

Introduction to Machine Learning

Unsupervised Learning: Dimension Reduction

Prof. Andreas Krause
Learning and Adaptive Systems (las.ethz.ch)

Typical approaches

- Assume $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$
- Obtain mapping $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^k$ where $k \ll d$
- Can distinguish
 - Linear dimension reduction: $\mathbf{f}(\mathbf{x}) = \underline{\mathbf{Ax}}$
 - Nonlinear dimension reduction
(parametric or non-parametric) $A \in \mathbb{R}^{k \times d}$
- **Key question:** Which mappings should we prefer?

Principal component analysis (PCA)

- Given centered data $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$, $1 \leq k \leq d$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \quad \mu = \frac{1}{n} \sum_i \mathbf{x}_i = 0$$

- The solution to the PCA problem

$$(\mathbf{W}, \mathbf{z}_1, \dots, \mathbf{z}_n) = \arg \min \sum_{i=1}^n \|\mathbf{W}\mathbf{z}_i - \mathbf{x}_i\|_2^2$$

where $\mathbf{W} \in \mathbb{R}^{d \times k}$ is orthogonal, $\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^k$

is given by $\mathbf{W} = (\mathbf{v}_1 | \dots | \mathbf{v}_k)$ and $\mathbf{z}_i = \mathbf{W}^T \mathbf{x}_i$

where

$$\Sigma = \sum_{i=1}^d \lambda_i \mathbf{v}_i \mathbf{v}_i^T \quad \lambda_1 \geq \dots \geq \lambda_d \geq 0$$

$= f(x_i)$

Side note: PCA vs. K-Means

PCA Problem:

$$(\mathbf{W}, \mathbf{z}_1, \dots, \mathbf{z}_n) = \arg \min \sum_{i=1}^n \|\mathbf{W}\mathbf{z}_i - \mathbf{x}_i\|_2^2$$

where $\mathbf{W} \in \mathbb{R}^{d \times k}$ is orthogonal, $\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^k$

k-Means problem: (equivalent formulation)

$$\mathbf{W} = (\mu_1 | \dots | \mu_k)$$

$$(\mathbf{W}, \mathbf{z}_1, \dots, \mathbf{z}_n) = \arg \min \sum_{i=1}^n \|\mathbf{W}\mathbf{z}_i - \mathbf{x}_i\|_2^2$$

$$\mathbf{W} \cdot e_j = \mu_j$$

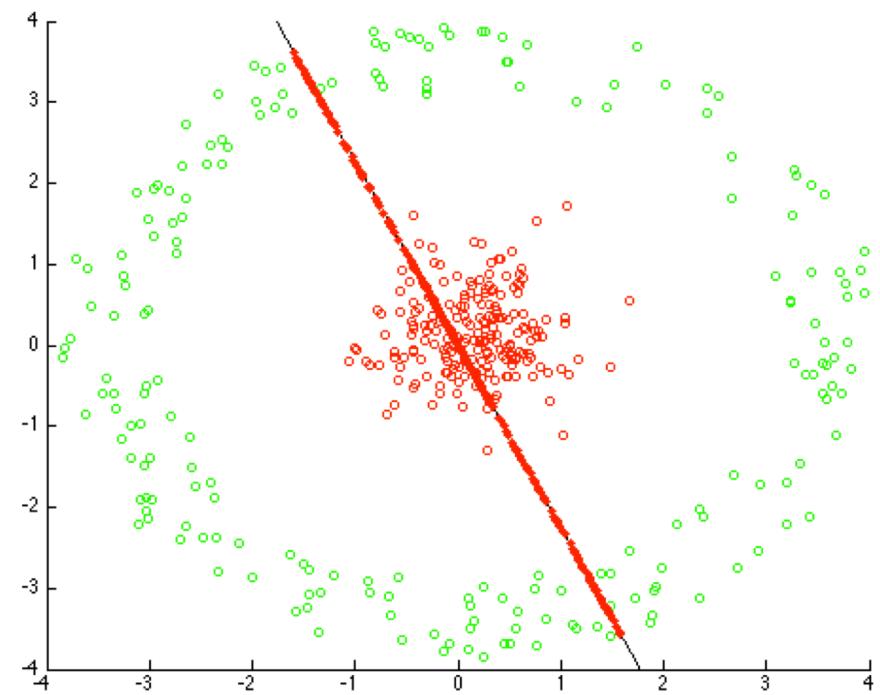
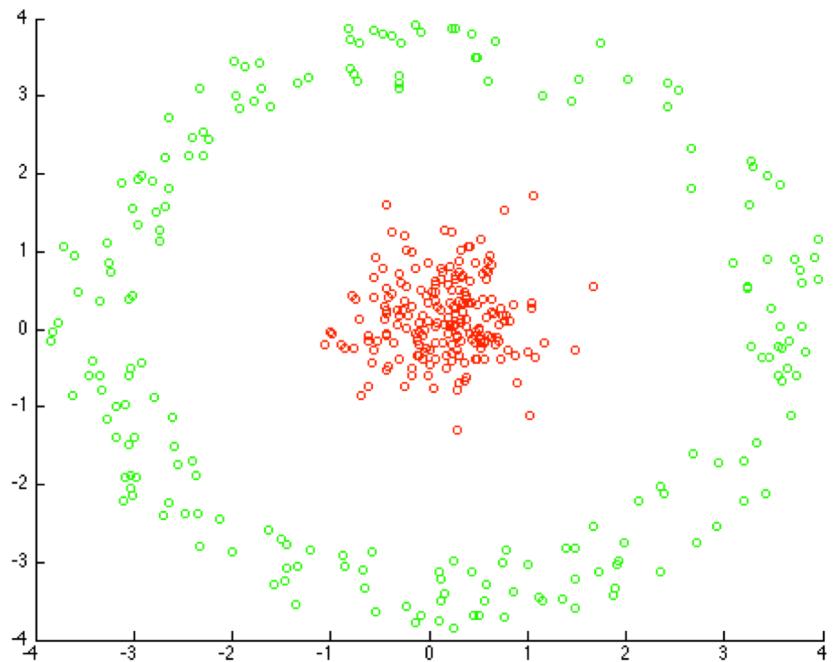
where $\mathbf{W} \in \mathbb{R}^{d \times k}$ arbitrary and $\mathbf{z}_1, \dots, \mathbf{z}_n \in E_k$

hereby, $E_k = \{[1, 0, \dots, 0], \dots, [0, \dots, 0, 1]\}$

is the set of unit vectors in \mathbb{R}^k

$$e_j = [0 \dots \underset{j\text{-th pos.}}{\underset{\sim}{1, 0, \dots, 0}}] \in E_k$$

Another example



Use Kernels!

- Recall: In supervised learning, kernels allowed us to solve non-linear problems by reducing them to linear ones in high-dimensional (implicitly represented) spaces
- Can take the same approach for unsupervised learning!

$$\text{Ansatz: } w = \sum_{i=1}^n \alpha_i x_i$$

$$x \in \mathbb{R}^d \xrightarrow{\phi(x)} \phi(x) \in \mathbb{R}^D \xrightarrow{z = w^T \phi(x)} z \in \mathbb{R}^k$$

Why: for PCA, w^\perp is leading eigenvector of $X^T X$

$$\Rightarrow \underbrace{X^T X w}_{\beta} = \lambda w \Rightarrow w = \frac{1}{\lambda} \underbrace{\sum_i \beta_i x_i}_{X^T \beta}$$

$$\Rightarrow \alpha_i = \frac{\beta_i}{\lambda}$$

Recall PCA for k=1

- Optimal solution to PCA problem solves, for $\Sigma = \frac{1}{n} X^T X$

$$\arg \max_{\substack{\|\mathbf{w}\|_2=1 \\ (\text{opt})}} \mathbf{w}^T X^T X \mathbf{w} = \arg \max_{\|\mathbf{w}\|_2=1} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i)^2$$

Using Ansatz $\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i$, show that both
Objective (A) and constraints (opt) can be expressed
in terms of α and inner products of $\mathbf{x}_i, \mathbf{x}_j$.

Derivation continued

Objective

$$\begin{aligned}
 w^T X^T X w &= \sum_{i=1}^n (w^T x_i)^2 = \sum_{i=1}^n \left(\left(\sum_{j=1}^n \alpha_j x_j \right)^T x_i \right)^2 \\
 &= \sum_{i=1}^n \left[\sum_j \alpha_j \underbrace{(x_j^T x_i)}_{k(x_i, x_j)} \right]^2 = \sum_{i=1}^n \left(\sum_j \alpha_j k(x_i, x_j) \right)^2 = \sum_{i=1}^n (\alpha^T k_i)^2 \\
 &= \alpha^T k^T k \alpha \rightarrow \max
 \end{aligned}$$

$k = \begin{pmatrix} k_{11} \\ \vdots \\ k_{nn} \end{pmatrix}$

Constraints

$$\begin{aligned}
 \|w\|_2^2 &= w^T w = \left(\sum_{i=1}^n \alpha_i x_i \right)^T \left(\sum_{j=1}^n \alpha_j x_j \right) = \sum_{i,j} \alpha_i \alpha_j \underbrace{x_i^T x_j}_{k(x_i, x_j)} \\
 &= \alpha^T k \alpha \stackrel{!}{=} 1
 \end{aligned}$$

Thus, kernel PCA requires solving

$$\max_{\alpha} \alpha^T k^T k \alpha \quad \text{s.t. } \alpha^T k \alpha = 1$$

regular PCA : $\underset{w}{\arg \max} \quad w^T X^T X w \quad \text{s.t. } w^T v = 1$

Recall PCA for k=1

- Optimal solution to PCA problem solves, for $\Sigma = \frac{1}{n} \mathbf{X}^T \mathbf{X}$

$$\arg \max_{\|\mathbf{w}\|_2=1} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} = \arg \max_{\|\mathbf{w}\|_2=1} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i)^2$$

- Applying feature maps, using $\mathbf{w} = \sum_{j=1}^n \alpha_j \phi(\mathbf{x}_j)$ and observing $\|\mathbf{w}\|_2^2 = \alpha^T \mathbf{K} \alpha$

$$\begin{aligned} & \arg \max_{\|\mathbf{w}\|_2=1} \sum_{i=1}^n (\mathbf{w}^T \phi(\mathbf{x}_i))^2 = \arg \max_{\alpha^T \mathbf{K} \alpha=1} \sum_{i=1}^n \left(\sum_{j=1}^n \alpha_j \phi(\mathbf{x}_j)^T \phi(\mathbf{x}_i) \right)^2 \\ &= \arg \max_{\alpha^T \mathbf{K} \alpha=1} \sum_{i=1}^n \left(\sum_{j=1}^n \alpha_j k(\mathbf{x}_j, \mathbf{x}_i) \right)^2 = \arg \max_{\alpha^T \mathbf{K} \alpha=1} \sum_{i=1}^n \left(\alpha^T \mathbf{K}_i \right)^2 \\ &= \arg \max_{\alpha^T \mathbf{K} \alpha=1} \alpha^T \mathbf{K}^T \mathbf{K} \alpha \end{aligned}$$

Kernel PCA (k=1)

- The Kernel-PCA problem (k=1) requires solving

$$\alpha^* = \arg \max_{\alpha^T \mathbf{K} \alpha = 1} \alpha^T \mathbf{K}^T \mathbf{K} \alpha$$

- The optimal solution is obtained in **closed form** from the eigendecomposition of \mathbf{K} :

$$\alpha^* = \frac{1}{\sqrt{\lambda_1}} \mathbf{v}_1$$

Why scaling?

$$\begin{aligned}\alpha^T \mathbf{K} \alpha &= \frac{1}{\sqrt{\lambda_1}} \mathbf{v}_1^T \mathbf{K} \frac{1}{\sqrt{\lambda_1}} \mathbf{v}_1 \\ &= \frac{1}{\lambda_1} \cdot \mathbf{v}_1^T \mathbf{K} \mathbf{v}_1 = \frac{1}{\lambda_1} \cdot \mathbf{v}_1^T (\lambda_1 \mathbf{v}_1) = \mathbf{v}_1^T \mathbf{v}_1 = 1\end{aligned}$$

$$\underbrace{\mathbf{K}}_{\mathbf{K} = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^T} \quad \lambda_1 \geq \dots \geq \lambda_d \geq 0$$

Kernel PCA (general k)

- For general $k \geq 1$, the Kernel Principal Components are given by $\alpha^{(1)}, \dots, \alpha^{(k)} \in \mathbb{R}^n$

where $\alpha^{(i)} \in \mathbb{R}^n$

$$\alpha^{(i)} = \frac{1}{\sqrt{\lambda_i}} \mathbf{v}_i$$

is obtained from: $\mathbf{K} = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^T \quad \lambda_1 \geq \dots \geq \lambda_d \geq 0$

- Given this, a new point x is projected as $z \in \mathbb{R}^k$

For $1 \leq i \leq k$: $z_i = w_i^T x = \left(\sum_{j=1}^n \alpha_j^{(i)} x_j \right)^T x = \sum_{j=1}^n \alpha_j^{(i)} \underbrace{x_j^T x}_{k(x_j, x)} = \sum_{j=1}^n \alpha_j^{(i)} k(x_j, x)$

Kernel PCA (general k)

- For general $k \geq 1$, the Kernel Principal Components are given by $\alpha^{(1)}, \dots, \alpha^{(k)} \in \mathbb{R}^n$

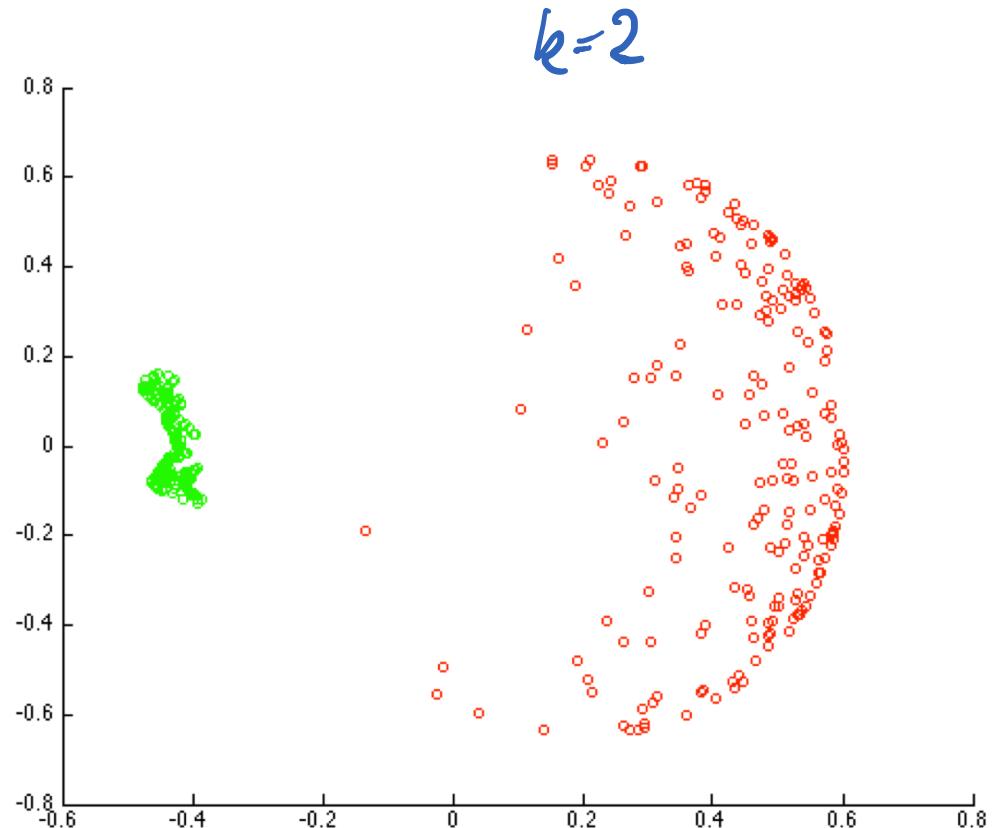
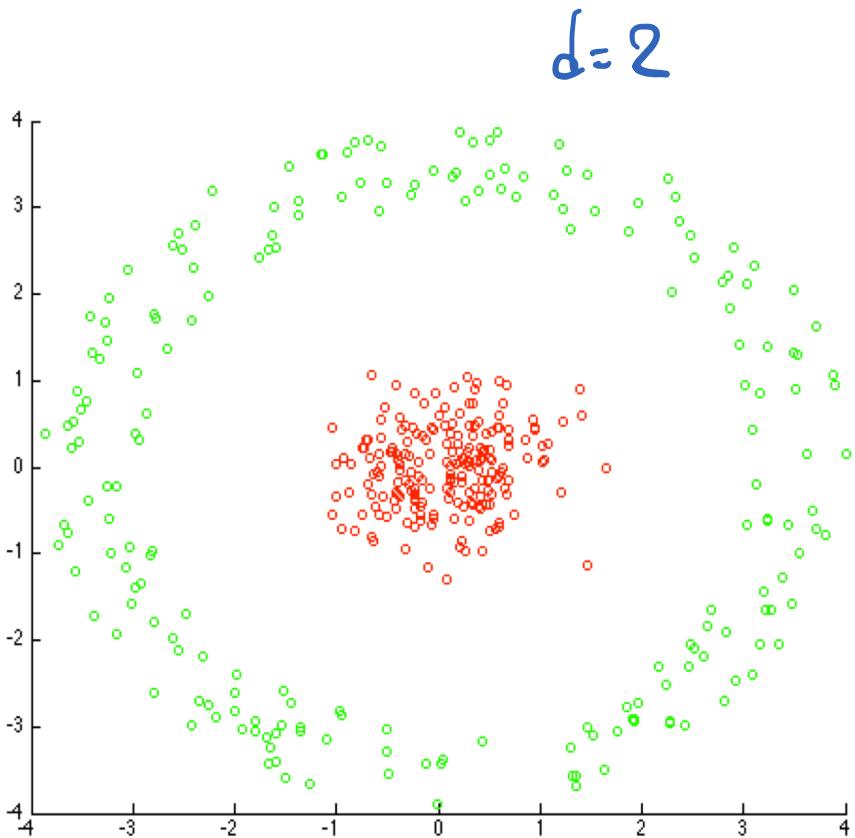
where $\alpha^{(i)} = \frac{1}{\sqrt{\lambda_i}} \mathbf{v}_i$

is obtained from: $\mathbf{K} = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^T \quad \lambda_1 \geq \dots \geq \lambda_d \geq 0$

- Given this, a new point \mathbf{x} is projected as $\mathbf{z} \in \mathbb{R}^k$

$$z_i = \sum_{j=1}^n \alpha_j^{(i)} k(\mathbf{x}, \mathbf{x}_j)$$

Illustration



RBF
Side note: Applying k-means on kernel principal components
is sometimes called Kernel-k-means or Spectral Clustering

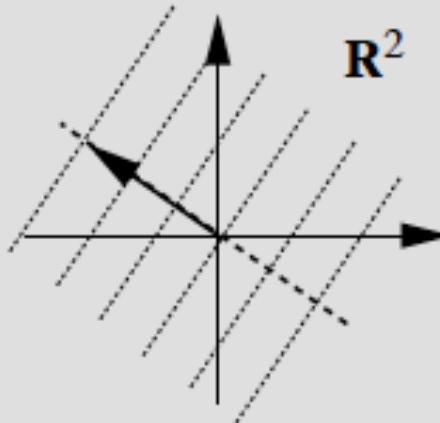
Kernel-PCA/k-means demo

Notes on Kernel-PCA

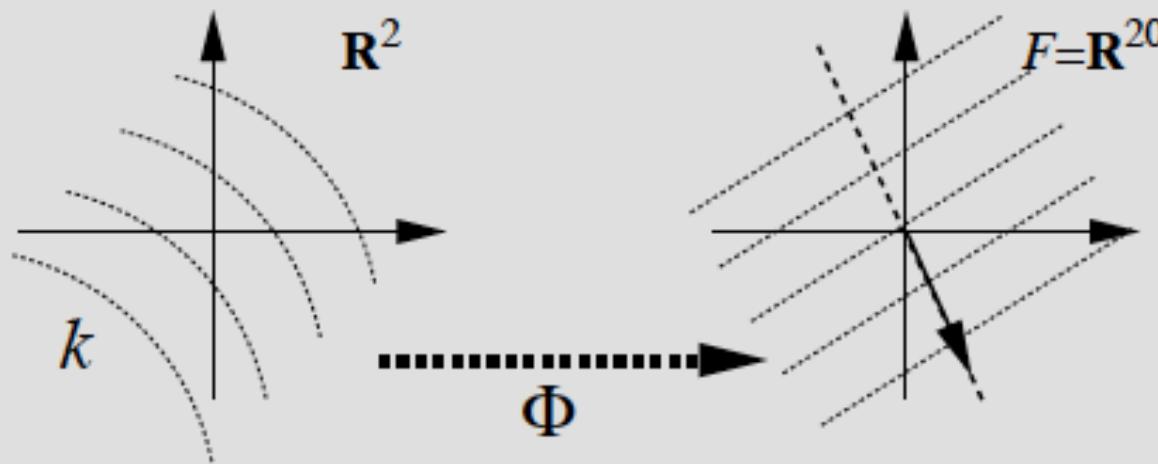
- Kernel-PCA corresponds to applying PCA in the feature space induced by the kernel k
- Can be used to discover non-linear feature maps in closed form
- This can be used as inputs, e.g., to SVMs given „multi-layer support vector machines“
- May want to „center“ the kernel:
$$\mathbf{E} = \frac{1}{n}[\mathbf{1}, \dots, \mathbf{1}][\mathbf{1}, \dots, \mathbf{1}]^T$$
$$\mathbf{K}' = \mathbf{K} - \mathbf{K}\mathbf{E} - \mathbf{E}\mathbf{K} + \mathbf{E}\mathbf{K}\mathbf{E}$$

Illustration: Kernel-PCA [Schölkopf et al.]

linear PCA



kernel PCA



Other non-linear methods

- Kernel-PCA requires data specified as kernel
 - Complexity grows with the number of data points
 - Cannot easily „explicitly“ embed high-dimensional data (unless we have an appropriate kernel)
- Alternatives
 - [Autoencoders](#)
 - [Locally linear embedding \(LLE\)](#)
 - [Multi-dimensional scaling](#)
 - [t-SNE](#)
 - ...

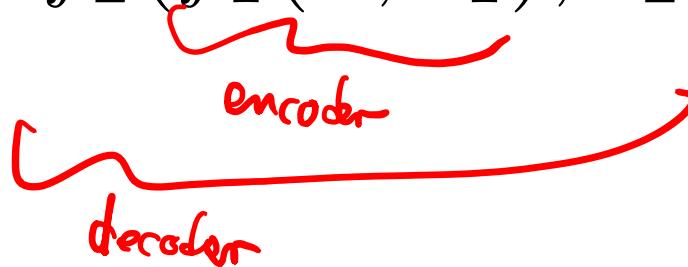
Dimension reduction via Autoencoders

- Key idea: Try to learn the identity function!

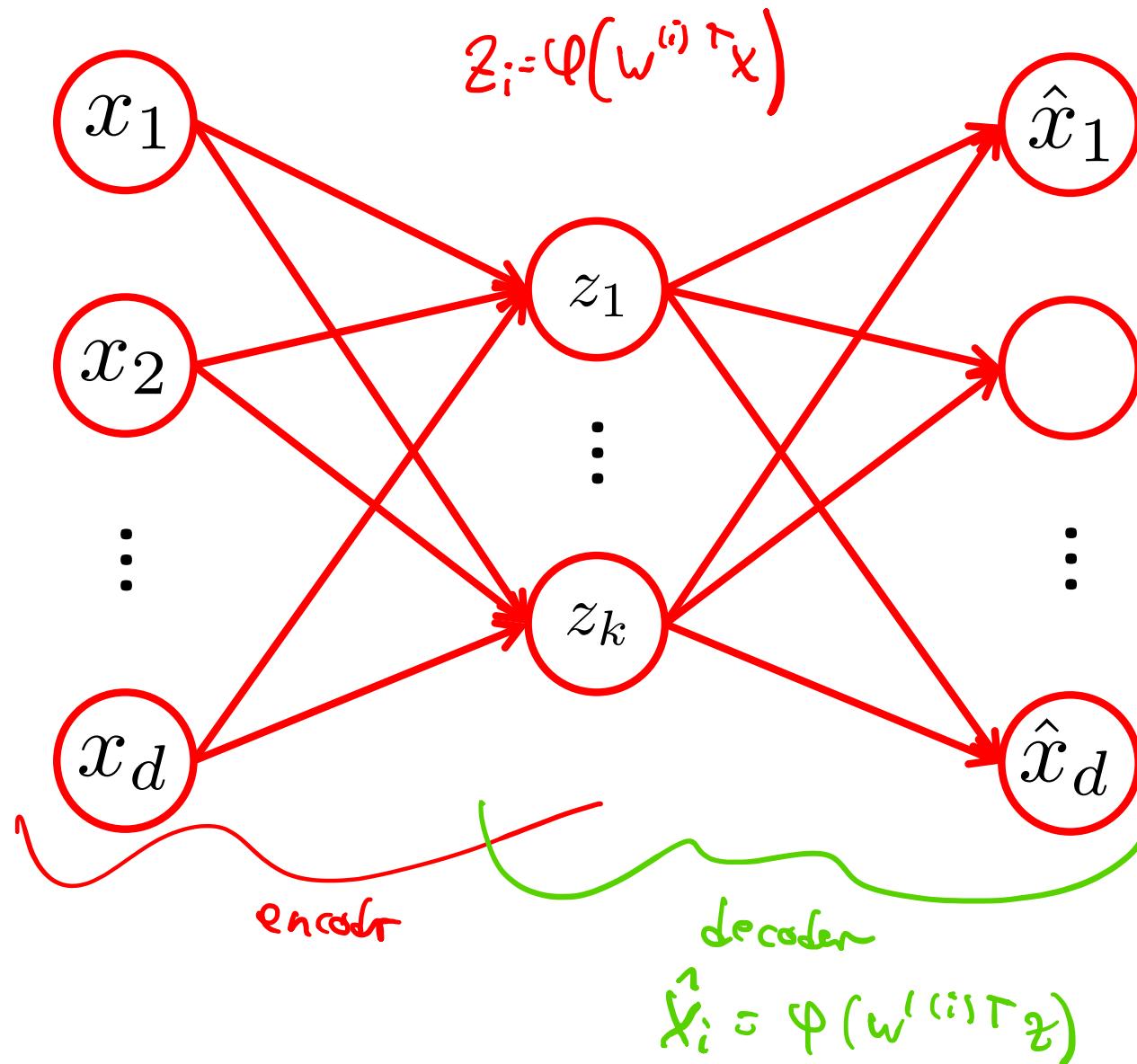
$$\mathbf{x} \approx f(\mathbf{x}; \theta)$$

- What function f should we pick?

$$f(\mathbf{x}; \theta) = f_2(f_1(\mathbf{x}; \theta_1); \theta_2)$$



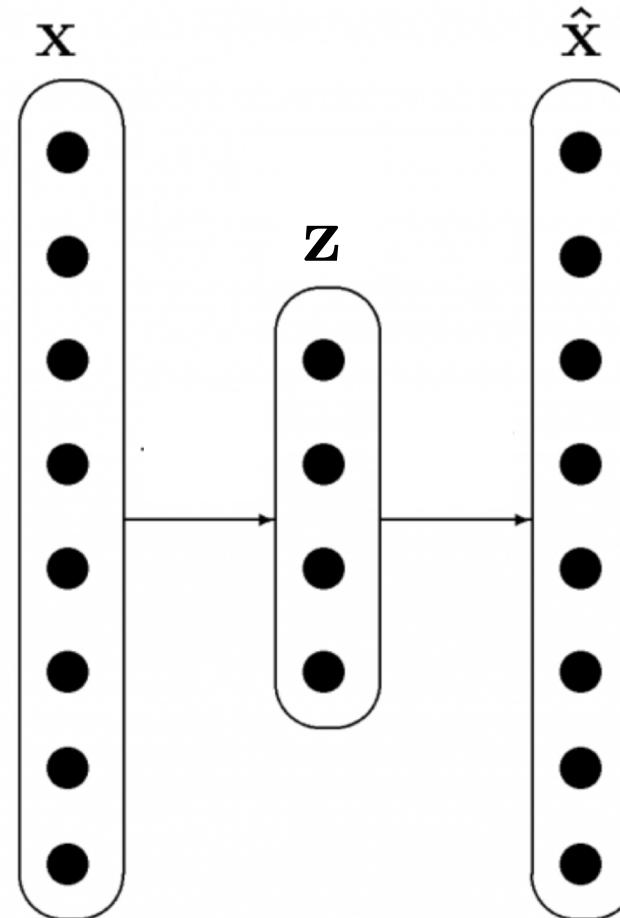
Autoencoders via Neural Nets



Neural network autoencoders

- Neural network Autodecoders are ANNs where
 - There is one output unit for each of the d input units
 - The number k of hidden units is usually smaller than the number of inputs
- The goal is to optimize the weights such that the output agrees with the input

Example



Training autoencoders

- The goal is to optimize the weights such that the **output agrees** with the input
 - For example, minimize the square loss

$$\min_{\mathbf{W}} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{f}(\mathbf{x}_i; \mathbf{W})\|_2^2$$

$= \sum_{j=1}^d (x_{i,j} - f_j(x_i; \mathbf{w}))^2$

- Find local minimum via SGD (backpropagation)
- Initialization matters and is challenging
 - C.f., work on „pretraining“ (e.g., layerwise training of restricted Boltzmann machines)

Autoencoders vs. PCA

- The internal representation $z = \varphi(W^{(1)}x)$ is the „dimensionality reduced“ input x
- If the activation function is the identity $\varphi(z) = z$ then in fact fitting a NN autoencoder is equivalent to **PCA!**

$$f(x) = f_z(f_1(x; W_1); W_2) = W_2 W_1 \cdot x$$

$W_1 \in \mathbb{R}^{k \times d}$
 $W_2 \in \mathbb{R}^{d \times k}$

The optimal solution to $\min_{W_1, W_2} \sum_i \|x_i - W_2 W_1 x_i\|_2^2$

is given by $W_1 = W^\top$ and $W_2 = V$, where W^\top is the PCA solution.
 $f(x) = W W^\top x$

Reconstruction comparison



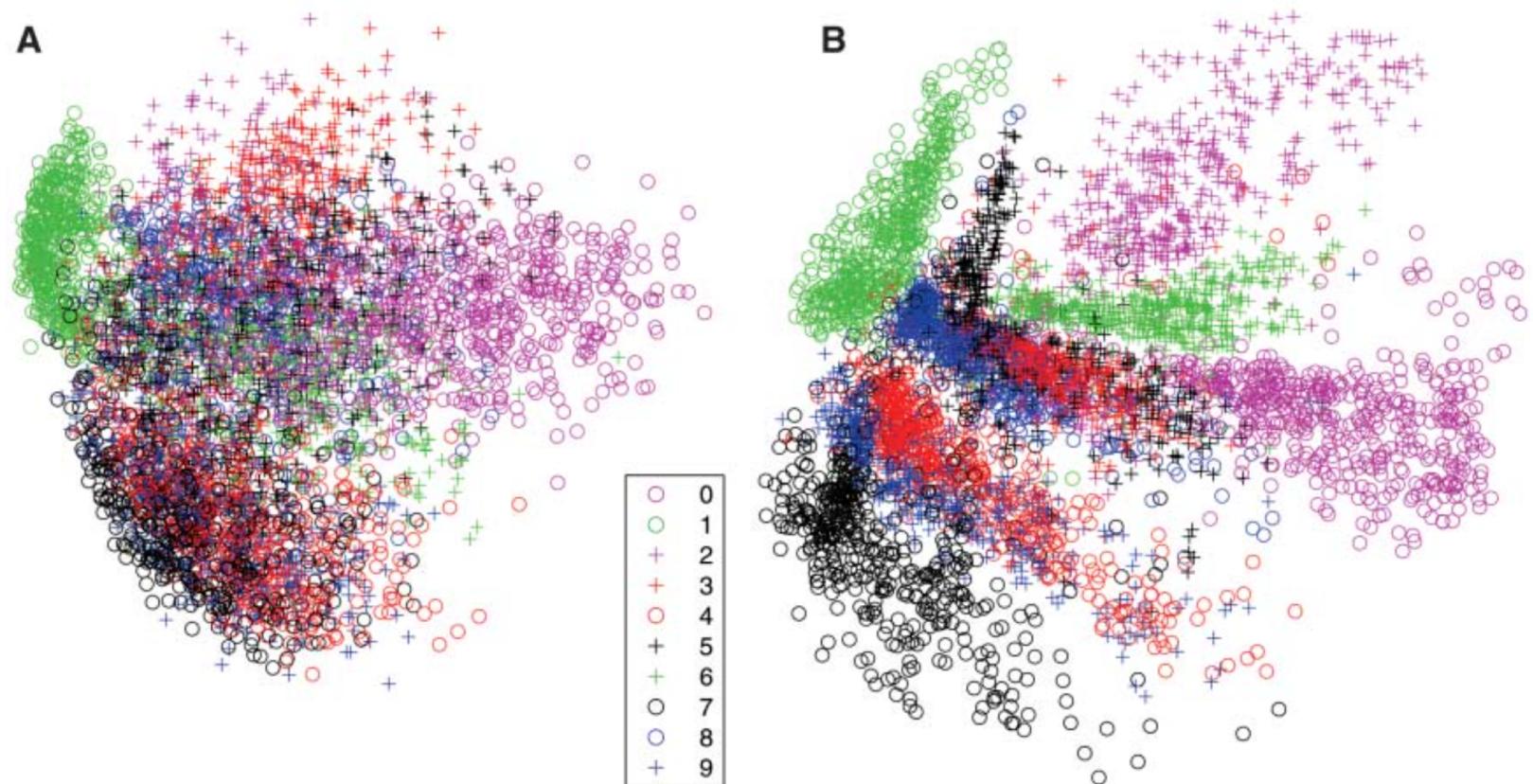
Top: original image

Middle: 30-dim Autoencoder

Bottom: 30-dim PCA

Dimension reduction

Fig. 3. (A) The two-dimensional codes for 500 digits of each class produced by taking the first two principal components of all 60,000 training images. (B) The two-dimensional codes found by a 784-1000-500-250-2 autoencoder. For an alternative visualization, see (8).

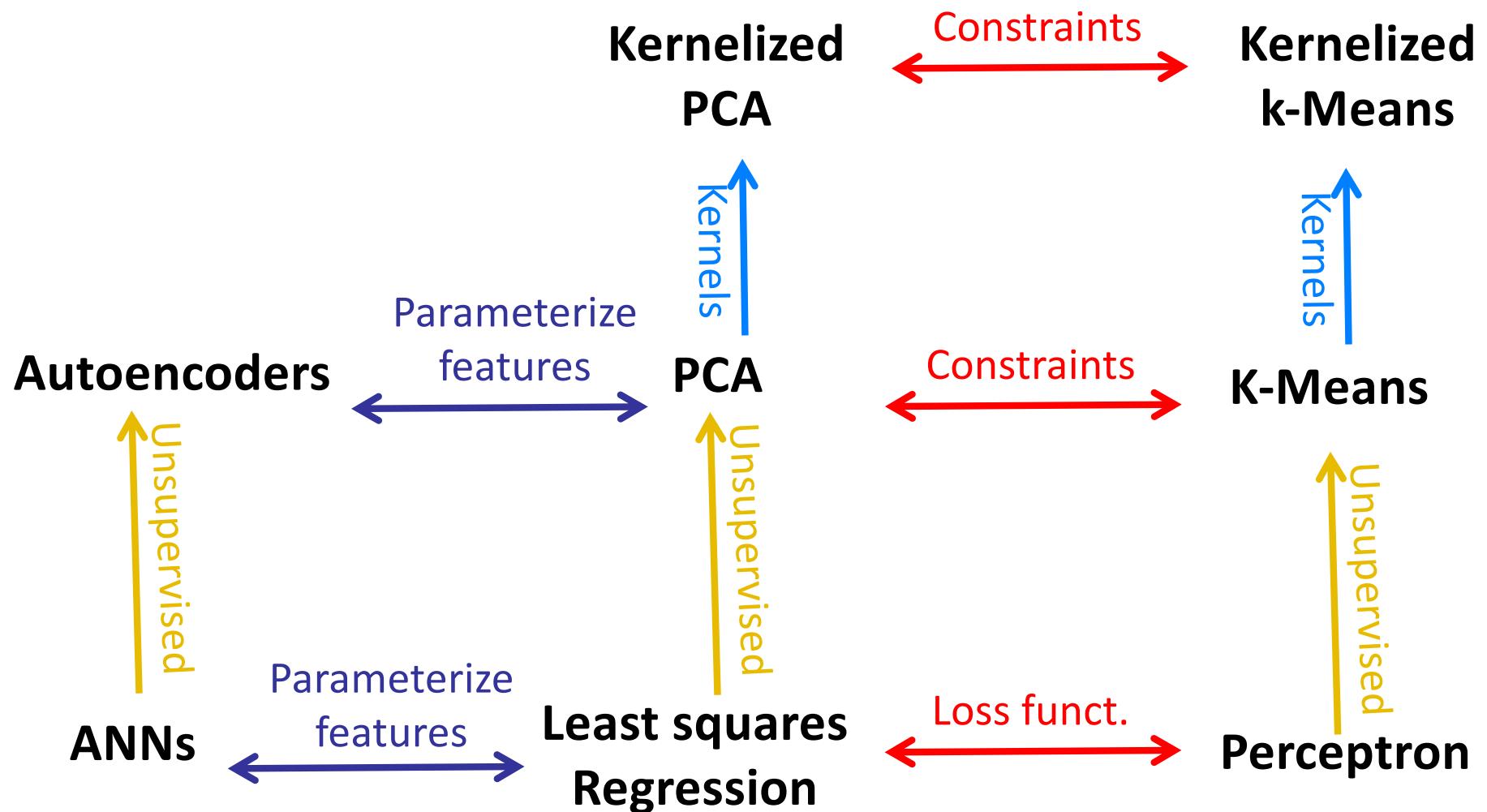


Autoencoder Demo

What you need to be able to do

- Apply Principal Component Analysis for linear dimension reduction
- Apply kernel PCA for nonlinear dimension reduction
- Understand how to train an neural network autoencoder via backpropagation

Supervised vs. unsupervised learning



Machine learning summary so far

Representation/ features	Linear hypotheses; nonlinear hypotheses with nonlinear feature transforms, kernels, learn nonlinear features via neural nets		
Model/ objective:	Loss-function	+	Regularization
	Squared loss, 0/1 loss, Perceptron loss, Hinge loss, cost sensitive losses, multi-class hinge loss, reconstruction error		L^2 norm, L^1 norm, early stopping, dropout
Method:	Exact solution, Gradient Descent, (mini-batch) SGD, Reductions, Lloyd's heuristic		
Evaluation metric:	Mean squared error, Accuracy, F1 score, AUC, Confusion matrices, compression performance		
Model selection:	K-fold Cross-Validation, Monte Carlo CV		