

ECE 595: Machine Learning I

Lecture 7.2: Kernel-PCA

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Outline

Feature Analysis

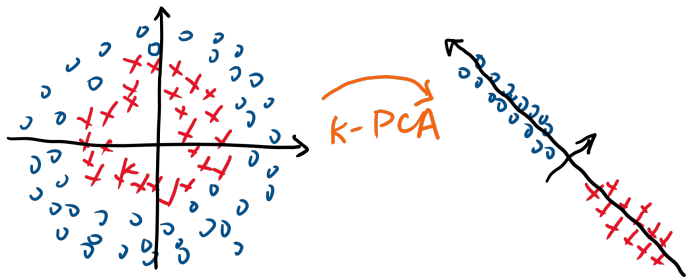
- Lecture 7 Principal Component Analysis (PCA)
- Lecture 8 Hand-Crafted and Deep Features

This Lecture

- PCA
 - Low-dimensional Representation
 - Geometric Interpretation
 - Eigen-Face Problem
- Kernel-PCA
 - Adding kernels to PCA
 - Algorithm
 - Examples

Motivation of Kernel PCA

- Data is originally difficult for PCA
- Find a nonlinear transform
- Idea: Leverage the kernel trick: $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle$
- Example: Left is hard for PCA. After K-PCA, right has a clear principal component.



Kernel for Covariance Matrix

- Assume $\phi(\mathbf{x}^{(n)})$ has zero mean. Then consider the covariance matrix

$$\Sigma = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (\mathbf{x}^{(n)})^T.$$

- Replacing the outer products by feature transforms

$$\mathbf{x}^{(n)} \rightarrow \phi(\mathbf{x}^{(n)}),$$

for some nonlinear transformation ϕ .

- If this can be done, then the covariance will become

$$\Sigma = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}^{(n)}) \phi(\mathbf{x}^{(n)})^T.$$

- But this is not enough because a kernel needs an inner product

$$k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) = \phi(\mathbf{x}^{(n)})^T \phi(\mathbf{x}^{(m)}).$$

Kernel Trick

- Recall: PCA solves the eigen-decomposition problem:

$$\mathbf{\Sigma} \mathbf{u} = \lambda \mathbf{u}$$

So we also need to consider \mathbf{u} .

- How about this candidate? (Recall: In Kernel Method we express the model parameter as a linear combination of the samples):

$$\mathbf{u} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}^{(n)}).$$

- Substitute this into the equation $\mathbf{\Sigma} \mathbf{u} = \lambda \mathbf{u}$:

$$\underbrace{\left(\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}^{(n)}) \phi(\mathbf{x}^{(n)})^T \right)}_{\mathbf{\Sigma}} \underbrace{\left(\sum_{m=1}^N \alpha_m \phi(\mathbf{x}^{(m)}) \right)}_{\mathbf{u}} = \lambda \underbrace{\left(\sum_{n=1}^N \alpha_n \phi(\mathbf{x}^{(n)}) \right)}_{\lambda \mathbf{u}}$$

Kernel Trick

- This means

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}^{(n)}) \left(\sum_{m=1}^N \alpha_m \phi(\mathbf{x}^{(n)})^T \phi(\mathbf{x}^{(m)}) \right) = \lambda \sum_{n=1}^N \alpha_n \phi(\mathbf{x}^{(n)})$$

- Recognizing $\phi(\mathbf{x}^{(n)})^T \phi(\mathbf{x}^{(m)}) = k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)})$:

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}^{(n)}) \left(\sum_{m=1}^N \alpha_n k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) \right) = \lambda \sum_{n=1}^N \alpha_n \phi(\mathbf{x}^{(n)})$$

- Multiply $\phi(\mathbf{x}^{(\ell)})^T$ on both sides.

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}^{(\ell)}, \mathbf{x}^{(n)}) \left(\sum_{m=1}^N \alpha_n k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) \right) = \lambda \sum_{n=1}^N \alpha_n k(\mathbf{x}^{(\ell)}, \mathbf{x}^{(n)})$$

- This is $\frac{1}{N} \mathbf{K}(\mathbf{K}\alpha) = \lambda \mathbf{K}\alpha$.

Eigenvectors of K-PCA

- Rearrange the terms we have that $\mathbf{K}^2\boldsymbol{\alpha} = N\lambda\mathbf{K}\boldsymbol{\alpha}$.
- We can remove one of the \mathbf{K} 's since it only causes issues with zero-eigenvalues which are not important to us anyway. So we have

$$\mathbf{K}\boldsymbol{\alpha} = N\lambda\boldsymbol{\alpha}. \quad (2)$$

- This is just another eigen-decomposition problem. We moved from $\boldsymbol{\Sigma}\mathbf{u} = \lambda\mathbf{u}$ to $\mathbf{K}\boldsymbol{\alpha} = N\lambda\boldsymbol{\alpha}$. Note that $\boldsymbol{\alpha}$ is the coefficients for \mathbf{u} :

$$\mathbf{u} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}^{(n)}) = \boldsymbol{\Phi}\boldsymbol{\alpha},$$

where $\boldsymbol{\Phi} = [\phi(\mathbf{x}^{(1)}), \dots, \phi(\mathbf{x}^{(N)})]$ is the transformed data matrix. Recall $\boldsymbol{\Phi}\boldsymbol{\Phi}^T = \mathbf{K}$ is the kernel matrix where

$$[\mathbf{K}]_{ij} = \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)}).$$

Representation in Kernel Space

- If you run eigen-decomposition on \mathbf{K} , you will get p eigen-vectors $\alpha_1, \dots, \alpha_p$ where p is the number you choose.
- When a new sample \mathbf{x} comes, the j -th representation coefficient is

$$\beta_j = \phi(\mathbf{x})^T \mathbf{u} = \phi(\mathbf{x})^T \sum_{n=1}^N \alpha_{jn} \phi(\mathbf{x}^{(n)}) = \sum_{n=1}^N \alpha_{jn} k(\mathbf{x}, \mathbf{x}^{(n)}). \quad (3)$$

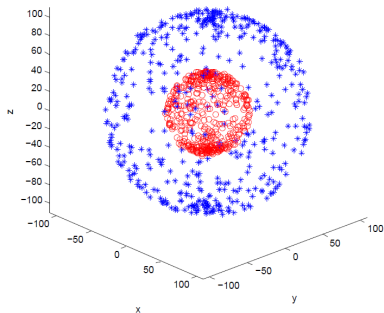
- For the entire representation $\beta \in \mathbb{R}^p$, we have

$$\beta = \begin{bmatrix} \text{---} \alpha_1^T \text{---} \\ \vdots \\ \text{---} \alpha_p^T \text{---} \end{bmatrix} \begin{bmatrix} k(\mathbf{x}, \mathbf{x}^{(1)}) \\ k(\mathbf{x}, \mathbf{x}^{(2)}) \\ \vdots \\ k(\mathbf{x}, \mathbf{x}^{(N)}) \end{bmatrix} \quad (4)$$

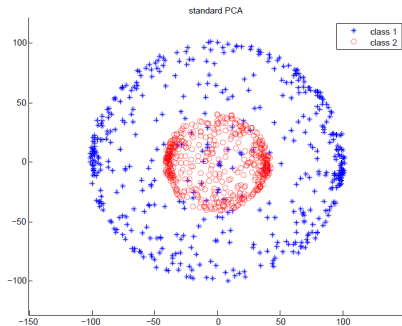
where $\alpha_j = [\alpha_{j1}, \dots, \alpha_{jN}]^T$.

Example

Here is an example taken from Wang (2012) Kernel Principal Component Analysis and its Applications <https://arxiv.org/abs/1207.3538>



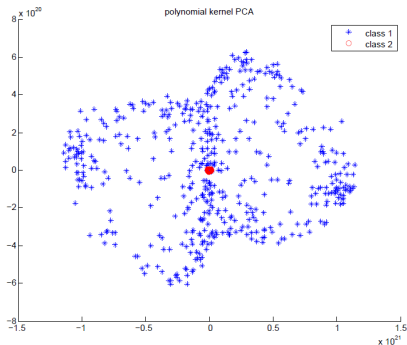
Original Data



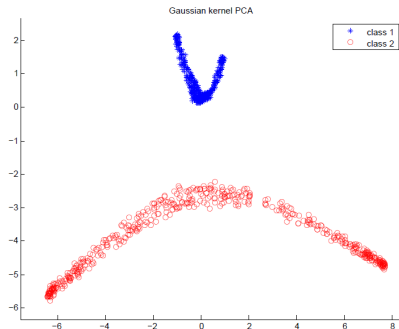
Linear PCA

Example

Here is an example taken from Wang (2012) Kernel Principal Component Analysis and its Applications <https://arxiv.org/abs/1207.3538>



K-PCA with polynomial



K-PCA with Gaussian

Reading List

PCA Tutorial

- Jonathon Shlens “A Tutorial on Principal Component Analysis”, <https://arxiv.org/pdf/1404.1100.pdf>

PCA: Should We Remove Mean?

- Paul Honeine, “An eigenanalysis of data centering in machine learning”, <https://arxiv.org/pdf/1407.2904.pdf>
- Does mean centering or feature scaling affect a Principal Component Analysis?
<https://sebastianraschka.com/faq/docs/pca-scaling.html>

K-PCA

- Quan Wang (2012), “Kernel Principal Component Analysis and its Applications”, <https://arxiv.org/abs/1207.3538>
- Schölkopf et al. (2005), “Kernel Principal Component Analysis”, <https://link.springer.com/chapter/10.1007/BFb0020217>