

Unsupervised Learning Methods

Problem Set II –

PCA and KPCA

Due: 24.05.2021

Guidelines

- Answer all questions (PDF + Jupyter notebook).
- You must type your solution manual (handwriting is not allowed).
- Submission in pairs (use the forum if needed).
- You **may** submit the entire solution in a single ipynb file (or in PDF + ipynb files).
- You **may** (and should) use the forum if you have any questions.
- Good luck!

1 PCA

1.1 Eigendecomposition

Trace

- Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ be a diagonalizable matrix, that is, $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{-1}$ where $\mathbf{\Lambda}$ is a diagonal matrix.

1.1.1

Prove that:

$$\text{Tr} \{ \mathbf{A} \} = \sum_{i=1}^d \lambda_i(\mathbf{A})$$

where $\lambda_i(\mathbf{A}) = \mathbf{\Lambda}[i, i]$ is the i th eigenvalue of \mathbf{A} .

Similarity

- Two (square) matrices $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{B} \in \mathbb{R}^{d \times d}$ are called similar, namely, $\mathbf{A} \sim \mathbf{B}$, if exists an (invertible) matrix $\mathbf{P} \in \mathbb{R}^{d \times d}$ such that:

$$\mathbf{B} = \mathbf{P}\mathbf{A}\mathbf{P}^{-1}$$

1.1.2

Prove that if \mathbf{A} is diagonalizable and $\mathbf{A} \sim \mathbf{B}$, then, \mathbf{A} and \mathbf{B} share the same set of eigenvalues, namely:

$$\mathbf{A} \sim \mathbf{B} \implies \{\lambda_i(\mathbf{A})\}_{i=1}^d = \{\lambda_i(\mathbf{B})\}_{i=1}^d$$

SPD matrices

A symmetric matrix $\mathbf{A} = \mathbf{A}^T$ is an Symmetric Positive Definite (SPD), namely $\mathbf{A} \succ 0$ if either:

1. $\lambda_i(\mathbf{A}) > 0$ for all i .
2. $\mathbf{v}^T \mathbf{A} \mathbf{v} > 0$ for all $\mathbf{v} \neq \mathbf{0}$.

1.1.3

Prove that the two conditions are equivalent, that is:

$$\lambda_i(\mathbf{A}) > 0 \iff \mathbf{v}^T \mathbf{A} \mathbf{v} > 0 \quad \forall \mathbf{v} \neq \mathbf{0}$$

1.2 PCA

Full PCA

- Consider the data $\mathcal{X} = \{\mathbf{x}_i \in \mathbb{R}^D\}_{i=1}^N$ with mean $\boldsymbol{\mu}_x \in \mathbb{R}^D$ and covariance $\boldsymbol{\Sigma}_x \in \mathbb{R}^{D \times D}$.
- Let $\boldsymbol{\Sigma}_x = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T$ be the eigendecomposition of $\boldsymbol{\Sigma}_x$.
- Let $\mathbf{z}_i = \mathbf{U}^T(\mathbf{x}_i - \boldsymbol{\mu}_x)$

1.2.1

Prove that:

1. The mean of $\mathcal{Z} = \{\mathbf{z}_i\}_{i=1}^N$ is zero, that is, $\boldsymbol{\mu}_z = \frac{1}{N} \sum_{i=1}^N \mathbf{z}_i = \mathbf{0}$.
 2. The covariance of \mathcal{Z} is diagonal, that is $\boldsymbol{\Sigma}_z$ is diagonal.
 3. $\|\mathbf{x}_i - \mathbf{x}_j\|_2 = \|\mathbf{z}_i - \mathbf{z}_j\|_2$ for all i and j .
-

Geometric PCA

- Let $\mathbf{U}_d \in \mathbb{R}^{D \times d}$ be a full rank matrix (with $d \leq D$).

1.2.2

Show that exists an invertible matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$ such that $\mathbf{O} = \mathbf{U}_d \mathbf{M} \in \mathbb{R}^{D \times d}$ is semi-orthogonal, that is:

$$\mathbf{O}^T \mathbf{O} = \mathbf{I}_d$$

-
- Consider the data $\mathbf{X} \in \mathbb{R}^{D \times N}$ with zero mean $\mathbf{X} \mathbf{1}_N = \mathbf{0} \in \mathbb{R}^D$ and covariance $\Sigma_x = \frac{1}{N} \mathbf{X} \mathbf{X}^T \in \mathbb{R}^{D \times D}$.
 - Consider the following optimization problems:

1. Reconstruction error minimization:

$$\begin{cases} \arg \min_{\mathbf{U}_d \in \mathbb{R}^{D \times d}} \|\mathbf{X} - \mathbf{U}_d \mathbf{U}_d^T \mathbf{X}\|_F^2 \\ \text{s.t. } \mathbf{U}_d^T \mathbf{U}_d = \mathbf{I}_d \end{cases}$$

2. Variance maximization:

$$\begin{cases} \arg \max_{\mathbf{U}_d \in \mathbb{R}^{D \times d}} \text{Tr} \{ \mathbf{U}_d^T \Sigma_x \mathbf{U}_d \} \\ \text{s.t. } \mathbf{U}_d^T \mathbf{U}_d = \mathbf{I}_d \end{cases}$$

1.2.3

Prove that both problems have the same optimal solution \mathbf{U}_d^* .

PCA analysis

- Consider the data $\{\mathbf{x}_i \in \mathbb{R}^D\}_{i=1}^N$ with mean $\boldsymbol{\mu}_x \in \mathbb{R}^D$ and covariance $\Sigma_x \in \mathbb{R}^{D \times D}$.
- Let $\mathbf{U}_d \in \mathbb{R}^{D \times d}$ be a semi-orthogonal matrix, that is, $\mathbf{U}_d^T \mathbf{U}_d = \mathbf{I}_d$.
- Let $\mathbf{z}_i = \mathbf{U}_d^T (\mathbf{x}_i - \boldsymbol{\mu}_x) \in \mathbb{R}^d$.
- Let $\hat{\mathbf{x}}_i = \mathbf{U}_d \mathbf{z}_i + \boldsymbol{\mu}_x \in \mathbb{R}^D$.
- Let $\boldsymbol{\epsilon}_i = \mathbf{x}_i - \hat{\mathbf{x}}_i \in \mathbb{R}^D$.

1.2.4

Prove that:

$$\text{Tr} \{ \Sigma_x \} = \text{Tr} \{ \Sigma_z \} + \text{Tr} \{ \Sigma_\epsilon \}$$

where:

- $\Sigma_z \in \mathbb{R}^{d \times d}$ is the covariance of $\{\mathbf{z}_i\}_{i=1}^N$.
- $\Sigma_\epsilon \in \mathbb{R}^{D \times D}$ is the covariance of $\{\boldsymbol{\epsilon}_i\}_{i=1}^N$.

1.2.5

- Let $\mathbf{U}_d \in \mathbb{R}^{D \times d}$ be the top d eigenvectors corresponding to the d largest eigenvalues of Σ_x .
- Show that:

$$\text{Tr} \{ \Sigma_\epsilon \} = \sum_{i=d+1}^D \lambda_i (\Sigma_x)$$

where we assume $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$

High-dimensional data PCA

- Consider the data $\mathbf{X} \in \mathbb{R}^{D \times N}$ where $D > N$.

1.2.6

- Provide a (tight) upper bound on the number of non-zero eigenvalues.
 - Consequently, can you apply PCA to $\mathbf{X} \in \mathbb{R}^{D \times N}$ to obtain $\mathbf{Z} \in \mathbb{R}^{d \times N}$ with $d < D$ such that there is no loss of information?
Explain your answer.
-

Rank minimization

- Let $\mathbf{A} \in \mathbb{R}^{D \times N}$.
- Consider the following rank minimization problem:

$$\begin{cases} \min_{\mathbf{M} \in \mathbb{R}^{D \times N}} \|\mathbf{A} - \mathbf{M}\|_F^2 \\ \text{s.t. rank}(\mathbf{M}) \leq d \end{cases}$$

1.2.7

- Solve the optimization problem.
- Write your final solution using the (truncated) matrices obtained by the SVD decomposition of \mathbf{A} , namely, $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$

Hints:

1. Any matrix $\mathbf{M} \in \mathbb{R}^{D \times N}$ with $\text{rank}(\mathbf{M}) = d$ can be written as $\mathbf{M} = \mathbf{B}\mathbf{C}$ where:

(a) $\mathbf{B} \in \mathbb{R}^{D \times d}$

(b) $\mathbf{C} \in \mathbb{R}^{d \times N}$

use this result to formulate (and solve) an equivalent unconstrained problem.

2. There is a strong connection to PCA.
-

1.3 Implementation and applications



Solve this section in the attached notebook.



2 KPCA

2.1

Centering matrix

- Let $\mathbf{J} = \mathbf{I} - \frac{1}{N}\mathbf{1}\mathbf{1}^T \in \mathbb{R}^{N \times N}$ be the centering matrix.

2.1.1

Prove that \mathbf{J} is idempotent, that is, $\mathbf{J}^2 = \mathbf{J}$.

In words, after applying centering once, the second centering has no effect.

Kernel matrix

- Let $\mathbf{X} \in \mathbb{R}^{D \times N}$ and let:

$$\Sigma_x := \mathbf{X}\mathbf{X}^T$$

$$\mathbf{K}_x := \mathbf{X}^T\mathbf{X}$$

Let $(\mathbf{u}_i, \lambda_i)$ be an eigen pair of Σ_x such that $\Sigma_x \mathbf{u}_i = \lambda_i \mathbf{u}_i$ with $\lambda_i > 0$.

2.1.2

1. Show that λ_i is an eigenvalue of \mathbf{K}_x as well.
 2. Find its corresponding eigenvector such that $\mathbf{K}_x \mathbf{w}_i = \lambda_i \mathbf{w}_i$.
-

Kernel functions

- Let $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ and consider $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^N$.

2.1.3

Show that if k can be written as an inner product, that is

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

for some ϕ , then, the matrix defined by:

$$\mathbf{K}_x[i, j] = k(\mathbf{x}_i, \mathbf{x}_j)$$

is an SPSP matrix, namely, $\mathbf{K}_x \succeq 0$.

-
- Let \mathbf{A} be an SPD matrix, and let:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{A} \mathbf{x}_j$$

2.1.4

Prove or disprove:
 k is a kernel function.

-
- Let $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^d$ and consider:

$$k(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$$

Prove or disprove:
 k is a kernel function.

-
- Consider $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^N$, and consider the kernel:

$$k(\mathbf{x}_i, \mathbf{x}_j) := \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

for some ϕ .

- Let:

$$\tilde{k}(\mathbf{x}_i, \mathbf{x}_j) := \langle \phi(\mathbf{x}_i) - \boldsymbol{\mu}_\phi, \phi(\mathbf{x}_j) - \boldsymbol{\mu}_\phi \rangle$$

be the centered version, where:

$$\boldsymbol{\mu}_\phi = \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{x}_i)$$

2.1.5

Show that \tilde{k} can be written using only k , and without using ϕ and μ_ϕ explicitly.

- Let $\mathbf{K}_x \in \mathbb{R}^{N \times N}$ be a kernel matrix, that is:

$$\mathbf{K}_x[i, j] = k(\mathbf{x}_i, \mathbf{x}_j)$$

for some kernel function k .

- Let $\widetilde{\mathbf{K}}_x$ be the centered version, that is:

$$\widetilde{\mathbf{K}}_x = \mathbf{J} \mathbf{K}_x \mathbf{J}$$

where $\mathbf{J} = \mathbf{I} - \frac{1}{N} \mathbf{1} \mathbf{1}^T$.

2.1.6

Prove or disprove:
 $\widetilde{\mathbf{K}}_x$ is an SPD matrix.

Out of sample extension

- Let \mathbf{K}_x be the kernel matrix obtained from the training set $\mathcal{X} = \{\mathbf{x}_i \in \mathbb{R}^D\}_{i=1}^N$.
- Let $\mathbf{Z} \in \mathbb{R}^{d \times N}$ be the low-dimensional representation obtained by applying KPCA, that is:

$$\mathbf{Z} = \Sigma_d \mathbf{V}_d^T$$

where $\mathbf{V} \Sigma \mathbf{V}^T = \mathbf{J} \mathbf{K}_x \mathbf{J}$ is an eigendecomposition (see lecture notes).

- Let $\mathbf{X}^* \in \mathbb{R}^{D \times N^*}$ be a set of new unseen data-points.

2.1.7


Write an expression (in a matrix form) for $\mathbf{Z}^* \in \mathbb{R}^{d \times N^*}$, the KPCA out of sample extension applied to \mathbf{X}^* .

- Let $\mathcal{X}^* = \{\mathbf{x}_i^*\}_{i=1}^{N^*} \subseteq \mathcal{X}$ be a subset of the training set \mathcal{X} .
- Let $\mathbf{X}^* \in \mathbb{R}^{D \times N^*}$ be the matrix from of \mathcal{X}^* .
- Let $\mathbf{Z}^* \in \mathbb{R}^{d \times N}$ be the low-dimensional representation obtained by the training encoding.

2.1.8

Prove that the out of sample encoding applied to \mathbf{X}^* coincide with the training encoding \mathbf{Z}^* .

2.2 Implementation and applications

 Solve this section in the attached notebook. 