

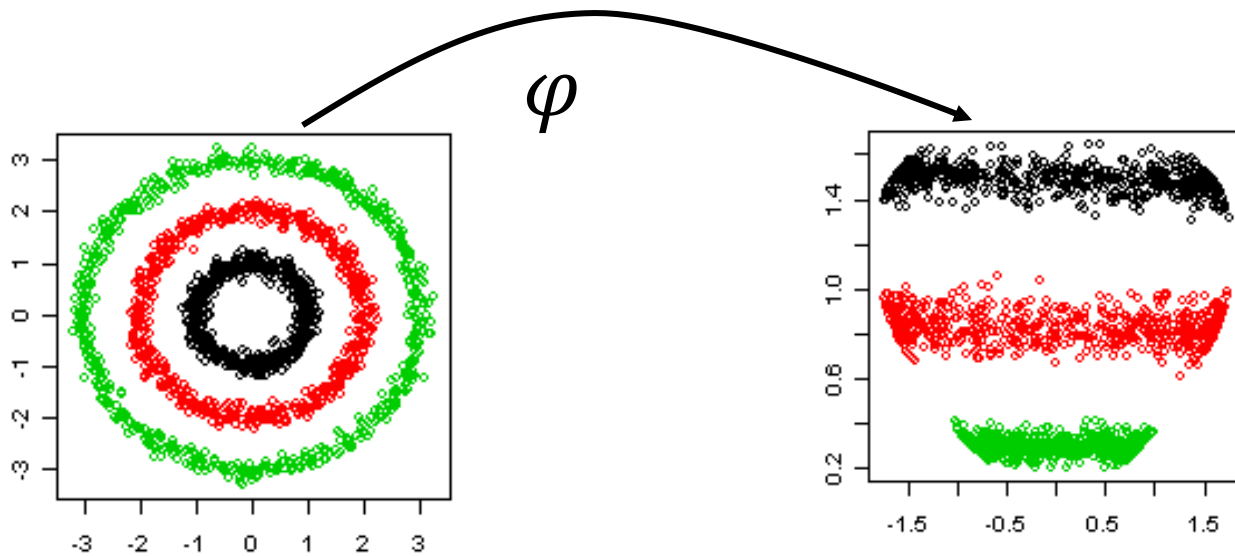
Course 495: Advanced Statistical Machine Learning/Pattern Recognition

- Goal (Lecture): To present Kernel Principal Component Analysis (KPCA) (and give a small flavour of Auto-encoders).

Materials

- Pattern Recognition & Machine Learning by C. Bishop Chapter 12
- **KPCA:** Schölkopf, Bernhard, Alexander Smola, and Klaus-Robert Müller. "Nonlinear component analysis as a kernel eigenvalue problem." *Neural computation* 10.5 (1998): 1299-1319.
- **Auto-Encoder:** Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *Science* 313.5786 (2006): 504-507.

Non-linear Component Analysis



$$\mathbf{x}_i \in \mathbb{R}^F \longrightarrow \varphi(\mathbf{x}_i) \in H$$

H can be of arbitrary dimensionality
(could be even infinite)

Kernel Principal Component Analysis

- $\varphi(\cdot)$ may not be explicitly known or is extremely expensive to compute and store.
- What is explicitly known is the dot product in H (also known as kernel k)

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

$$(\mathbf{x}_i, \mathbf{x}_j) \in R^F \times R^F \longrightarrow k(\cdot, \cdot) \in R$$

- All positive (semi)-definite functions can be used as kernels

KPCA-Kernel Matrix

- Given a training population of n samples $[\mathbf{x}_1, \dots, \mathbf{x}_n]$ we compute the training kernel matrix (also called Gram matrix).

$$\mathbf{K} = [\varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)] = [k(\mathbf{x}_i, \mathbf{x}_j)]$$

- All the computations are performed via the use of the kernel or the centralized kernel matrix.

$$\bar{\mathbf{K}} = (\varphi(\mathbf{x}_i) - \mathbf{m}^\Phi)^T (\varphi(\mathbf{x}_j) - \mathbf{m}^\Phi) \quad \mathbf{m}^\Phi = \frac{1}{n} \sum_{i=1}^n \varphi(\mathbf{x}_i)$$

KPCA- Popular Kernels

Gaussian Radial Basis Function (RBF) kernel:

$$k(\mathbf{x}_i, \mathbf{x}_j) = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / r^2}$$

Polynomial kernel: $k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + b)^n$

Hyperbolic Tangent kernel:

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

KPCA-kernel Matrix

Input space: $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$

Feature space: $\mathbf{X}^\Phi = [\varphi(\mathbf{x}_1), \dots, \varphi(\mathbf{x}_n)]$

Centralised: $\bar{\mathbf{X}}^\Phi = [\varphi(\mathbf{x}_1) - \mathbf{m}^\Phi, \dots, \varphi(\mathbf{x}_n) - \mathbf{m}^\Phi]$
 $= \mathbf{X}^\Phi (\mathbf{I} - \mathbf{E}) = \mathbf{X}^\Phi \mathbf{M}, \quad \mathbf{E} = \frac{1}{n} \mathbf{1} \mathbf{1}^T$

Kernel: $\mathbf{K} = [\varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)] = [k(\mathbf{x}_i, \mathbf{x}_j)] = \mathbf{X}^{\Phi T} \mathbf{X}^\Phi$

Centralised Kernel:

$$\begin{aligned} \bar{\mathbf{K}} &= [(\varphi(\mathbf{x}_i) - \mathbf{m}^\Phi)^T (\varphi(\mathbf{x}_j) - \mathbf{m}^\Phi)] = (\mathbf{I} - \mathbf{E}) \mathbf{X}^{\Phi T} \mathbf{X}^\Phi (\mathbf{I} - \mathbf{E}) \\ &= (\mathbf{I} - \mathbf{E}) \mathbf{K} (\mathbf{I} - \mathbf{E}) = \mathbf{K} - \mathbf{E} \mathbf{K} - \mathbf{K} \mathbf{E} + \mathbf{E} \mathbf{K} \mathbf{E} \end{aligned}$$

KPCA-Optimization problem

- KPCA cost function

$$\begin{aligned} \mathbf{U}^{\Phi}_o &= \arg \max_{\mathbf{U}^{\Phi}} \text{tr}[\mathbf{U}^{\Phi T} \mathbf{S}_t^{\Phi} \mathbf{U}^{\Phi}] \\ &= \arg \max_{\mathbf{U}^{\Phi}} \text{tr}[\mathbf{U}^{\Phi T} \bar{\mathbf{X}}^{\Phi} \bar{\mathbf{X}}^{\Phi T} \mathbf{U}^{\Phi}] \end{aligned}$$

subject to $\mathbf{U}^{\Phi T} \mathbf{U}^{\Phi} = \mathbf{I}$

- The solution is given by the d eigenvectors that correspond to the d largest eigenvalues

$$\mathbf{S}_t^{\Phi} \mathbf{U}^{\Phi}_o = \mathbf{U}^{\Phi}_o \mathbf{\Lambda}$$

KPCA- Computing Principal Components

- Do you see any problem with that?
- How can we compute the eigenvectors of \mathbf{S}_t^Φ ?
We do not even know φ !!!!
- Remember our Lemma that links the eigenvectors and eigenvalues of matrices of the form $\mathbf{A}\mathbf{A}^T$ and $\mathbf{A}^T\mathbf{A}$

$$\bar{\mathbf{K}} = \bar{\mathbf{X}}^\Phi{}^T \bar{\mathbf{X}}^\Phi = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T \quad \text{then } \mathbf{U}^\Phi_o = \bar{\mathbf{X}}^\Phi \mathbf{V} \mathbf{\Lambda}^{-\frac{1}{2}}$$

- All computations are performed via the use of $\bar{\mathbf{K}}$ (so-called kernel trick)

KPCA- Computing Principal Components

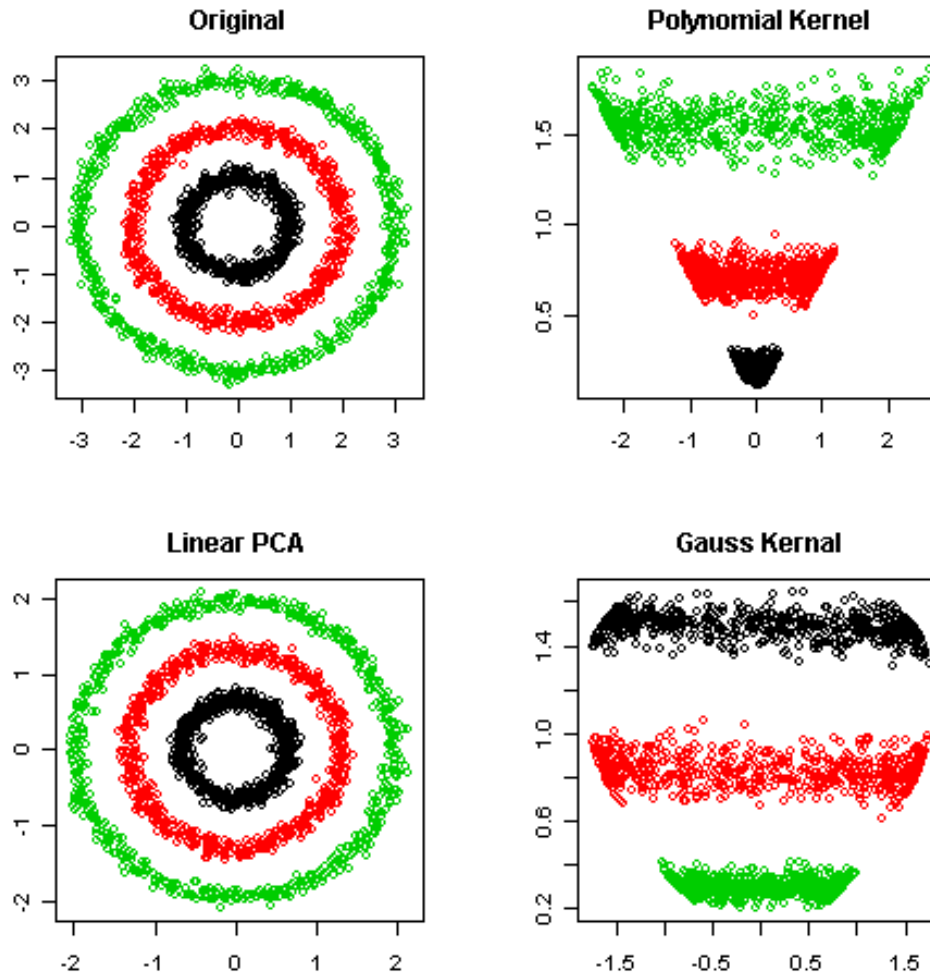
- Still $\mathbf{U}^{\Phi}_o = \bar{\mathbf{X}}^{\Phi} \mathbf{V} \mathbf{\Lambda}^{-\frac{1}{2}}$ cannot be analytically computed.
- But we do not want to compute \mathbf{U}^{Φ}_o .
- What we want is to compute latent features.
- That is, given a test sample \mathbf{x}_t we want to compute $\mathbf{y} = \mathbf{U}^{\Phi}_o{}^T \varphi(\mathbf{x}_t)$ (this can be performed via the kernel trick)

KPCA- Extracting Latent Features

$$\begin{aligned} \mathbf{y} &= \mathbf{U}^{\Phi}{}^T (\varphi(\mathbf{x}_t) - \mathbf{m}^{\Phi}) \\ &= \Lambda^{-\frac{1}{2}} \mathbf{V}^T \bar{\mathbf{X}}^{\Phi T} (\varphi(\mathbf{x}_t) - \mathbf{m}^{\Phi}) \\ &= \Lambda^{-\frac{1}{2}} \mathbf{V}^T (\mathbf{I} - \mathbf{E}) \mathbf{X}^{\Phi T} \left(\varphi(\mathbf{x}_t) - \frac{1}{n} \mathbf{X}^{\Phi} \mathbf{1} \right) \\ &= \Lambda^{-\frac{1}{2}} \mathbf{V}^T (\mathbf{I} - \mathbf{E}) \left(\mathbf{X}^{\Phi T} \varphi(\mathbf{x}_t) - \frac{1}{n} \mathbf{X}^{\Phi T} \mathbf{X}^{\Phi} \mathbf{1} \right) \\ &= \Lambda^{-\frac{1}{2}} \mathbf{V}^T (\mathbf{I} - \mathbf{E}) \left(g(\mathbf{x}_t) - \frac{1}{n} \mathbf{K} \mathbf{1} \right) \end{aligned}$$

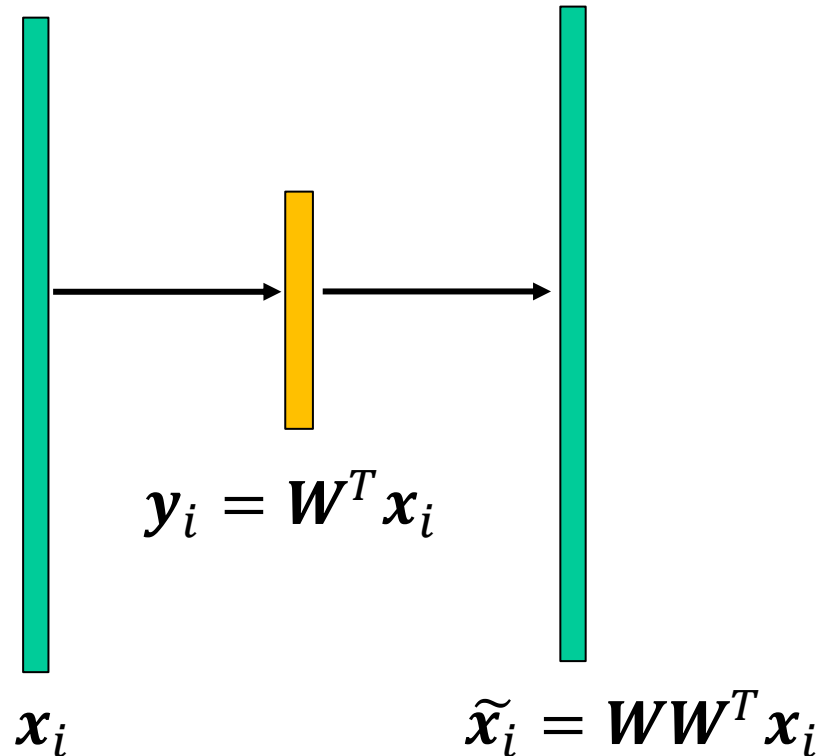
$$g(\mathbf{x}_t) = \mathbf{X}^{\Phi T} \varphi(\mathbf{x}_t) = \begin{bmatrix} \varphi(\mathbf{x}_1)^T \varphi(\mathbf{x}_t) \\ \dots \\ \varphi(\mathbf{x}_n)^T \varphi(\mathbf{x}_t) \end{bmatrix} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_t) \\ \dots \\ k(\mathbf{x}_n, \mathbf{x}_t) \end{bmatrix}$$

KPCA- Example

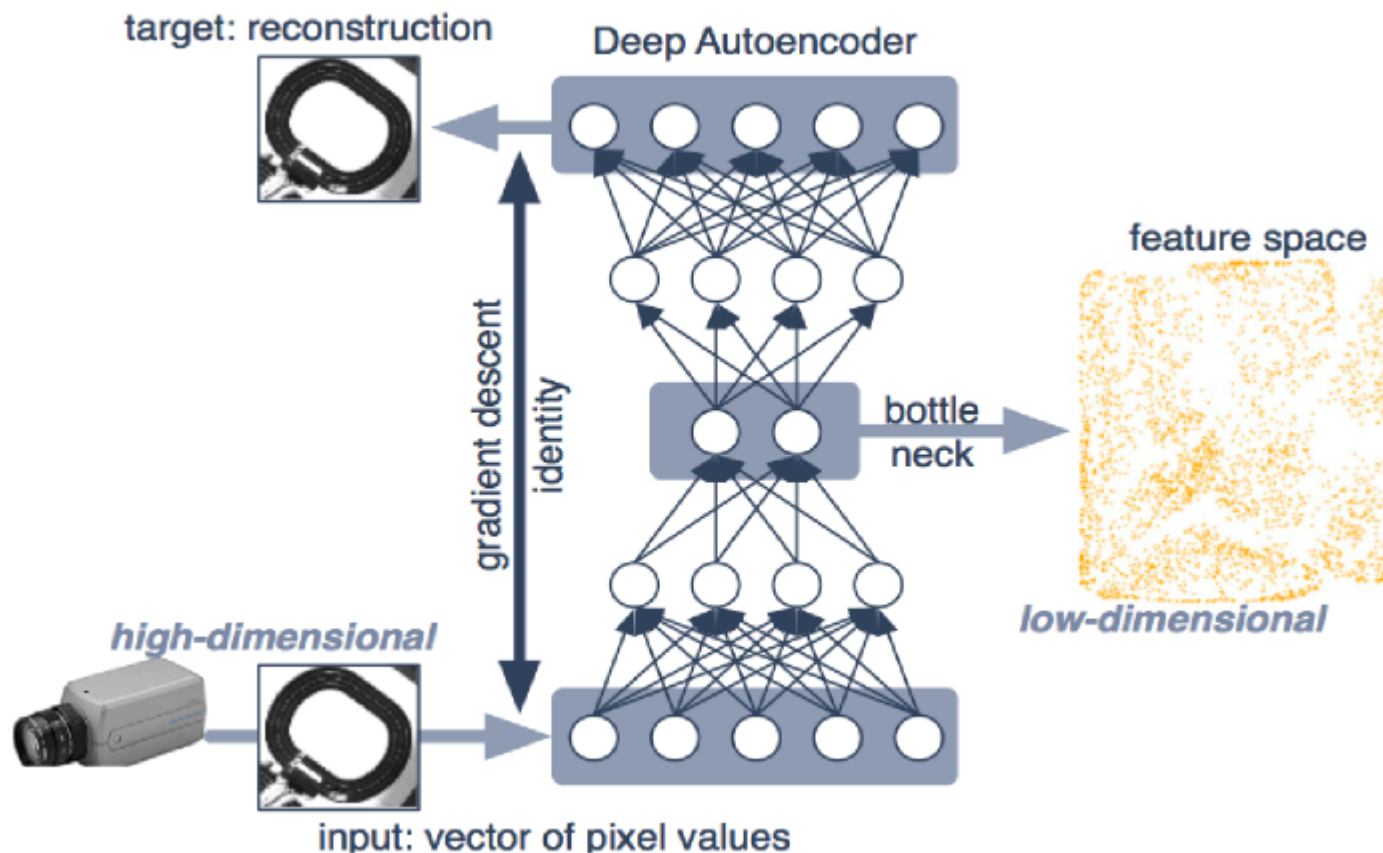


Latent Feature Extraction with Neural Networks

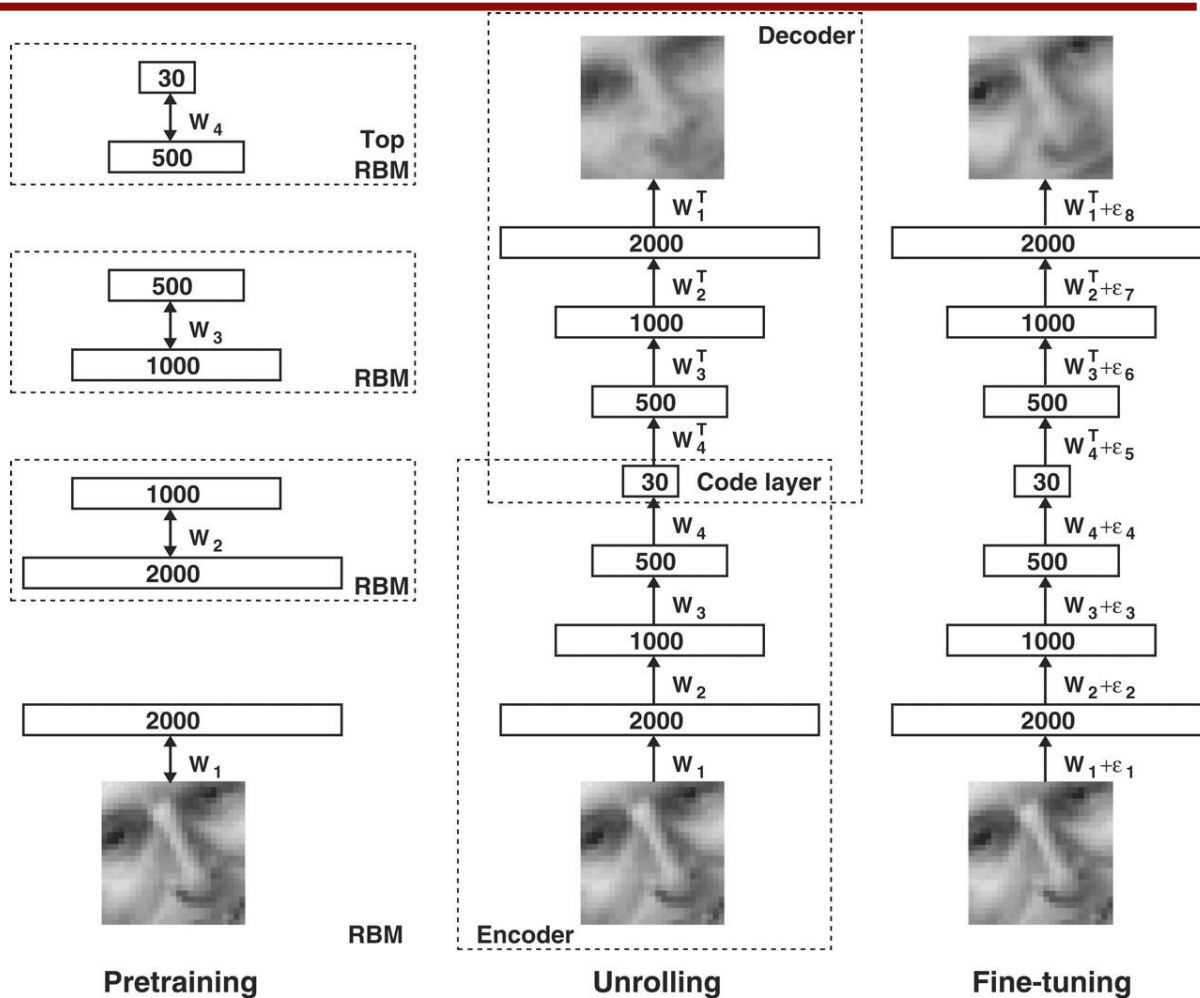
Remember the PCA model?



Latent Feature Extraction with Neural Networks



Latent Feature Extraction with Neural Networks



•G E Hinton, and R R Salakhutdinov Science 2006;313:504-507

Toolboxes on Component Analysis

Matlab Toolbox for Dimensionality Reduction

http://homepage.tudelft.nl/19j49/Matlab_Toolbox_for_Dimensionality_Reduction.html

Matlab 2013b has PPCA implemented

<http://www.mathworks.co.uk/help/stats/ppca.html>

Toolboxes on Component Analysis

ICA toolboxes for image and signal processing:

<http://www.bsp.brain.riken.jp/ICALAB/>

ICA for EEG Analysis:

<http://mialab.mrn.org/software/>

FastICA

<http://research.ics.aalto.fi/ica/fastica/>