TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Statistical Machine Learning

## Lecture 10: Linear Dimensionality Reduction & Statistical Learning Theory

**Kristian Kersting**
**TU Darmstadt**

Summer Term 2020

# Today's Objectives

- Make you understand how to reduce dimensionality and how to understand the power of function approximators!

- Covered Topics:
  - Principal Component Analysis

  - Statistical Learning Theory

# **Outline**

**1. Linear Dimensionality Reduction**

**2. Statistical Learning Theory**

**3. Wrap-Up**

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# **Outline**

## **1. Linear Dimensionality Reduction**

## 2. Statistical Learning Theory

## 3. Wrap-Up

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Towards Principal Component Analysis

- We have covered a number of linear methods
  - For classification and regression

  - Both are supervised learning settings, i.e., pairs of input/output training data points

- Sometimes, it is quite helpful to analyze the data points themselves
  - Unsupervised learning

  - Particularly: Reduce the dimensionality of the data

  - Possible application: Visualization of the data

# Motivation from Linear Least-squares Regression

- In straightforward least-squares linear regression the parameters are computed as
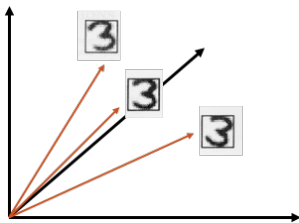
$$\hat{\mathbf{w}} = \left( \hat{\mathbf{X}} \hat{\mathbf{X}}^\mathsf{T} \right)^{-1} \hat{\mathbf{X}} \mathbf{y}$$

  where $\hat{\mathbf{X}} \in \mathbb{R}^{d \times n}$ and $\mathbf{y} \in \mathbb{R}^{n \times 1}$

- We need to invert a $d \times d$ matrix, which naively costs $O\left(d^3\right)$

- Hence, it would be helpful to find a new $d^{\text{new}} << d$ to gain computational advantage while not loosing prediction performance

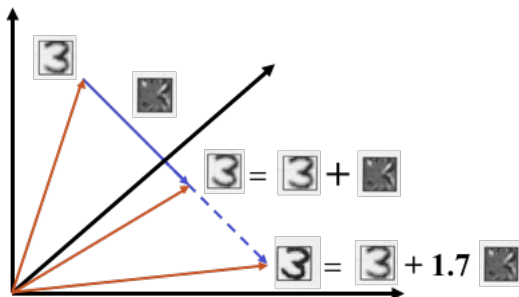TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Dimensionality Reduction

- How can we find more efficient representations for our data? How can we capture the "essence" of the data?

- Example: images of the digit 3



- The images can be represented as points in a high-dimensional space (e.g., with one dimension per pixel, in a 4k image there are around 9 million dimensions!)
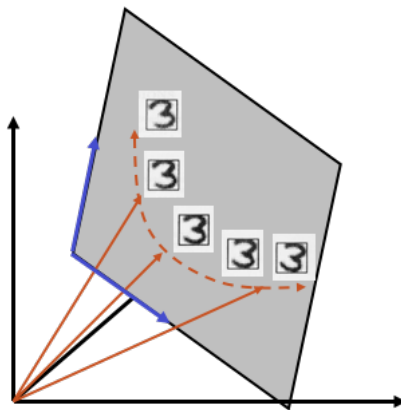
# Linear Dimensionality Reduction

- To make things easier, we will once again assume linear models. A data point (here: one image) can be written as a linear combination of bases (here: basis images)

# Linear Dimensionality Reduction

- What linear transformations of the data can be used to define a lower-dimensional subspace that captures most of the structure?

# Goal of Linear Dimensionality Reduction

- Original data point $n$: $\boldsymbol{x}^n \in \mathbb{R}^M$

- Low-dimensional representation of data point $n$: $\boldsymbol{a}^n \in \mathbb{R}^D$ with $D << M$

- Goal: find a mapping

$$\boldsymbol{x}^n \to \boldsymbol{a}^n$$

- Restrict this mapping to be a linear function

$$\boldsymbol{a}^n = \boldsymbol{B}\boldsymbol{x}^n, \quad \boldsymbol{B} \in \mathbb{R}^{D \times M}$$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Simple Observation

- We can always write a vector as

$$\mathbf{x} = \sum_{i=1}^{M} a_i \mathbf{u}_i, \text{ where } \mathbf{u}_i^\mathsf{T} \mathbf{u}_j = \delta_{ij}$$

  - $\delta_{ij} = 1$ if $i = j$, 0 otherwise. $\mathbf{u}_i$ and $\mathbf{u}_j$ are **orthonormal** to each other

- Example

$$\left[ \begin{array}{c} 3 \\ 7 \end{array} \right] = 3 \left[ \begin{array}{c} 1 \\ 0 \end{array} \right] + 7 \left[ \begin{array}{c} 0 \\ 1 \end{array} \right]$$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Simple Observation

$$
\begin{aligned}
\boldsymbol{x} &= a_1\boldsymbol{u}_1 + a_2\boldsymbol{u}_2 \\
a_1 &= \boldsymbol{u}_1^{\mathsf{T}}\left(\boldsymbol{x} - a_2\boldsymbol{u}_2\right) \\
&= \boldsymbol{u}_1^{\mathsf{T}}\boldsymbol{x} - \boldsymbol{u}_1^{\mathsf{T}}a_2\boldsymbol{u}_2 \\
&= \boldsymbol{u}_1^{\mathsf{T}}\boldsymbol{x} - a_2\underbrace{\boldsymbol{u}_1^{\mathsf{T}}\boldsymbol{u}_2}_{=0} \\
&= \boldsymbol{u}_1^{\mathsf{T}}\boldsymbol{x}
\end{aligned}
$$

■ More generally

$$
\underbrace{a_i}_{\text{scalar coefficient}} = \underbrace{\boldsymbol{u}_i^{\mathsf{T}}\boldsymbol{x}}_{\text{projection}}
$$

# Decomposition

$$\boldsymbol{x}^n = \sum_{i=1}^{D} a_i \boldsymbol{u}_i + \underbrace{\sum_{j=D+1}^{M} b_j \boldsymbol{u}_j}_{\text{error}} \approx \tilde{\boldsymbol{x}}^n$$

- We want the $D$ bases that minimize the mean squared error over the training data

$$\boldsymbol{u}_1, \ldots, \boldsymbol{u}_D = \arg \min_{\boldsymbol{u}_1, \ldots, \boldsymbol{u}_D} E\left(\boldsymbol{u}_1, \ldots, \boldsymbol{u}_D\right) = \arg \min_{\boldsymbol{u}_1, \ldots, \boldsymbol{u}_D} \sum_{n=1}^{N} \|\boldsymbol{x}^n - \tilde{\boldsymbol{x}}^n\|^2$$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Minimizing the error

- Rewrite the error (assuming a single basis vector)

$$
\begin{aligned}
E(\boldsymbol{u}) &= \sum_{n=1}^{N} \|\boldsymbol{x}^n - \tilde{\boldsymbol{x}}^n\|^2 \\
&= \sum_{n=1}^{N} \|\boldsymbol{x}^n - (\boldsymbol{u}^\mathsf{T}\boldsymbol{x}^n)\,\boldsymbol{u}\|^2 \\
&= \sum_{n=1}^{N} \|\boldsymbol{x}^n\|^2 - 2\,(\boldsymbol{u}^\mathsf{T}\boldsymbol{x}^n)^2 + (\boldsymbol{u}^\mathsf{T}\boldsymbol{x}^n)^2 \cdot \boldsymbol{u}^\mathsf{T}\boldsymbol{u} \\
&= \sum_{n=1}^{N} \|\boldsymbol{x}^n\|^2 - (\boldsymbol{u}^\mathsf{T}\boldsymbol{x}^n)^2 \\
&= \sum_{n=1}^{N} \|\boldsymbol{x}^n\|^2 - (a_n)^2
\end{aligned}
$$

# Minimizing the error

- Rewrite the error (assuming a single basis vector)

$$E(\boldsymbol{u}) = \sum_{n=1}^{N} \|\boldsymbol{x}^n\|^2 - (\boldsymbol{u}^\mathsf{T}\boldsymbol{x}^n)^2 = \sum_{n=1}^{N} \|\boldsymbol{x}^n\|^2 - (a_n)^2$$

- Minimizing the error is equivalent to maximizing the variance of the projection. Assuming a zero mean on the data

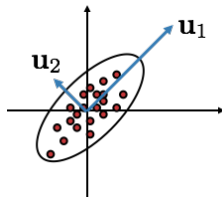$$\max \frac{1}{N} \sum_{n=1}^{N} a_n^2$$

- We can ensure a zero mean projection by subtracting the mean from every data point

$$\boldsymbol{x}^n - \bar{\boldsymbol{x}}$$

# Intuition

$$\boldsymbol{x}^n - \bar{\boldsymbol{x}} = \sum_{i=1}^{D} a_i \boldsymbol{u}_i + \sum_{j=D+1}^{M} b_j \boldsymbol{u}_j$$

$$\tilde{\boldsymbol{x}}^n = \sum_{i=1}^{D} a_i \boldsymbol{u}_i + \bar{\boldsymbol{x}}$$
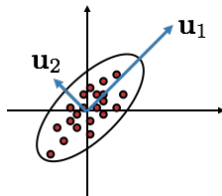


- Projecting onto $\boldsymbol{u}_1$ captures the majority of the variance and hence projecting onto it minimizes the error

- What is that and how do we find the axis of largest variance?

# Intuition

$$\boldsymbol{x}^n - \bar{\boldsymbol{x}} = \sum_{i=1}^{D} a_i \boldsymbol{u}_i + \sum_{j=D+1}^{M} b_j \boldsymbol{u}_j$$

$$\tilde{\boldsymbol{x}}^n = \sum_{i=1}^{D} a_i \boldsymbol{u}_i + \bar{\boldsymbol{x}}$$



- Note that these axes are orthogonal and decorrelate the data; i.e. in the coordinate frame of these axes, the data is uncorrelated (side note: this only works for Gaussians)

# Principal Component Analysis

- Goal: find the so-called principal directions, and the variance of the data along each principal direction



- $\lambda_i$ is the marginal variance along the principal direction $\boldsymbol{u}_i$

# Principal Component Analysis

- The first principal direction $u_1$ is the direction along which the variance of the data is maximal

$$u_1 = \arg\max_u u^\top C u \qquad \text{s.t. } u^\top u = 1$$

- The second principal direction maximizes the variance of the data in the orthogonal complement of the first principal direction

# Problem Formulation

- Let $\boldsymbol{X} = \left[\boldsymbol{x}^1, \ldots, \boldsymbol{x}^n\right] \in \mathbb{R}^{M \times N}$ be a matrix of $N$ vectors in a $M$-dimensional input space, with $\boldsymbol{x}^j \in \mathbb{R}^M$

- Let $\boldsymbol{u} \in \mathbb{R}^M$ be a direction (a vector of length 1) in the input space

- The projection of the $j$–th vector $\boldsymbol{x}^j$ onto the vector $\boldsymbol{u}$ can be computed as

$$a_j = \boldsymbol{u}^\mathsf{T} \boldsymbol{x}^j = \sum_{i=1}^{M} X_{ij} u_i$$

- The goal is to find a direction $\boldsymbol{u}$ that maximizes the variance of the projections of all input vectors

$$\boldsymbol{x}^j \quad \forall j = 1, \ldots, N$$

# Variance of the Projection

$$\bar{a} = \frac{1}{N} \sum_{j=1}^{N} a_j = \frac{1}{N} \sum_{j=1}^{N} \sum_{i=1}^{M} \mathbf{X}_{ij} u_i = \sum_{i=1}^{M} u_i \mu_i$$

where $\mu_i = \frac{1}{N} \sum_{j=1}^{N} X_{ij}$ and $\boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}_j$

$$
\begin{aligned}
\sigma^2 &= \frac{1}{N} \sum_{j=1}^{N} (a_j - \bar{a}) = \frac{1}{N} \sum_{j=1}^{N} \left( \sum_{i=1}^{M} u_i \mathbf{X}_{ij} - \sum_{i=1}^{M} u_i \mu_i \right)^2 \\
&= \frac{1}{N} \sum_{j=1}^{N} \left( \sum_{i=1}^{M} u_i \hat{\mathbf{X}}_{ij} \right)^2 = \frac{1}{N} \sum_{j=1}^{N} \sum_{i=1}^{M} \sum_{k=1}^{M} u_i \hat{\mathbf{X}}_{ij} u_k \hat{\mathbf{X}}_{kj} \\
&= \sum_{i=1}^{M} \sum_{k=1}^{M} u_i u_k \underbrace{\frac{1}{N} \left\langle \hat{\boldsymbol{X}}_{i:}, \hat{\boldsymbol{X}}_{k:} \right\rangle}_{\boldsymbol{C} = \frac{1}{N} \hat{\boldsymbol{X}} \hat{\boldsymbol{X}}^\top} = \sum_{i=1}^{M} \sum_{k=1}^{M} u_i C_{ik} u_k = \boldsymbol{u}^\top \boldsymbol{C} \boldsymbol{u}
\end{aligned}
$$

- $\boldsymbol{C}$ is the covariance matrix

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Maximizing the Variance

- Maximize $\sigma^2$ under the constraint that $\|\boldsymbol{u}\| = 1$

- We can translate into an optimization with Lagrangian multipliers

$$F\left(\boldsymbol{u}; \lambda\right) = \boldsymbol{u}^{\mathsf{T}} \boldsymbol{C} \boldsymbol{u} - \lambda \left( \sum_{k=1}^{M} u_k^2 - 1 \right) = \sum_{i=1}^{M} \sum_{j=1}^{M} u_i C_{ij} u_j - \lambda \left( \sum_{k=1}^{M} u_k^2 - 1 \right)$$

$$\frac{\partial F}{\partial u_l} = \sum_{j=1}^{M} C_{lj} u_j + \sum_{i=1}^{M} u_i C_{il} - \lambda 2 u_l = 0 \quad \forall l = 1, \ldots, M$$

$$\sum_{j=1}^{M} C_{lj} u_j = \lambda u_l \quad \forall l = 1, \ldots, M \Longleftrightarrow \mathbf{C} \mathbf{u} = \lambda \mathbf{u}$$

# Maximizing the Variance

$$\sum_{j=1}^{M} C_{lj} u_j = \lambda u_l \quad \forall l = 1, \ldots, M \Longleftrightarrow \mathbf{C} \mathbf{u} = \lambda \mathbf{u}$$

- Equation of Eigenvalues-Eigenvectors

- The largest Eigenvalue gives us the maximal variance

- The corresponding Eigenvector gives us the direction with maximal variance

# Eigendecomposition - Reminder

- Let the Eigenvectors and Eigenvalues of **C** be $\mathbf{u}_k$ and $\lambda_k$ for $k \leq M$, i.e., $\mathbf{C}\mathbf{u}_k = \lambda_k \mathbf{u}_k$ with $\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_M$

- In matrix form: $\mathbf{C}\mathbf{U} = \mathbf{U}\Lambda$, where $\Lambda = \mathrm{diag}\,(\lambda_1, \ldots, \lambda_M)$ and $\mathbf{U} = [\mathbf{u}_1, \ldots, \mathbf{u}_M]$

- Because **U** is orthonormal (we assume that the Eigenvectors have unit norm) we know that $\mathbf{U}\mathbf{U}^\mathsf{T} = \mathbf{I}$

- This means that we can decompose **C** as

$$\mathbf{C} = \mathbf{U}\Lambda\mathbf{U}^{-1}$$
$$\mathbf{C} = \mathbf{U}\Lambda\mathbf{U}^\mathsf{T}$$

## Principal Component Analysis

- By definition of a covariance matrix, **C** is **real**, **symmetric** and **positive-definite**. Thus we can decompose it in its Eigendecomposition as

$$
\mathbf{C} = \mathbf{U}\Lambda\mathbf{U}^\mathsf{T} = \underbrace{\left[\begin{array}{cccc} \mathbf{u}_1 & \mathbf{u}_2 & \ldots & \mathbf{u}_M \end{array}\right]}_{\text{Eigenvectors}} \underbrace{\left[\begin{array}{cccc} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_2 \end{array}\right]}_{\text{Eigenvalues}} \left[\begin{array}{c} \mathbf{u}_1^\mathsf{T} \\ \mathbf{u}_2^\mathsf{T} \\ \vdots \\ \mathbf{u}_M^\mathsf{T} \end{array}\right]
$$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Principal Component Analysis

- **Observation**: If $\lambda_k \approx 0$ for $k > D$ for some $D << M$, then we can use the subset of the first $D$ eigenvectors to define a basis for approximating the data vectors

$$\mathbf{x}^n - \bar{\mathbf{x}} = \sum_{i=1}^{D} a_i \mathbf{u}_i + \sum_{j=D+1}^{M} b_j \mathbf{u}_j$$

$$\mathbf{x}^n \approx \tilde{\mathbf{x}}^n = \bar{\mathbf{x}} + \sum_{i=1}^{D} a_i \mathbf{u}_i \quad \text{where } a_i = \mathbf{u}_i^{\mathsf{T}} (\mathbf{x}^n - \bar{\mathbf{x}})$$

- This representation has the minimal mean squared error (MSE) of all linear representations of dimension $D$

$$\min E(\mathbf{u}_1, \ldots, \mathbf{u}_D) = \sum_{n=1}^{N} \|\mathbf{x}^n - \tilde{\mathbf{x}}^n\|^2$$

# Principal Component Analysis

- Now we know how we can represent our data in a lower dimensional space in a principled way
    - Center the data around the mean (compute the mean of the data and subtract it)

    - Compute the covariance matrix, decompose it, and choose the first $D$ largest Eigenvalues and corresponding Eigenvectors

    - This gives us an (Eigen)basis for representing the data

    $$\mathbf{a}^n = \mathbf{B}^\mathsf{T} \left( \mathbf{x}^n - \bar{\mathbf{x}} \right)$$
    $$\tilde{\mathbf{x}}^n = \bar{\mathbf{x}} + \mathbf{B}\mathbf{a}^n$$

    where $\mathbf{B} = \begin{bmatrix} \mathbf{u}_1 & \dots & \mathbf{u}_D \end{bmatrix}$

    - It is also common to normalize the variance of each dimension

- But how to choose D?

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Choosing *D*

- A larger *D* leads to a better approximation. In the limit, when *D* = *M* we stay in the initial data dimensions

- There are at least 2 good possibilities for choosing *D*

# Choosing $D$

- A larger $D$ leads to a better approximation. In the limit, when $D = M$ we stay in the initial data dimensions

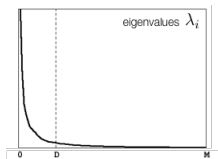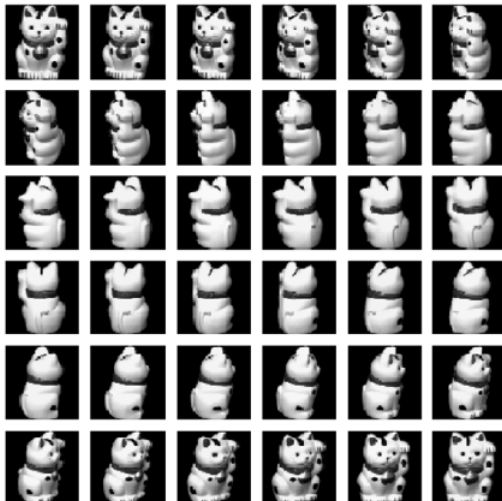- There are at least 2 good possibilities for choosing $D$
  1. Choose $D$ based on application performance, i.e. choose the smallest $D$ that makes the application work *well enough*

# **Choosing** $D$

- A larger $D$ leads to a better approximation. In the limit, when $D = M$ we stay in the initial data dimensions

- There are at least 2 good possibilities for choosing $D$
  1. Choose $D$ based on application performance, i.e. <span style="color:red">choose the smallest $D$ that makes the application work *well enough*</span>

  2. Choose $D$ so that the <span style="color:red">eigenbasis captures some fraction of the variance</span> (for example $\eta = 0.9$)
     The eigenvalue $\lambda_i$ describes the marginal variance captured by $\mathbf{u}_i$

Choose $D$ s.t. $\displaystyle\sum_{i=1}^{D} \lambda_i \geq \eta \sum_{i=1}^{M} \lambda_i$



eigenvalues $\lambda_i$

# Image Representation

# Image Representation with PCA



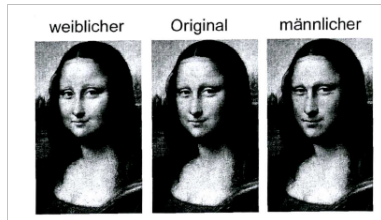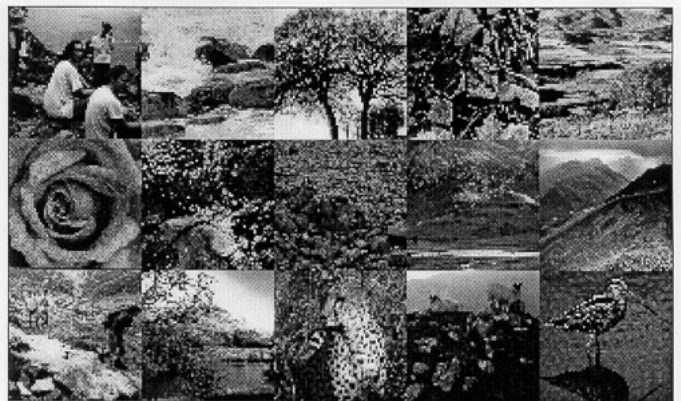$$= \quad + a_1 \quad + a_2 \quad + a_3$$

# EigenFaces

- The first popular use of PCA for object recognition was for the detection and recognition of faces [Turk and Pentland, 1991]

- Collect a face ensemble

- Normalize for contrast, scale, & orientation

- Remove backgrounds

- Apply PCA & choose the first *D* eigen-images that account for most of the variance of the data

# Morphing with EigenFaces

weiblicher    Original    männlicher
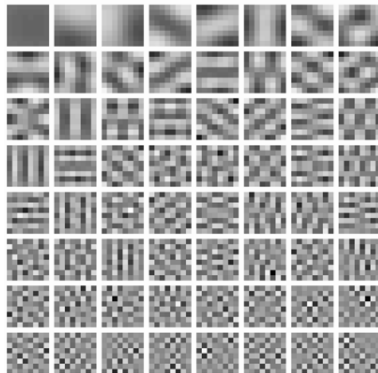
# Generic Image Ensembles

# Generic Image Ensembles



Is there a low dimensional model describing natural images?
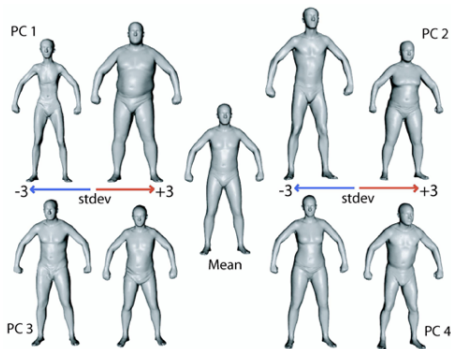
# PCA of Natural Image Patches

8x8 image patches



What do these bases look like?

# PCA Model of Body Shapes

PCA on a detailed triangle model of human bodies [Anguelov et al. 05]

# PCA Summary and Applications

- **Summary**
  - PCA projects the data into a linear subspace

  - PCA maximizes the variance of the projection

  - PCA minimizes the error of the reconstruction

- **Applications**
  - PCA allows us to transform a high-dimensional input space to a low-dimensional feature space, while capturing the essence of the data

  - PCA finds a more *natural* coordinate system for the data

  - PCA is a very common preprocessing step for high-dimensional input data
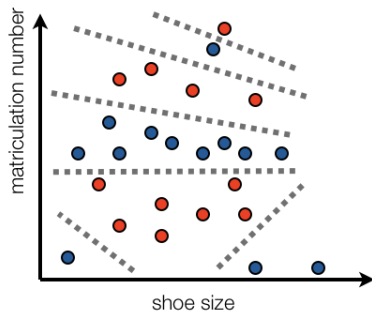
# **Outline**

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Our approach to Machine Learning until now

- We have mostly dealt with *Classical Statistical Learning*
  - Estimate a parameter vector **w** of a fixed model (*learning machine*)
  - Examples
    - Maximum Likelihood estimation: determine the parameters that model a density function
    - Discriminant functions: determine the orientation of the separating hyperplane

- Learning occurs by optimizing the model parameters

- In this setting we assume that we already know in advance the correct model (Bayesian learning is the exception)
  - For instance we model the density function with a linear combination of features, a neural network, a Gaussian, ...

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Generalization Abilities

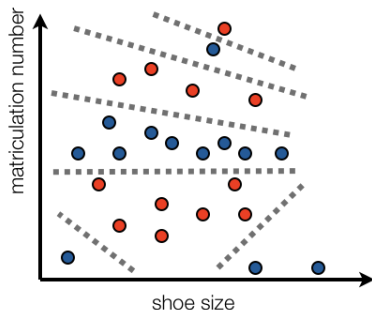- Example: Classify students who will pass the exam. What are some common challenges?

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Generalization Abilities

■ Example: Classify students who will pass the exam. What are some common challenges?



■ Problem 1: Possibly irrelevant features

■ Problem 2: Overfitting from overly complex model

■ We are really interested in the generalization ability and the corresponding risk

# Statistical Learning Theory

- Statistical Learning Theory
  - Does not assume that the correct model is known in advance

  - The goal is to choose an *optimal* model from a set of specified models

  - It is a form of *model selection*

- Optimality here means
  - Ability of the model to generalize, i.e. to have the lowest error probability on all data, and not only on the training data

# Statistical Learning Theory



Vladimir Vapnik [1936-]

- Statistical Learning Theory (from Vladimir Vapnik)
  - Is concerned with the question how one can control the generalization abilities of a learning machine

  - Aims at a formal description of the generalization ability

  - The goal is to develop a rigorous theory as opposed to commonly used heuristics

- We will see later that this may be a noble goal, but the theory unfortunately does not say that much about real problems...

# Supervised Learning

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Statistical Learning Theory

- Supervised learning as an example
    - The environment is stationary, i.e., the data points $\mathbf{x}_i$ have an unknown, but fixed probability density

        $$\mathbf{x}_i \sim p_X$$

    - The supervisor returns the intended classification label $y$ for every data point $\mathbf{x}$, possibly with some added noise $\upsilon$

        $$y = g\left(\mathbf{x}, \upsilon\right)$$

    - The learning machine is represented through a class of functions (with parameters $\mathbf{w}$) that return an output $y$ for every input $\mathbf{x}$

        $$y = f\left(\mathbf{x}, \mathbf{w}\right)$$

# Statistical Learning Theory

- Supervised learning from the learning machine view
  - Choose a particular function

  $$y = f(\mathbf{x}, \mathbf{w})$$

  - Given a set of training examples

  $$\{\mathbf{x}_i, y_i\}_{i=1}^N$$

  - Goal: the desired output $y$ should be approximated *optimally*

- Assessment of *optimality*
  - A loss function, for example the quadratic loss

  $$L(y, f(\mathbf{x}, \mathbf{w})) = (y - f(\mathbf{x}, \mathbf{w}))^2$$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Assessment of optimality - Risk

- Loss function

$$L\left(y, f\left(\mathbf{x}, \mathbf{w}\right)\right)$$

- Empirical risk

$$R_{\text{emp}}\left(\mathbf{w}\right) = \frac{1}{N} \sum_{i=1}^{N} L\left(y_i, f\left(\mathbf{x}_i, \mathbf{w}\right)\right)$$

where $N$ is the number of data points

- Example with the quadratic loss function

$$R_{\text{emp}}\left(\mathbf{w}\right) = \frac{1}{N} \sum_{i=1}^{N} \left(y_i - f\left(\mathbf{x}_i, \mathbf{w}\right)\right)^2$$

# Risk

- In reality, we are instead interested in the True Risk

$$R\left(\mathbf{w}\right) = \int L\left(y, f\left(\mathbf{x}, \mathbf{w}\right)\right) p\left(\mathbf{x}, y\right) d\mathbf{x}dy$$
$$= \mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)}\left[L\left(y, f\left(\mathbf{x}, \mathbf{w}\right)\right)\right]$$

  where $p\left(\mathbf{x}, y\right)$ is the true joint probability density of $\mathbf{x}$ and $y$

- The risk is the expected error over all data sets

- The risk is the expectation of the generalization error

- What is the problem here?

# Risk

- In reality, we are instead interested in the True Risk

$$R\left(\mathbf{w}\right) = \int L\left(y, f\left(\mathbf{x}, \mathbf{w}\right)\right) p\left(\mathbf{x}, y\right) d\mathbf{x}dy$$

$$= \mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)}\left[L\left(y, f\left(\mathbf{x}, \mathbf{w}\right)\right)\right]$$

  where $p\left(\mathbf{x}, y\right)$ is the true joint probability density of $\mathbf{x}$ and $y$

- The risk is the expected error over all data sets

- The risk is the expectation of the generalization error
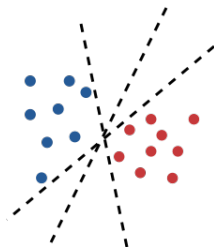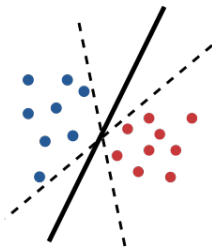
- What is the problem here?
  - $p\left(\mathbf{x}, y\right)$ is fixed but usually unknown

  - We cannot compute the (actual) risk directly

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Empirical vs True Risk



- All 3 decision boundaries have zero empirical risk

- Which one is preferable?

- Which one generalizes best?

# Empirical vs True Risk



- All 3 decision boundaries have zero empirical risk

- Which one is preferable?

- Which one generalizes best?

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Empirical vs True Risk

$$R\left(\mathbf{w}\right) = \mathbb{E}_{\mathbf{x},y\sim p(\mathbf{x},y)}\left[L\left(y,f\left(\mathbf{x},\mathbf{w}\right)\right)\right], \quad R_{\text{emp}}\left(\mathbf{w}\right) = \frac{1}{N}\sum_{i=1}^{N} L\left(y_i,f\left(\mathbf{x}_i,\mathbf{w}\right)\right)$$

- True risk
  - Advantage: Measure for the generalization ability
  - Disadvantage: In general, we do not know $p\left(\mathbf{x},y\right)$

- Empirical risk
  - Disadvantage: No direct measure for the generalization ability
  - Advantage: Does not depend on $p\left(\mathbf{x},y\right)$
  - Learning algorithms often minimize the empirical risk

- We are interested in the dependencies between these two risks
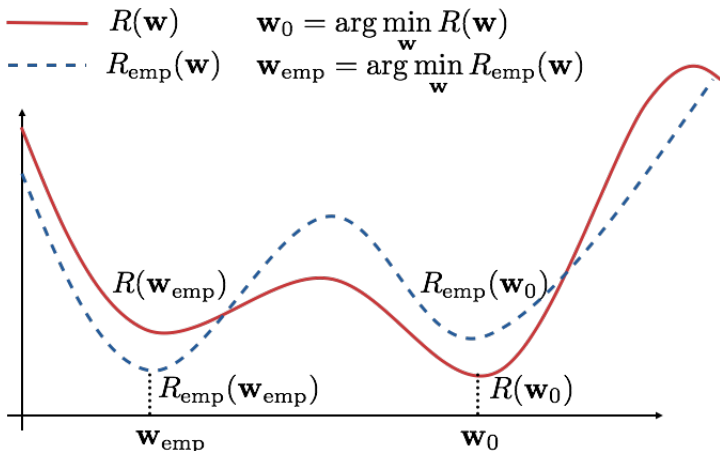
TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Empirical vs True Risk

- You may be asking yourself: we learned that with Monte Carlo we can estimate an expectation by sampling and averaging. So, isn't the empirical risk the same as the true risk?

$$R\left(\mathbf{w}\right) = \mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)}\left[L\left(y, f\left(\mathbf{x}, \mathbf{w}\right)\right)\right] \approx \frac{1}{N}\sum_{i=1}^{N} L\left(y_i, f\left(\mathbf{x}_i, \mathbf{w}\right)\right) = R_{\text{emp}}\left(\mathbf{w}\right)$$

- It is just an approximation, which works very well if our distribution is very concentrated

- The approximation is very good if we have infinitely many samples from the joint distribution

- We will see now how the two risks relate

# Convergence Properties

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Convergence Properties - Derivation

- We assume that the empirical risk converges to the true risk as we supply more and more training data

$$\inf_{\mathbf{w}} R_{\text{emp}}(\mathbf{w}) \xrightarrow{P} \inf_{\mathbf{w}} R(\mathbf{w}) \quad \text{as } N \to \infty$$

where inf stands for *infimum*

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Convergence Properties - Derivation

- We assume that the empirical risk converges to the true risk as we supply more and more training data

$$\inf_{\mathbf{w}} R_{\mathsf{emp}}(\mathbf{w}) \xrightarrow{P} \inf_{\mathbf{w}} R(\mathbf{w}) \quad \text{as } N \to \infty$$

where inf stands for *infimum*

- We further assume that the convergence has to be uniform

$$P\left(\sup_{\mathbf{w}} |R(\mathbf{w}) - R_{\mathsf{emp}}(\mathbf{w})| > \epsilon\right) = 0 \quad \text{as } N \to \infty$$

where sup stands for *supremum*

- Intuitively, "the learning machine gets better as we give it more training data"

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Convergence Properties - Derivation

- If the convergence is uniform

$$P\left(\sup_{\mathbf{w}} |R\left(\mathbf{w}\right) - R_{\text{emp}}\left(\mathbf{w}\right)| > \epsilon\right) < p^{\star} \quad \text{for some } p^{\star} > 0$$

- Then with probability $(1 - p^{\star})$ it holds that

$$
\begin{aligned}
|R\left(\mathbf{w}_{\text{emp}}\right) - R_{\text{emp}}\left(\mathbf{w}_{\text{emp}}\right)| &< \epsilon \\
|R\left(\mathbf{w}_0\right) - R_{\text{emp}}\left(\mathbf{w}_0\right)| &< \epsilon
\end{aligned}
$$

- Hence we have

$$P\left(|R\left(\mathbf{w}_0\right) - R_{\text{emp}}\left(\mathbf{w}_{\text{emp}}\right)| > 2\epsilon\right) < p^{\star}$$

# Empirical Risk Minimizing Principle

- Necessary and sufficient condition: uniform convergence

$$P\left(\sup_{\mathbf{w}} |R(\mathbf{w}) - R_{\text{emp}}(\mathbf{w})| > \epsilon\right) = 0 \quad \text{as } N \to \infty$$

- Compute the empirical risk from the available training data $\{\mathbf{x}_i, y_i\}_{i=1}^{N}$

$$R_{\text{emp}} = \frac{1}{N} \sum_{i=1}^{N} L(y_i, f(\mathbf{x}_i, \mathbf{w}))$$

- Minimizing the empirical risk guarantees that we minimize the true risk in the limit of $N \to \infty$

# Empirical Risk Minimizing Principle

- Advantages
    - We have a formal criterion as to what we can expect in terms of generalization
    - Necessary and sufficient condition is the uniform convergence
- Disadvantages
    - In reality, the amount of training data is very limited
    - We can't easily "take the limit" of $N \to \infty$
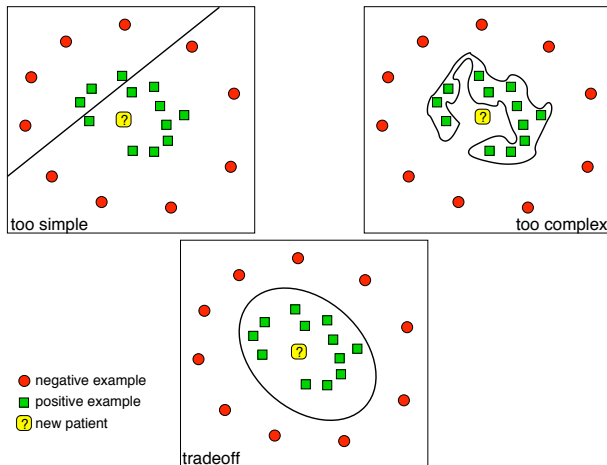    - Nonetheless, what we have essentially been assuming so far in many cases

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Risk Bound

- Idea: determine an upper bound on the true risk based on the empirical risk

$$R\left(\mathbf{w}\right) \leq R_{\text{emp}}\left(\mathbf{w}\right) + \epsilon\left(N, p^\star, h\right)$$

  - $N$ is the number of training examples

  - $p^\star$ is the probability that the bound is met

  - $h$ is the *learning power* of the learning machine, formally called the VC-dimension

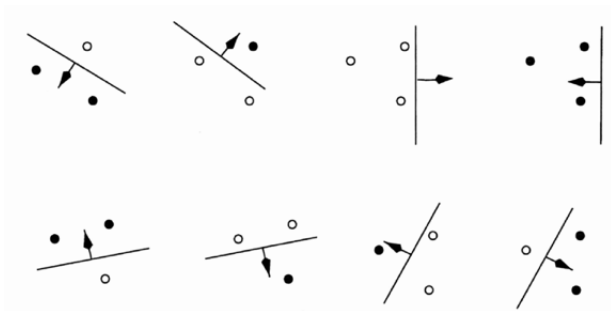# Learning Power



[Florian Markowetz]

# VC-Dimension

- VC stands for Vapnik-Chervonenkis

- Informal definition of the VC-dimension
    - The VC-dimension of a family of functions is the maximal number of data points that can be correctly classified by a function from that family (no matter which label configuration the data points have)

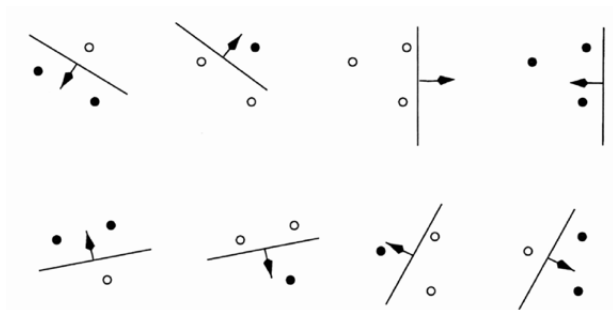    - The VC-dimension is a measure for the capacity, or "learning power", of a classifier

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# VC-Dimension Example

- The VC-dimension of linear classifiers (lines) in $\mathbb{R}^2$ is 3

TECHNISCHE
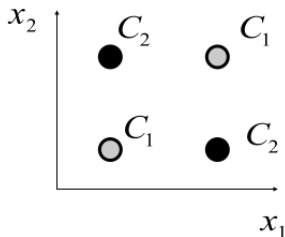UNIVERSITÄT
DARMSTADT

# VC-Dimension Example

- The VC-dimension of linear classifiers (lines) in $\mathbb{R}^2$ is 3



- More generally, the VC-dimension of linear classifiers (hyperplanes) in $\mathbb{R}^n$ is $n + 1$.

# VC-Dimension Example

■ The VC-dimension of linear classifiers (lines) in $\mathbb{R}^2$ is 3. Why is it not bigger than 3?

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# VC-Dimension

- Intuition (but unfortunately not quite right)
  - The VC-dimension is directly related to the number of parameters

- Counter example
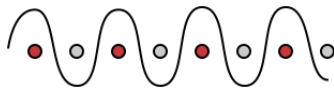
$$
\begin{aligned}
f(x, \mathbf{w}) &= g\left(\sin\left(w_1 x + w_0\right)\right) \\
g(x) &= \left\{ \begin{array}{ll} 1 & x > 0 \\ -1 & x \leq 0 \end{array} \right.
\end{aligned}
$$

  - Has only two parameters but infinite VC-dimension

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# VC-Dimension Counter Example

$$f(x, \mathbf{w}) = g(\sin(w_1 x + w_0))$$

$$g(x) = \begin{cases} 1 & x > 0 \\ -1 & x \leq 0 \end{cases}$$



- For other patterns of data, there is always a setting that correctly classifies any labeling of a fixed number of points

# Risk Bound Example

- Loss function

$$L\left(y, f\left(\mathbf{x}, \mathbf{w}\right)\right) = \frac{1}{2}\left|y - f\left(\mathbf{x}, \mathbf{w}\right)\right|$$

- True risk

$$R\left(\mathbf{w}\right) = \int \frac{1}{2}\left|y - f\left(\mathbf{x}, \mathbf{w}\right)\right| p\left(\mathbf{x}, y\right) d\mathbf{x}dy$$

- Empirical risk

$$R_{\text{emp}}\left(\mathbf{w}\right) = \frac{1}{2N}\sum_{i=1}^{N}\left|y_i - f\left(\mathbf{x}_i, \mathbf{w}\right)\right|$$

# Risk Bound Example

- With probability $p^\star$ it holds that

$$R\left(\mathbf{w}\right) \leq R_{\text{emp}}\left(\mathbf{w}\right) + \sqrt{\frac{h\left(\log\left(2N/h\right) + 1\right) - \log\left(\left(1 - p^\star\right)/4\right)}{N}}$$

- The upper bound is independent of $p\left(\mathbf{x}, y\right)$

- Typically we cannot compute the true risk

- But if we know the VC-dimension, we can compute a bound of the type

$$R\left(\mathbf{w}\right) \leq R_{\text{emp}}\left(\mathbf{w}\right) + \epsilon\left(N, p^\star, h\right)$$

- Note, however, that in practice this bound is very loose

# Risk Bound

# Structural Risk Minimization

- Given a family of $n$ models $f_i(\mathbf{x}_i; \mathbf{w}_i)$ with non-decreasing VC-dimension

$$h_1 \leq h_2 \leq \ldots \leq h_n$$

- Minimize the empirical risk for every model

- Choose the model that minimizes the risk bound (Right Hand Side of the true/empirical risk inequality)

- In general, this is not the same model that minimizes the empirical risk

- This is formally justified by the bound on the true risk

- The result is only sensible, however, if the upper bound on the true risk is a tight bound

# **Outline**

1. **Linear Dimensionality Reduction**

2. **Statistical Learning Theory**

## **3. Wrap-Up**

# 3. Wrap-Up

You know now:

- What dimensionality reduction is and why do we need it

- What the intuition behind PCA is

- Why maximizing the variance of the projection is a good idea

- How PCA relates to Eigenvectors and Eigenvalues

- What Statistical Learning Theory is

- What empirical and true risk are and how they relate

- Why looking just at the empirical risk is not enough

- What the VC-Dimension is

- How to relate the VC-Dimension with the model complexity

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Self-Test Questions

- What does dimensionality reduction mean?

- What is PCA? What are the three things that it does?

- What are the roles of the Eigenvectors and Eigenvalues in PCA?

- Can you describe applications of PCA?

- What does risk in statistical learning theory mean?

- How is the true risk different from the empirical risk?

- What is the learning power of a function approximator?

- What is expressed by a VC-Dimension?

- Is the VC-Dimension always correlated with the number of parameters?

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# **Homework**

- Reading Assignment for next lecture
  - Bishop 5.1, 5.2, 5.3