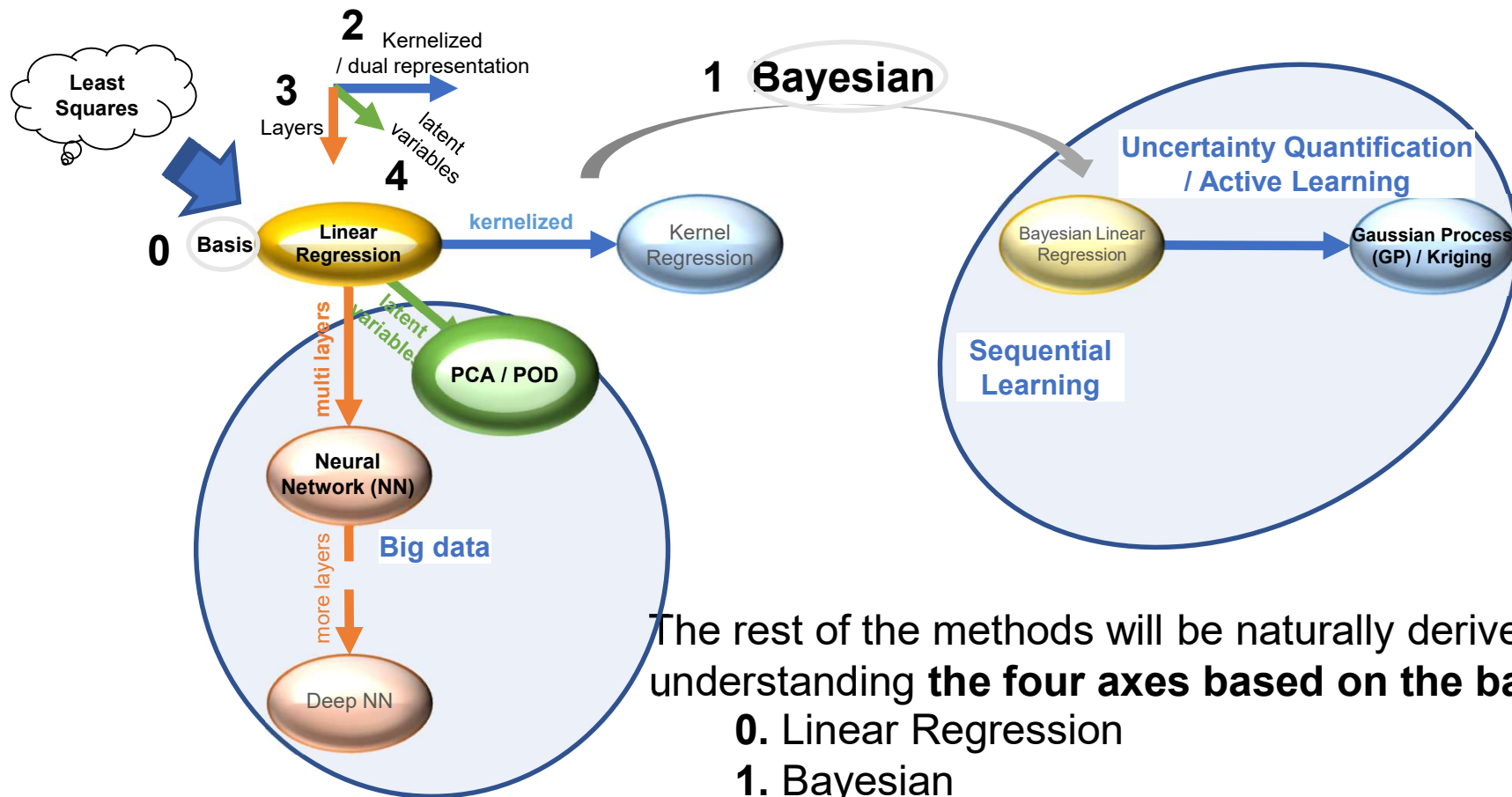# Scientific Machine Learning
## *Lecture 11: Unsupervised Learning*

Dr. Daigo Maruyama

Prof. Dr. Ali Elham

# Key Components (Repost from Lecture 1)



The rest of the methods will be naturally derived after understanding **the four axes based on the basis.**

**0.** Linear Regression

**1.** Bayesian

**2.** Kernel / Dual Representation

**3.** Layers

**4.** Latent Variables (Unsupervised Learning)

Finishing all the four axes here

# Lecture content

- Unsupervised Learning (Concepts / Applications)

- Principal Component Analysis (PCA) / Proper Orthogonal Decomposition (POD) – Linear dimensionality reduction

- Nonlinear dimensionality reduction / Other methods

The lecture of this time partially follows the Chapters 12 of the book:
Christopher M. Bishop "Pattern Recognition And Machine Learning" Springer-Verlag (2006)
The name of this book is shown as "PRML" when it is referred in the slides.
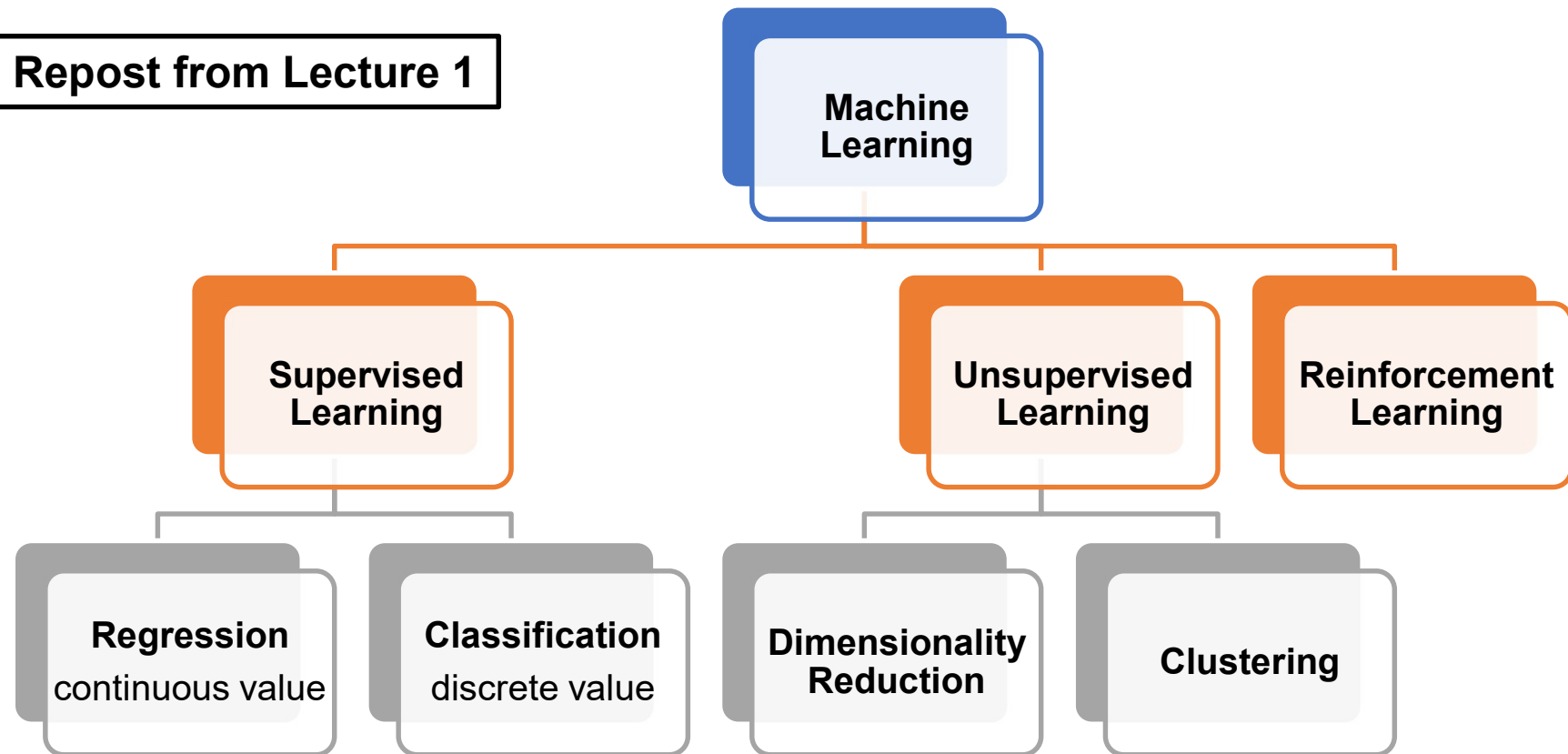
The lecture slides contains original topics in addition to the above.

Technische
Universität
Braunschweig

IFL

# Lecture content

- Unsupervised Learning (Concepts / Applications)

IFL

# Machine Learning Classification by Use/Application

**Repost from Lecture 1**

**Machine Learning**

**Supervised Learning**

**Unsupervised Learning**

**Reinforcement Learning**

**Regression**
continuous value

**Classification**
discrete value

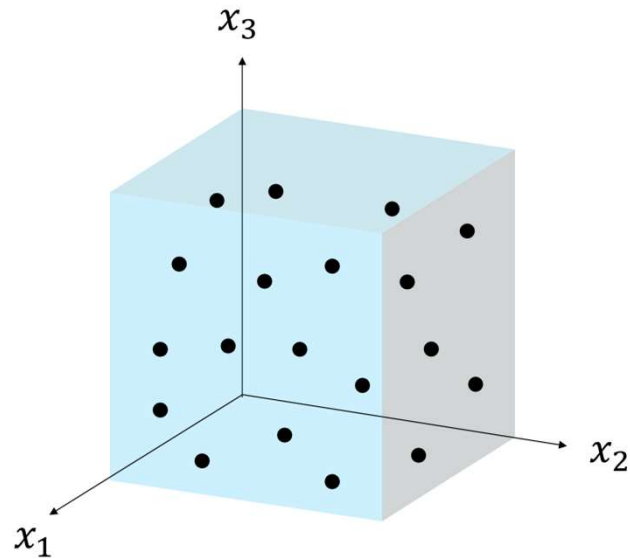**Dimensionality Reduction**

**Clustering**

In this course, machine learning classification is done by **methods and their concepts**.

➡ Then the use/application is naturally derived/understood.

Technische
Universität
Braunschweig

IFL

# Concept / Motivation of Unsupervised Learning

3D input space

collected pictures that look numbers

e.g. a sample generated by DoE

e.g. a sample selected <u>for some reasons</u>

Looks like a plane



Looks nothing to do

Looks possible to be projected on a 2D space

# Concept / Motivation of Unsupervised Learning



MNIST database



https://en.wikipedia.org/wiki/Pixel#/media/File:Closeup_of_pixels.JPG

Image data
= color (or brightness of black) data at each pixel
= **value at each coordinate**

$$x = (x_1, x_2, \cdots, x_{100 \times 100})$$

10,000 dimensional input!

Technische
Universität
Braunschweig

IFL

# Concept / Motivation of Unsupervised Learning

Example of data



only two input parameters actually:
- Rotation
- Translation

$$x = (x_1, x_2)$$

**Inherently 2 dimensional input**

**Latent variables**

Image data
= color (or brightness of black) data at each pixel
= **value at each coordinate**

$$x = (x_1, x_2, \cdots, x_{100 \times 100})$$

10,000 dimensional input!
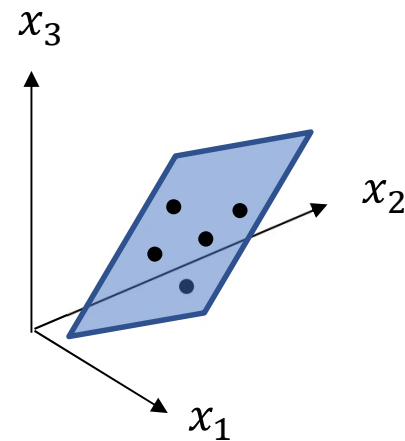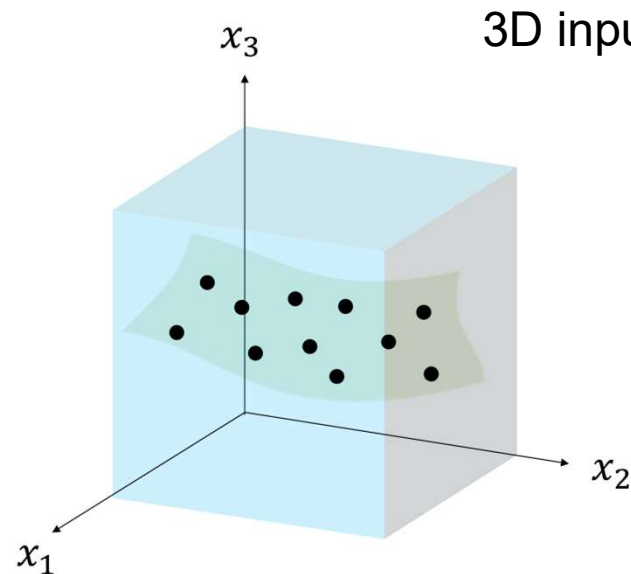
Technische
Universität
Braunschweig

IFL

# Use / Applications

- **Dimensionality reduction**
  - various benefits on computational expense
    - GP models (lower sample size $N$ as a result)
    - NN models

- **Feature extraction**
  - can possibly clarify the essence of the data
    - e.g. rotation and translation in the previous example
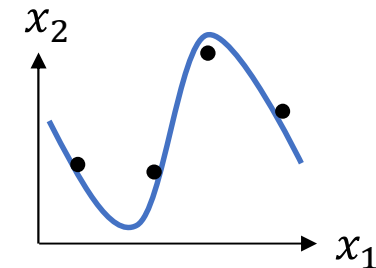    - e.g. the contours of the objects in pictures

- **Visualization (2D or 3D)**

Examples of benchmark testcases
- MNIST
- Iris dataset

Technische
Universität
Braunschweig

IFL

# How can we extract the Latent Variables?

3D input space

2D input space

Have you ever seen something similar to this?

➡ the regression models (in Lecture 4)

**Supervised Learning** in general

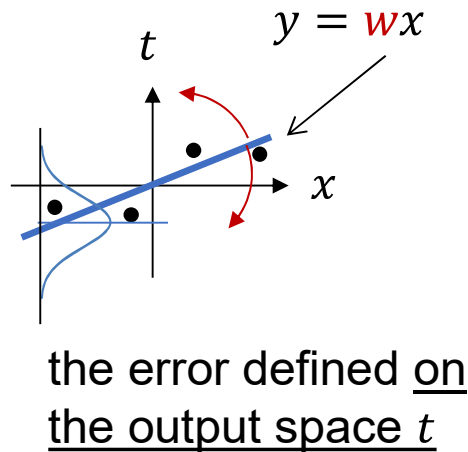Technische
Universität
Braunschweig

IFL

# Lecture content

- Principal Component Analysis (PCA) / Proper Orthogonal Decomposition (POD) – Linear dimensionality reduction

Technische
Universität
Braunschweig

IFL

# Extension from Supervised Learning

Input-output space

$y = wx$



the error defined <u>on the output space $t$</u>

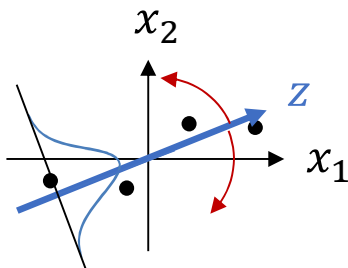**Least Squares**

$$\min \quad E = \sum_{i=1}^{N} \{t_i - wx_i\}^2$$

w.r.t. $w$

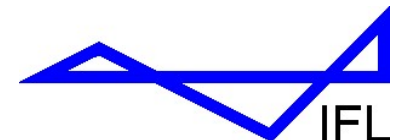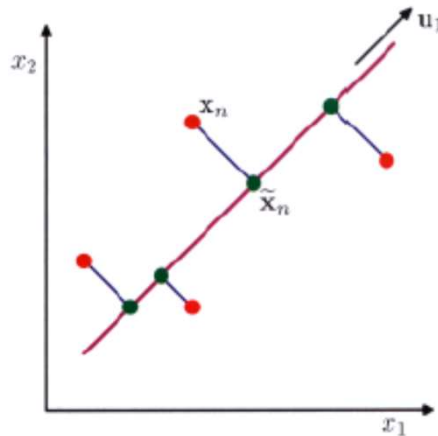$\widehat{w}$

Input space

$z$ is the new axis.



How we define the error (the probabilistic model)? Normal to $z$?

Technische Universität Braunschweig

IFL

# Extension from Supervised Learning



PRML, Fig. 12.2

consider the error on the axis normal to $z$

$$\min \qquad E = \sum_{n=1}^{N}\{(x_{1n} - \tilde{x}_{1n})^2 + (x_{2n} - \tilde{x}_{2n})^2\}$$

w.r.t. $\tilde{x}_1, \tilde{x}_2$

s.t. $\qquad \tilde{x}_2 = w\tilde{x}_1$



$z$ (as $\boldsymbol{u}_1$) has been determined.

$\min E$ is realized on the direction of:

the eigenvector $\boldsymbol{u}_2$ of the minimum eigenvalue of the covariance matrix $\mathbf{S}$ of the data $\mathbf{X}$

$$\mathbf{S} = \frac{1}{N}\mathbf{X}\mathbf{X}^{\mathrm{T}}$$

$$\mathbf{X} = \mathbf{X}_{org} - \boldsymbol{m}_{\mathrm{X}}$$

centered

Technische Universität Braunschweig

IFL

**Principal Component Analysis (PCA) / Proper Orthogonal Decomposition (POD)**

Two commonly used definition:

- **Maximum variance formulation**
  - PRML, 12.1.1

- **Minimum-error formulation**
  - PRML, 12.1.2
  - Another explanation: introduced in the previous slides as an extension of the supervised learning techniques (simple linear regression models)

$\Rightarrow$ the eigenvalue problem     Linear algebra

of the covariance matrix $\mathbf{S}$ of the data $\mathbf{X}$     $\mathbf{S} = \dfrac{1}{N}\mathbf{X}\mathbf{X}^{\mathrm{T}}$

Technische Universität Braunschweig

IFL

# Image of maximum variance formulation

$$\mathbf{S} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^{-1}$$

$D$: dimensionality of the input parameter $x$
$D = 2$ in the previous example

$$\boldsymbol{\Lambda} = \begin{pmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_D \end{pmatrix} \Big\} \ 1, \cdots, D$$

$1, \cdots, M$

$\lambda_i$: eigenvalue in **descending order**

$$\mathbf{U} = \begin{pmatrix} \boldsymbol{u}_1 \\ \vdots \\ \boldsymbol{u}_D \end{pmatrix} \Big\} \ 1, \cdots, D$$

$1, \cdots, M$

$\boldsymbol{u}_i$: eigenvector in the respective order of $\lambda_i$
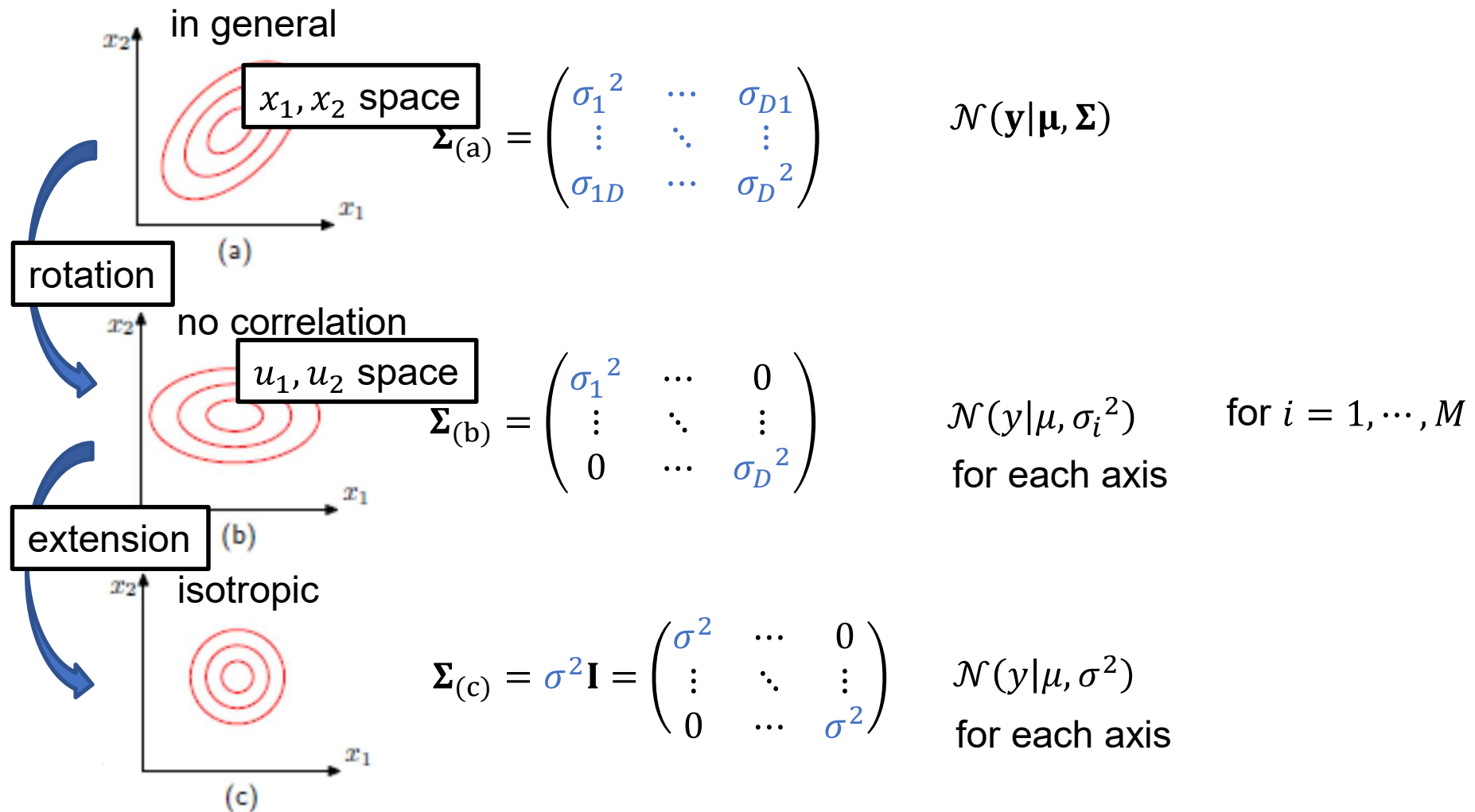
$M < D$ ➡ **Dimensionality reduction**

usually $M$ is judged by the eigenvalues

All the equations can be naturally
understood from the linear algebra.

$$x = \sum_{i=1}^{M} (x^{\mathrm{T}}\boldsymbol{u}_i)\boldsymbol{u}_i$$

# Probability Distributions (Repost from Lecture 3)

in general



$x_1, x_2$ space

$$\Sigma_{(a)} = \begin{pmatrix} \sigma_1{}^2 & \cdots & \sigma_{D1} \\ \vdots & \ddots & \vdots \\ \sigma_{1D} & \cdots & \sigma_D{}^2 \end{pmatrix}$$

$\mathcal{N}(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$

(a)

rotation

no correlation



$u_1, u_2$ space

$$\Sigma_{(b)} = \begin{pmatrix} \sigma_1{}^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_D{}^2 \end{pmatrix}$$

$\mathcal{N}(y|\mu, \sigma_i{}^2)$    for $i = 1, \cdots, M$

for each axis

extension   (b)

isotropic



$$\Sigma_{(c)} = \sigma^2 \mathbf{I} = \begin{pmatrix} \sigma^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma^2 \end{pmatrix}$$

$\mathcal{N}(y|\mu, \sigma^2)$

for each axis

(c)

Technische Universität Braunschweig
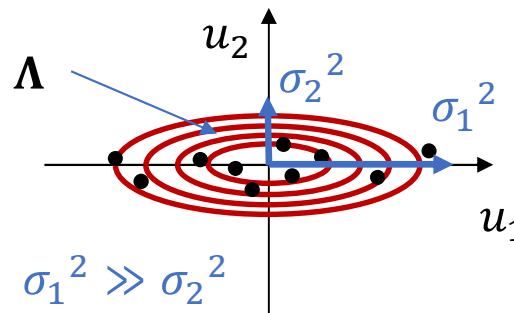
IFL

# Image of maximum variance formulation

connection between the linear algebra and the probability theory



How the data **X** was generated.

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \xrightarrow{\mathbf{S}} \mathbf{S}\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \qquad \mathbf{S} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^{-1}$$

$\mathbf{U}^{-1} \downarrow \qquad\qquad\qquad\qquad \uparrow \mathbf{U}$

$$\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \xrightarrow[\text{scaled by } \sigma_i{}^2 \text{ in this space}]{\boldsymbol{\Lambda}} \boldsymbol{\Lambda}\begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$

The data info on $u_2$ can be neglected.



$\sigma_1{}^2 \gg \sigma_2{}^2$

$$\boldsymbol{\Lambda} = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} = \begin{pmatrix} \sigma_1{}^2 & 0 \\ 0 & \sigma_2{}^2 \end{pmatrix}$$

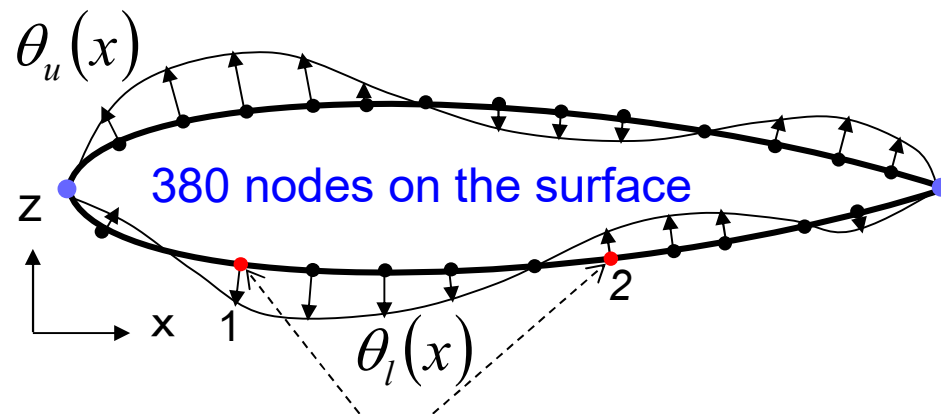e.g. $\dfrac{\lambda_1}{\lambda_1 + \lambda_2} \geq 90\%$

Technische Universität Braunschweig

IFL

# Example

$$x = (x_1, x_2, \cdots, x_{380})$$

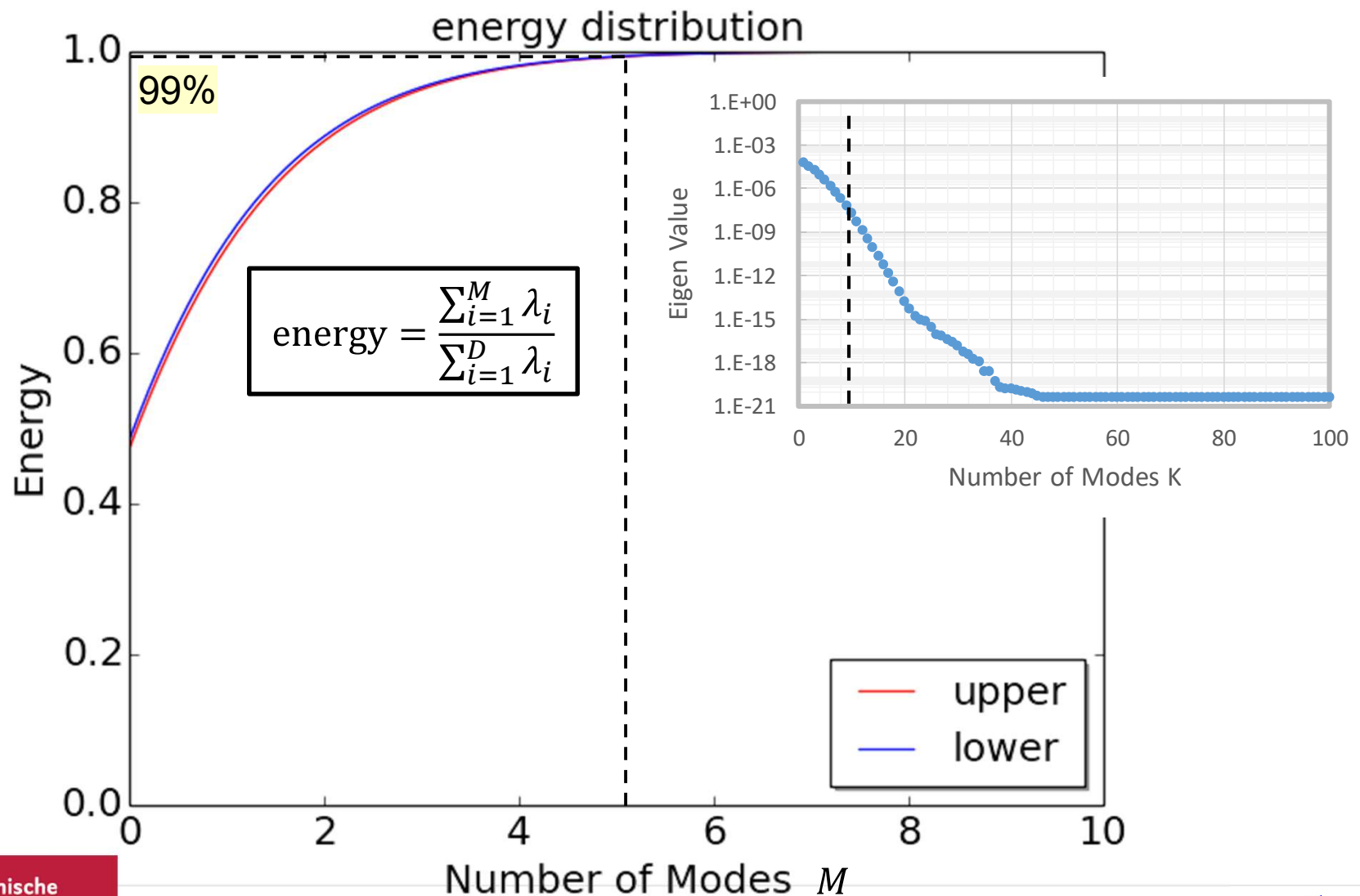describing the airfoil shape by using <u>380</u> original parameters

objective: reduced



The correlation was intentionally assumed to achieve plausible (smooth like the left sketch) surfaces.

= can be also regarded as a GP model

$$\mathrm{cov}[\theta_l(x_1), \theta_l(x_2)] = \sigma_\theta(x_1)\sigma_\theta(x_2)\exp\left(-\frac{(x_1 - x_2)^2}{l^2}\right)$$
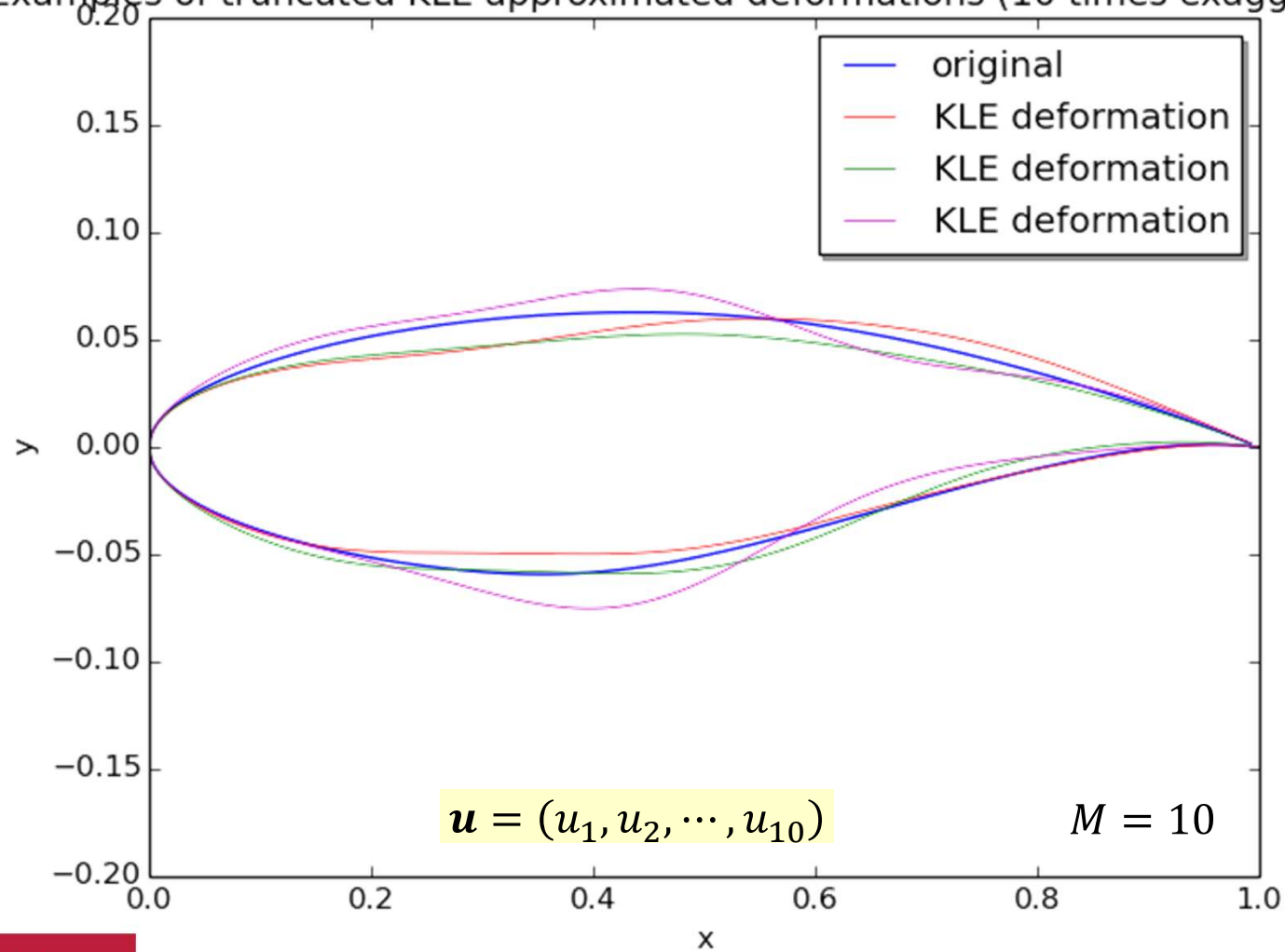
Maruyama, D., Liu D., and Görtz S., "An efficient aerodynamic shape optimization framework for robust design of airfoils using surrogate models," in: Proceedings of the VII European Congress on Computational Methods in Applied Sciences and Engineering (ECCOMAS Congress 2016), 2016.
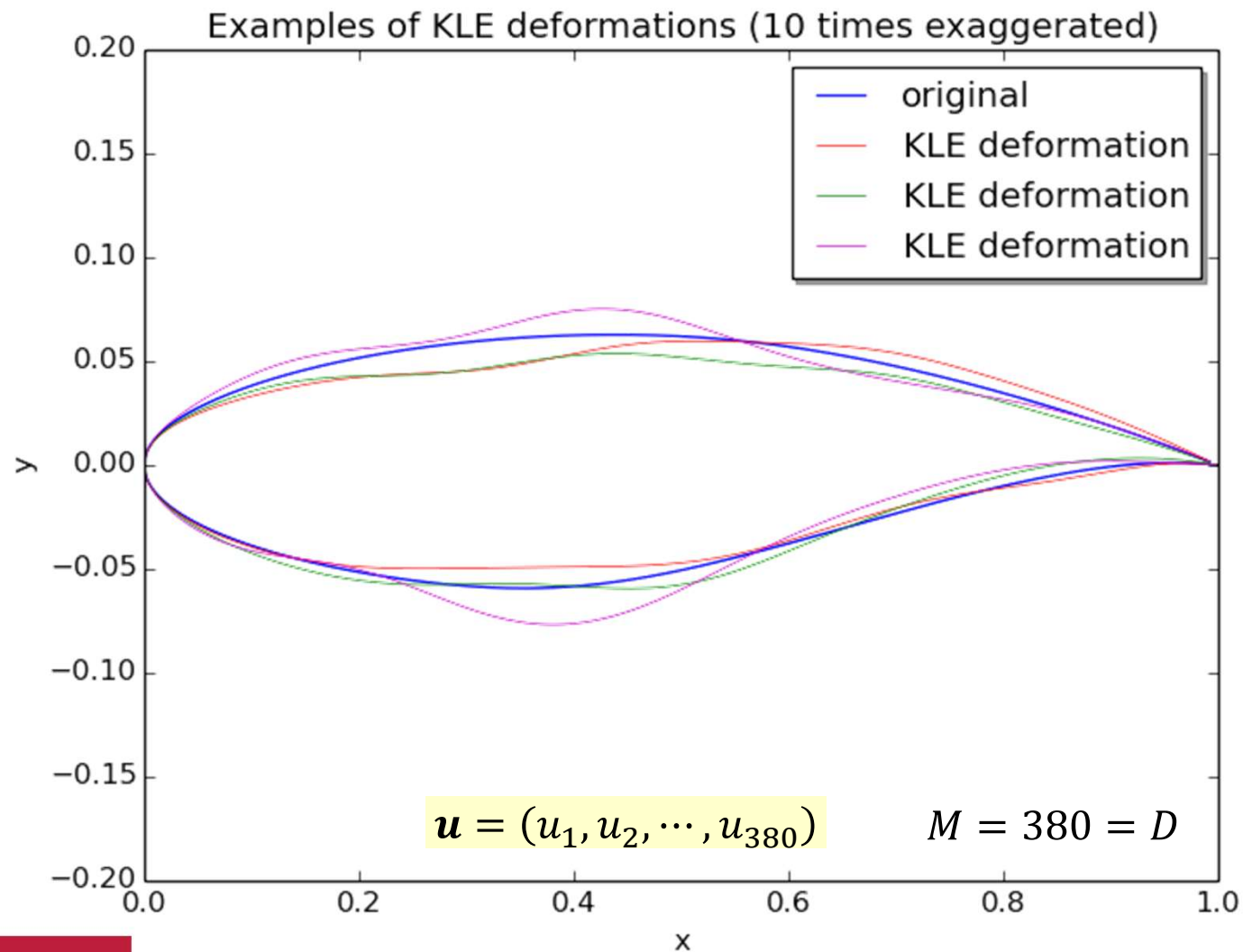
Technische Universität Braunschweig

IFL

# Example



energy distribution

$$\text{energy} = \frac{\sum_{i=1}^{M} \lambda_i}{\sum_{i=1}^{D} \lambda_i}$$

99%

# Example



Examples of truncated KLE approximated deformations (10 times exaggerated)

$$\boldsymbol{u} = (u_1, u_2, \cdots, u_{10})$$

$$M = 10$$

Legend: original, KLE deformation, KLE deformation, KLE deformation

# Example



Examples of KLE deformations (10 times exaggerated)

legend:
- original
- KLE deformation
- KLE deformation
- KLE deformation

$$\boldsymbol{u} = (u_1, u_2, \cdots, u_{380})$$

$$M = 380 = D$$

Technische Universität Braunschweig

IFL

# Lecture content

- Nonlinear dimensionality reduction / Other methods
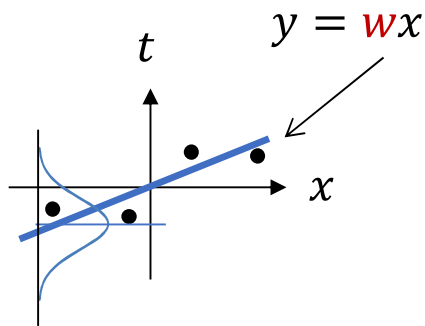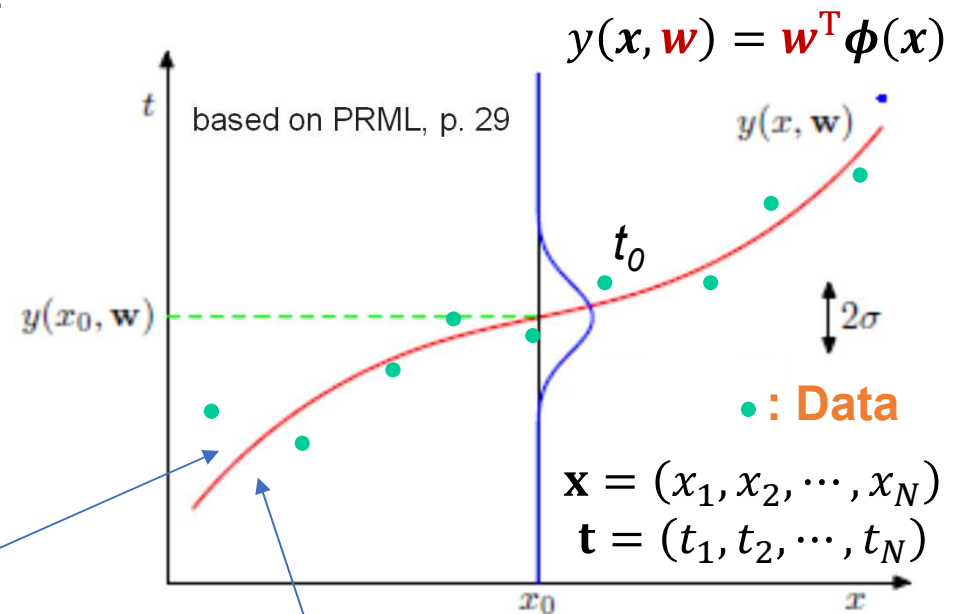
Technische
Universität
Braunschweig

IFL

# Extension to Nonlinear Dimensionality Reduction

**PCA/POD** is a **linear** dimensionality reduction method.

We introduced the PCA/POD from this:

$$y(x, w) = w^{\mathrm{T}} \phi(x)$$

$t$    $y = wx$

based on PRML, p. 29

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

$t_0$

$y(x_0, \mathbf{w})$

$2\sigma$

● : **Data**

in the previous slides.

$$\mathbf{x} = (x_1, x_2, \cdots, x_N)$$
$$\mathbf{t} = (t_1, t_2, \cdots, t_N)$$

We know the left figure's case is one special (simple) case of the generalized **linear regression models**.

**Manifold\***

Technische Universität Braunschweig

IFL

# Kernel PCA (Extension using the Linear Regression Models)

In PCA/POD, we considered the original data $\mathbf{X}$:

$$\mathbf{S} = \frac{1}{N}\mathbf{X}\mathbf{X}^{\mathrm{T}}$$

We first map $x$ into a feature space $\boldsymbol{\phi}(x)$:

(exactly what we have been doing to consider the linear regression models since lecture 4)

Then, the covariance matrix is described in this situation:

$$\mathbf{S} = \frac{1}{N}\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}}$$

In the process of solving the eigenvalue problem, the terms including $\boldsymbol{\phi}(x)$ are always expressed by $\boldsymbol{\phi}(x)^{\mathrm{T}}\boldsymbol{\phi}(x')$.

$$k(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{\phi}(\boldsymbol{x})^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}')$$

The same as when the kernel trick in Lecture 7 was introduced.

Technische
Universität
Braunschweig

IFL

# Linear Regression (Repost from Lecture 4)

Shifting to the simplest linear regression by mapping $\boldsymbol{\phi}: x \to s$ $\qquad s = \boldsymbol{\phi}(x)$

$$y(x, \boldsymbol{w}) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 = \sum_{i=0}^{3} w_i x^i$$



from 1D to 4D
in this example

Linear Regression

In this example:

$$\boldsymbol{\phi}(x) = \left(\phi_0(x), \phi_1(x), \phi_2(x), \phi_3(x)\right)^{\mathrm{T}}$$
$$= (x^0, x^1, x^2, x^3)^{\mathrm{T}}$$

$$y(\boldsymbol{s}, \boldsymbol{a}) = a_0 s_0 + a s_1 + a_2 s_2 + a_3 s_3 = \sum_{i=0}^{3} a_i s_i$$

Technische
Universität
Braunschweig

IFL

# Kernel PCA

The original data space $(x_1, x_2)$

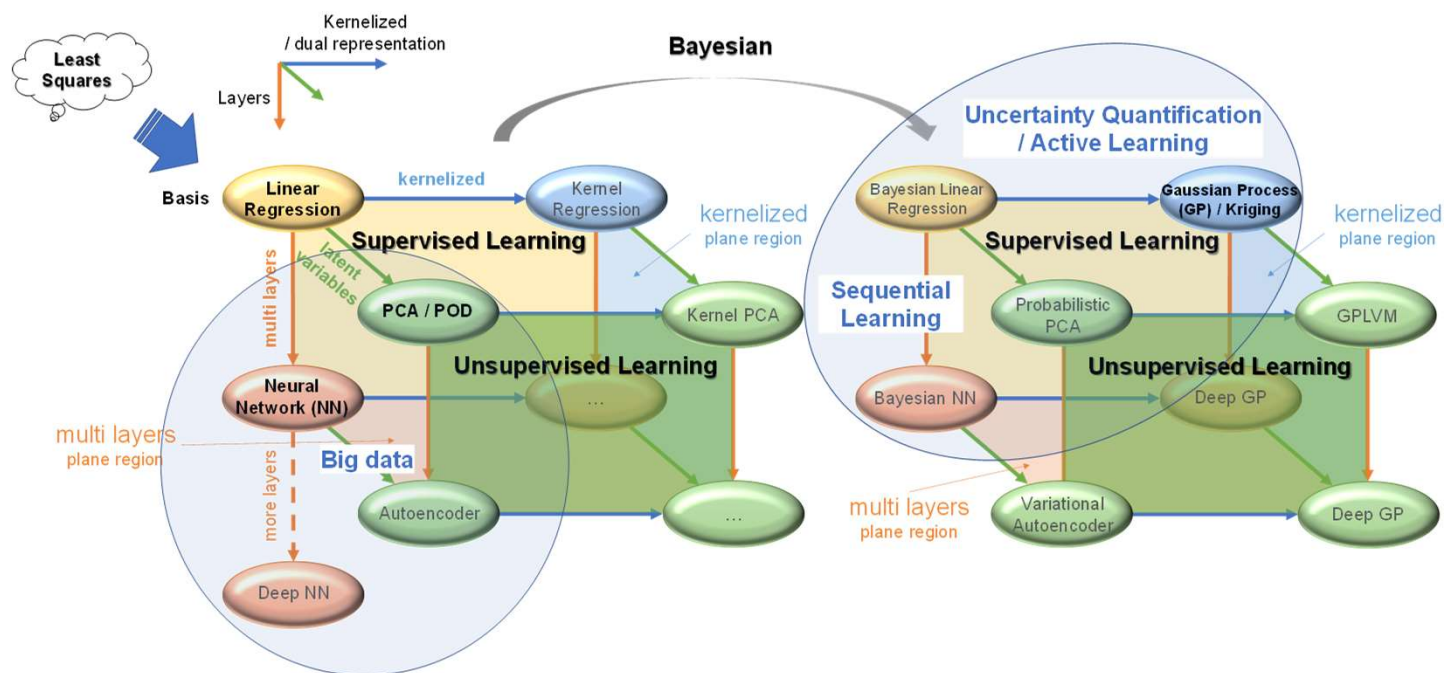A feature space $(\phi_1(x), \phi_2(x))$

# Extension to Nonlinear Dimensionality Reduction

Eventually you may think all the **supervised learning techniques** learned so far can be **extended to unsupervised learning techniques**.

Caution:
Simple linear regression    PCA/POD

Technische
Universität
Braunschweig

IFL

# Extension to Nonlinear Dimensionality Reduction

**supervised**          **unsupervised**

Linear regression          Kernel PCA

Neural network      ➡      Autoencoder

Gaussian process          GPLVM

$$\mathbf{y}(\boldsymbol{x}, \mathbf{W}) = \mathbf{W}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x})$$

The tuned parameter $\widehat{w}_i$ ($i$=1,2,3) can be obtained at the same time as $\widehat{\mathbf{W}}$.

input    output

$\phi_1(\boldsymbol{x})$   $y_1$

$\phi_2(\boldsymbol{x})$

$\phi_3(\boldsymbol{x})$   $\mathbf{W}$   $y_2$

$\phi_4(\boldsymbol{x})$

$\phi_5(\boldsymbol{x})$   $y_3$

$\boldsymbol{\phi}(\boldsymbol{x})$   $\mathbf{y}$

Looks like a neural network?

$$\widehat{\mathbf{W}} = \left(\mathbf{\Phi}^{\mathrm{T}}\mathbf{\Phi}\right)^{-1}\mathbf{\Phi}^{\mathrm{T}}\mathbf{T}$$

$$= \begin{pmatrix} & \cdots & \\ \vdots & \ddots & \vdots \\ & \cdots & \end{pmatrix} \begin{pmatrix} t_1^{(1)} & \cdots & t_3^{(1)} \\ \vdots & \ddots & \vdots \\ t_1^{(N)} & \cdots & t_3^{(N)} \end{pmatrix} = \begin{pmatrix} \widehat{w}_{11} & \cdots & \widehat{w}_{13} \\ \vdots & \ddots & \vdots \\ \widehat{w}_{51} & \cdots & \widehat{w}_{53} \end{pmatrix}$$

$$\widehat{\mathbf{W}} = (\widehat{w}_1, \widehat{w}_2, \widehat{w}_3)$$

Consider a multiple output case:

$$\mathbf{T} = \left(\boldsymbol{t}^{(1)}, \boldsymbol{t}^{(2)}, \cdots, \boldsymbol{t}^{(N)}\right)^{\mathrm{T}}$$

This is still the **linear regression model**.

The objective is to obtain $\widehat{\mathbf{W}}$ or $p(\mathbf{W})$ (by using data).

Technische Universität Braunschweig

IFL

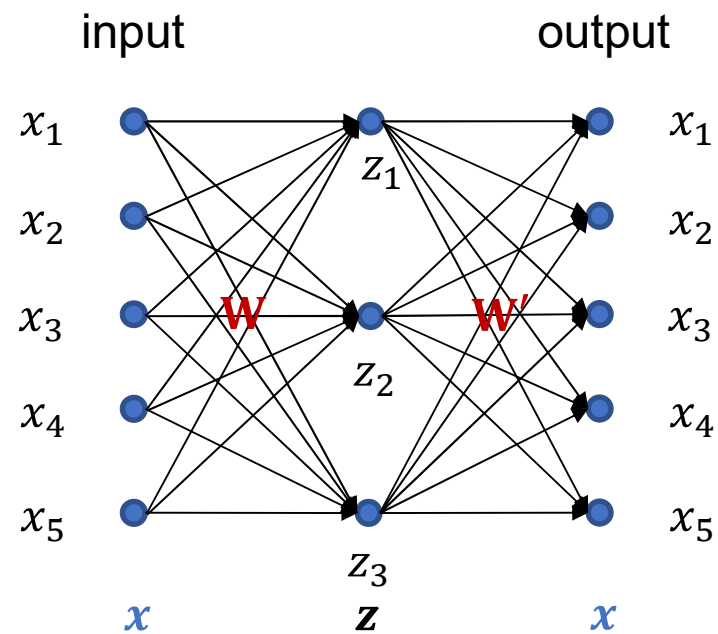# Autoencoder (Extension using the Neural Networks)

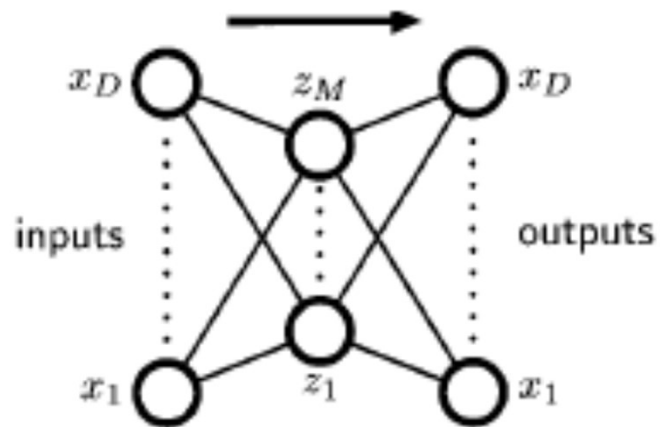**(Simple) Linear regression**

This is still the PCA/POD.

**PCA/POD**



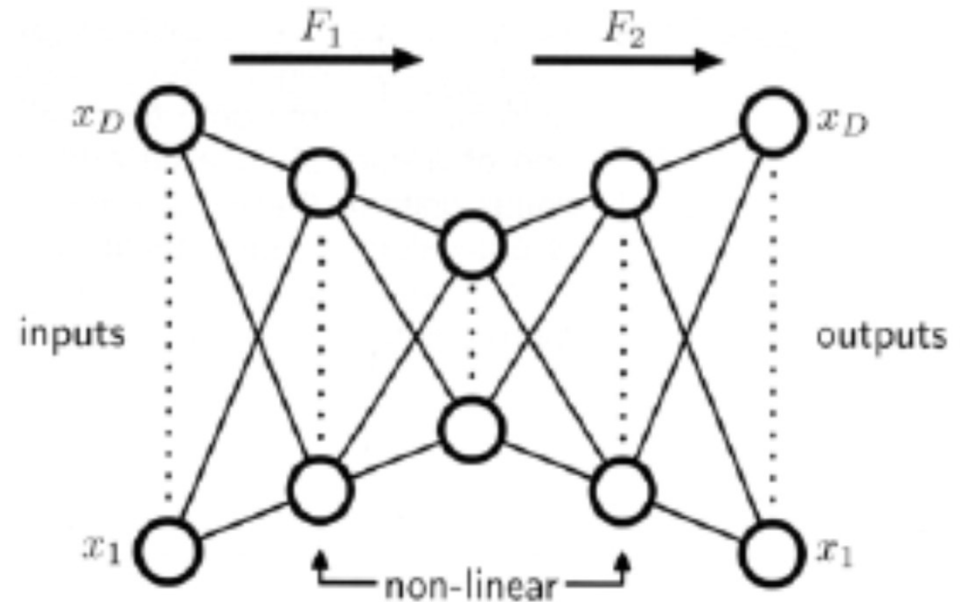Compare this with $E(\boldsymbol{w})$ of the NN models in Lecture 10.

$$E(\boldsymbol{w}) = \sum_{n=1}^{N}\left\{\boldsymbol{x}^{(n)} - \mathbf{y}\left(\boldsymbol{x}^{(n)}, \boldsymbol{w}\right)\right\}^2$$

Technische Universität Braunschweig

IFL

# Autoencoder (Extension using the Neural Networks)



PRML, Fig. 12.18

**PCA/POD**



PRML, Fig. 12.19

**Autoencoder**

compress

recover

Combination of the NN and the autoencoder is one of the standard models.

PRML, Fig. 12.20

Technische
Universität
Braunschweig
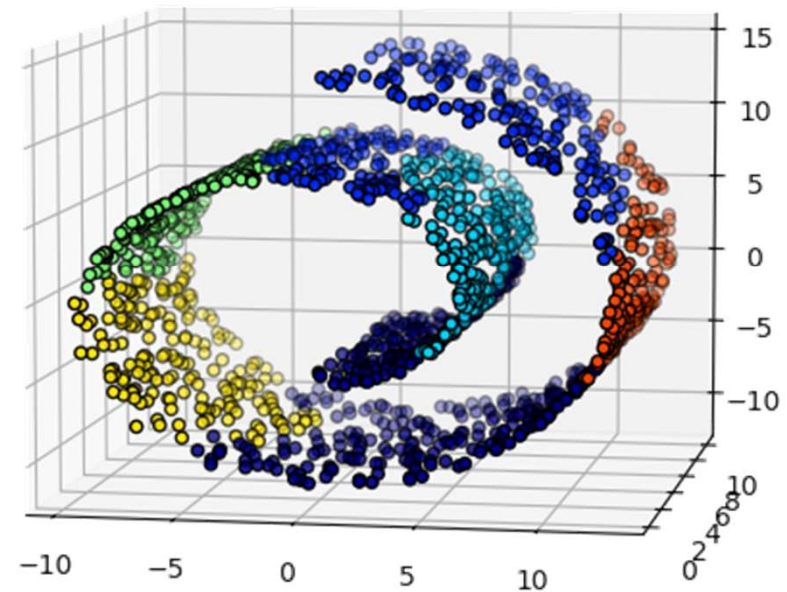
IFL

## Other Methods

Methods of identifying the **manifold explicitly**

- Locally linear embedding (LLE)
    - smooth joining local linear models

- Isomap
    - assuming an Euclidian space locally
    - Then, take the geodesic (shortest path)

Methods often used in 2D mapping (**visualization**)

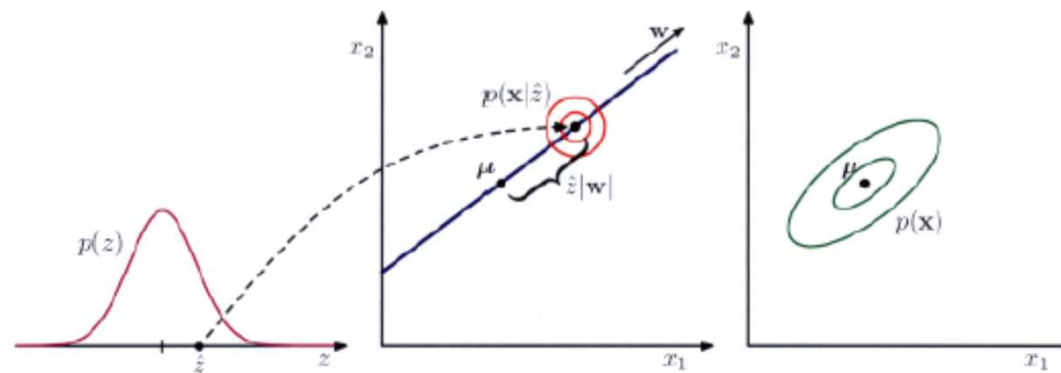- Self-organizing map (SOM)

- t-SNE

Swiss roll:
benchmark test case

# Modelling from the Probabilistic Approaches

The topics here are advanced but related to the topics in Lectures (12) and 13.

This brings the most generalized perspective of modelling.



PRML, Fig. 12.9

Since you may have thought that we started from the least squares, so there might exist a (Gaussian) probabilistic model behind, and also extended to Bayesian approaches.

**Probabilistic PCA, Bayesian PCA**

This is true. The topics are summarized with other techniques.

The keyword is the Latent variables.

# Summary

Unsupervised learning methods were introduced.

- The concepts / applications :
    - Dimensionality reduction
    - Feature extraction
    - 2D, 3D Visualization

- The linear dimensionality reduction (PCA/POD) comes down to the eigenvalue problem.

- Nonlinear dimensionality reduction methods can be naturally extended by the use of the supervised learning techniques learned until here.
    - by regarding the PCA/POD as a simple linear regression model

Technische
Universität
Braunschweig

IFL