

ASSIGNMENT 3: PCA AND KERNEL PCA



Institute for Machine Learning

Contact

Heads:

Thomas Adler,
Philipp Renz,
Andreas Radler

Institute for Machine Learning
Johannes Kepler University
Altenberger Str. 69
A-4040 Linz

E-Mail: {adler,renz,radler}@ml.jku.at

[Institute Homepage](#)

Copyright statement:

This material, no matter whether in printed or electronic form, may be used for personal and non-commercial educational use only. Any reproduction of this material, no matter whether as a whole or in parts, no matter whether in printed or in electronic form, requires explicit prior acceptance of the authors.

Agenda

- Recap of PCA
- Intuition of Kernel PCA
- Example of Kernel PCA
- Derivation of Kernel PCA

Lecture notes

Mathematics for Machine Learning, [2018, Marc Peter Deisenroth et. al.]

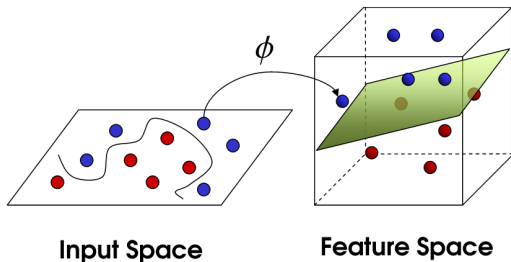
Notation

Notation as in the lecture, i.e. we are given a data matrix $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$, where each $\mathbf{x}_i \in \mathbb{R}^m$ for $1 \leq i \leq n$. Thus, \mathbf{X} is an $n \times m$ matrix.

Recap of Principal Component Analysis

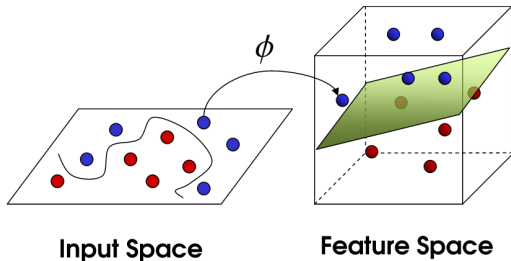
- Intuition: find directions of **largest variances** in X
- Assume $X \in \mathbb{R}^{n \times m}$ is centered, i.e. $\sum_{i=1}^n x_i = 0$
- Compute the **covariance matrix** $C = \frac{1}{n} X^\top X$
- Compute its **eigendecomposition** $C = U \Lambda U^\top$
 - Λ is the diagonal matrix of **eigenvalues** $\text{diag}(\lambda_1, \dots, \lambda_m)$ in descending order, i.e. $\lambda_1 \geq \dots \geq \lambda_m$
 - U is the matrix of respective **eigenvectors** (u_1, \dots, u_m)
 - U is **orthogonal** (i.e. $U^{-1} = U^\top$) because C is **symmetric**
- The eigenvector u_i is the i -th **principal component** (PC)
- The eigenvalue λ_i is the **variance** along the direction u_i
 - U uncorrelates X because $\frac{1}{n} (XU)^\top XU = \Lambda$
 - Λ is the (diagonal) covariance matrix of XU

Intuition of Kernel PCA (1/2)



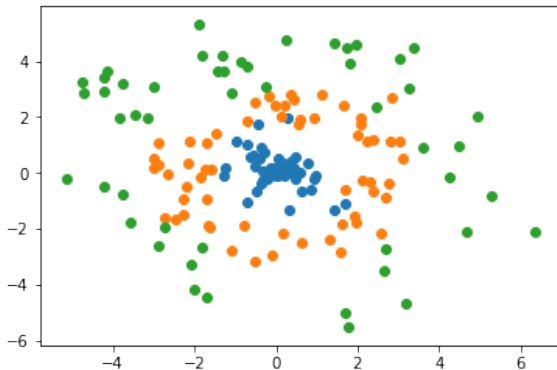
- PCA is a linear method
- However, data are often not linearly separable
- **Idea:** map data into high-dimensional space where it becomes linearly separable and apply PCA there

Intuition of Kernel PCA (2/2)



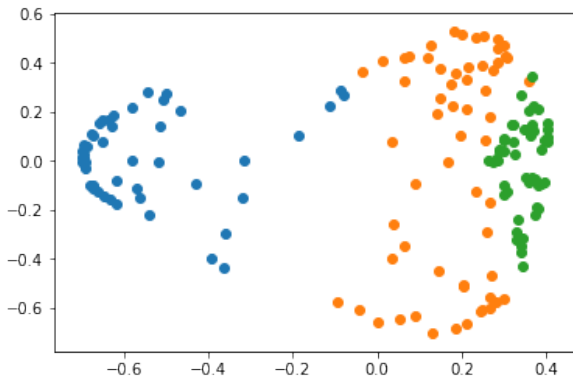
- This space is called the **feature space**
- The mapping Φ is called the **feature map**
- Use a **kernel function** instead of explicitly computing Φ
- This approach is called the **kernel trick**

Example of Kernel PCA (1/2)



■ Data cannot be separated linearly

Example of Kernel PCA (2/2)



■ Data can be well separated linearly using only the first PC

Definitions and Assumptions

- Let \mathcal{V} denote the **feature space**
- Let $\Phi : \mathbb{R}^m \rightarrow \mathcal{V}$ denote the **feature map**
- Let $k : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ denote the **kernel**

$$k(\mathbf{x}, \mathbf{y}) = \Phi^\top(\mathbf{x})\Phi(\mathbf{y})$$

- Assume the data is **centered** in \mathcal{V} , i.e. $\sum_{i=1}^n \Phi(\mathbf{x}_i) = \mathbf{0}$
- The covariance matrix $\mathbf{C} \in \mathbb{R}^{\dim \mathcal{V} \times \dim \mathcal{V}}$ is

$$\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)\Phi^\top(\mathbf{x}_i)$$

- The **Gram matrix** $\mathbf{K} \in \mathbb{R}^{n \times n}$ is defined as

$$K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = \Phi^\top(\mathbf{x}_i)\Phi(\mathbf{x}_j)$$

Eigenvector in Feature Space

- Note that possibly $\dim \mathcal{V} > n$
- In fact, even $\dim \mathcal{V} = \infty$ is possible
- Therefore, consider only eigenvectors $\mathbf{u} \in \mathcal{V}$ which are in the span of the mapped data $\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_n)$
- Thus, the considered eigenvectors can be written as a linear combination of the mapped samples

$$\mathbf{u} = \sum_{i=1}^n \alpha_i \Phi(\mathbf{x}_i)$$

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)^\top$

Eigenvalue Equation in Feature Space

$$Cu = \lambda u$$

$$C \sum_{i=1}^n \alpha_i \Phi(\mathbf{x}_i) = \lambda \sum_{i=1}^n \alpha_i \Phi(\mathbf{x}_i)$$

$$\Phi^\top(\mathbf{x}_l) C \sum_{i=1}^n \alpha_i \Phi(\mathbf{x}_i) = \lambda \sum_{i=1}^n \alpha_i \Phi^\top(\mathbf{x}_l) \Phi(\mathbf{x}_i)$$

$$\frac{1}{n} \Phi^\top(\mathbf{x}_l) \left(\sum_{j=1}^n \Phi(\mathbf{x}_j) \Phi^\top(\mathbf{x}_j) \right) \sum_{i=1}^n \alpha_i \Phi(\mathbf{x}_i) = \lambda [K\alpha]_l$$

$$K^2 \alpha = n\lambda K \alpha$$

Eigenvalue Equation of the Gram Matrix

- We can solve the eigenvalue equation in feature space by solving the eigenvalue equation of the Gram matrix

$$\mathbf{K}\boldsymbol{\alpha} = n\lambda\boldsymbol{\alpha}$$

- We have to normalize $\boldsymbol{\alpha} \leftarrow \left(\|\boldsymbol{\alpha}\|\sqrt{n\lambda}\right)^{-1}\boldsymbol{\alpha}$ to account for an orthonormal basis in feature space because

$$1 = \mathbf{u}^\top \mathbf{u} = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \boldsymbol{\Phi}^\top(\mathbf{x}_i) \boldsymbol{\Phi}(\mathbf{x}_j) = \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} = n\lambda \boldsymbol{\alpha}^\top \boldsymbol{\alpha}$$

- Project \mathbf{x} onto a principal component \mathbf{u} by

$$\mathbf{u}^\top \boldsymbol{\Phi}(\mathbf{x}) = \sum_{i=1}^n \alpha_i \boldsymbol{\Phi}^\top(\mathbf{x}_i) \boldsymbol{\Phi}(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

The Kernel Trick

- The **feature map** $\Phi(\cdot)$ need never be computed explicitly
- Instead, we reformulated everything which requires explicit computation using only the **kernel** $k(\cdot, \cdot)$
- This is known as the **kernel trick**
- It can be applied not only to PCA but also to various other methods (e.g. support vector machines)
- It is a popular device to make a linear method non-linear
- The linear method is applied in a high-dimensional space that has a non-linear relation to the input space

Centering the Gram Matrix

- We assumed that the data are centered in feature space
- We can center K_{ij} by

$$\begin{aligned} & \left(\Phi(\mathbf{x}_i) - \frac{1}{n} \sum_{p=1}^n \Phi(\mathbf{x}_p) \right)^\top \left(\Phi(\mathbf{x}_j) - \frac{1}{n} \sum_{q=1}^n \Phi(\mathbf{x}_q) \right) \\ &= \Phi^\top(\mathbf{x}_i) \Phi(\mathbf{x}_j) - \frac{1}{n} \sum_{q=1}^n \Phi^\top(\mathbf{x}_i) \Phi(\mathbf{x}_q) \\ &\quad - \frac{1}{n} \sum_{p=1}^n \Phi^\top(\mathbf{x}_p) \Phi(\mathbf{x}_j) + \frac{1}{n^2} \sum_{p=1}^n \sum_{q=1}^n \Phi^\top(\mathbf{x}_p) \Phi(\mathbf{x}_q) \end{aligned}$$

- In matrix form, this is

$$\mathbf{K} - \frac{1}{n} \mathbf{K} \mathbf{1} \mathbf{1}^\top - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \mathbf{K} + \frac{\mathbf{1}^\top \mathbf{K} \mathbf{1}}{n^2}$$

Summary

- Choose a kernel $k(\cdot, \cdot)$
- Compute the Gram matrix \mathbf{K} by $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$
- Center the Gram matrix in feature space by

$$\mathbf{K} \leftarrow \mathbf{K} - \frac{1}{n} \mathbf{K} \mathbf{1} \mathbf{1}^\top - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \mathbf{K} + \frac{\mathbf{1}^\top \mathbf{K} \mathbf{1}}{n^2}$$

- Compute eigenvectors α and eigenvalues $n\lambda$
- Normalize eigenvectors $\alpha \leftarrow \left(\|\alpha\| \sqrt{n\lambda} \right)^{-1} \alpha$
- Given a new point \mathbf{x} , define $\mathbf{k} = (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n))^\top$
- Center it by $\mathbf{k} \leftarrow \mathbf{k} - \frac{1}{n} \mathbf{K} \mathbf{1} - \frac{1}{n} \mathbf{1}^\top \mathbf{k} + \frac{1}{n^2} \mathbf{1}^\top \mathbf{K} \mathbf{1}$
- Project \mathbf{k} onto principal component by $\alpha^\top \mathbf{k}$