

Introduction to Machine Learning

Unsupervised Learning: Dimension Reduction

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

We will

- Introduce basic dimension reduction algorithms
 - Principal Component Analysis (PCA)
 - Kernel PCA
 - Neural network autoencoders
- Much more details in
 - Computational Intelligence Lab
 - Deep Learning

Basic challenge

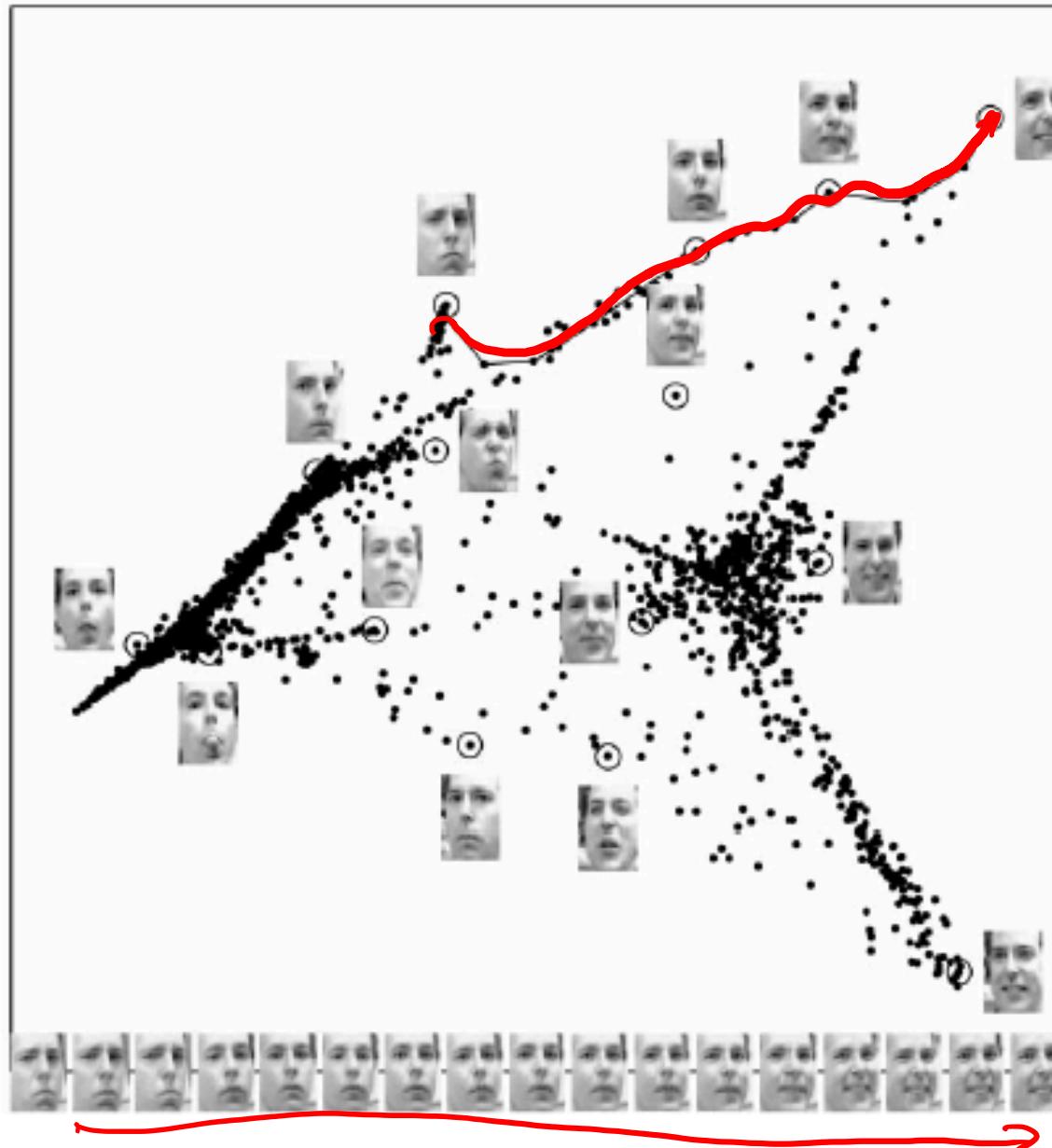
- Given data set $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \quad \mathbf{x}_i \in \mathbb{R}^d$
obtain „embedding“ (low-dimensional representation)

$$\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^k \quad k < d$$

- Motivation**
 - Visualization ($k=1,2,3$)
 - Regularization (model selection)
 - Unsupervised feature discovery
(i.e., determine features from data!)
 - ...

Example: Embedding of faces

[Saul & Roweis]

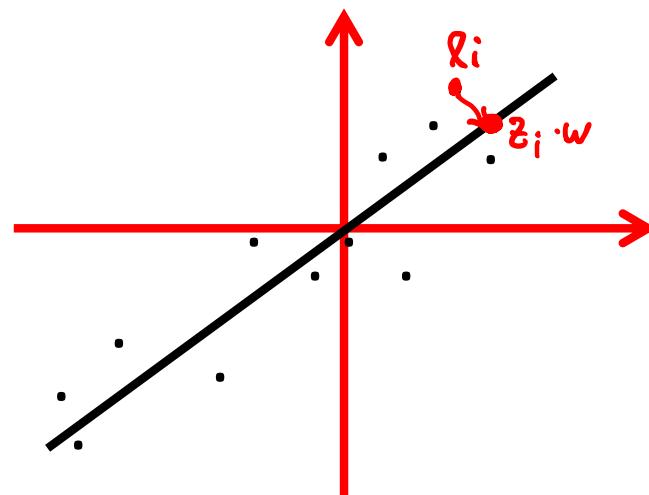


Typical approaches

- 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: $\mathbf{f}(\mathbf{x}) = \mathbf{Ax}$
 $A \in \mathbb{R}^{k \times d}$
 - Nonlinear dimension reduction
(parametric or non-parametric)
- **Key question:** Which mappings should we prefer?

Linear dim. reduction as *compression*

- **Motivation:** Low-dimensional representation should allow to *compress* original data (accurate reconstruction)
- **Example:** $k=1$
- Given data set $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$
- Want to represent data as points on a *line* $\underline{\mathbf{w}} \in \mathbb{R}^d$ with coefficients z_1, \dots, z_n
- I.e., want $z_i \mathbf{w} \approx \mathbf{x}_i$, assuming $\mu = \frac{1}{n} \sum_i \mathbf{x}_i = 0$



Linear dim. reduction for reconstruction

- Given data set $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$
- Want $z_i \mathbf{w} \approx \mathbf{x}_i$, e.g., minimizing $\|z_i \mathbf{w} - \mathbf{x}_i\|_2^2$
- To ensure uniqueness, normalize: $\|\mathbf{w}\|_2 = 1$

$$(z_i) \left(\frac{1}{\|\mathbf{w}\|_2} \mathbf{w} \right) = z_i \mathbf{w}$$

- Optimize over $\mathbf{w}, z_1, \dots, z_n$ jointly:

$$(\mathbf{w}^*, z_{1:n}^*) = \underset{\mathbf{w}, z_{1:n}}{\operatorname{arg\,min}} \sum_{i=1}^n \|x_i - z_i \cdot \mathbf{w}\|_2^2$$

Linear dim. reduction for reconstruction

- Given data set $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$
- Want $z_i \mathbf{w} \approx \mathbf{x}_i$, e.g., minimizing $\|z_i \mathbf{w} - \mathbf{x}_i\|_2^2$
- To ensure uniqueness, normalize: $\|\mathbf{w}\|_2 = 1$
- Optimize over $\mathbf{w}, z_1, \dots, z_n$ jointly:

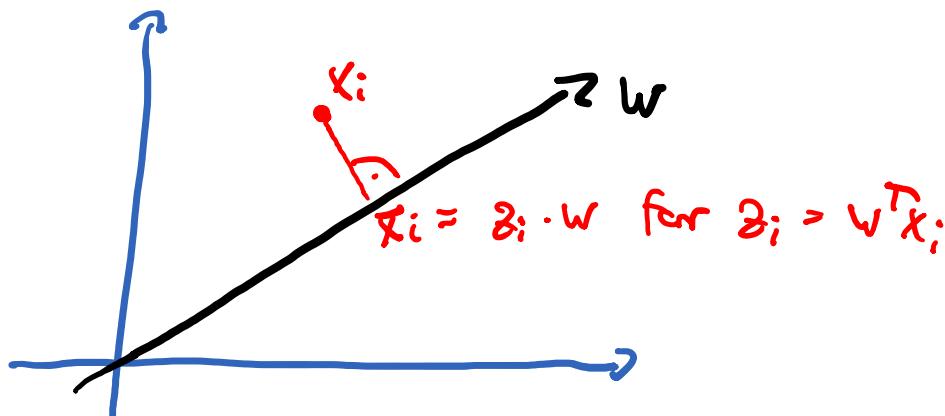
$$(\mathbf{w}^*, \mathbf{z}^*) = \arg \min_{\|\mathbf{w}\|_2=1, \mathbf{z}} \sum_{i=1}^n \|z_i \mathbf{w} - \mathbf{x}_i\|_2^2$$

$\hat{R}(\mathbf{w}, \mathbf{z}_{1:n})$

Solving for z given w

$$(\mathbf{w}^*, \mathbf{z}^*) = \arg \min_{\|\mathbf{w}\|_2=1, \mathbf{z}} \sum_{i=1}^n \|z_i \mathbf{w} - \mathbf{x}_i\|_2^2$$

- Suppose we consider some vector w . What can we say about the optimal z ?



Solving for z given \mathbf{w}

$$(\mathbf{w}^*, \mathbf{z}^*) = \arg \min_{\|\mathbf{w}\|_2=1, \mathbf{z}} \sum_{i=1}^n \|z_i \mathbf{w} - \mathbf{x}_i\|_2^2$$

- Suppose we consider some vector \mathbf{w} . What can we say about the optimal z ?

$$z_i^* = \mathbf{w}^T \mathbf{x}_i$$

- Thus, we effectively solve a *regression* problem, interpreting \mathbf{x} as features and z as labels!

Solving for z given w

- Want to solve

$$(\mathbf{w}^*, \mathbf{z}^*) = \arg \min_{\|\mathbf{w}\|_2=1, \mathbf{z}} \sum_{i=1}^n \|z_i \mathbf{w} - \mathbf{x}_i\|_2^2$$

- Note:** For any fixed $\|\mathbf{w}\|_2 = 1$ it holds that $z_i^* = \mathbf{w}^T \mathbf{x}_i$ therefore, only need

$$\mathbf{w}^* = \arg \min_{\|\mathbf{w}\|_2=1} \sum_{i=1}^n \|\mathbf{w} \mathbf{w}^T \mathbf{x}_i - \mathbf{x}_i\|_2^2$$

$$\begin{aligned} \hat{R}(\mathbf{w}) &= \sum_{i=1}^n (\mathbf{w} \mathbf{w}^T \mathbf{x}_i - \mathbf{x}_i)^T (\mathbf{w} \mathbf{w}^T \mathbf{x}_i - \mathbf{x}_i) \\ &= \sum_{i=1}^n \left[\underbrace{\mathbf{x}_i^T \mathbf{w} \mathbf{w}^T \mathbf{w}^T \mathbf{x}_i}_{1} - 2 \mathbf{x}_i^T \mathbf{w} \mathbf{w}^T \mathbf{x}_i + \underbrace{\mathbf{x}_i^T \mathbf{x}_i}_{2} \right] \\ &\quad - \underbrace{(\mathbf{x}_i^T \mathbf{w})(\mathbf{w}^T \mathbf{x}_i)}_{= -(\mathbf{w}^T \mathbf{x}_i)^2} \underbrace{\|\mathbf{x}_i\|_2^2}_{2} \\ &= \sum_{i=1}^n \|\mathbf{x}_i\|_2^2 - \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i)^2 \end{aligned}$$

Constant w.r.t w

$\Rightarrow \underset{\mathbf{w}}{\operatorname{argmin}} \hat{R}(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i)^2$

Solving for \mathbf{w}

- The objective

$$\mathbf{w}^* = \arg \min_{\|\mathbf{w}\|_2=1} \sum_{i=1}^n \|\mathbf{w}\mathbf{w}^T \mathbf{x}_i - \mathbf{x}_i\|_2^2$$

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

$$\begin{aligned} A &= \sum_{i=1}^n \mathbf{w}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{w} = \mathbf{w}^T \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{w} \\ &= n \cdot \underbrace{\sum}_{\substack{\rightarrow \text{empirical covariance of data} \\ (\text{assuming mean } 0)}} \end{aligned}$$

Solving for \mathbf{w}

- Further:

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

is equivalent to: $\mathbf{w}^* = \arg \max_{\|\mathbf{w}\|_2=1} \mathbf{w}^T \Sigma \mathbf{w}$

where $\Sigma = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$ is the *empirical covariance*

assuming the data is centered: $\mu = \frac{1}{n} \sum_i \mathbf{x}_i = 0$

Solving for \mathbf{w}

$$\underline{\Sigma} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$$

- The optimal solution to

$$\mathbf{w}^* = \arg \max_{\|\mathbf{w}\|_2=1} \mathbf{w}^T \Sigma \mathbf{w} \quad \mu = \frac{1}{n} \sum_i \mathbf{x}_i = 0$$

is given by the principal eigenvector of $\underline{\Sigma}$

i.e., $\mathbf{w}^* = \mathbf{v}_1$ where

$$\Sigma = \sum_{i=1}^d \lambda_i \mathbf{v}_i \mathbf{v}_i^T \quad \boxed{\lambda_1 \geq \dots \geq \lambda_d \geq 0 \quad (\text{H.H})}$$

Why? $\mathbf{w} = \sum_{i=1}^d \alpha_i \mathbf{v}_i$ Thus, objective (H) = $\left(\sum_i \alpha_i \mathbf{v}_i \right)^T \left(\sum_{j=1}^d \lambda_j \mathbf{v}_j \mathbf{v}_j^T \right) \left(\sum_k \alpha_k \mathbf{v}_k \right)$

$$= \sum_{i,j,k} \alpha_i \alpha_k \lambda_j \underbrace{(\mathbf{v}_i^T \mathbf{v}_j)}_{= \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}} \underbrace{(\mathbf{v}_j^T \mathbf{v}_k)}_{= \begin{cases} 1 & \text{if } j=k \\ 0 & \text{otherwise} \end{cases}} = \sum_{i=1}^d \alpha_i^2 \lambda_i$$

Constraint (H) = $\sum_{i=1}^d \alpha_i^2 = 1$

Thus: from (H.P): $\alpha_i = 1, \alpha_i = 0 \quad \forall i > 1$

How about $k>1$?

- Suppose we wish to project to more than one dimension. Thus we want:

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

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

$$\mathbf{W} = \begin{pmatrix} \mathbf{v}_1 & | & \mathbf{v}_k \\ \vdots & | & \vdots \end{pmatrix} \in \mathbb{R}^{d \times k}$$

- This is called the **Principal Component Analysis** problem
- Its solution can be obtained in closed form even for $k>1$

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

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

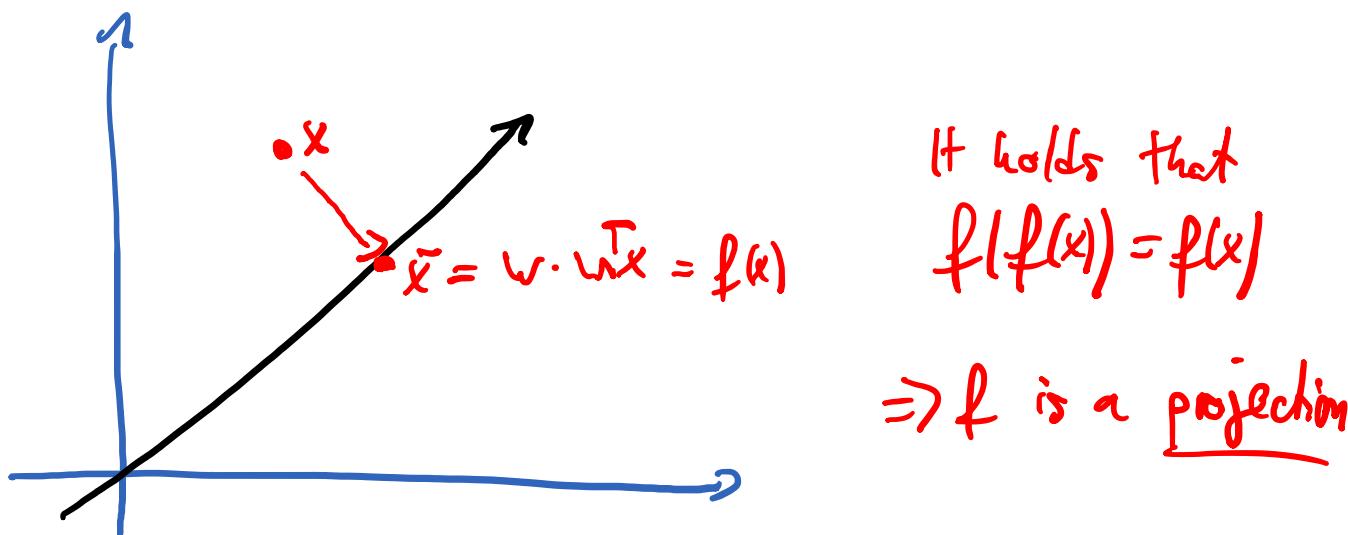
where $\mathbf{W} \in \mathbb{R}^{d \times k}$ is orthogonal, $\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^k$

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

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

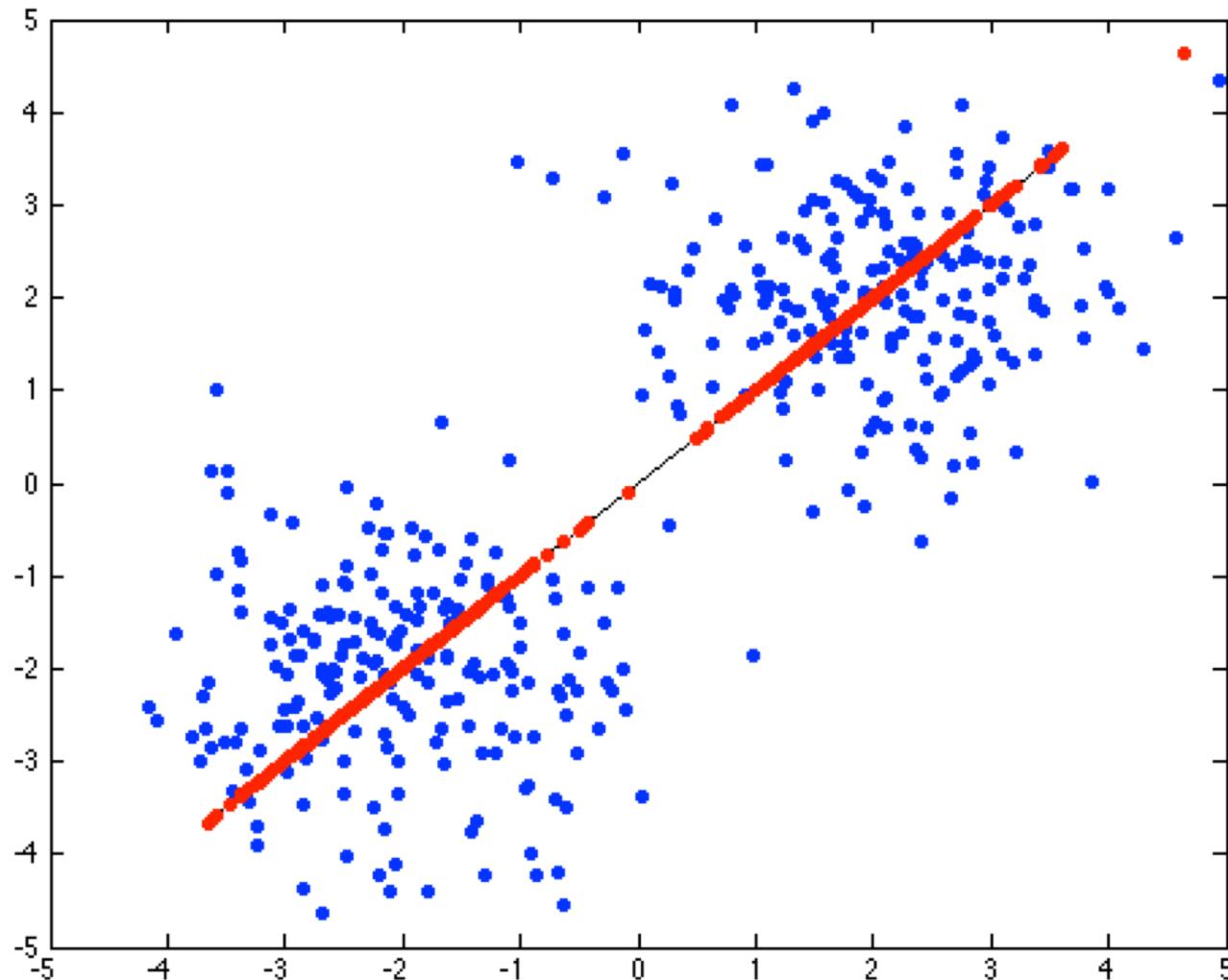
PCA is a projection

- The linear mapping $f(\mathbf{x}) = \mathbf{W}^T \mathbf{x}$ obtained from PCA projects vectors $\mathbf{x} \in \mathbb{R}^d$ into a k -dimensional subspace



- This projection is chosen to minimize the reconstruction error (measured in Euclidean norm)

PCA Illustration



PCA Demo

Connection to SVD

- Can obtain PCA through Singular-Value Decomposition
- Recall: Can represent any $\mathbf{X} \in \mathbb{R}^{n \times d}$ as $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ where $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\mathbf{V} \in \mathbb{R}^{d \times d}$ are orthogonal, and $\mathbf{S} \in \mathbb{R}^{n \times d}$ is diagonal (wlog in decreasing order). Its entries are called singular values

$$\left[\begin{array}{c} \\ \\ \end{array} \right] \in \mathbb{R}^{n \times d} = \left[\begin{array}{c} \mathbf{u} \\ \vdots \\ \mathbf{u} \end{array} \right] \in \mathbb{R}^{n \times n} \left[\begin{array}{c} \mathbf{S} \\ \vdots \\ \mathbf{S} \end{array} \right] \in \mathbb{R}^{n \times d} \left[\begin{array}{c} \mathbf{v}^\top \\ \vdots \\ \mathbf{v}^\top \end{array} \right] \in \mathbb{R}^{d \times d}$$

PCA via SVD

- Recall: Can represent any $\mathbf{X} \in \mathbb{R}^{n \times d}$ as $\underline{\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T}$ where $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\mathbf{V} \in \mathbb{R}^{d \times d}$ are orthogonal, and $\mathbf{S} \in \mathbb{R}^{n \times d}$ is diagonal (wlog in decreasing order).
- The top k principal components are exactly the first k columns of \mathbf{V}

$$\mathbf{n} \mathcal{L} = \mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{S}^T \underbrace{\mathbf{U}^T \mathbf{U}}_{\mathbf{I}} \mathbf{S} \mathbf{V}^T = \mathbf{V} \underbrace{\mathbf{S}^T \mathbf{S}}_{\mathbf{D}} \mathbf{V}^T$$

Common PCA Usecases

- Visualization ($k=1, 2, 3$)
- Feature learning
- Compression

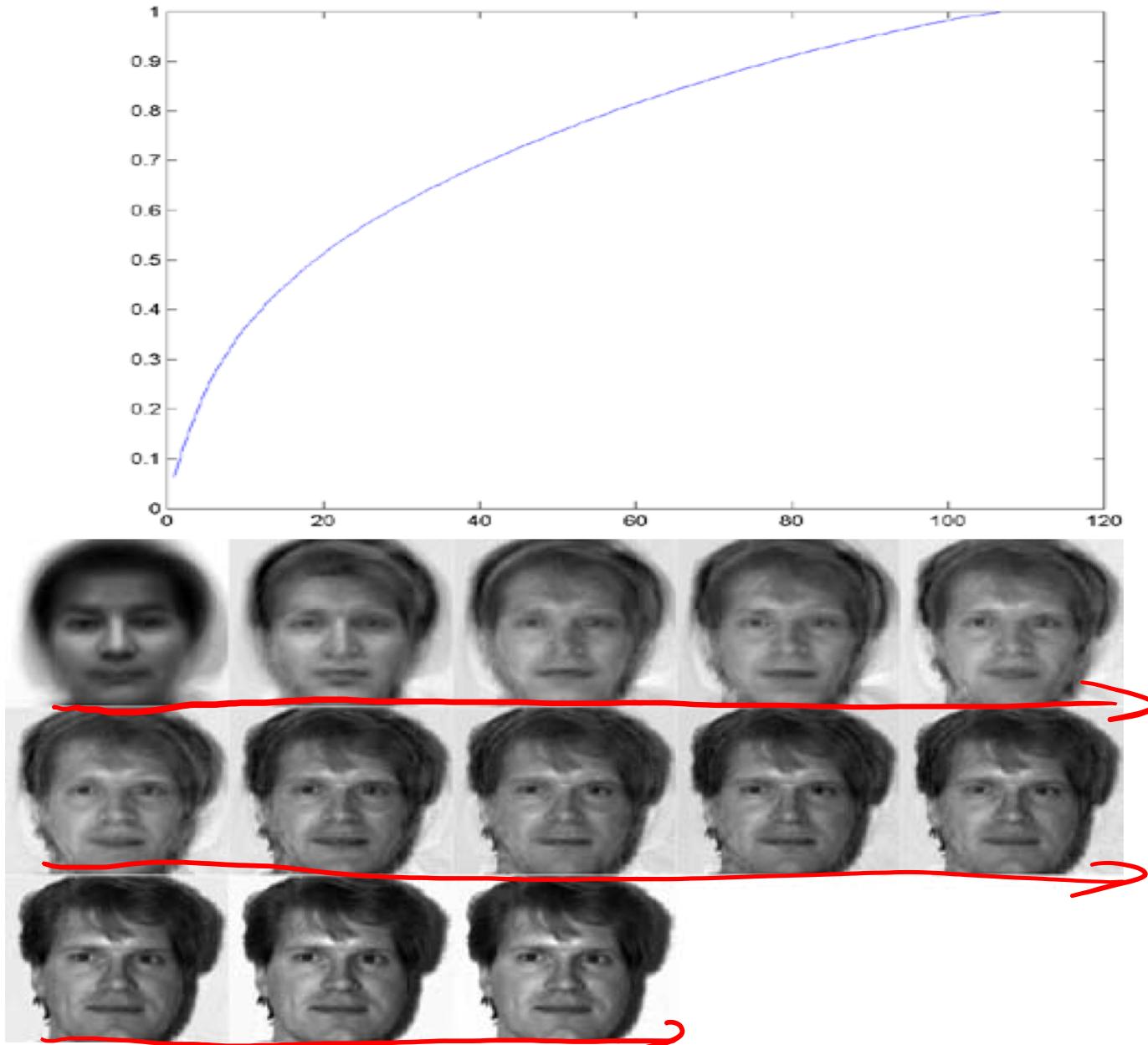
Example: Eigenfaces [de Coro]



Choosing k

- For **visualization**: by inspection ☺
- For **feature induction**: by cross-validation
- **Otherwise**: Pick k so that most of the variance is explained (similar to the choice in k -means)

Reconstruction performance



Now we will

- Discuss relationship of PCA to other methods
- Introduce nonlinear dimension reduction techniques

Side note: PCA vs. K-Means

$W \in \mathbb{R}^{d \times k}$
 $\boxed{\square}_{\mathbb{R}^{k \times m}}$

PCA Problem:

$$X \approx WZ$$

$$(W, z_1, \dots, z_n) = \arg \min \sum_{i=1}^n \|Wz_i - x_i\|_2^2$$

where $W \in \mathbb{R}^{d \times k}$ is orthogonal, $z_1, \dots, z_n \in \mathbb{R}^k$

k-Means problem: (equivalent formulation)

$$(W, z_1, \dots, z_n) = \arg \min \sum_{i=1}^n \|Wz_i - x_i\|_2^2$$

$W = [\mu_1 | \dots | \mu_k]$

where $W \in \mathbb{R}^{d \times k}$ arbitrary and $z_1, \dots, z_n \in E_k$

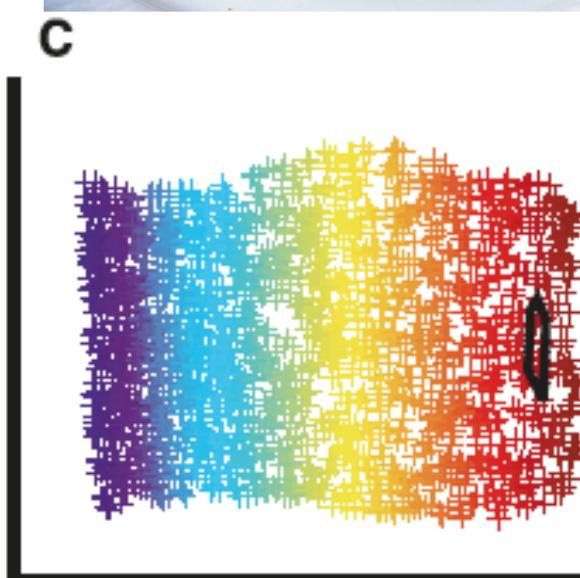
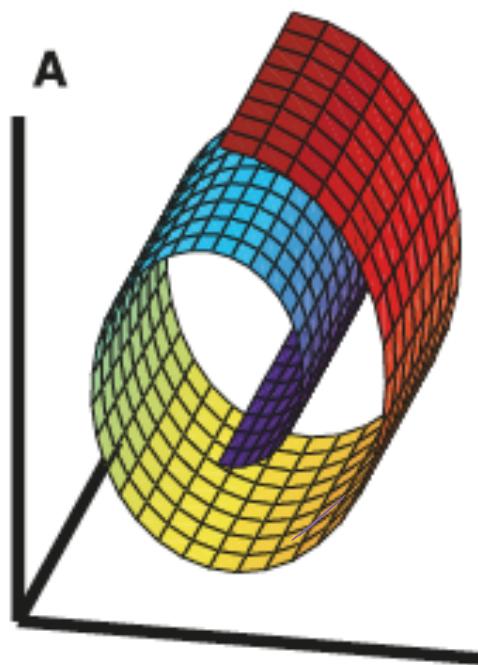
hereby, $E_k = \{[1, 0, \dots, 0], \dots, [0, \dots, 0, 1]\}$
 is the set of unit vectors in \mathbb{R}^k

PCA vs. k-Means

- Can think of PCA and k-Means to solve a similar unsupervised learning problem, with different constraints
- Both aim to compress the data with maximum fidelity under constraints on the model complexity
- This insight gives rise to a much broader class of techniques!
 - [Matrix factorization](#), see Computational Intelligence Lab

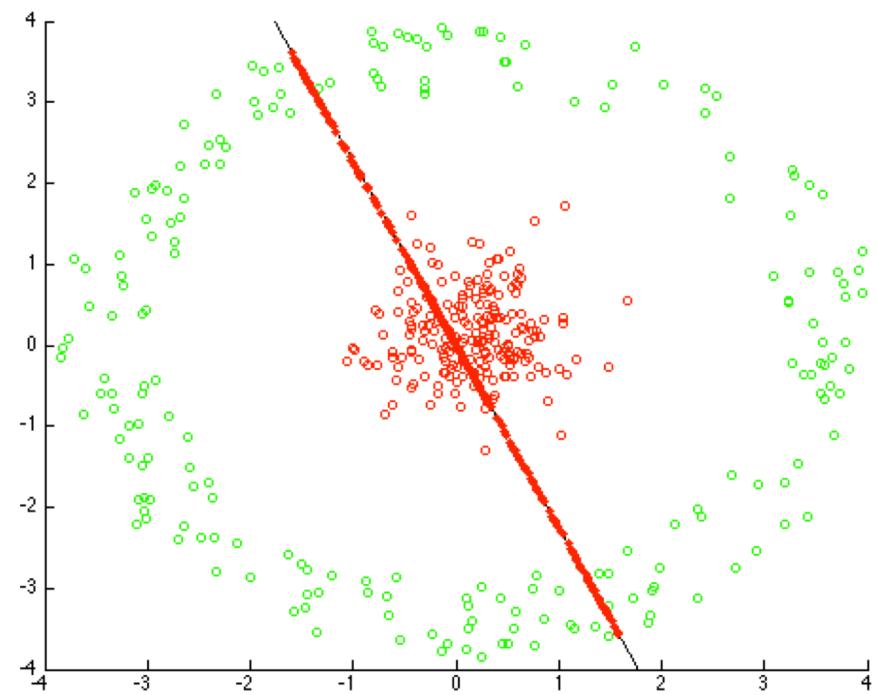
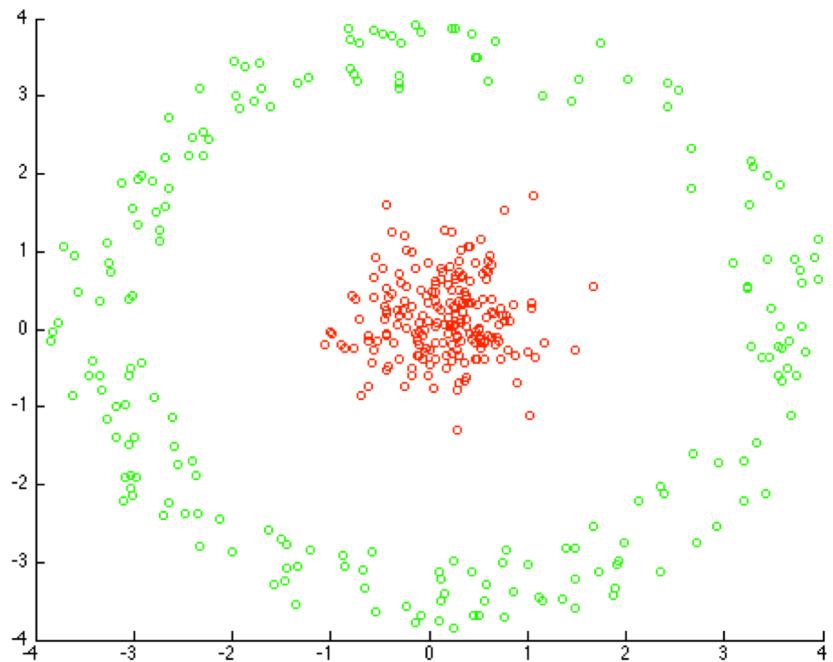
Nonlinear Dimension Reduction

- Motivating Example: „Swiss Roll“



- What is the result of a linear projection?

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!

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

Why: for PCA, w^\perp is [early] eigenvector of $X^T X$

$$\Rightarrow \underbrace{X^T X w}_{\beta} = 2w \Rightarrow w = \frac{1}{2} \underbrace{\sum_i \beta_i x_i}_{X^T \beta}$$

$$\Rightarrow \alpha_i = \frac{\beta_i}{2}$$