathbf{XV}_k\mathbf{V}_k^\top\|_F^2 = \sum_{i=k+1}^r s_i^2. \end{aligned}$$

Additionally,

$$\begin{aligned} \min_{\substack{D \in \mathbb{R}^{d \times k} \\ D^\top D = I}} \|\mathbf{X} - \mathbf{XDD}^\top\|_F^2 &= \|\mathbf{X}\|_F^2 - \max_{\substack{D \in \mathbb{R}^{d \times k} \\ D^\top D = I}} \|\mathbf{XD}\|_F^2 \\ &= \|\mathbf{X}\|_F^2 - \|\mathbf{XV}_k\|_F^2 = \|\mathbf{X}\|_F^2 - \sum_{i=1}^k s_i^2. \end{aligned}$$

**Remark 1.** SVD is not unique, but  $\sum_{i=1}^r s_i^2$  is identical across SVD choices.

# PCA properties

Let's try to capture “best low-rank approximation” and “best linear encoder/decoder”.

**Theorem.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  with SVD  $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$  and integer  $k \leq r$  be given.

$$\begin{aligned} \min_{\substack{\mathbf{D} \in \mathbb{R}^{k \times d} \\ \mathbf{E} \in \mathbb{R}^{d \times k}}} \|\mathbf{X} - \mathbf{X}\mathbf{E}\mathbf{D}\|_F^2 &= \min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X} - \mathbf{X}\mathbf{D}\mathbf{D}^\top\|_F^2 \\ &= \|\mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^\top\|_F^2 = \sum_{i=k+1}^r s_i^2. \end{aligned}$$

Additionally,

$$\begin{aligned} \min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X} - \mathbf{X}\mathbf{D}\mathbf{D}^\top\|_F^2 &= \|\mathbf{X}\|_F^2 - \max_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X}\mathbf{D}\|_F^2 \\ &= \|\mathbf{X}\|_F^2 - \|\mathbf{X}\mathbf{V}_k\|_F^2 = \|\mathbf{X}\|_F^2 - \sum_{i=1}^k s_i^2. \end{aligned}$$

**Remark 1.** SVD is not unique, but  $\sum_{i=1}^r s_i^2$  is identical across SVD choices.

**Remark 2.** As written, this is not a convex optimization problem!



# PCA properties

Let's try to capture “best low-rank approximation” and “best linear encoder/decoder”.

**Theorem.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  with SVD  $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$  and integer  $k \leq r$  be given.

$$\begin{aligned} \min_{\substack{\mathbf{D} \in \mathbb{R}^{k \times d} \\ \mathbf{E} \in \mathbb{R}^{d \times k}}} \|\mathbf{X} - \mathbf{X}\mathbf{E}\mathbf{D}\|_F^2 &= \min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X} - \mathbf{X}\mathbf{D}\mathbf{D}^\top\|_F^2 \\ &= \|\mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^\top\|_F^2 = \sum_{i=k+1}^r s_i^2. \end{aligned}$$

Additionally,

$$\begin{aligned} \min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X} - \mathbf{X}\mathbf{D}\mathbf{D}^\top\|_F^2 &= \|\mathbf{X}\|_F^2 - \max_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X}\mathbf{D}\|_F^2 \\ &= \|\mathbf{X}\|_F^2 - \|\mathbf{X}\mathbf{V}_k\|_F^2 = \|\mathbf{X}\|_F^2 - \sum_{i=1}^k s_i^2. \end{aligned}$$

**Remark 1.** SVD is not unique, but  $\sum_{i=1}^r s_i^2$  is identical across SVD choices.

**Remark 2.** As written, this is not a convex optimization problem!

**Remark 3.** The second form is interesting...

# PCA properties

Let's try to capture “best low-rank approximation” and “best linear encoder/decoder”.

**Theorem.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  with SVD  $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$  and integer  $k \leq r$  be given.

$$\begin{aligned} \min_{\substack{\mathbf{D} \in \mathbb{R}^{k \times d} \\ \mathbf{E} \in \mathbb{R}^{d \times k}}} \|\mathbf{X} - \mathbf{X}\mathbf{E}\mathbf{D}\|_F^2 &= \min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X} - \mathbf{X}\mathbf{D}\mathbf{D}^\top\|_F^2 \\ &= \|\mathbf{X} - \mathbf{X}\mathbf{V}_k\mathbf{V}_k^\top\|_F^2 = \sum_{i=k+1}^r s_i^2. \end{aligned}$$

Additionally,

$$\begin{aligned} \min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X} - \mathbf{X}\mathbf{D}\mathbf{D}^\top\|_F^2 &= \|\mathbf{X}\|_F^2 - \max_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X}\mathbf{D}\|_F^2 \\ &= \|\mathbf{X}\|_F^2 - \|\mathbf{X}\mathbf{V}_k\|_F^2 = \|\mathbf{X}\|_F^2 - \sum_{i=1}^k s_i^2. \end{aligned}$$

**Remark 1.** SVD is not unique, but  $\sum_{i=1}^r s_i^2$  is identical across SVD choices.

**Remark 2.** As written, this is not a convex optimization problem!

**Remark 3.** The second form is interesting. . .

**Remark 4.** Easy to show “best rank- $k$  approximation” from here.

## Centered PCA

Some treatments replace  $\mathbf{X}$  with  $\mathbf{X} - \mathbf{1}\boldsymbol{\mu}^\top$ ,  
with mean  $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1} \mathbf{x}_i$ .

## Centered PCA

Some treatments replace  $\mathbf{X}$  with  $\mathbf{X} - \mathbf{1}\boldsymbol{\mu}^\top$ ,  
with mean  $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1} \mathbf{x}_i$ .

$\frac{1}{n} \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{d \times d}$  is data covariance;

# Centered PCA

Some treatments replace  $\mathbf{X}$  with  $\mathbf{X} - \mathbf{1}\boldsymbol{\mu}^\top$ ,  
with mean  $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1} \mathbf{x}_i$ .

$\frac{1}{n} \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{d \times d}$  is data covariance;

$\frac{1}{n} (\mathbf{X}\mathbf{D})^\top (\mathbf{X}\mathbf{D})$  is data covariance after projection;

# Centered PCA

Some treatments replace  $\mathbf{X}$  with  $\mathbf{X} - \mathbf{1}\boldsymbol{\mu}^\top$ ,  
with mean  $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1} \mathbf{x}_i$ .

$\frac{1}{n} \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{d \times d}$  is data covariance;

$\frac{1}{n} (\mathbf{X}\mathbf{D})^\top (\mathbf{X}\mathbf{D})$  is data covariance after projection;  
lastly

$$\frac{1}{n} \|\mathbf{X}\mathbf{D}\|_F^2 = \frac{1}{n} \operatorname{tr} \left( (\mathbf{X}\mathbf{D})^\top (\mathbf{X}\mathbf{D}) \right) = \frac{1}{n} \sum_{i=1}^k (\mathbf{X}\mathbf{D}\mathbf{e}_i)^\top (\mathbf{X}\mathbf{D}\mathbf{e}_i).$$

Since PCA is also solving  $\max_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X}\mathbf{D}\|_F^2$ , therefore centered PCA is  
maximizing the resulting per-coordinate variances!

# PCA example

- ▶ Image data; e.g., “eigenfaces” .

- ▶ Image data; e.g., “eigenfaces” .

Weirdness: negative faces?

This motivates *non-negative matrix factorization*.



# PCA example

- ▶ Image data; e.g., “eigenfaces” .  
Weirdness: negative faces?  
This motivates *non-negative matrix factorization*.
- ▶ *LSI (Latent Semantic Indexing)*:  
collect many documents into  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ,  
where  $x_i$  is a normalized bag-of-words vector (plus nonlinear mappings).  
Can interpret new representation as weighting over “topics” .

# Application: digit data

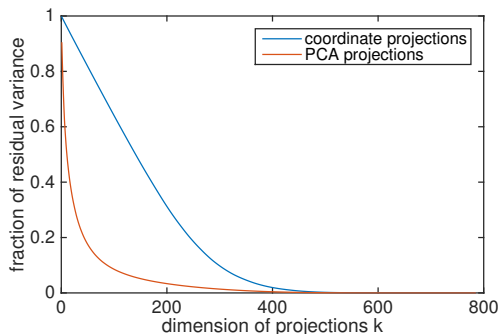
Data  $(\mathbf{x}_i)_{i=1}^n$  with  $\mathbf{x}_i \in \mathbb{R}^{784}$ .

- ▶ Residual variance left by rank- $k$  PCA projection:

$$1 - \frac{\sum_{j=1}^k \text{variance in direction } \mathbf{v}_j}{\text{total variance}} = 1 - \frac{\|\mathbf{X}\mathbf{V}_k\|_F^2}{\|\mathbf{X}\|_F^2}.$$

- ▶ Residual variance left by best  $k$  coordinate projections:

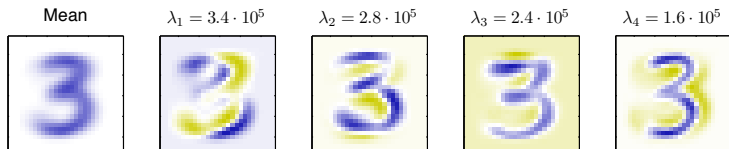
$$1 - \frac{\sum_{j=1}^k \text{variance in direction } \mathbf{e}_j}{\text{total variance}} = 1 - \frac{\sum_{j=1}^k (\mathbf{X}\mathbf{e}_j)^\top (\mathbf{X}\mathbf{e}_j)}{\|\mathbf{X}\|_F^2}.$$



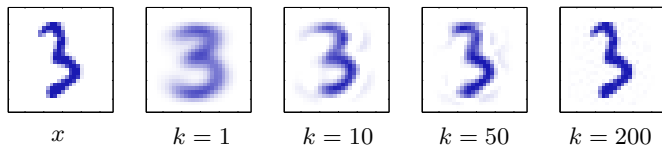
# Application: digit data

$16 \times 16$  pixel images of handwritten 3s (as vectors in  $\mathbb{R}^{256}$ )

**Mean  $\mu$  and right singular vectors  $v_1, v_2, v_3, v_4$**



**Reconstructions:**



Only have to store  $k$  numbers per image,  
along with the mean  $\mu$  and  $k$  eigenvectors ( $256(k + 1)$  numbers).

# Proof of PCA theorem.

**Fact.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and  $k \leq r$  be given.

$$\min_{\substack{\mathbf{M} \in \mathbb{R}^{d \times d} \\ \text{rank}(\mathbf{M})=k}} \|\mathbf{X} - \mathbf{X}\mathbf{M}\|_F^2 = \min_{\substack{\mathbf{D} \in \mathbb{R}^{k \times d} \\ \mathbf{E} \in \mathbb{R}^{d \times k}}} \|\mathbf{X} - \mathbf{X}\mathbf{E}\mathbf{D}\|_F^2 = \min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X} - \mathbf{X}\mathbf{D}\mathbf{D}^\top\|_F^2.$$

# Proof of PCA theorem.

**Fact.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and  $k \leq r$  be given.

$$\min_{\substack{\mathbf{M} \in \mathbb{R}^{d \times d} \\ \text{rank}(\mathbf{M})=k}} \|\mathbf{X} - \mathbf{X}\mathbf{M}\|_{\text{F}}^2 = \min_{\substack{\mathbf{D} \in \mathbb{R}^{k \times d} \\ \mathbf{E} \in \mathbb{R}^{d \times k}}} \|\mathbf{X} - \mathbf{X}\mathbf{E}\mathbf{D}\|_{\text{F}}^2 = \min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^{\text{T}}\mathbf{D}=\mathbf{I}}} \|\mathbf{X} - \mathbf{X}\mathbf{D}\mathbf{D}^{\text{T}}\|_{\text{F}}^2.$$

**Proof.** Since

$$\begin{aligned} \left\{ \mathbf{M} \in \mathbb{R}^{d \times d} : \text{rank}(\mathbf{M}) = k \right\} &\supseteq \left\{ \mathbf{D}\mathbf{E} : \mathbf{D} \in \mathbb{R}^{k \times d}, \mathbf{E} \in \mathbb{R}^{d \times k} \right\} \\ &\supseteq \left\{ \mathbf{D}\mathbf{D}^{\text{T}} : \mathbf{D} \in \mathbb{R}^{d \times k}, \mathbf{D}^{\text{T}}\mathbf{D} = \mathbf{I} \right\}, \end{aligned}$$

it suffices to show

$$\min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^{\text{T}}\mathbf{D}=\mathbf{I}}} \|\mathbf{X} - \mathbf{X}\mathbf{D}\mathbf{D}^{\text{T}}\|_{\text{F}}^2 \leq \min_{\substack{\mathbf{M} \in \mathbb{R}^{d \times d} \\ \text{rank}(\mathbf{M})=k}} \|\mathbf{X} - \mathbf{X}\mathbf{M}\|_{\text{F}}^2.$$

**Proof (continued).** For any  $M = USV^T \in \mathbb{R}^{d \times d}$  with  $\text{rank}(M) \leq k$  (whereby  $M = U_k S_k V_k^T$ ),

$$\begin{aligned}\|X - XM\|_F^2 &= \left\| X - XV_k V_k^T + XV_k V_k^T - XM \right\|_F^2 \\ &= \left\| X - XV_k V_k^T \right\|_F^2 + 2 \text{tr} \left( \left( X - XV_k V_k^T \right)^T \left( XV_k V_k^T - XM \right) \right) \\ &\quad + \left\| XV_k V_k^T - XM \right\|_F^2.\end{aligned}$$

We'll show the middle (trace) term is 0, and therefore

$$\|X - M\|_F^2 = \left\| X - XV_k V_k^T \right\|_F^2 + \left\| XV_k V_k^T - XM \right\|_F^2 \geq \left\| X - XV_k V_k^T \right\|_F^2.$$

**Proof (continued).** Note

$$\begin{aligned} & \text{tr} \left( \left( \mathbf{X} - \mathbf{XV}_k \mathbf{V}_k^\top \right)^\top \left( \mathbf{XV}_k \mathbf{V}_k^\top - \mathbf{XM} \right) \right) \\ &= \text{tr} \left( \left( \mathbf{I} - \mathbf{V}_k \mathbf{V}_k^\top \right)^\top \mathbf{X}^\top \mathbf{X} \left( \mathbf{X} - \mathbf{XU}_k \mathbf{S}_k \mathbf{V}_k^\top \right) \mathbf{V}_k \mathbf{V}_k^\top \right) \\ &= \text{tr} \left( \mathbf{X}^\top \mathbf{X} \left( \mathbf{X} - \mathbf{XU}_k \mathbf{S}_k \mathbf{V}_k^\top \right) \mathbf{V}_k \mathbf{V}_k^\top \left( \mathbf{I} - \mathbf{V}_k \mathbf{V}_k^\top \right)^\top \right), \end{aligned}$$

and

$$\begin{aligned} \left( \mathbf{I} - \mathbf{V}_k \mathbf{V}_k^\top \right)^\top \left( \mathbf{V}_k \mathbf{V}_k^\top \right) &= \left( \sum_{i=1}^d \mathbf{v}_i \mathbf{v}_i^\top - \sum_{i=1}^k \mathbf{v}_i \mathbf{v}_i^\top \right)^\top \sum_{j=1}^k \mathbf{v}_j \mathbf{v}_j^\top \\ &= \left( \sum_{i=k+1}^d \mathbf{v}_i \mathbf{v}_i^\top \right) \sum_{j=1}^k \mathbf{v}_j \mathbf{v}_j^\top = 0. \end{aligned}$$

Therefore

$$\text{tr} \left( \left( \mathbf{X} - \mathbf{XV}_k \mathbf{V}_k^\top \right)^\top \left( \mathbf{XV}_k \mathbf{V}_k^\top - \mathbf{XM} \right) \right) = 0.$$

□

**Fact.** Let  $X \in \mathbb{R}^{n \times d}$  be given along with  $D \in \mathbb{R}^{d \times k}$  with  $D^T D = I$ . Then  $\|X - XDD^T\|_F^2 = \|X\|_F^2 - \|XD\|_F^2$ , and

$$\min_{\substack{D \in \mathbb{R}^{d \times k} \\ D^T D = I}} \|X - XDD^T\|_F^2 = \|X\|_F^2 - \max_{\substack{D \in \mathbb{R}^{d \times k} \\ D^T D = I}} \|XD\|_F^2,$$

$$\arg \min_{\substack{D \in \mathbb{R}^{d \times k} \\ D^T D = I}} \|X - XDD^T\|_F^2 = \arg \max_{\substack{D \in \mathbb{R}^{d \times k} \\ D^T D = I}} \|XD\|_F^2.$$



**Fact.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be given along with  $\mathbf{D} \in \mathbb{R}^{d \times k}$  with  $\mathbf{D}^\top \mathbf{D} = \mathbf{I}$ . Then  $\|\mathbf{X} - \mathbf{XDD}^\top\|_F^2 = \|\mathbf{X}\|_F^2 - \|\mathbf{XD}\|_F^2$ , and

$$\min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X} - \mathbf{XDD}^\top\|_F^2 = \|\mathbf{X}\|_F^2 - \max_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{XD}\|_F^2,$$

$$\arg \min_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X} - \mathbf{XDD}^\top\|_F^2 = \arg \max_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{XD}\|_F^2.$$

**Proof.** Since

$$\begin{aligned} \|\mathbf{XDD}^\top\|_F^2 &= \text{tr} \left( (\mathbf{XDD}^\top)^\top (\mathbf{XDD}^\top) \right) = \text{tr} \left( (\mathbf{XD})^\top (\mathbf{XDD}^\top \mathbf{D}) \right) \\ &= \text{tr} \left( (\mathbf{XD})^\top (\mathbf{XD}) \right) = \|\mathbf{XD}\|_F^2, \end{aligned}$$

therefore

$$\begin{aligned} \|\mathbf{X} - \mathbf{XDD}^\top\|_F^2 &= \|\mathbf{X}\|_F^2 - 2 \text{tr} \left( (\mathbf{XDD}^\top)^\top \mathbf{X} \right) + \|\mathbf{XDD}^\top\|_F^2 \\ &= \|\mathbf{X}\|_F^2 - \|\mathbf{XD}\|_F^2. \end{aligned}$$

□

**Fact.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be given with SVD  $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ . Then

$$\max_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X}\mathbf{D}\|_F^2 = \|\mathbf{X}\mathbf{V}_k\|_F^2 = \sum_{i=1}^k s_i^2.$$

**Fact.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be given with SVD  $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ . Then

$$\max_{\substack{\mathbf{D} \in \mathbb{R}^{d \times k} \\ \mathbf{D}^\top \mathbf{D} = \mathbf{I}}} \|\mathbf{X}\mathbf{D}\|_F^2 = \|\mathbf{X}\mathbf{V}_k\|_F^2 = \sum_{i=1}^k s_i^2.$$

**Proof.** Define

$$S_1 := \{\mathbf{D} \in \mathbb{R}^{d \times k} : \mathbf{D}^\top \mathbf{D} = \mathbf{I}\}, \quad S_2 := \{\mathbf{V}\mathbf{D} : \mathbf{D} \in S_1\}.$$

Note  $S_1 = S_2$ :

- ▶  $S_1 \subseteq S_2$ , since  $\mathbf{D} \in S_1$  implies  $(\mathbf{V}^\top \mathbf{D})^\top \mathbf{V}^\top \mathbf{D} = \mathbf{I}$ , thus  $\mathbf{D} = \mathbf{V}(\mathbf{V}^\top \mathbf{D}) \in S_2$ .
- ▶  $S_2 \subseteq S_1$ , since  $\mathbf{V}\mathbf{D} \in S_2$  implies  $(\mathbf{V}\mathbf{D})^\top (\mathbf{V}\mathbf{D}) = \mathbf{I}$  thus  $\mathbf{V}\mathbf{D} \in S_1$ .

Therefore

$$\begin{aligned} \max_{\mathbf{D} \in S_1} \|\mathbf{X}\mathbf{D}\|_F^2 &= \max_{\mathbf{M} \in S_2} \|\mathbf{X}\mathbf{M}\|_F^2 = \max_{\mathbf{D} \in S_1} \|\mathbf{X}\mathbf{V}\mathbf{D}\|_F^2 = \max_{\mathbf{D} \in S_1} \|\mathbf{U}\mathbf{S}\mathbf{V}^\top \mathbf{V}\mathbf{D}\|_F^2 \\ &= \max_{\mathbf{D} \in S_1} \text{tr} \left( (\mathbf{U}\mathbf{S}\mathbf{D})^\top (\mathbf{U}\mathbf{S}\mathbf{D}) \right) = \max_{\mathbf{D} \in S_1} \text{tr} \left( \mathbf{D}\mathbf{D}^\top \mathbf{S}^\top \mathbf{S} \right) \\ &= \max_{\mathbf{D} \in S_1} \sum_{j=1}^r s_j^2 \sum_{i=1}^k D_{ij}^2. \end{aligned}$$

**Proof (continued).** We've reduced the proof to showing

$$\max_{\substack{D \in \mathbb{R}^{d \times k} \\ D^\top D = I}} \sum_{j=1}^r s_j^2 \sum_{i=1}^k D_{ij}^2 = \|\mathbf{XV}_k\|_F^2,$$

and note moreover  $\|\mathbf{XV}_k\|_F = \text{tr} \left( (\mathbf{U}_k \mathbf{S}_k)^\top (\mathbf{U}_k \mathbf{S}_k) \right) = \sum_{i=1}^k s_i^2$ . Lastly:

- ▶ Since  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$  and  $\mathbf{V}_k^\top \mathbf{V}_k = \mathbf{I}$ ,

$$\max_{\substack{D \in \mathbb{R}^{d \times k} \\ D^\top D = I}} \sum_{j=1}^r s_j^2 \sum_{i=1}^k D_{ij}^2 \geq \|\mathbf{XV}_k\|_F^2.$$

- ▶ For any feasible  $D \in \mathbb{R}^{d \times k}$ , extend it to orthonormal  $M \in \mathbb{R}^{d \times d}$ ; since  $M^\top M = I = M M^\top$ , then  $M^\top$  is orthonormal as well, and  $\sum_{i=1}^k D_{ij}^2 \leq \sum_{i=1}^d M_{ij}^2 = 1$ . Moreover,  $\sum_{i,j} D_{ij}^2 \leq k$ , so

$$\max_{\substack{D \in \mathbb{R}^{d \times k} \\ D^\top D = I}} \sum_{j=1}^r s_j^2 \left( \sum_{i=1}^k D_{ij}^2 \right) \leq \max_{\substack{\mathbf{w} \in [0,1]^d \\ \sum_i w_i \leq k}} \sum_{j=1}^r s_j^2 w_j^2 \leq \sum_{j=1}^k s_j^2.$$

□

# Summary for today

- ▶ **Unsupervised learning?**
- ▶ PCA.
- ▶ PCA example.
- ▶ PCA proofs.