

ECE 595: Machine Learning I

Lecture 07 Feature Analysis via PCA

Spring 2020

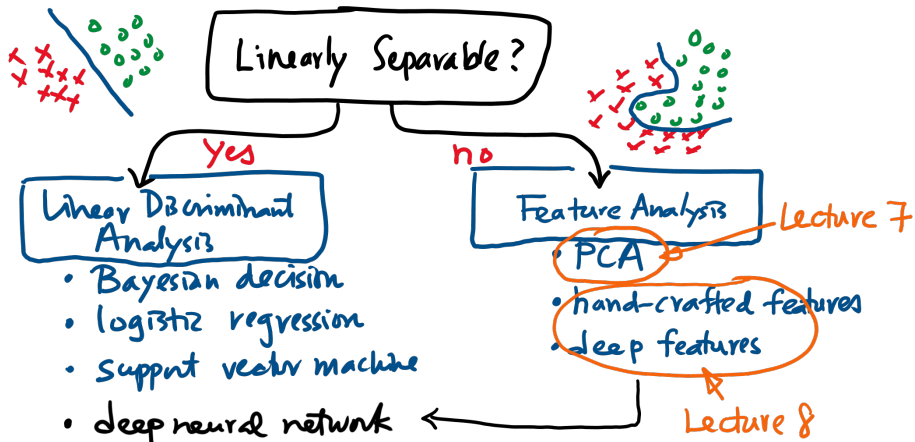
Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Overview

Supervised Learning for Classification



Outline

Feature Analysis

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

This Lecture

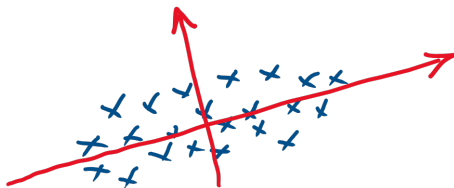
- PCA
 - Low-dimensional Representation
 - Geometric Interpretation
 - Eigen-Face Problem
- Kernel-PCA
 - Adding kernels to PCA
 - Algorithm
 - Examples

Low-Dimensional Representation

- Consider a set of data point $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$
- These data points are living in a **high dimensional space** $\mathbf{x}^{(n)} \in \mathbb{R}^d$
- Find a **low dimensional representation** in \mathbb{R}^p where $p < d$
- Equivalent to finding the **principal components** $\mathbf{v}_1, \dots, \mathbf{v}_p$ such that

$$\mathbf{x}^{(n)} \approx \sum_{i=1}^p \alpha_i^{(n)} \mathbf{v}_i$$

- Then every $\mathbf{x}^{(n)} \in \mathbb{R}^d$ can be represented using $\alpha^{(n)} \in \mathbb{R}^p$.



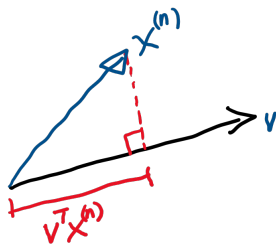
One Sample Analysis

- Consider a simpler problem: One data point \mathbf{x} and one direction \mathbf{v} .
- We want to find a direction $\hat{\mathbf{v}}$ and a scalar $\hat{\alpha}$ such that

$$(\hat{\mathbf{v}}, \hat{\alpha}) = \underset{\|\mathbf{v}\|_2=1, \alpha}{\operatorname{argmin}} \left\| \begin{bmatrix} | \\ | \\ \mathbf{x} \\ | \\ | \end{bmatrix} - \alpha \begin{bmatrix} | \\ | \\ \mathbf{v} \\ | \\ | \end{bmatrix} \right\|^2$$

- First assume \mathbf{v} is available. Then take derivative w.r.t. α :

$$2\mathbf{v}^T(\mathbf{x} - \alpha\mathbf{v}) = 0 \quad \Rightarrow \quad \alpha = \mathbf{v}^T\mathbf{x}.$$



One Sample Analysis

- Substitute $\alpha = \mathbf{x}^T \mathbf{v}$ into the optimization
- Then the optimization becomes

$$\begin{aligned} \operatorname{argmin}_{\|\mathbf{v}\|_2=1} \|\mathbf{x} - \alpha \mathbf{v}\|^2 &= \operatorname{argmin}_{\|\mathbf{v}\|_2=1} \left\{ \mathbf{x}^T \mathbf{x} - 2\alpha \mathbf{x}^T \mathbf{v} + \alpha^2 \mathbf{v}^T \mathbf{v} \right\} \\ &= \operatorname{argmin}_{\|\mathbf{v}\|_2=1} \left\{ -2\alpha \mathbf{x}^T \mathbf{v} + \alpha^2 \right\} \\ &= \operatorname{argmin}_{\|\mathbf{v}\|_2=1} \left\{ -2(\mathbf{x}^T \mathbf{v}) \mathbf{x}^T \mathbf{v} + (\mathbf{x}^T \mathbf{v})^2 \right\} \\ &= \operatorname{argmax}_{\|\mathbf{v}\|_2=1} \left\{ \mathbf{v}^T \mathbf{x} \mathbf{x}^T \mathbf{v} \right\} \end{aligned}$$

- Take expectation on both sides:

$$\operatorname{argmin}_{\|\mathbf{v}\|_2=1} \mathbb{E}_{\mathbf{x}} \|\mathbf{x} - \alpha \mathbf{v}\|^2 = \operatorname{argmax}_{\|\mathbf{v}\|_2=1} \mathbf{v}^T \mathbb{E}_{\mathbf{x}} \left\{ \mathbf{x} \mathbf{x}^T \right\} \mathbf{v}$$

Eigenvalue Problem

- Let $\Sigma \stackrel{\text{def}}{=} \mathbb{E}[\mathbf{x}\mathbf{x}^T]$.
- Then the optimization problem is

$$\operatorname{argmax}_{\|\mathbf{v}\|_2=1} \mathbf{v}^T \Sigma \mathbf{v}.$$

- The solution to this problem is the eigenvalue and eigenvectors of Σ .

Theorem

Let Σ be a $d \times d$ matrix with eigen-decomposition $\Sigma = \mathbf{U}\mathbf{S}\mathbf{U}^T$. Then, the optimization

$$\hat{\mathbf{v}} = \operatorname{argmax}_{\|\mathbf{v}\|_2=1} \mathbf{v}^T \Sigma \mathbf{v}.$$

has a solution $\hat{\mathbf{v}} = \mathbf{u}_i$ for any $i = 1, \dots, d$.

Proof: See Appendix.

Finite Samples

- When there are N training samples, the optimization is

$$\operatorname{argmin}_{\|\mathbf{v}\|_2=1} \underbrace{\frac{1}{N} \sum_{n=1}^N \|\mathbf{x}^{(n)} - \alpha^{(n)} \mathbf{v}\|^2}_{=\mathbb{E}[\|\mathbf{x} - \alpha \mathbf{v}\|^2], N \rightarrow \infty} = \operatorname{argmax}_{\|\mathbf{v}\|_2=1} \mathbf{v}^T \underbrace{\left\{ \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (\mathbf{x}^{(n)})^T \right\}}_{=\mathbb{E}[\mathbf{x} \mathbf{x}^T], N \rightarrow \infty} \mathbf{v}$$

- In practice, given $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$, we approximate Σ by its empirical estimate

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

- You can also remove the mean vectors: $\boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)}$:

$$\Sigma \approx \frac{1}{N} \sum_{n=1}^N (\mathbf{x}^{(n)} - \boldsymbol{\mu})(\mathbf{x}^{(n)} - \boldsymbol{\mu})^T$$

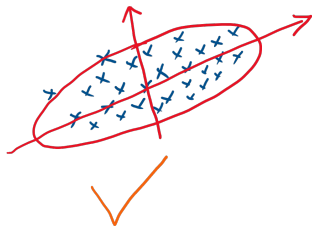
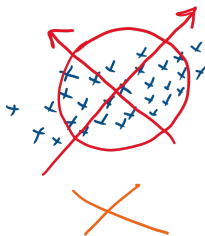
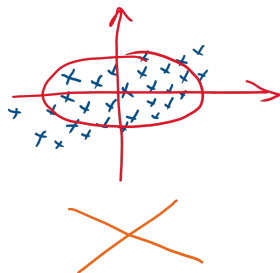
Statistical Interpretation

- The optimization

$$\operatorname{argmax}_{\|\mathbf{v}\|_2=1} \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}.$$

asks us to find a principal direction that maximizes the **variance**.

- Belief: Large variance = “signal”, small variance = “noise”



The Eigenface Problem



Figure: The extended Yale Face Database B.

- Dataset: $\{\mathbf{x}^{(n)}\}_{n=1}^N$.
- Each $\mathbf{x}^{(n)} \in \mathbb{R}^d$ is a vector representation of a $\sqrt{d} \times \sqrt{d}$ image.
- Task 1: Find a low-dimensional representation (This lecture)
- Task 2: Classify faces for a new image (Later)

Low Dimensional Representation

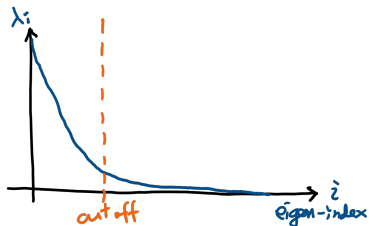
- Estimate the mean vector $\boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)}$.
- Estimate the covariance matrix

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

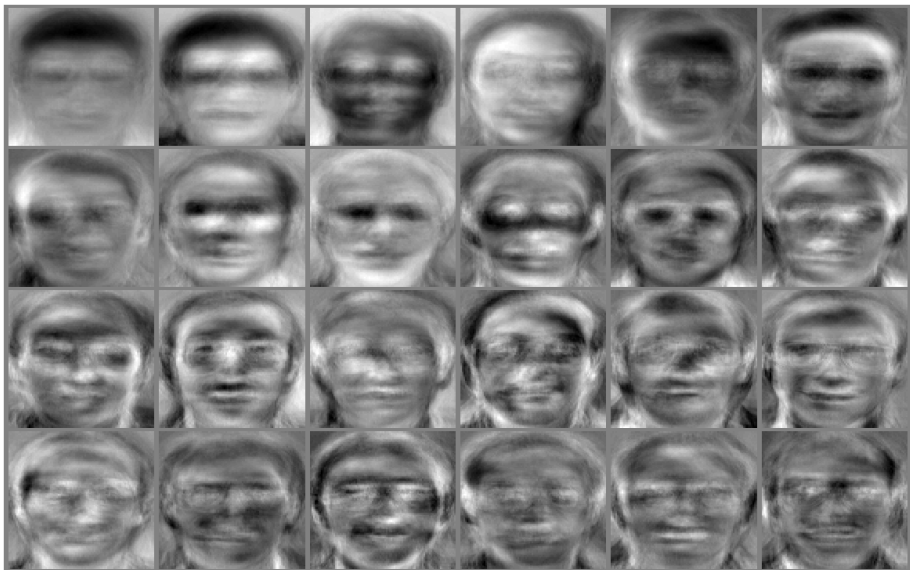
- Eigen-decomposition: $\boldsymbol{\Sigma} = \mathbf{U}\mathbf{S}\mathbf{U}^T$.
- When a new image \mathbf{y} comes, estimate the coefficients:

$$\alpha_j = \mathbf{u}_j^T \mathbf{y}$$

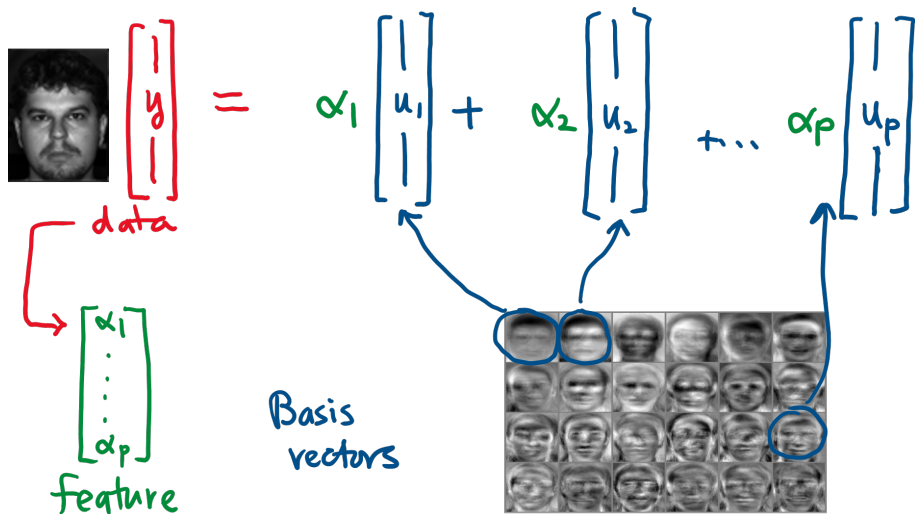
- How many coefficients to use?



The Basis Vectors u_i



Representing Faces



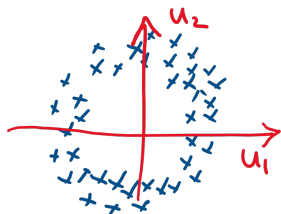
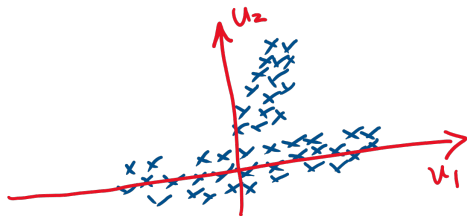
Discussion

What does PCA do?

- PCA is a tool for **dimension reduction**.
- It compresses a raw data vector $\mathbf{y} \in \mathbb{R}^d$ into a smaller feature vector $\alpha \in \mathbb{R}^p$.
- You can now do classification in \mathbb{R}^p instead of \mathbb{R}^d .

When will PCA fail?

- When data intrinsically does not have orthogonal projections
- For example, the distributions below



Outline

Feature Analysis

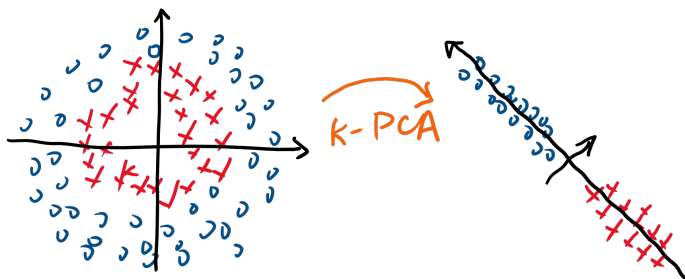
- Lecture 7 Principal Component Analysis (PCA)
- Lecture 8 Hand-Crafted and Deep Features

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

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

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

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

- 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

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

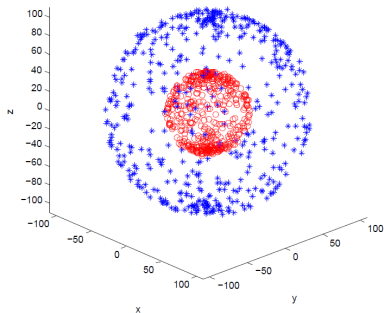
- For the entire representation $\beta \in \mathbb{R}^p$, we have

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

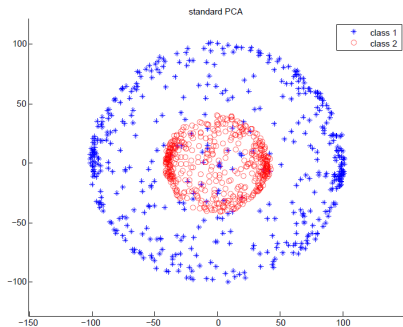
where $\alpha_j = [\alpha_{j1}, \dots, \alpha_{jN}]^T$.

Example

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



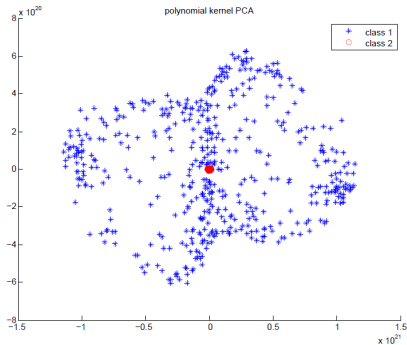
Original Data



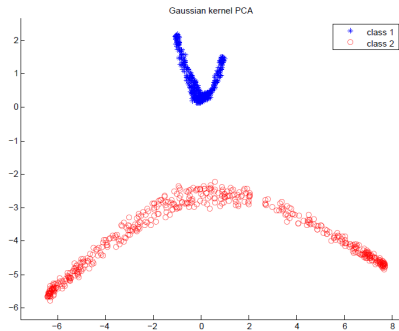
Linear PCA

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>

Appendix

Proof of Eigenvalue Problem

We want to prove that the solution to the problem

$$\hat{\mathbf{v}} = \operatorname{argmax}_{\|\mathbf{v}\|_2=1} \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}.$$

is the eigenvector of the matrix $\boldsymbol{\Sigma}$. To show that, we first write down the Lagrangian:

$$L(\mathbf{v}, \lambda) = \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} - \lambda(\|\mathbf{v}\|^2 - 1)$$

Take derivative w.r.t. \mathbf{v} and setting to zero yields

$$\nabla_{\mathbf{v}} L(\mathbf{v}, \lambda) = 2\boldsymbol{\Sigma} \mathbf{v} - 2\lambda \mathbf{v} = \mathbf{0}.$$

This is equivalent to $\boldsymbol{\Sigma} \mathbf{v} = \lambda \mathbf{v}$. So if $\boldsymbol{\Sigma} = \mathbf{U} \mathbf{S} \mathbf{U}^T$, then by letting $\mathbf{v} = \mathbf{u}_i$ and $\lambda = s_i$ we can satisfy the condition since $\boldsymbol{\Sigma} \mathbf{u}_i = \mathbf{U} \mathbf{S} \mathbf{U}^T \mathbf{u}_i = \mathbf{U} \mathbf{S} \mathbf{e}_i = s_i \mathbf{u}_i$.