


PCA and kernel PCA

Kaiwen Chang
01/03/2017

A decorative light blue triangle is located in the bottom right corner of the slide, pointing towards the top right.

Outline

I. PCA

1. Definitions
2. Solution
3. Example

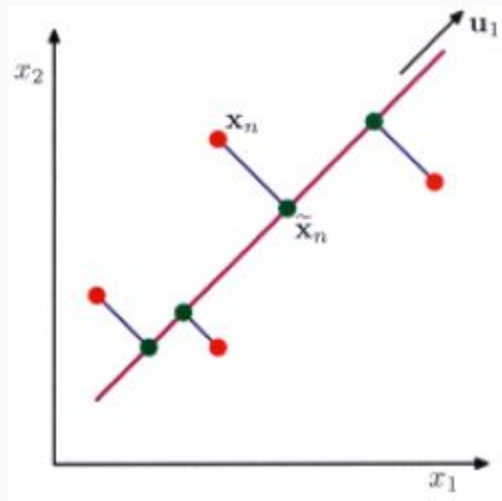
II. Kernel PCA

1. Why kernel
2. Solution
3. Example

I. PCA

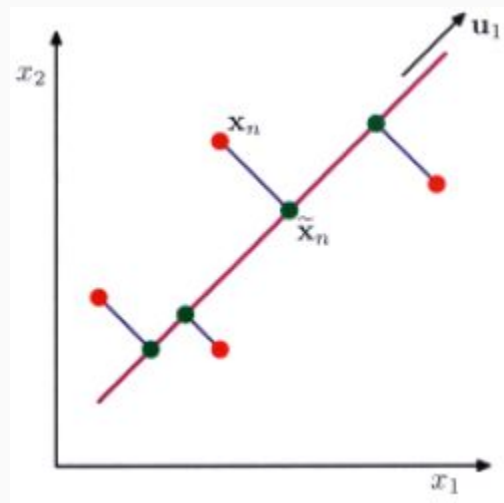
1. Definitions

1. The orthogonal projection of the data onto a lower dimensional linear space (principal subspace), such that the **variance** of the projected data is maximized



Magenta line: principal subspace

2. The linear projection that minimizes the average projection cost (the mean squared distance between the data points and their projections)



Blue line: projection error

2. Solution

Data set $\{\mathbf{x}_n\}_{n=1,2,\dots,N}$

\mathbf{x}_n : column vector with dimension D

Aim: project data onto a space with dimension $M < D$, while maximizing the variance of projected points

Projection of \mathbf{x}_n : $\mathbf{u}_1^T \mathbf{x}_n$

(\mathbf{u}_1 : the direction of space, dimension D)

Mean of projected data: $\mathbf{u}_1^T \bar{\mathbf{x}} \quad \bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$

Variance:

$$\frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}}\}^2 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

Covariance matrix:

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

Maximize variance $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ with constraint $\mathbf{u}_1^T \mathbf{u}_1 = 1$.

Maximize $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$

λ_1 : lagrange multiplier

Setting the derivative to zero:

\mathbf{u}_1 : an eigenvector of S

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1$$

Largest eigenvalue \longrightarrow maximum variance

\mathbf{u}_1 : first principal component

Additional principal component \mathbf{u}_n : choosing a new direction among all possible directions orthogonal to $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{n-1}$, that maximize the projected variance

Optimal linear projection: eigenvectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M$ of covariance matrix S corresponding to M largest eigenvalues

Summary

1. Evaluating the mean and covariance matrix **S** of data set
2. Finding the M eigenvectors of S corresponding to M largest eigenvalues

Computation cost: $O(MD^2)$

Example: http://sebastianraschka.com/Articles/2014_pca_step_by_step.html

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$

$$J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$$

3. Example

Applications:

Dimensionality reduction

Lossy data compression

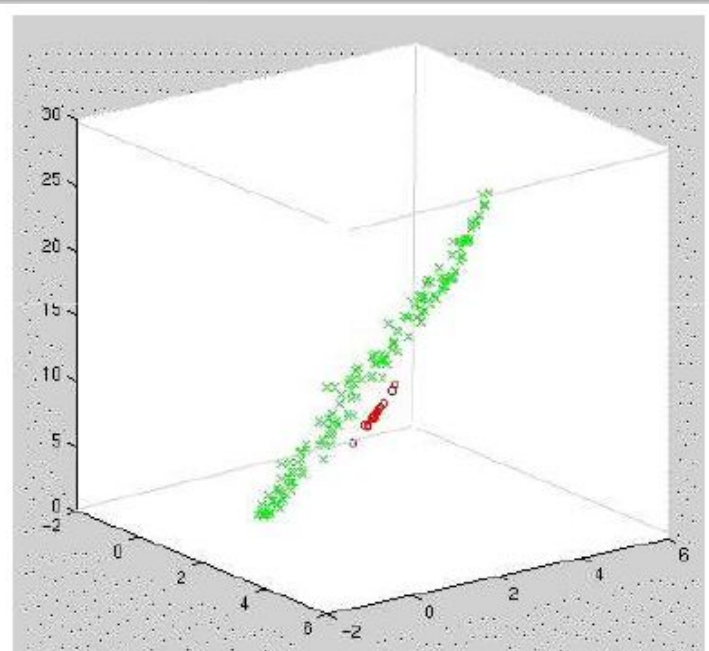
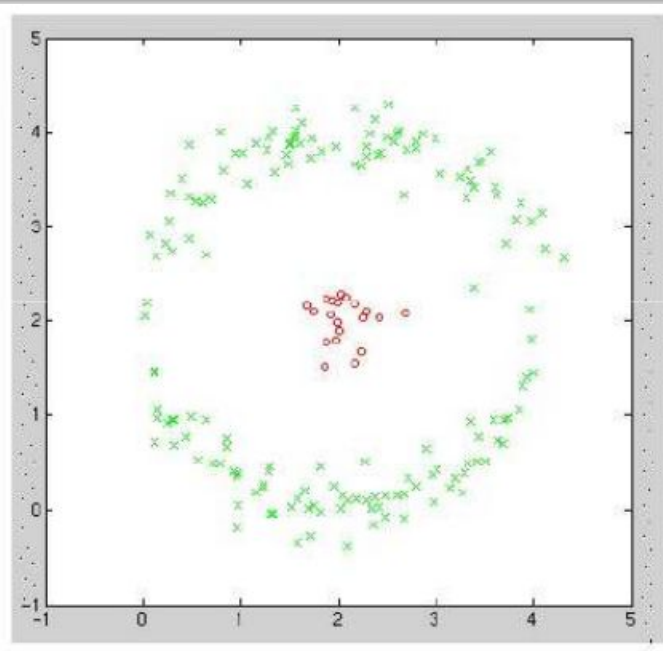
Feature extraction

Data visualization

http://scikit-learn.org/stable/auto_examples/decomposition/plot_pca_iris.html#sphx-glr-auto-examples-decomposition-plot-pca-iris-py

II. Kernel PCA

1. Why kernel



Not linear

Kernel trick $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$

2. Solution

$$\sum_n \mathbf{x}_n = \mathbf{0}.$$

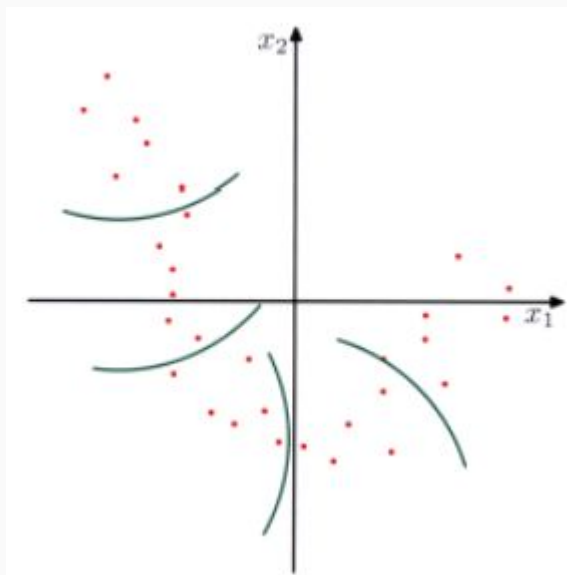
Covariance matrix

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$

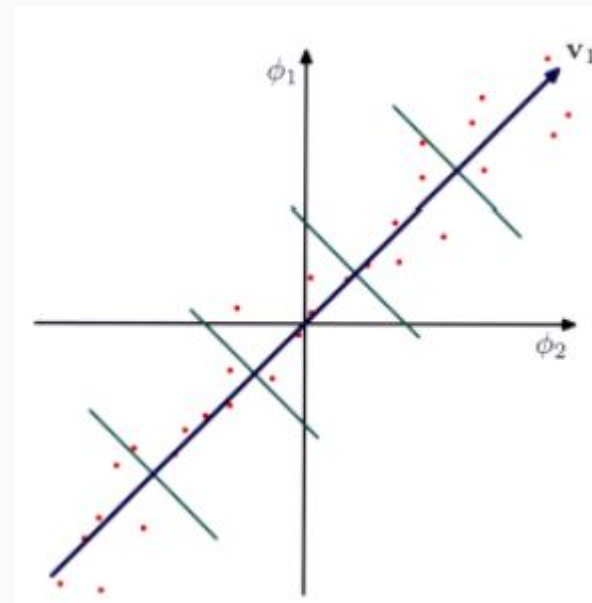
$\phi(\mathbf{x})$: a nonlinear transformation into an M-dimensional feature space ($M > D$)

$$\mathbf{x}_n \longrightarrow \phi(\mathbf{x}_n)$$

Perform standard PCA in feature space



Original data space
Nonlinear projection



Feature space
Principal components

Assume $\sum_n \phi(\mathbf{x}_n) = \mathbf{0}$.

Covariance matrix in feature space:

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

$$\mathbf{C} \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \{ \phi(\mathbf{x}_n)^T \mathbf{v}_i \} = \lambda_i \mathbf{v}_i$$

$$\mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n)$$

Goal: solve the problem without having to work explicitly in the feature space

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \sum_{m=1}^N a_{im} \phi(\mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n)$$

Kernel function: $k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$.

Multiply by $\phi(\mathbf{x}_l)^T$

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_l, \mathbf{x}_n) \sum_{m=1}^m a_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} k(\mathbf{x}_l, \mathbf{x}_n)$$

$$\mathbf{K}^2 \mathbf{a}_i = \lambda_i N \mathbf{K} \mathbf{a}_i$$

\mathbf{a}_i : N dimensional column vector with a_{in}

Kernel matrix

$$\mathbf{K} = \phi(\mathbf{X}) \phi(\mathbf{X})^T = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \kappa(\mathbf{x}_1, \mathbf{x}_2) & \cdots & \kappa(\mathbf{x}_1, \mathbf{x}_N) \\ \kappa(\mathbf{x}_2, \mathbf{x}_1) & \kappa(\mathbf{x}_2, \mathbf{x}_2) & \cdots & \kappa(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \kappa(\mathbf{x}_N, \mathbf{x}_1) & \kappa(\mathbf{x}_N, \mathbf{x}_2) & \cdots & \kappa(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

$$\mathbf{K}\mathbf{a}_i = \lambda_i N \mathbf{a}_i$$

Normalization

$$1 = \mathbf{v}_i^T \mathbf{v}_i = \sum_{n=1}^N \sum_{m=1}^N a_{in} a_{im} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = \mathbf{a}_i^T \mathbf{K} \mathbf{a}_i = \lambda_i N \mathbf{a}_i^T \mathbf{a}_i$$

Projection

$$y_i(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x})^T \phi(\mathbf{x}_n) = \sum_{n=1}^N a_{in} k(\mathbf{x}, \mathbf{x}_n)$$

Centralizing projected data:

$$\tilde{\phi}(\mathbf{x}_n) = \phi(\mathbf{x}_n) - \frac{1}{N} \sum_{l=1}^N \phi(\mathbf{x}_l)$$

Element of the Gram matrix

$$\tilde{K}_{nm} = \tilde{\phi}(\mathbf{x}_n)^T \tilde{\phi}(\mathbf{x}_m)$$

$$\begin{aligned} &= \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) - \frac{1}{N} \sum_{l=1}^N \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_l) - \frac{1}{N} \sum_{l=1}^N \phi(\mathbf{x}_l)^T \phi(\mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N \phi(\mathbf{x}_j)^T \phi(\mathbf{x}_l) \\ &= k(\mathbf{x}_n, \mathbf{x}_m) - \frac{1}{N} \sum_{l=1}^N k(\mathbf{x}_l, \mathbf{x}_m) - \frac{1}{N} \sum_{l=1}^N k(\mathbf{x}_n, \mathbf{x}_l) + \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N k(\mathbf{x}_j, \mathbf{x}_l) \end{aligned}$$

In matrix notation $\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{1}_N \mathbf{K} - \mathbf{K} \mathbf{1}_N + \mathbf{1}_N \mathbf{K} \mathbf{1}_N$

$\mathbf{1}_N$: NxN matrix with each element 1/N

1. Pick a kernel

Gaussian

$$K(\vec{x}, \vec{x}') = \exp(-\beta \|\vec{x} - \vec{x}'\|^2)$$

Polynomial

$$K(\vec{x}, \vec{x}') = (1 + \vec{x} \cdot \vec{x}')^p$$

2. Construct the normalized kernel matrix of data (NxN)

$$\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{1}_N \mathbf{K} - \mathbf{K} \mathbf{1}_N + \mathbf{1}_N \mathbf{K} \mathbf{1}_N$$

3. Solve an eigenvalue problem:

$$\tilde{\mathbf{K}} \mathbf{a}_i = \lambda_i N \mathbf{a}_i$$

4. Represent the data point i as

$$y_i(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x})^T \phi(\mathbf{x}_n) = \sum_{n=1}^N a_{in} k(\mathbf{x}, \mathbf{x}_n)$$

3. Example

Example: http://sebastianraschka.com/Articles/2014_kernel_pca.html

RBF: radial basis function

http://scikit-learn.org/stable/auto_examples/plot_digits_pipe.html#sphx-glr-auto-examples-plot-digits-pipe-py

Reference

1. Pattern recognition and machine learning. Christopher M. Bishop. 2006
2. Lecture: kernel PCA. Unsupervised learning 2011. Rita Osadchy.
3. Implementing a Principal Component Analysis
http://sebastianraschka.com/Articles/2014_pca_step_by_step.html
4. Kernel tricks and nonlinear dimensionality reduction via RBF kernel PCA
http://sebastianraschka.com/Articles/2014_kernel_pca.html