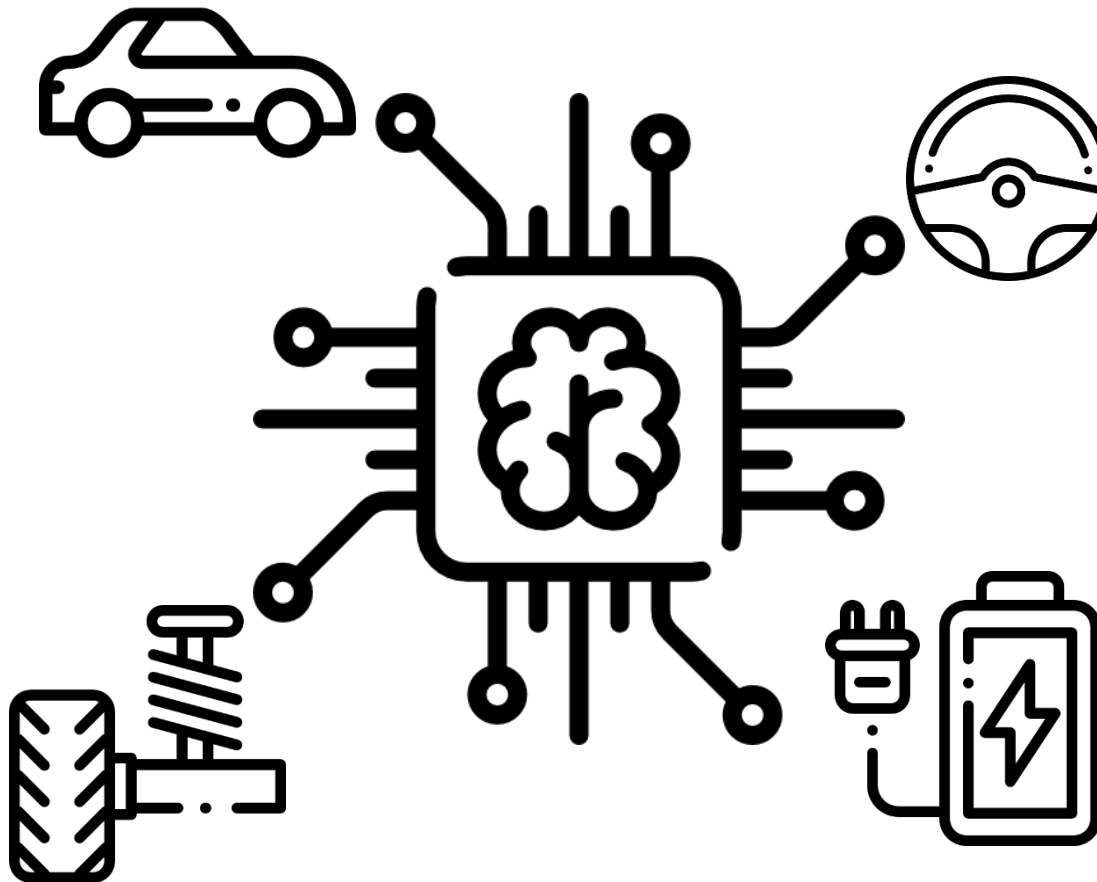


# Artificial Intelligence in Automotive Technology

Johannes Betz / Prof. Dr.-Ing. Markus Lienkamp / Prof. Dr. Boris Lohmann



## Lecture Overview

Lecture 16:15 – 17:45	Practice 17:45 – 18:30
<b>1 Introduction: Artificial Intelligence</b> 17.10.2019 – Johannes Betz	<b>Practice 1</b> 17.10.2019 – Johannes Betz
<b>2 Perception</b> 24.10.2019 – Johannes Betz	<b>Practice 2</b> 24.10.2019 – Johannes Betz
<b>3 Supervised Learning: Regression</b> 31.10.2019 – Alexander Wischnewski	<b>Practice 3</b> 31.10.2019 – Alexander Wischnewski
<b>4 Supervised Learning: Classification</b> 7.11.2019 – Jan Cedric Mertens	<b>Practice 4</b> 7.11.2019 – Jan Cedric Mertens
<b>5 Unsupervised Learning: Clustering</b> 14.11.2019 – Jan Cedric Mertens	<b>Practice 5</b> 14.11.2019 – Jan Cedric Mertens
<b>6 Pathfinding: From British Museum to A*</b> 21.11.2019 – Lennart Adenaw	<b>Practice 6</b> 21.11.2019 – Lennart Adenaw
<b>7 Introduction: Artificial Neural Networks</b> 28.11.2019 – Lennart Adenaw	<b>Practice 7</b> 28.11.2019 – Lennart Adenaw
<b>8 Deep Neural Networks</b> 5.12.2019 – Jean-Michael Georg	<b>Practice 8</b> 5.12.2019 – Jean-Michael Georg
<b>9 Convolutional Neural Networks</b> 12.12.2019 – Jean-Michael Georg	<b>Practice 9</b> 12.12.2019 – Jean-Michael Georg
<b>10 Recurrent Neural Networks</b> 19.12.2019 – Christian Dengler	<b>Practice 10</b> 19.12.2019 – Christian Dengler
<b>11 Reinforcement Learning</b> 09.01.2020 – Christian Dengler	<b>Practice 11</b> 09.01.2020 – Christian Dengler
<b>12 AI-Development</b> 16.01.2020 – Johannes Betz	<b>Practice 12</b> 16.01.2020 – Johannes Betz
<b>13 Guest Lecture: VW Data:Lab</b> 23.01.2020 –	

# Objectives for Lecture 3: Regression

After the lecture you are able to...

... understand the definition of a Regression problem and know the relevant vocabulary.

.... understand in which areas of Automotive Technology Regression can be applied.

... know the strengths and weaknesses of different basis functions for applying a regression model

... choose the right loss function for training your model based on your dataset

.... analyze the result of the training process

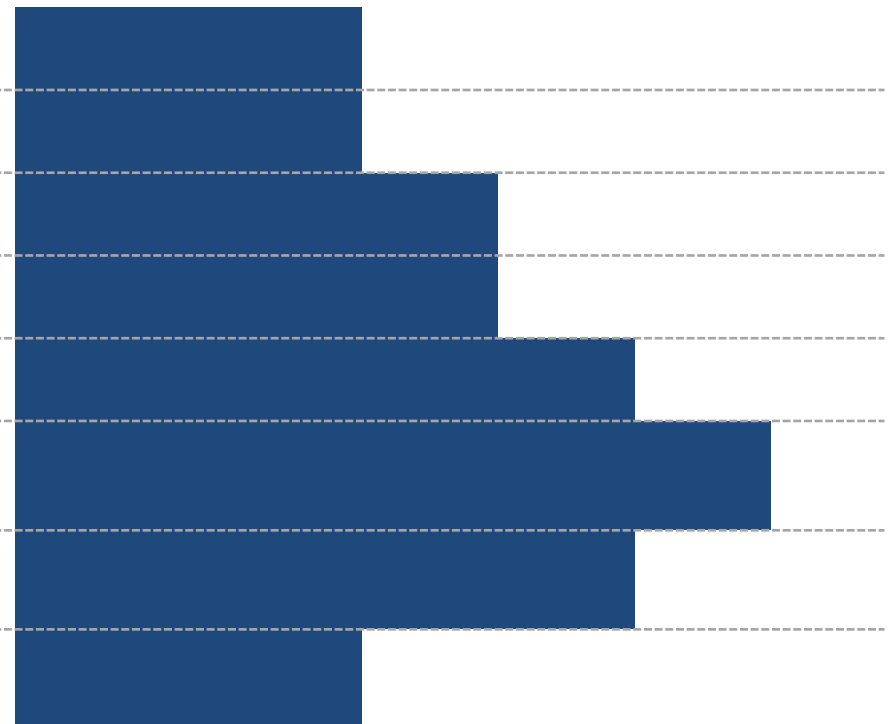
.... apply and evaluate the usefulness of regularization techniques to control model complexity

.... analyze potential problems arising from real world datasets used in regression

.... understand the basic mathematical principles behind regression

Depth of understanding

Remember	Understand	Apply	Analyze	Evaluate	Develop
----------	------------	-------	---------	----------	---------



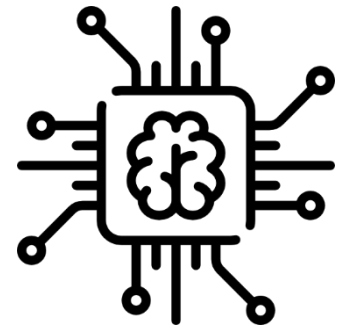
# Supervised Learning: Regression

Johannes Betz / Prof. Dr. Markus Lienkamp / Prof. Dr. Boris Lohmann  
(Alexander Wischnewski, M. Sc.)

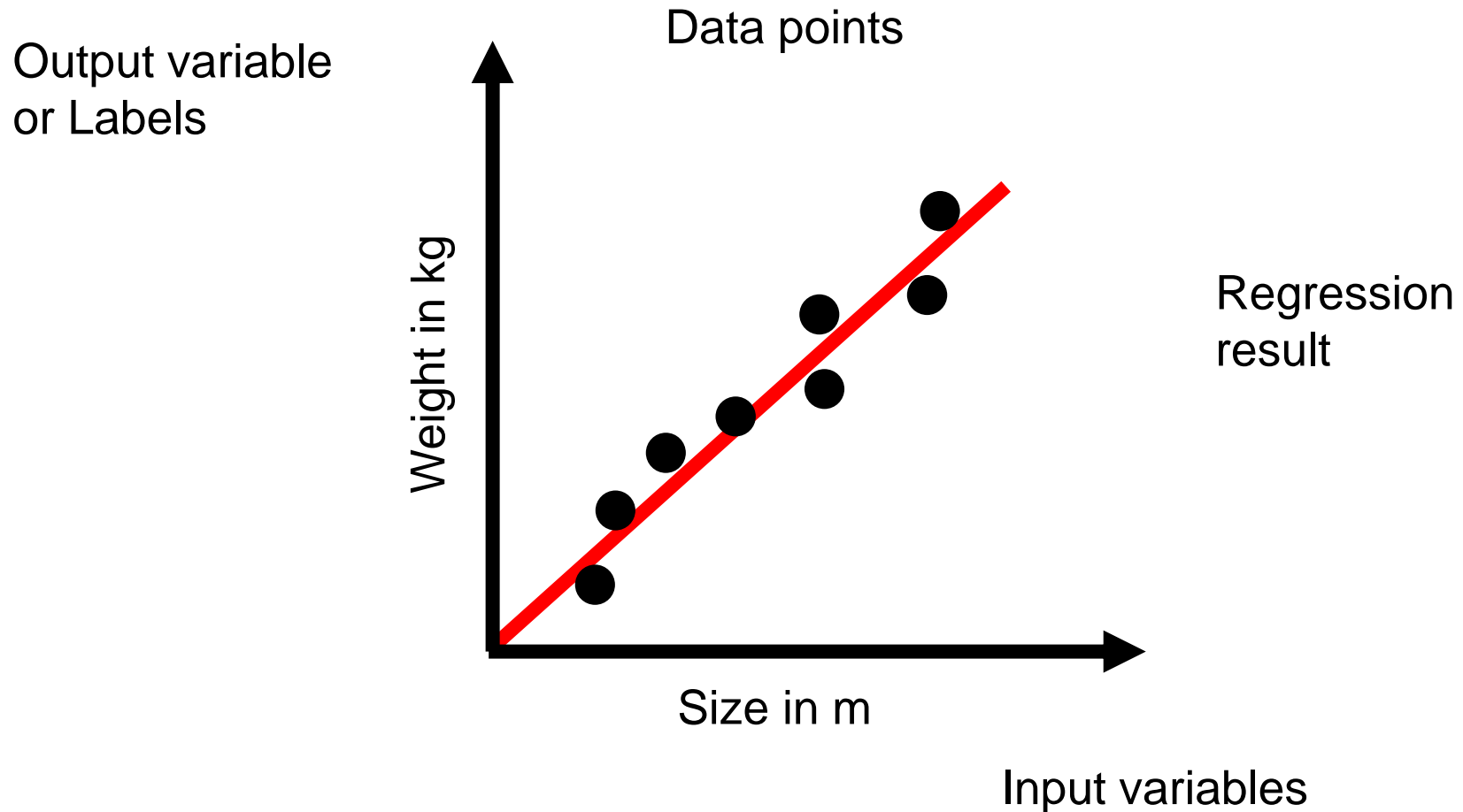
## Agenda

---

1. **Chapter: Motivation**
2. Chapter: Linear Models
3. Chapter: Loss Functions
4. Chapter: Regularization & Validation
5. Chapter: Practical Considerations
6. Chapter: Summary



## Motivation – Regression Example

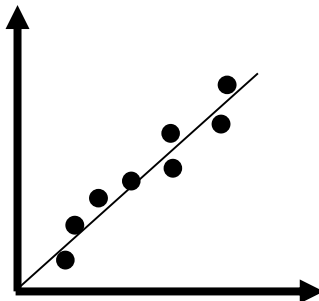


# Motivation – Algorithms in Machine Learning

## Pattern recognition

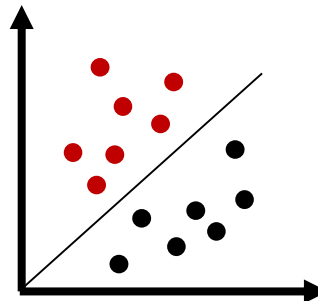
### Regression

- Predict **continuous** values
- Supervised



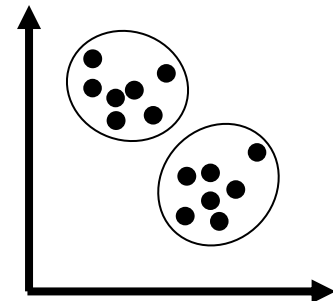
### Classification

- Predict **discrete** values
- Supervised



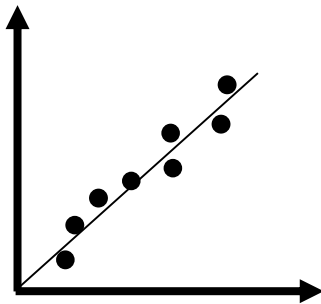
### Clustering

- Predict discrete values
- **Unsupervised**



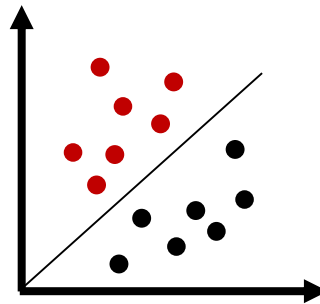
# Motivation – Algorithms in Machine Learning

## Regression



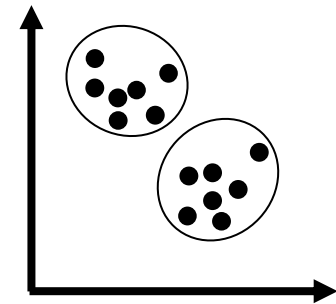
- House pricing
- Sales
- Persons weight

## Classification



- Object detection
- Spam detection
- Cancer detection

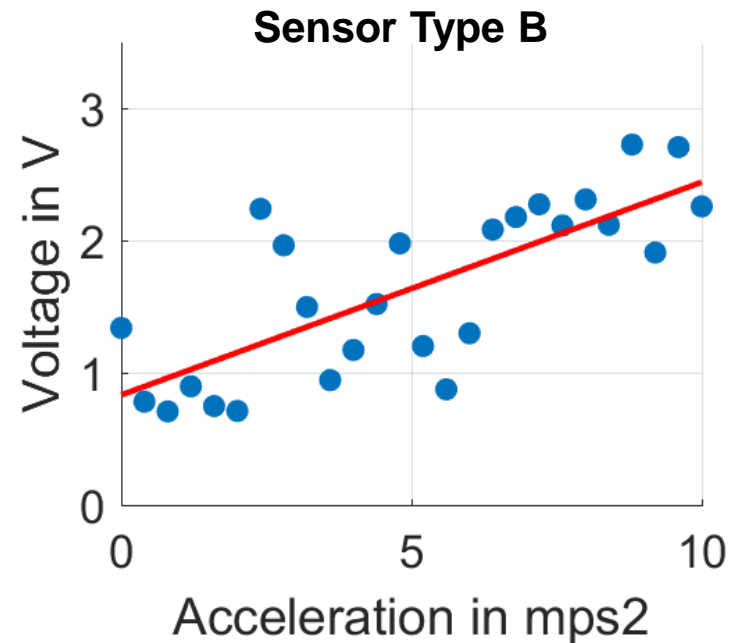
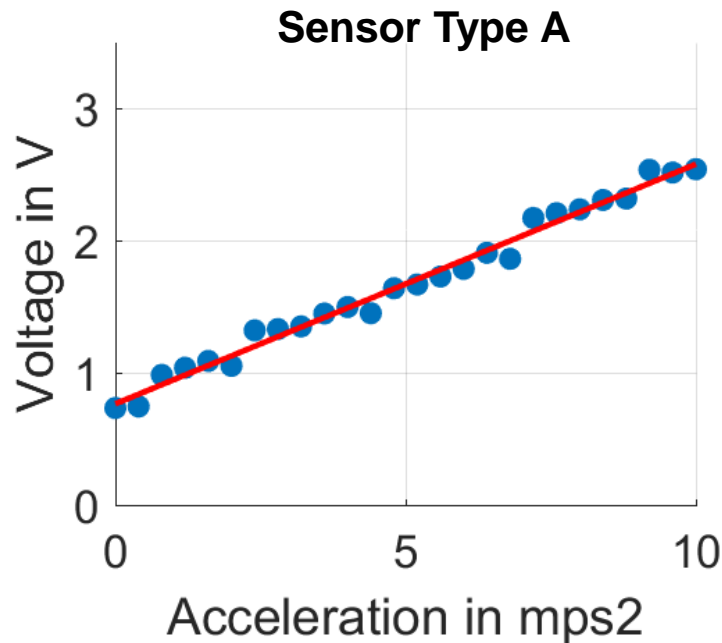
## Clustering



- Genome patterns
- Google news
- Pointcloud (Lidar) processing

# Motivation – Regression in Automotive Technology

## Sensor calibration

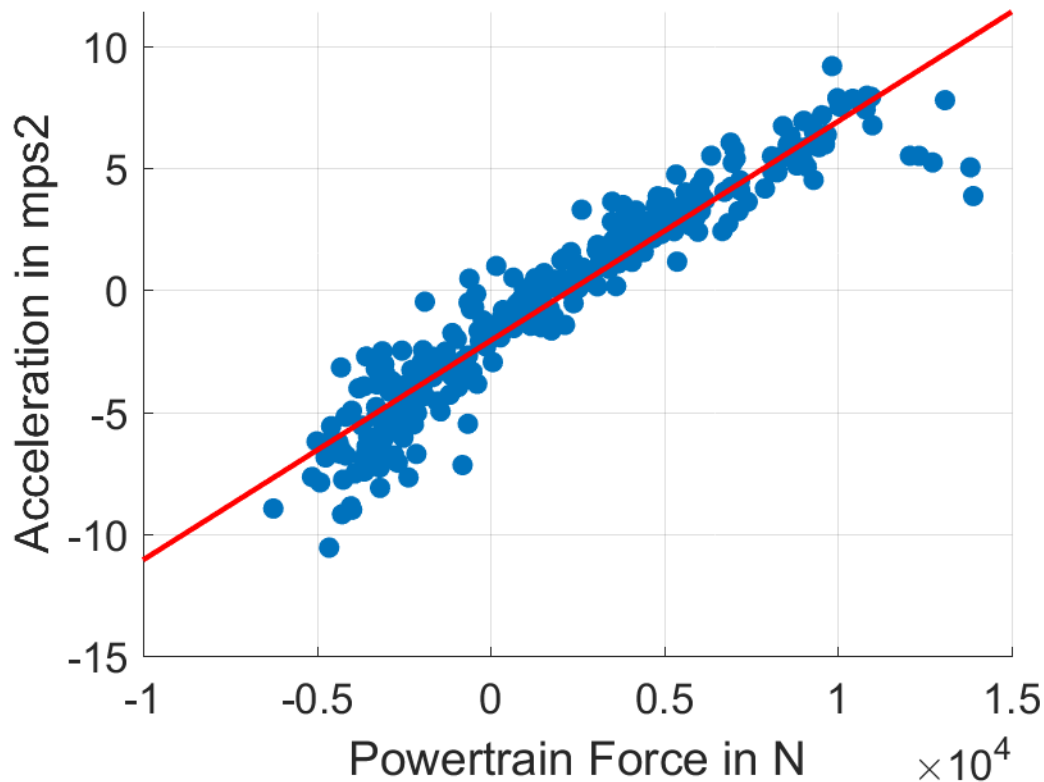


- Usually, electric quantities are measured
- Necessary to convert them to physical quantities
- Examples:
  - Accelerometers
  - Gyroscopes
  - Displacement sensors



# Motivation – Regression in Automotive Technology

## Parameter estimation

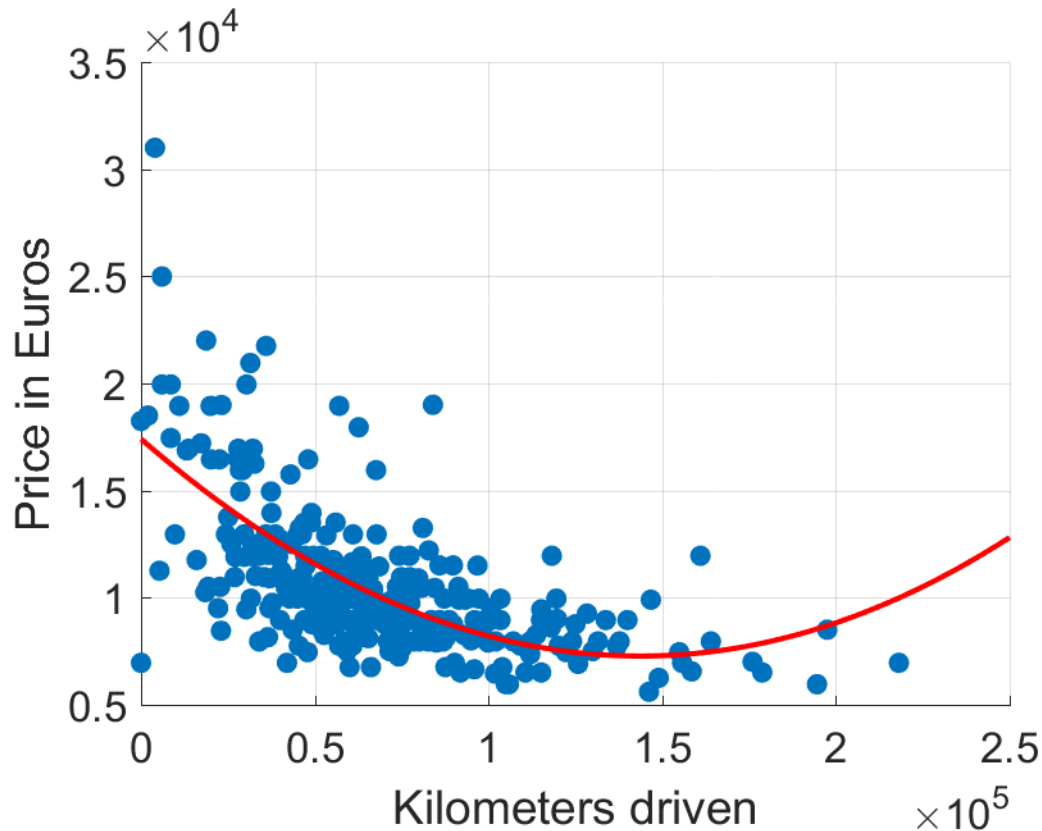


- Vehicle parameters are often only roughly known
- Estimation by regression techniques

$$a_x = \frac{1}{m} F_{PT}$$

# Motivation – Regression in Automotive Technology

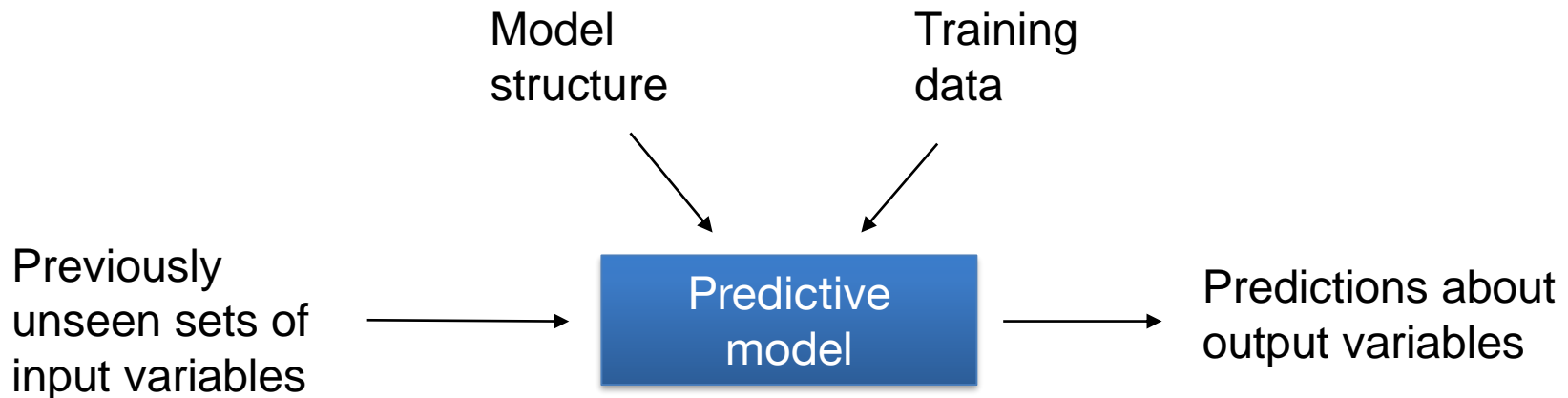
## Vehicle pricing



- Regression is widely used for financial relations
- Allows to compress data into a simple model and evaluate derivatives

## Motivation – Why should you use Regression?

- Based on the combination of data and model structure, it is possible to **predict the outcome** of a process or system
- Training dataset is usually only a representation at sparse points and contains lots of noise
- Allows usage of information in simulation, optimization, etc.



# Relation of Statistics and Machine learning

***How can we extract information from data and use them to reason and predict in beforehand unseen cases? (learning)***

- Nearly all classical machine learning methods can be reinterpreted in terms of statistics
- Focus in machine learning is mainly on prediction
- Statistics often focusses on relation analysis
- Lots of advanced regression techniques build upon a statistical interpretation of regression

## Additional Information

- Definition according to Wikipedia:

In statistical modeling, **regression analysis** is a set of statistical processes for estimating the relationships among variables. It includes many techniques for modeling and analyzing several variables, when the focus is on the relationship between a dependent variable and one or more independent variables (or 'predictors'). More specifically, regression analysis helps one understand how the typical value of the dependent variable (or 'criterion variable') changes when any one of the independent variables is varied, while the other independent variables are held fixed.

Source: [https://en.wikipedia.org/wiki/Regression\\_analysis](https://en.wikipedia.org/wiki/Regression_analysis)

- Definition according to a machine learning blog:

Regression is a method of modelling a target value based on independent predictors. This method is mostly used for forecasting and finding out cause and effect relationship between variables. Regression techniques mostly differ based on the number of independent variables and the type of relationship between the independent and dependent variables.

Source: <https://towardsdatascience.com/introduction-to-machine-learning-algorithms-linear-regression-14c4e325882a>

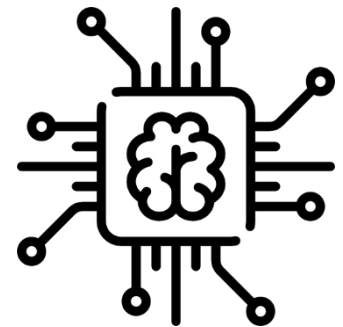
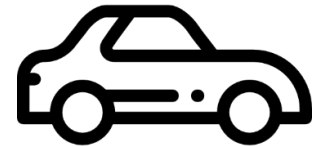
## Supervised Learning: Regression

Johannes Betz / Prof. Dr. Markus Lienkamp / Prof. Dr. Boris Lohmann  
(Alexander Wischnewski, M. Sc.)

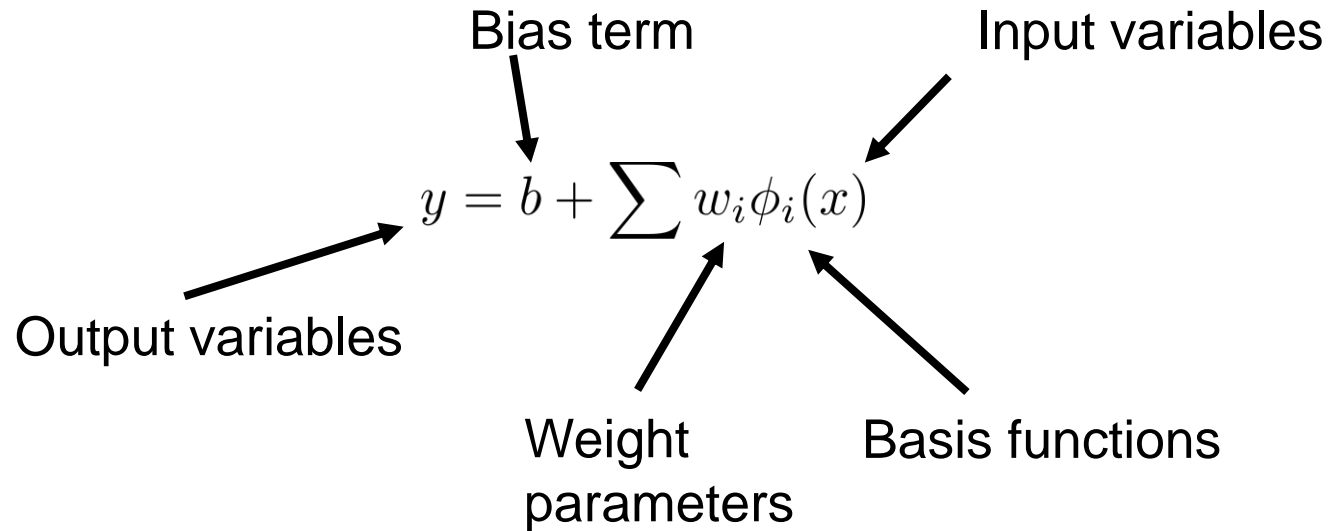
### Agenda

---

1. Chapter: Motivation
2. **Chapter: Linear Models**
3. Chapter: Loss Functions
4. Chapter: Regularization & Validation
5. Chapter: Practical Considerations
6. Chapter: Summary



# Linear Basis Function Model



$$y = \begin{bmatrix} 1 & \phi_1(x) & \dots & \phi_k(x) \end{bmatrix} \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_k \end{bmatrix}$$

Labels and arrows pointing to the matrix equation components:

- Basis functions** points to the row vector  $\begin{bmatrix} 1 & \phi_1(x) & \dots & \phi_k(x) \end{bmatrix}$ .
- Weight parameters** points to the column vector  $\begin{bmatrix} b \\ w_1 \\ \vdots \\ w_k \end{bmatrix}$ .

# Representing the Dataset as a Matrix

$$\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} 1 & \phi_1(x_1) & \dots & \phi_k(x_1) \\ \vdots & \ddots & \ddots & \vdots \\ 1 & \phi_1(x_N) & \dots & \phi_k(x_N) \end{bmatrix} \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_k \end{bmatrix}$$

The diagram illustrates the relationship between the matrix representation of a dataset and the vector equation  $\vec{y} = X\vec{w}$ . The matrix equation at the top shows the output vector  $\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$  as the product of the design matrix  $\begin{bmatrix} 1 & \phi_1(x_1) & \dots & \phi_k(x_1) \\ \vdots & \ddots & \ddots & \vdots \\ 1 & \phi_1(x_N) & \dots & \phi_k(x_N) \end{bmatrix}$  and the weight vector  $\begin{bmatrix} b \\ w_1 \\ \vdots \\ w_k \end{bmatrix}$ . Below this, the vector equation  $\vec{y} = X\vec{w}$  is shown. Four arrows point from the labels 'Output vector', 'Design matrix', and 'Weight vector' to their respective components in the vector equation. A fifth arrow points from the vector equation up to the matrix equation, indicating that the matrix equation is a detailed representation of the vector equation.

Output vector

Design matrix

Weight vector



## Additional Information

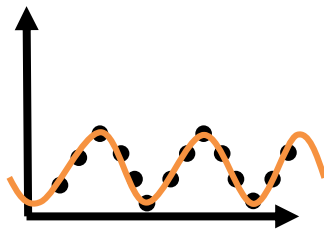
### On the differences between linear regression and nonlinear regression

The term linear is derived from the fact that the problem is 'linear in the parameters'. This is reflected by the fact that the output of a regression can be written as a vector or matrix vector multiplication as depicted on the previous slide. The basis functions are very often nonlinear in the independent variable domain.

There are many extensions to this basic regression technique out there, a few examples can be found below:

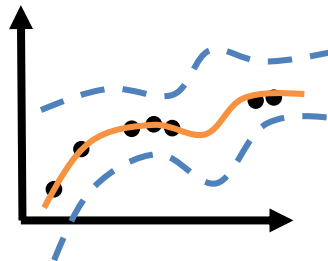
#### Nonlinear optimization

- Specify set of parameters
- Solve nonlinear optimization problem



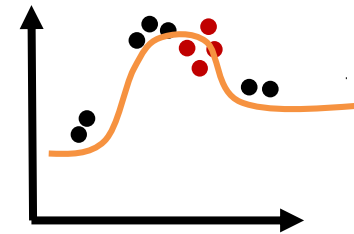
#### Gaussian processes

- Uses kernel methods
- No parameters
- Bayesian interpretation

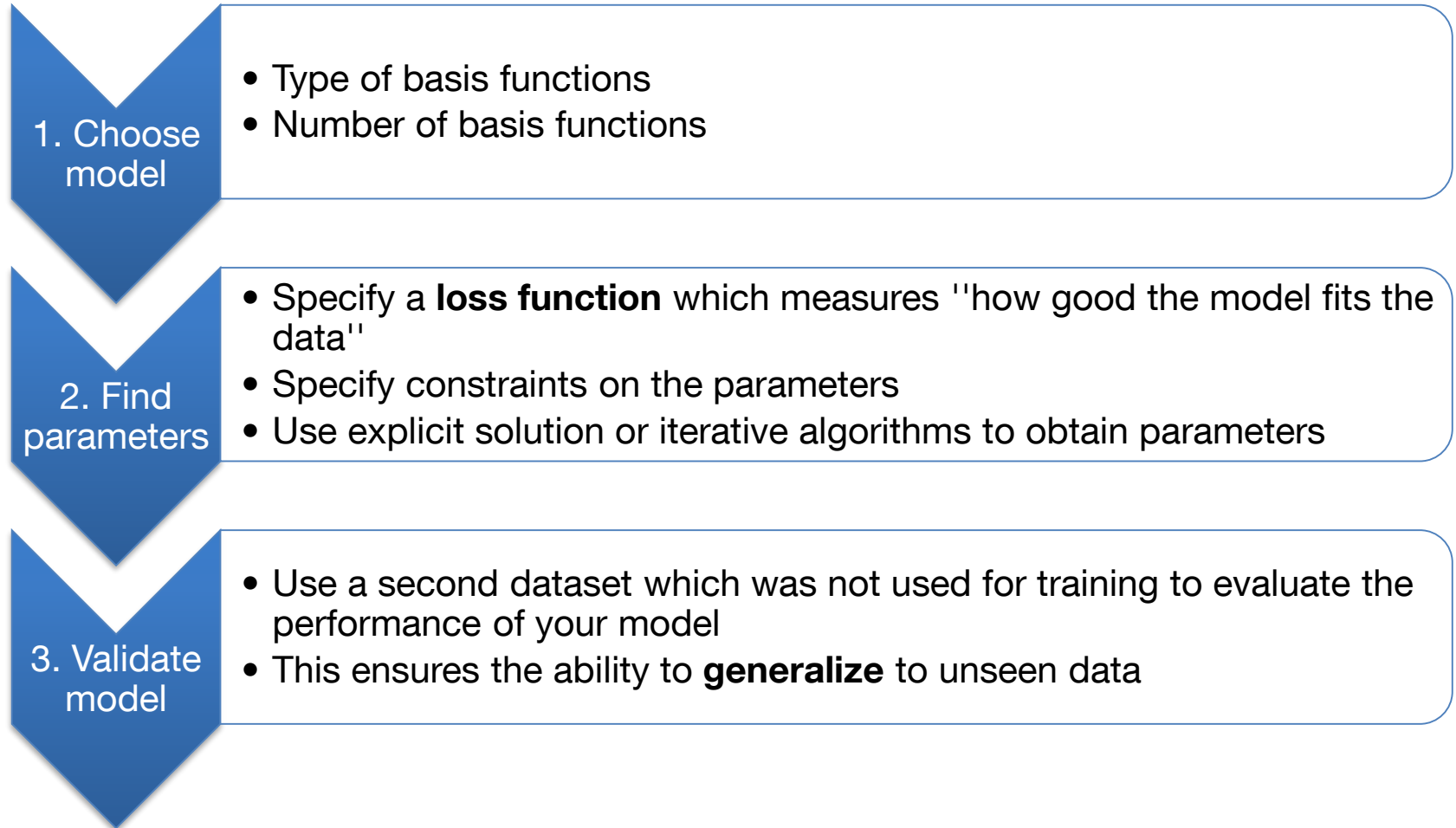


#### Support vector machines

- Uses kernel methods
- No parameters
- Mostly used for classification

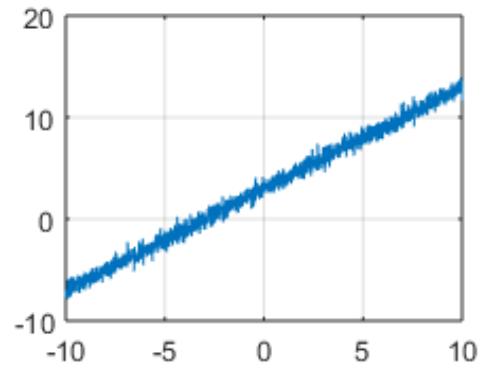


# Workflow – How do we Obtain Model Parameters?

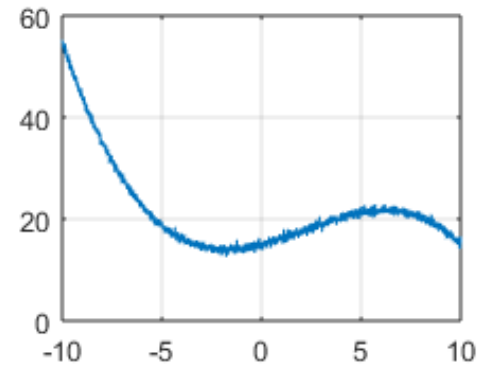


# Choose Model – Examples

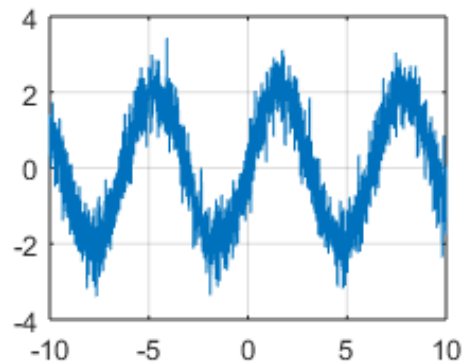
Linear  
function



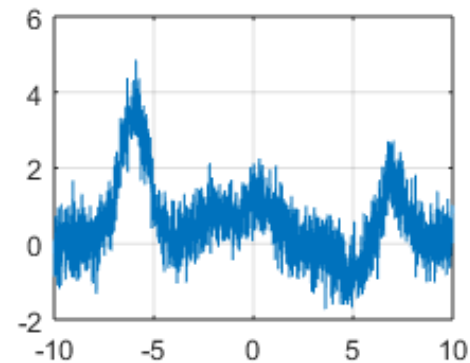
Polynomial  
function



Sinusoidal  
function

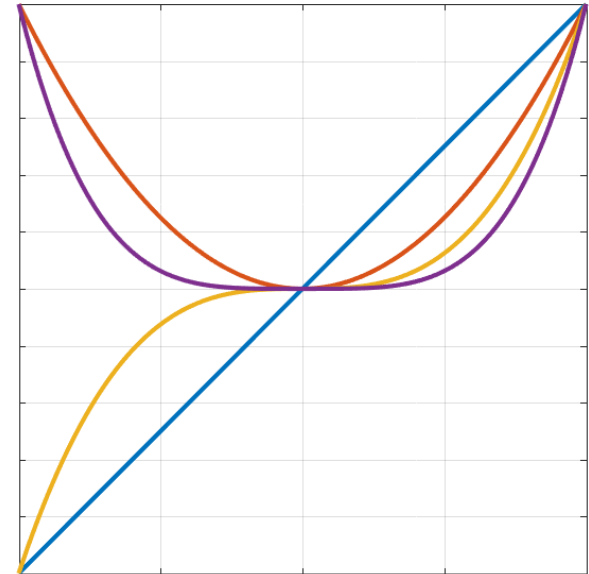


Gaussian  
basis function



## Choose Model – Polynomials

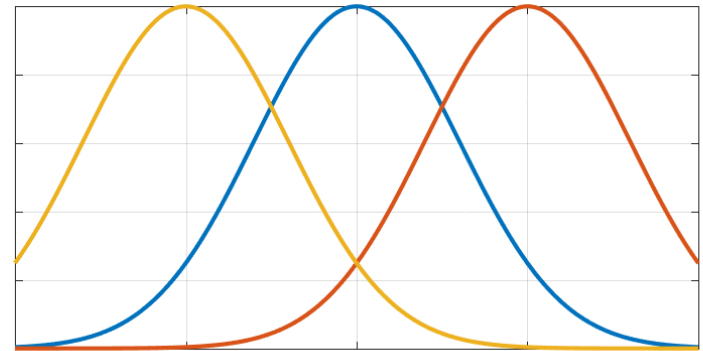
- Globally defined on independent variable domain
- Design matrix becomes ill-conditioned for large input domain variables for standard polynomials
- Hyperparameter:
  - Polynomial degree



$$\begin{array}{ll} \phi_0(x) = 1 & \phi_3(x) = x^3 \\ \phi_1(x) = x & \dots \\ \phi_2(x) = x^2 & \phi_i(x) = x^i \end{array}$$

## Choose Model – Gaussians

- Locally defined on independent variable domain
- Sparse design matrix
- Infinitely differentiable
- Hyperparameter:
  - Number of Gaussian functions
  - Width  $s$  of each basis function
  - Mean  $\mu$  of each basis function



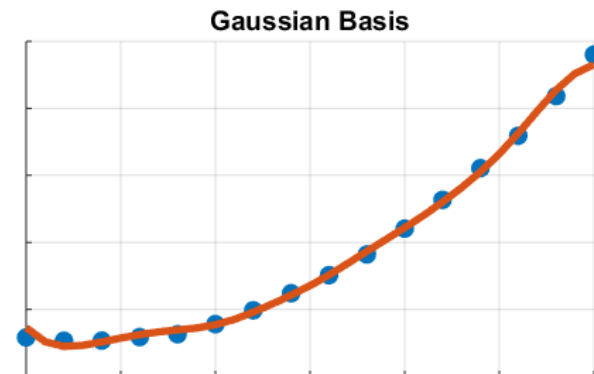
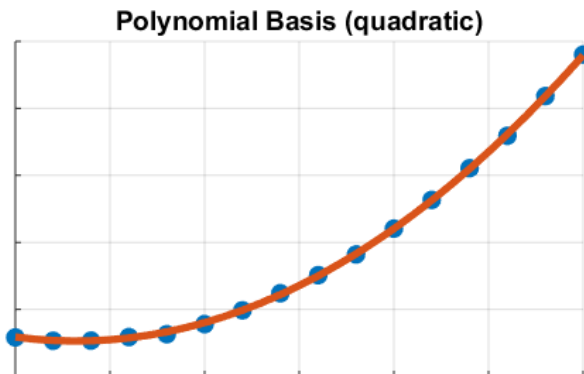
$$\phi(x) = \exp\left(\frac{-(x - \mu)^2}{2s^2}\right)$$

# Choose Model – Comparison of Local and Global BF

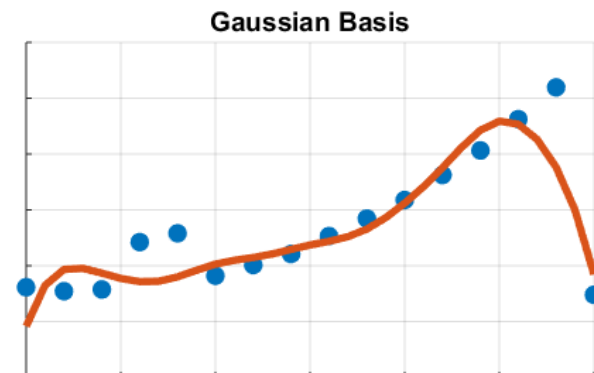
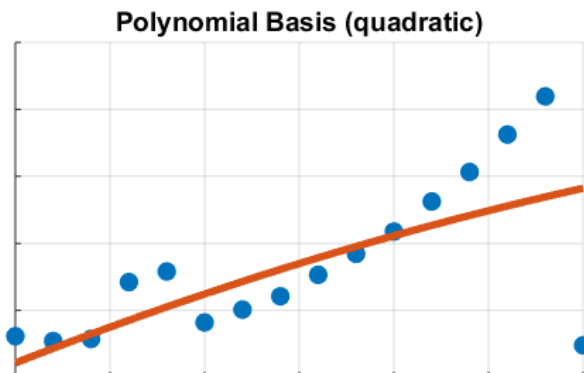
## Global basis function

## Local basis function Spread parameter: 0.3

Dataset A



Dataset B



## Additional Information

The previous slide compares the reaction of regression to outliers if different basis functions are used. All examples are calculated using a quadratic loss function. Details on other possibilities follow in the upcoming slides.

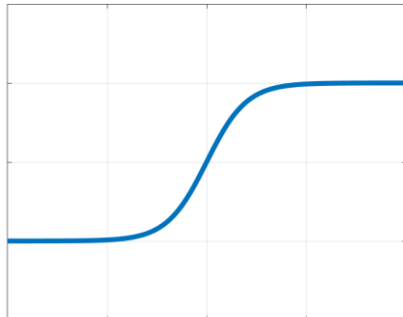
The upper plots are constructed based on slightly noisy observations of a quadratic polynomial function. The ground truth is  $y = 3x^2 - 2x + 3$ . The upper left plot uses a quadratic polynomial as regression model. Since this exactly matches the underlying generation model, the coefficients are estimated to 2.99, -1.96 and 2.96. The estimation quality is really good which is also reflected in the plot. The upper right plot uses seven gaussian basis function equally spaced along the relevant interval of independent variables. The prediction quality is only slightly worse than that of the polynomial model. However, note the little oscillations which are induced by the fact that the regression model does not perfectly fit the underlying model. This leads to local overfitting.

The lower plots contain four outliers. This heavily influences the polynomial model (right[RT]{Nicht left?}). Its characteristic does not fit the rest of the data at all. This is also reflected in the coefficients, which are estimated to -0.38, 5.55 and 0.97. Note that not even the coefficient signs are estimated correctly. As a consequence of the global nature of the basis functions, the model does not match the data over the complete interval. The gaussian model performs way better. It is only affected locally by the outliers and performs well on the rest of the interval.

## Choose Model – Other Basis Functions

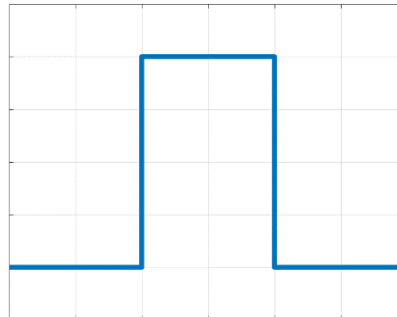
### Sigmoid functions

- Bounded
- Globally defined



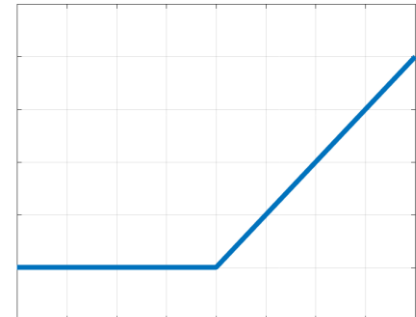
### Bin functions

- Bounded
- Locally defined
- No continuity



### Piecewise linear

- Not bounded
- Globally defined





## Additional Information

Background details on different basis functions:

- <http://madrury.github.io/jekyll/update/statistics/2017/08/04/basis-expansions.html>

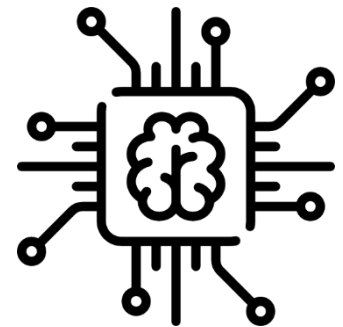
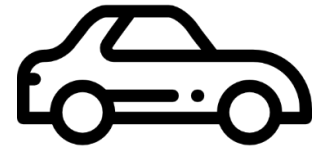
## Supervised Learning: Regression

Johannes Betz / Prof. Dr. Markus Lienkamp / Prof. Dr. Boris Lohmann  
(Alexander Wischnewski, M. Sc.)

### Agenda

---

1. Chapter: Motivation
2. Chapter: Linear Models
- 3. Chapter: Loss Functions**
4. Chapter: Regularization & Validation
5. Chapter: Practical Considerations
6. Chapter: Summary



# Loss Functions

- The loss functions measures the accuracy of a model based on the training dataset
- The best model we can obtain is the minimum loss model
  - Choice of loss function is fundamental in the regression problem

$$\min_{\vec{w}} L(\vec{x}, \vec{y}, \vec{w})$$

- Minimize loss function  $L$  for the training dataset consisting of independent variables  $\vec{x}$  and target variables  $\vec{y}$  by variation of the basis function weights  $\vec{w}$ .

# Loss Functions – Mean Squared Error (MSE or L2)

Pros:

- Very important in practical applications
- Solution can be easily obtained analytically

Cons:

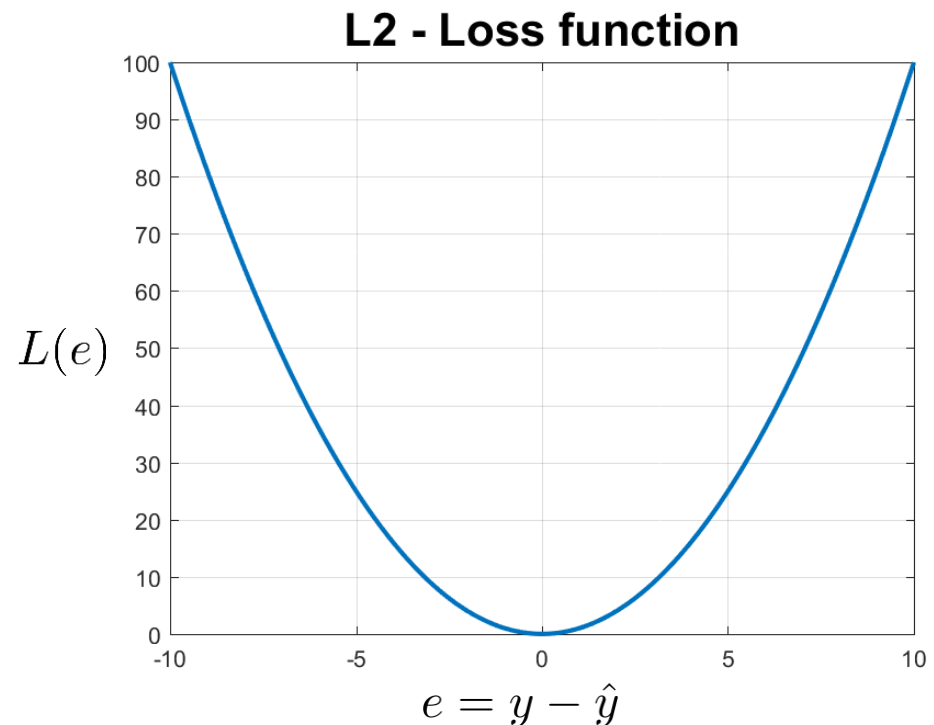
- Not robust to outliers

Examples:

- Basic regression
- Energy optimization
- Control applications

Number of datapoints

$$L = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$



# Loss Functions – Mean Absolute Error (MAE or L1)

Pros:

- Robust to outliers

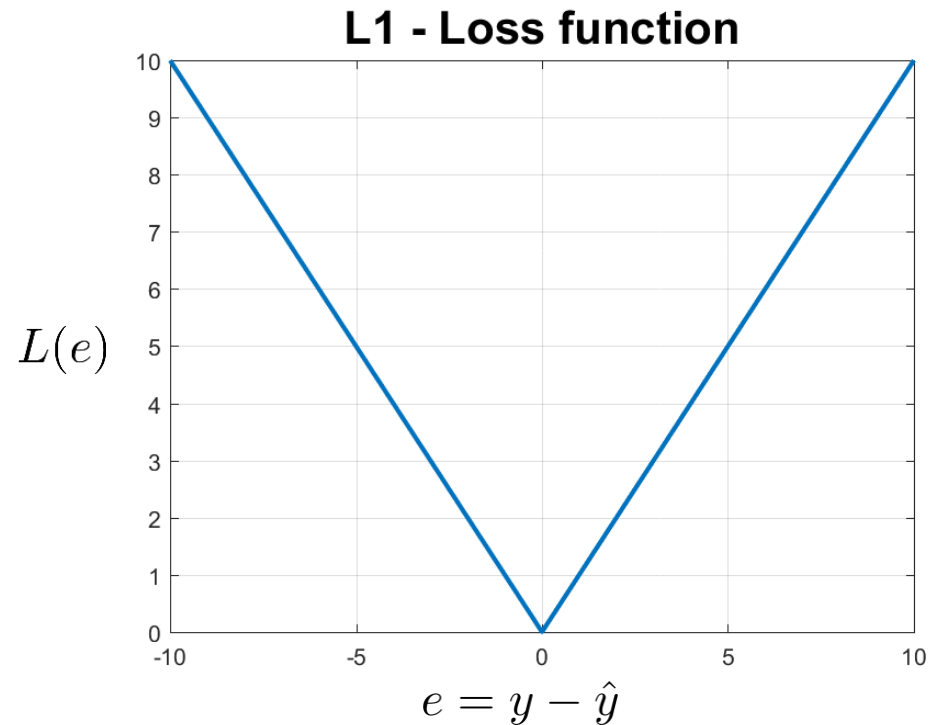
$$L = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

Cons:

- No analytical solution
- Non-differentiable at the origin

Examples:

- Financial applications



# Loss Functions – Huber Loss

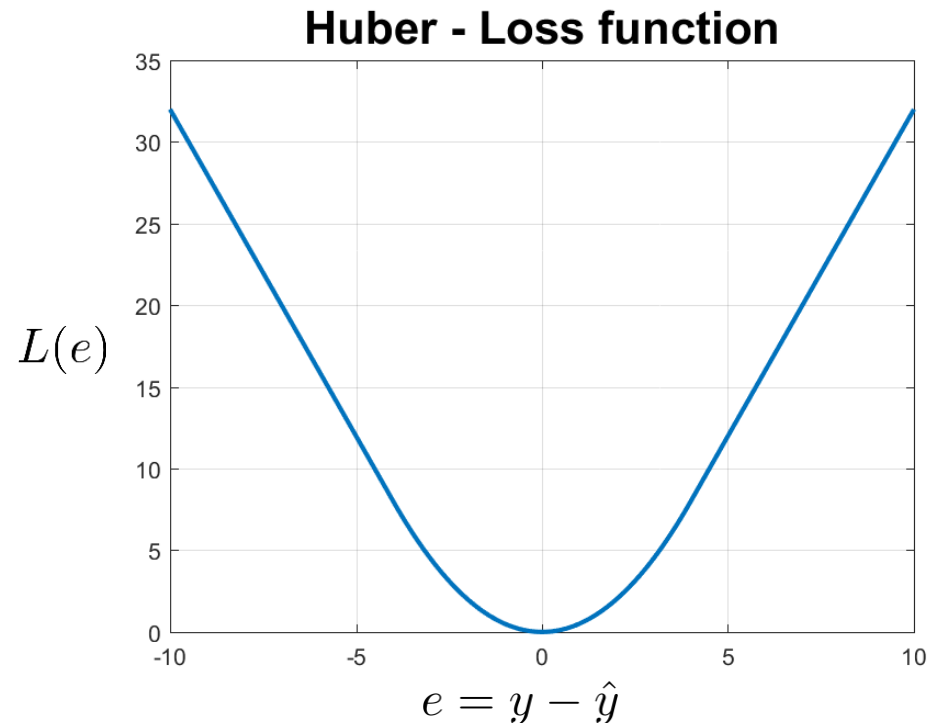
Pros:

- Combines strengths and weaknesses of L1 and L2 loss functions
- Robust + differentiable

Cons:

- More hyperparameters
- No analytical solution

$$L = \frac{1}{N} \sum_{i=1}^N \begin{cases} 0.5 (y_i - \hat{y}_i)^2, & \text{for } |y_i - \hat{y}_i| \leq \delta \\ \delta |y_i - \hat{y}_i| - 0.5\delta^2, & \text{otherwise} \end{cases}$$

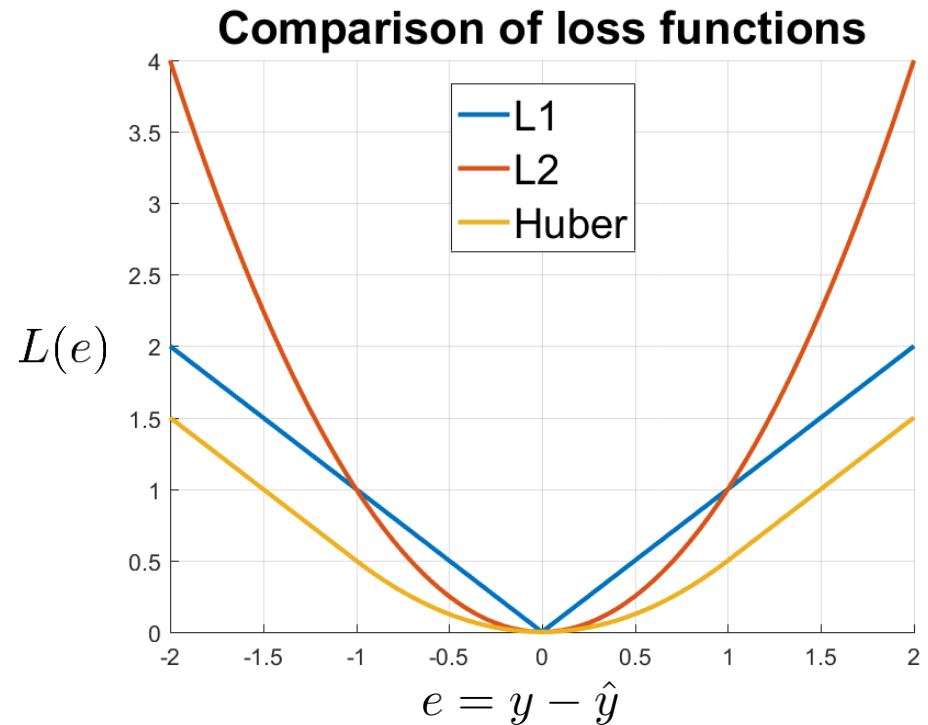


# Loss Functions – Comparison

- L2 loss is differentiable
- L1 loss is more intuitive
- Huber Loss combines theoretical strengths of both

Practical hints:

- Start with L2 loss whenever possible
- Think about physical insights and your intent!



## Additional Information

Additional material on loss functions:

- Good overview and basic discussion of pros and cons of different loss functions

<https://heartbeat.fritz.ai/5-regression-loss-functions-all-machine-learners-should-know-4fb140e9d4b0>

In detail comparison of MSE and MAE loss functions and their robustness

<http://rishy.github.io/ml/2015/07/28/l1-vs-l2-loss/>



## Analytic Solution – Low Dimensional Example

Solve optimization problem

$$\min_{\vec{w}} \frac{1}{2} \sum (y - \hat{y})^2$$

with model  $\hat{y} = w_1 x + b$

Insert model and data points  $\{x_1, y_1\}, \{x_2, y_2\}, \{x_3, y_3\}$

$$\min_{w_1, b} \frac{1}{2} \left[ (y_1 - w_1 x_1 - b)^2 + (y_2 - w_1 x_2 - b)^2 + (y_3 - w_1 x_3 - b)^2 \right]$$

In general, optimal solutions are obtained at the points where the gradient vanishes to zero.

# Analytic Solution – Low Dimensional Example

Calculate gradient

$$\begin{aligned}
 \Delta L(x, y, w) &= \begin{bmatrix} \frac{\partial L(x, y, w)}{\partial w_1} \\ \frac{\partial L(x, y, w)}{\partial b} \end{bmatrix} \\
 &= \begin{bmatrix} -(y_1 - w_1 x_1 - b) x_1 - (y_2 - w_1 x_2 - b) x_2 - (y_3 - w_1 x_3 - b) x_3 \\ -(y_1 - w_1 x_1 - b) - (y_2 - w_1 x_2 - b) - (y_3 - w_1 x_3 - b) \end{bmatrix} \\
 &= \begin{bmatrix} x_1^2 + x_2^2 + x_3^2 & x_1 + x_2 + x_3 \\ x_1 + x_2 + x_3 & 3 \end{bmatrix} \begin{bmatrix} w_1 \\ b \end{bmatrix} - \begin{bmatrix} y_1 x_1 + y_2 x_2 + y_3 x_3 \\ y_1 + y_2 + y_3 \end{bmatrix}
 \end{aligned}$$

and set it equal to zero

$$\begin{bmatrix} x_1^2 + x_2^2 + x_3^2 & x_1 + x_2 + x_3 \\ x_1 + x_2 + x_3 & 3 \end{bmatrix} \begin{bmatrix} w_1 \\ b \end{bmatrix} - \begin{bmatrix} y_1 x_1 + y_2 x_2 + y_3 x_3 \\ y_1 + y_2 + y_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

## Analytic Solution – Low Dimensional Example

Solve resulting equation (also called normal equation):

$$\begin{bmatrix} x_1^2 + x_2^2 + x_3^2 & x_1 + x_2 + x_3 \\ x_1 + x_2 + x_3 & 3 \end{bmatrix} \begin{bmatrix} w_1 \\ b \end{bmatrix} - \begin{bmatrix} y_1 x_1 + y_2 x_2 + y_3 x_3 \\ y_1 + y_2 + y_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} w_1 \\ b \end{bmatrix} = \begin{bmatrix} x_1^2 + x_2^2 + x_3^2 & x_1 + x_2 + x_3 \\ x_1 + x_2 + x_3 & 3 \end{bmatrix}^{-1} \begin{bmatrix} y_1 x_1 + y_2 x_2 + y_3 x_3 \\ y_1 + y_2 + y_3 \end{bmatrix}$$

## **Additional Information**

The previous derivation of the low-dimensional example is a demonstration of the mathematical principles behind the regression problem. The derivation can be conducted in a similar way for more complex cases and will be done in a general form in the next slides.

Analytical solutions are mainly important for small-scale problems and for theoretical analysis. In many cases, it is possible to precompute several parts of the solution, which enables applications on hardware systems with low-computational power.

However, from a practical point of view, it is an important fact that the problem can be cast as an optimization problem. This allows for several modifications, e.g., regularization or adding constraints on physically meaningful variables. Some of them are discussed later in this lecture. The price to pay is a demanding solution procedure in form of a numerical optimization problem. This requires more background knowledge and computational power than the analytical solution.

## Analytic Solution – General Form

- Minimizing MSE loss function can be rewritten in matrix form

$$\min_{\vec{w}} \frac{1}{2} \sum (y - \hat{y})^2$$
$$\min_{\vec{w}} \frac{1}{2} (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w})$$

- Optimum value for  $w$  is equal to setting the gradient to zero and solve for

$$\vec{w} = (X^T X)^{-1} X^T \vec{y}$$

- The importance of this loss function is tightly related to the fact that the analytical solution is available and can be calculated explicitly for low- to medium sized datasets!

## Additional Information

Detailed derivation of least squares solution

$$\begin{aligned} \min_{\vec{w}} L(\vec{x}, \vec{y}, \vec{w}) \\ \min_{\vec{w}} \frac{1}{2} (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w}) \\ \min_w \frac{1}{2} (\vec{y}^T \vec{y} - 2\vec{w}^T X^T \vec{y} + \vec{w}^T X^T X \vec{w}) \end{aligned}$$

Taking the gradient with respect to  $w$

$$\Delta_{\vec{w}} L(\vec{x}, \vec{y}, \vec{w}) = -2X^T \vec{y} + 2X^T X \vec{w}$$

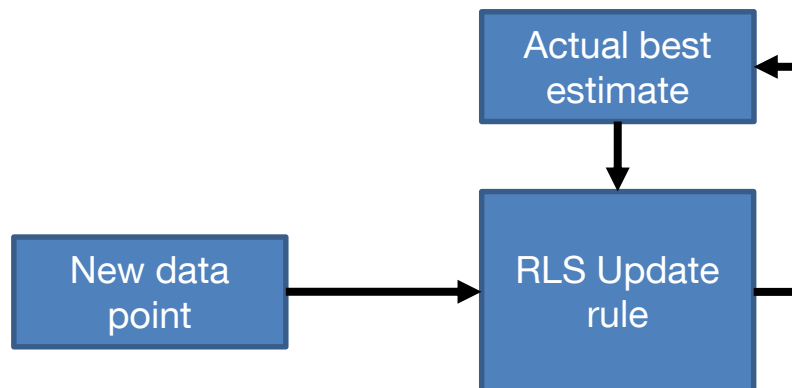
From the necessary condition for a minimum, it follows that

$$\begin{aligned} -2X^T \vec{y} + 2X^T X \vec{w} &= 0 \\ \vec{w} &= (X^T X)^{-1} X^T \vec{y} \end{aligned}$$

The term  $(X^T X)^{-1} X^T$  is also called Moore-Penrose pseudo-inverse. It is also the MSE solution of an overdetermined system of equations of the form  $\vec{b} = A\vec{x}$  with  $A$  being a tall, full-rank matrix. Note the equivalence of the problem formulation using the Data Matrix before.

# Sequential Analytic Solution - Motivation

- Consider the following cases:
  - Apply regression during operation of the product
  - There is not enough memory to store all data points
- A possible solution is given by recursive least squares (RLS)



# Sequential Analytic Solution – The Algorithm

- New data point:  $\{\bar{x}, \bar{y}\}$
- Update parameters

$$\vec{w}(k+1) = \underbrace{\vec{w}(k)}_{\text{Old parameter estimate}} + \underbrace{P(k)\bar{x}^T (I + \bar{x}^T P(k)\bar{x})^{-1}}_{\text{Correction gain}} \underbrace{(\bar{y} - \bar{x}^T \vec{w}(k))}_{\text{Residual}}$$

Prediction based on old parameters

The diagram shows the equation for updating the parameter vector  $\vec{w}$  at step  $k+1$ . The equation is  $\vec{w}(k+1) = \vec{w}(k) + P(k)\bar{x}^T (I + \bar{x}^T P(k)\bar{x})^{-1} (\bar{y} - \bar{x}^T \vec{w}(k))$ . Annotations include: an arrow pointing to  $\vec{w}(k)$  labeled 'Old parameter estimate'; a bracket under  $P(k)\bar{x}^T (I + \bar{x}^T P(k)\bar{x})^{-1}$  labeled 'Correction gain'; a bracket under  $(\bar{y} - \bar{x}^T \vec{w}(k))$  labeled 'Residual'; and an arrow pointing to  $\bar{x}^T \vec{w}(k)$  labeled 'Prediction based on old parameters'.

- And memory matrix  $P$

$$P(k+1) = \left( I - P(k)\bar{x}^T (I + \bar{x}^T P(k)\bar{x})^{-1} \bar{x} \right) P(k)$$

- with  $I$  being the identity matrix of appropriate dimension



## Additional Information

Please note the following facts:

The initialization of this algorithm is a crucial part in its application. There are two major possibilities:

- Initialize by using the first datapoints available
- Initialize using values which you assume to be right, e.g., obtained from basic physical laws

The first technique leads to the solution of the ordinate least squares problem. The second technique however, is tightly related to the regularization techniques already discussed. If one initializes the parameters as follows

$$\begin{aligned}P(0) &= \lambda^{-1}I \\ w(0) &= 0\end{aligned}$$

the resulting solution  $w(k)$  is equivalent to the solution of the L2 regularized least squares problem. Moreover, you have the possibility to set a *starting value*. This can be helpful if a first principles model is available which should be refined online.

Moreover, the presented algorithm is tightly connected to the Kalman filter algorithm. In this case,  $P$  takes the role of the estimation covariance matrix and  $w$  is the corresponding state estimate. More details on this correspondence can be understood from the careful comparison of the update equations of both algorithms. This constitutes the fact, that the Kalman Filter delivers the least squares state estimate.

## Sequential Analytic Solution – Forgetting Factor

- Some applications show slowly varying conditions in long term, but can be considered stationary on short to medium time periods
  - Aging of products leads to slight parameter changes
  - Vehicle mass is usually constant over a significant period of time
- The RLS algorithm can deal with this by introduction of a forgetting factor  $\gamma$ . This leads to a reduction of weights for old samples.

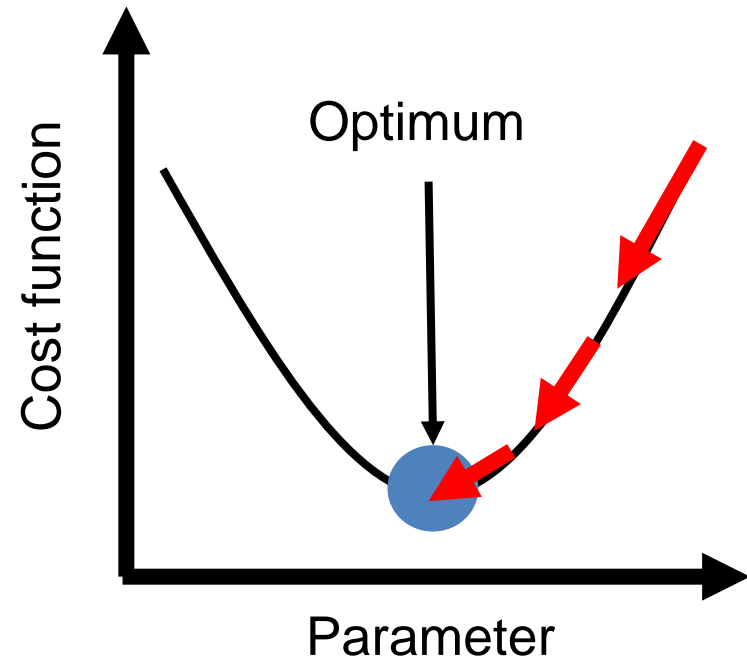
$$\begin{aligned}\vec{w}(k+1) &= \vec{w}(k) + P(k)\bar{x}^T (\gamma I + \bar{x}^T P(k)\bar{x})^{-1} (\bar{y} - \bar{x}^T \vec{w}(k)) \\ P(k+1) &= \gamma^{-1} \left( I - P(k)\bar{x}^T (\gamma I + \bar{x}^T P(k)\bar{x})^{-1} \bar{x} \right) P(k)\end{aligned}$$

# Numerical Iterative Solutions

- Regression can be solved numerically
- Important for large-scale problems and for non-quadratic loss functions
- Popular methods:
  - Gradient descent
  - Gauss-Newton
  - Levenberg-Marquardt

Pros:

- Very generic



Cons:

- Knowledge about numeric optimization necessary

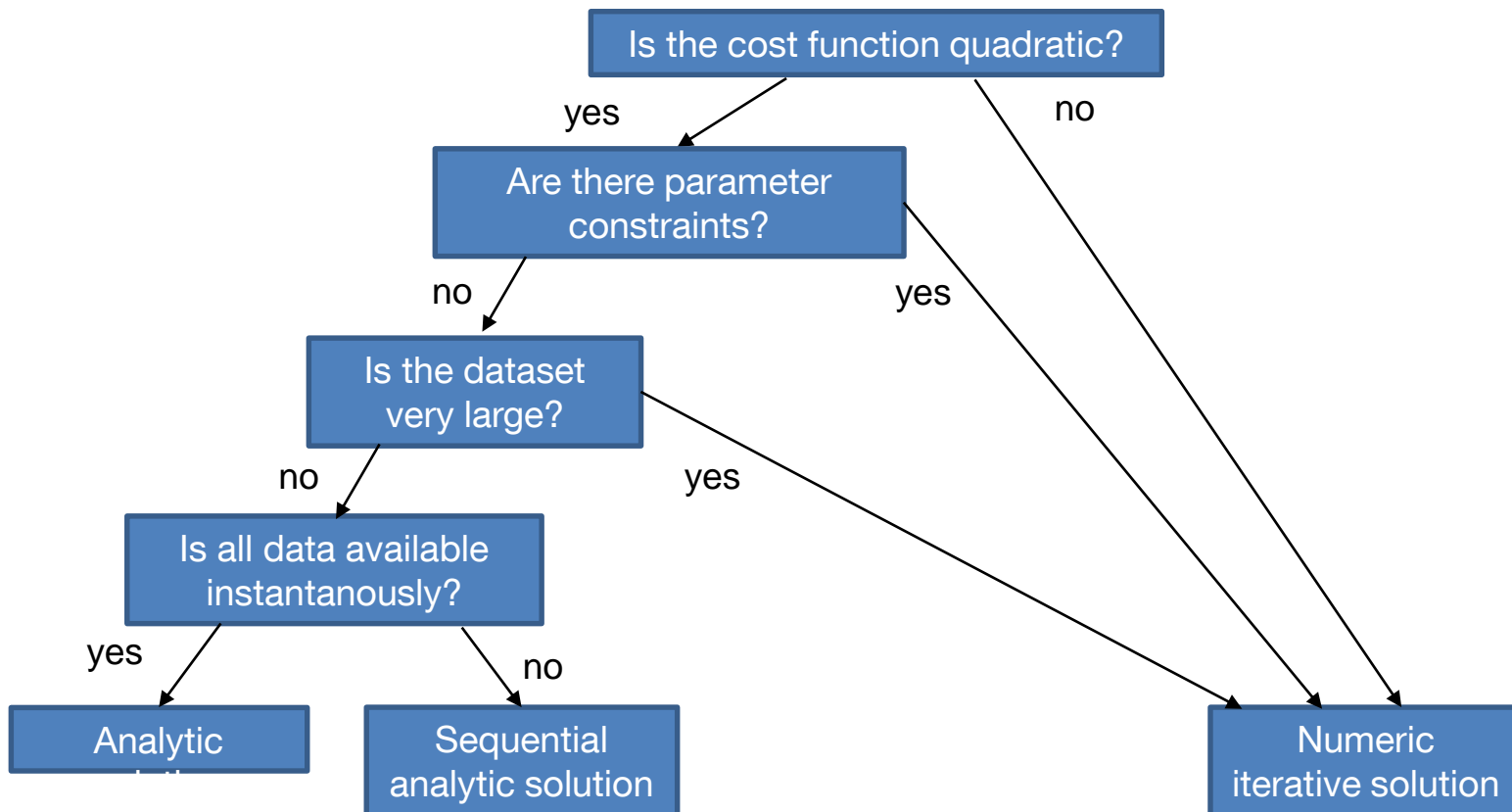
# Constraints on Weights

- Weights can be interpreted as physical quantities
  - Temperature (non-negative)
  - Spring constants (non-negative)
  - Mass (non-negative)
- A valid range is known for the weights
  - Tire and other friction models
  - Efficiency (0 – 100 %)

$$\begin{aligned} & \underset{\vec{w}}{\text{minimize}} \quad L(\vec{x}, \vec{y}, \vec{w}) \\ & \text{subject to} \quad \vec{c}_1 \leq \vec{w} \leq \vec{c}_2 \end{aligned}$$

- Improves robustness
- More difficult to solve

# How to Solve the Regression Problem?



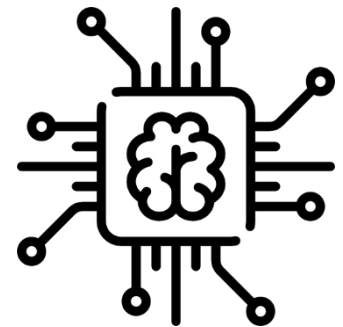
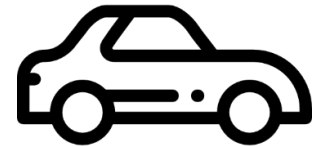
## Supervised Learning: Regression

Johannes Betz / Prof. Dr. Markus Lienkamp / Prof. Dr. Boris Lohmann  
(Alexander Wischnewski, M. Sc.)

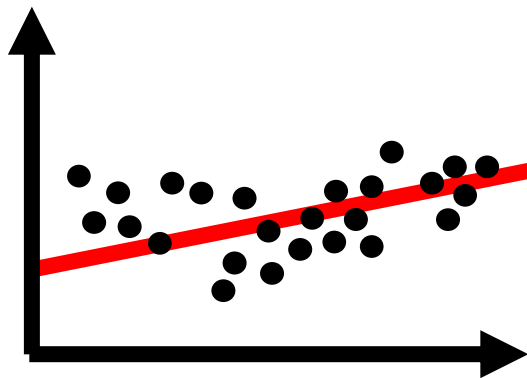
### Agenda

---

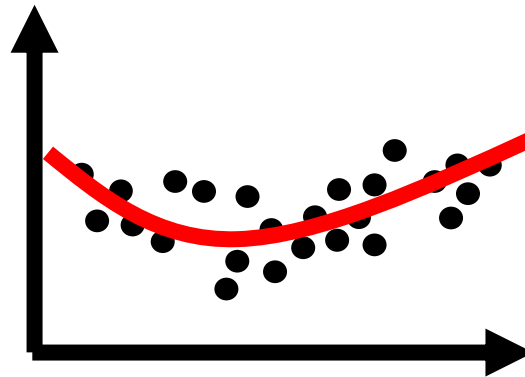
1. Chapter: Motivation
2. Chapter: Linear Models
3. Chapter: Loss Functions
4. **Chapter: Regularization & Validation**
5. Chapter: Practical Considerations
6. Chapter: Summary



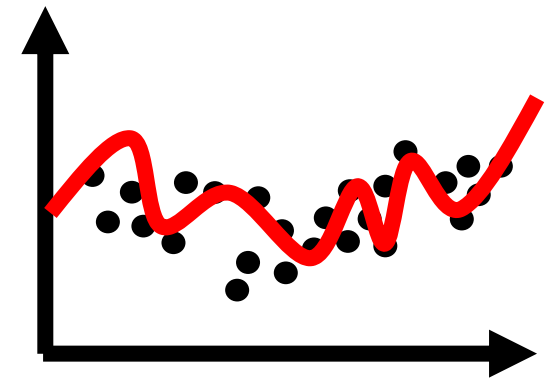
## How to Choose the Model?



Underfitted



Well done



Overfitted

- ←
- Not enough features
  - Wrong structure

- 
- Too many features
  - Irrelevant features

# Overfitting – Choice of Hyperparameters

- Overfitting is the failure to generalize properly between data points
- Cost function decreases with increased model complexity
- Noise and irrelevant effects become too important

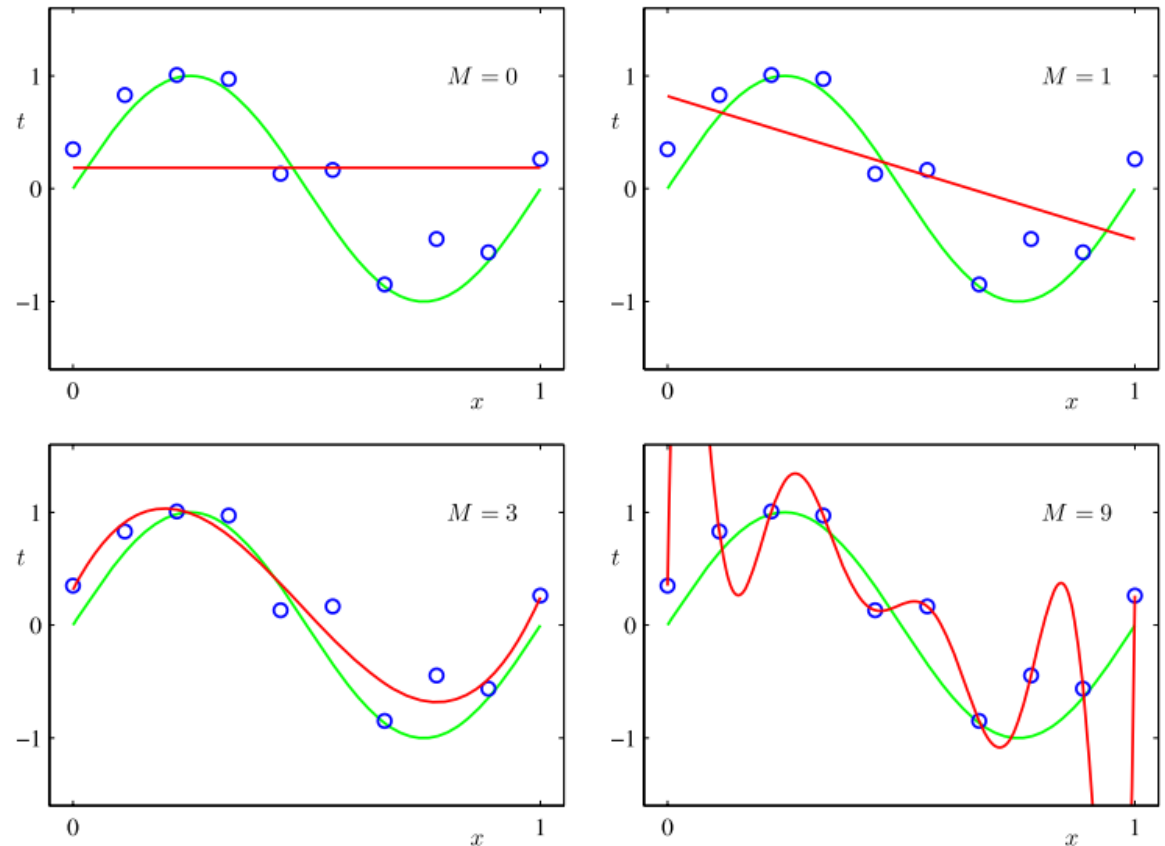


Figure source: Bishop – Pattern Recognition and Machine Learning



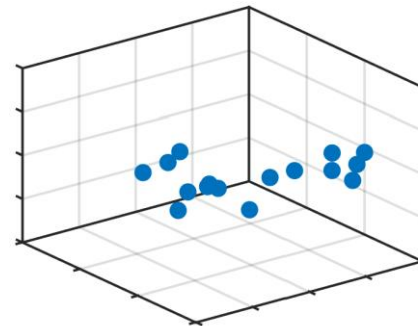
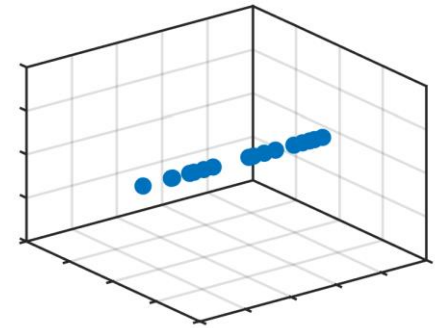
## Additional Information

Overfitting due to large model complexity is one of the most common errors in practical applications. In general, it is advisable to look for the least complex model with a descent performance on the training and test dataset. Its results are more likely to generalize to new data and less prone to numerical issues and outliers.

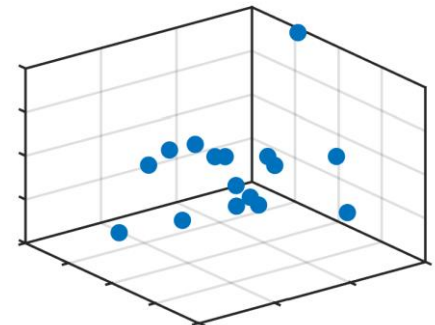
It is difficult to quantify whether a model is overfitted or not. This depends on the application and the structure of the dataset. If you use a model for predictive purposes only, some oscillations in the output domain may be acceptable. If you plan to use the model for analytic purposes, e.g., sensitivity analysis, preventing of overfitting is crucial and you should aim for a less complex model to capture the relevant effects.

# Overfitting – Curse of dimensionality

- Overfitting occurs if
  - Data points are sparse
  - Model complexity is high
- Sparsity of data points is difficult to grasp
- Sparsity increases fast with increased input dimension



16 samples  
in one, two and  
three dimensional  
space



## Additional Information

The curse of dimensionality is a widespread problem along all data based algorithms. The most intuitive is the combinatorial aspect mentioned on the slide before. It affects the required sample size as well as computational efforts needed. Beside the mathematical implications, it is therefore of big impact regarding feasible implementations in a descent time and can be considered critical for the application of systems with real time constraints.

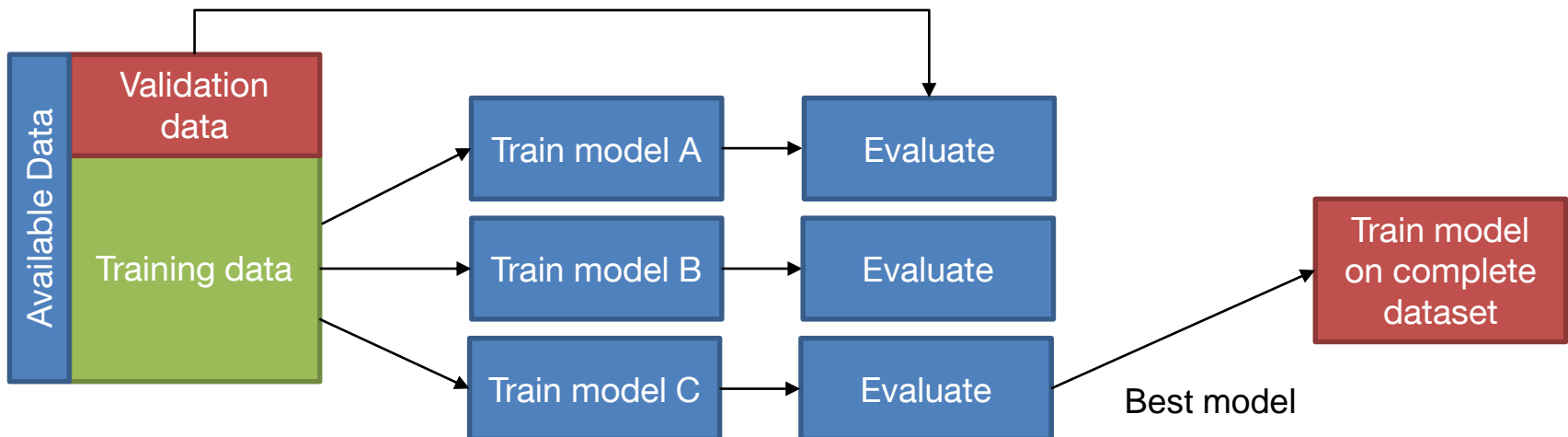
Despite the raw number of samples, it is helpful for regression if samples are "nicely" located in sample space. This usually transforms to well-spread samples along the domain of interest. It is often possible to design sample points in advance to the experiment carried out. In this case, techniques like design of experiment can help to optimize the location of sample points. This helps increasing the model quality with a limited sample size. (Further details can be found here:

[https://en.wikipedia.org/wiki/Design\\_of\\_experiments](https://en.wikipedia.org/wiki/Design_of_experiments))

There are other effects related to the dimensionality of the spaces as well. A fundamental one which is more difficult to grasp is the fact, that distances become less meaningful in high dimensional space. (Details can be found here: [https://en.wikipedia.org/wiki/Curse\\_of\\_dimensionality](https://en.wikipedia.org/wiki/Curse_of_dimensionality))

# Validation Datasets

- Difficult to judge overfitting in high-dimensional domains and autonomous systems
- A standard technique is to separate the data into training and validation data



## Validation Datasets

- Different hyperparameters can be used to tune the model
- Validation technique works for all of them

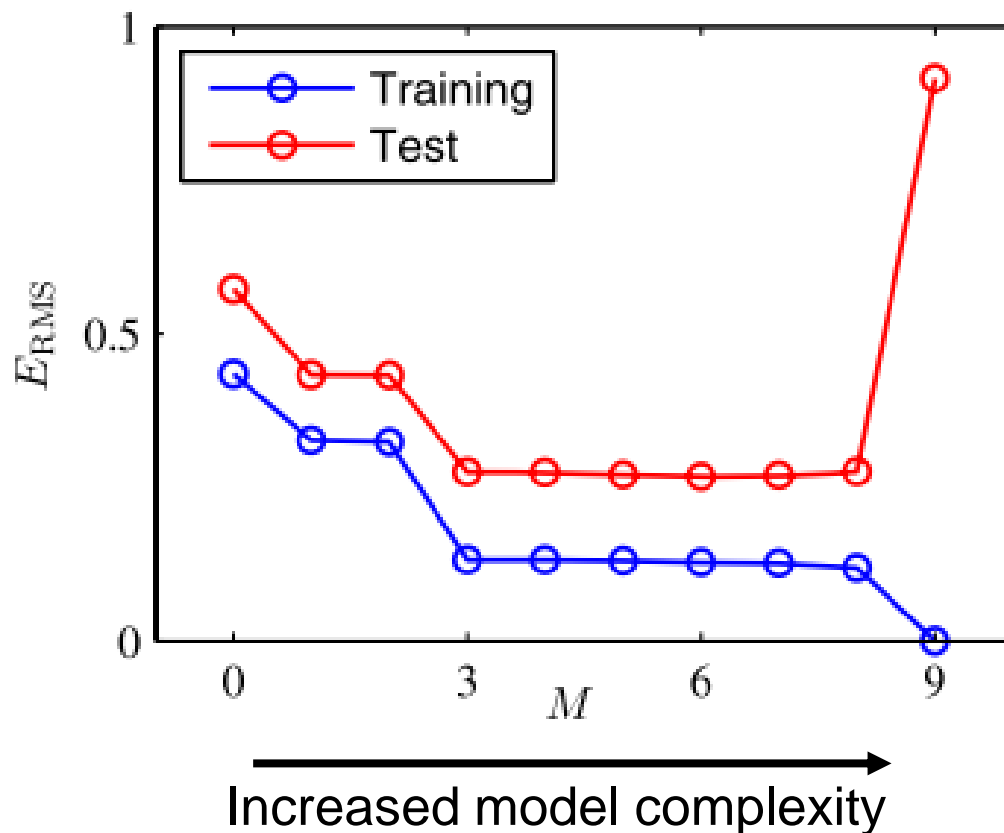


Figure source: Bishop –  
Pattern Recognition and  
Machine Learning

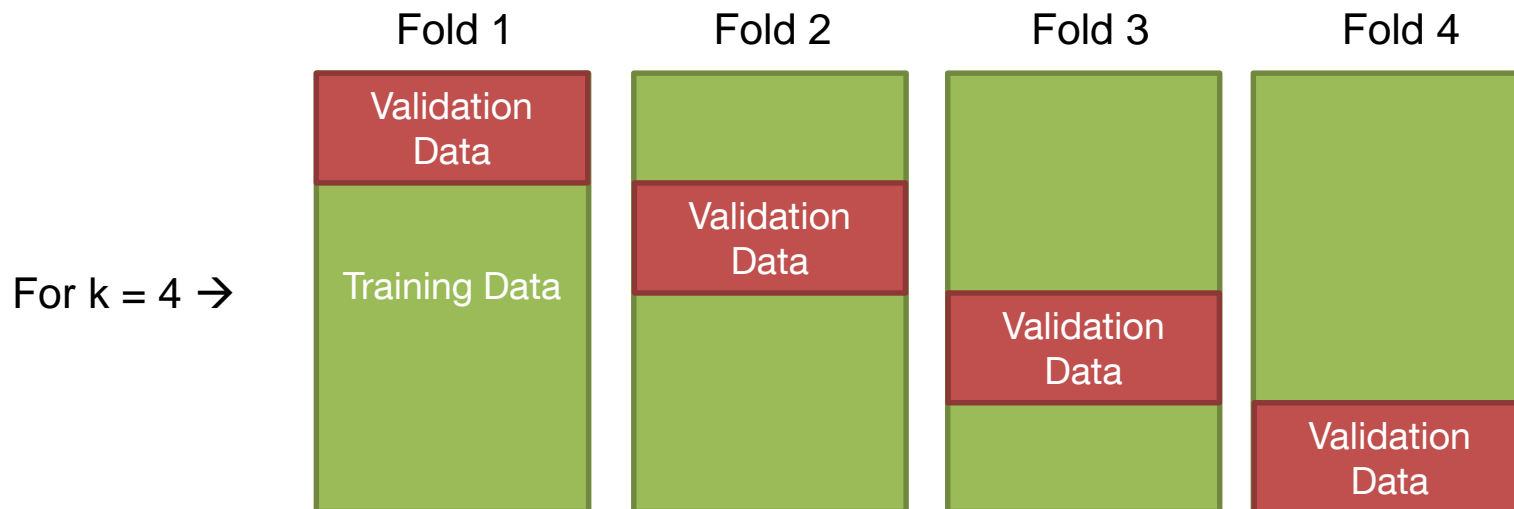
## Common Pitfalls with Validation Datasets

- Be aware that your validation dataset must **reflect the future properties** of the underlying physical relationship.
- **Do not reuse validation datasets.** If the same validation set is used again and again for testing the model performance it is somehow incorporated into the modelling process and does not give the expected results anymore!
- Split the data before fitting the model is therefore essential. Taking  $\frac{2}{3}$  of the data as training data is a good starting value and keep  $\frac{1}{3}$  as test dataset for final performance evaluation

***Visualize your data as much as possible!***

## k-Fold Cross-Validation

- In case of limited data size sets, one may not want to remove a substantial part of the data for the fitting process
- One can use smaller validation sets to estimate the true prediction error by splitting the data into multiple "folds"
- Variance of the estimation error is an indicator for model stability



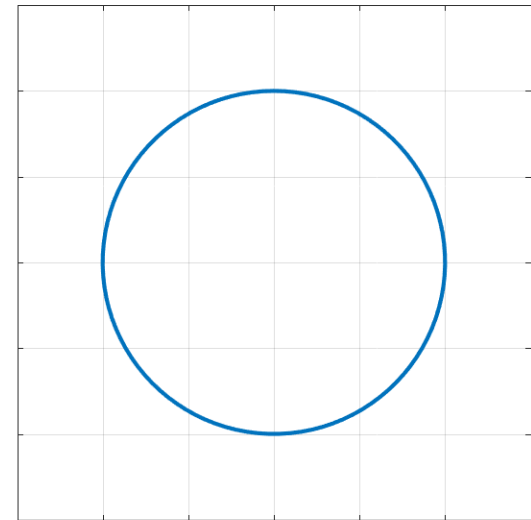
# Regularization

- From a design point of view, we want to choose the model structure based on underlying physical principles and not on characteristics of the dataset
  - Polynomial basis functions tend to have large coefficients for sparse datasets
  - Gaussian basis functions tend to overfit locally leading to single, large coefficients
- A technique to circumvent this is regularization
  - Penalize high coefficients in the optimization prevents these effects
  - Weighting of the penalty term gives an intuitive hyperparameter to control model complexity



# Typical Regularization – Ridge Regression

- Other names: L2 regularization, Thikonov regularization
- Prevents overfitting well
- Analytic solution is available as an extension to the MSE problem
- Difficult to apply and tune in high-dimensional feature spaces



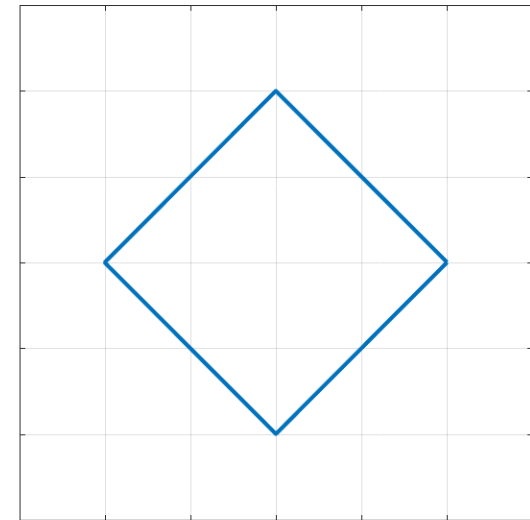
$$\min_{\vec{w}} L(\vec{x}, \vec{y}, \vec{w}) + \lambda \vec{w}^T \vec{w}$$

Regularization term

# Typical Regularization – Lasso Regression

- Other names: L1 regularization
- Tends to produce sparse solutions and can therefore be applied for feature selection
- Sparse solution means that several coefficients go to zero:

$$\vec{w} = [0 \quad w_1 \quad 0 \quad 0 \quad 0 \quad w_2]$$



$$\min_{\vec{w}} L(\vec{x}, \vec{y}, \vec{w}) + \lambda \sum_i |w_i|$$

Regularization term

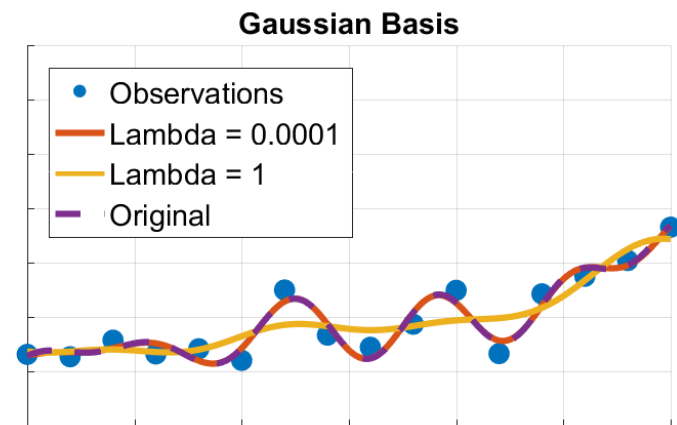
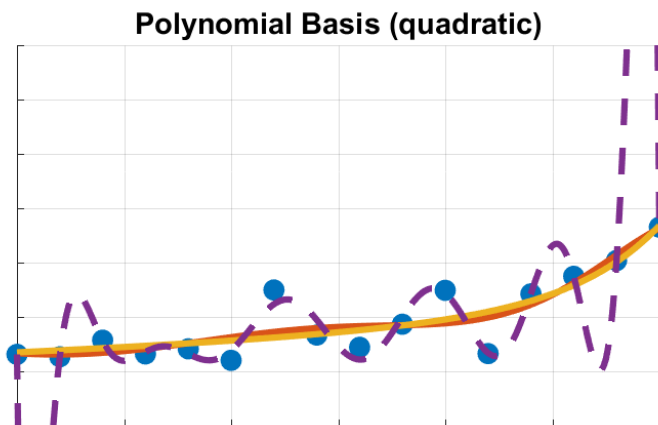
## Additional Information

Comparison between ridge and lasso regularization:

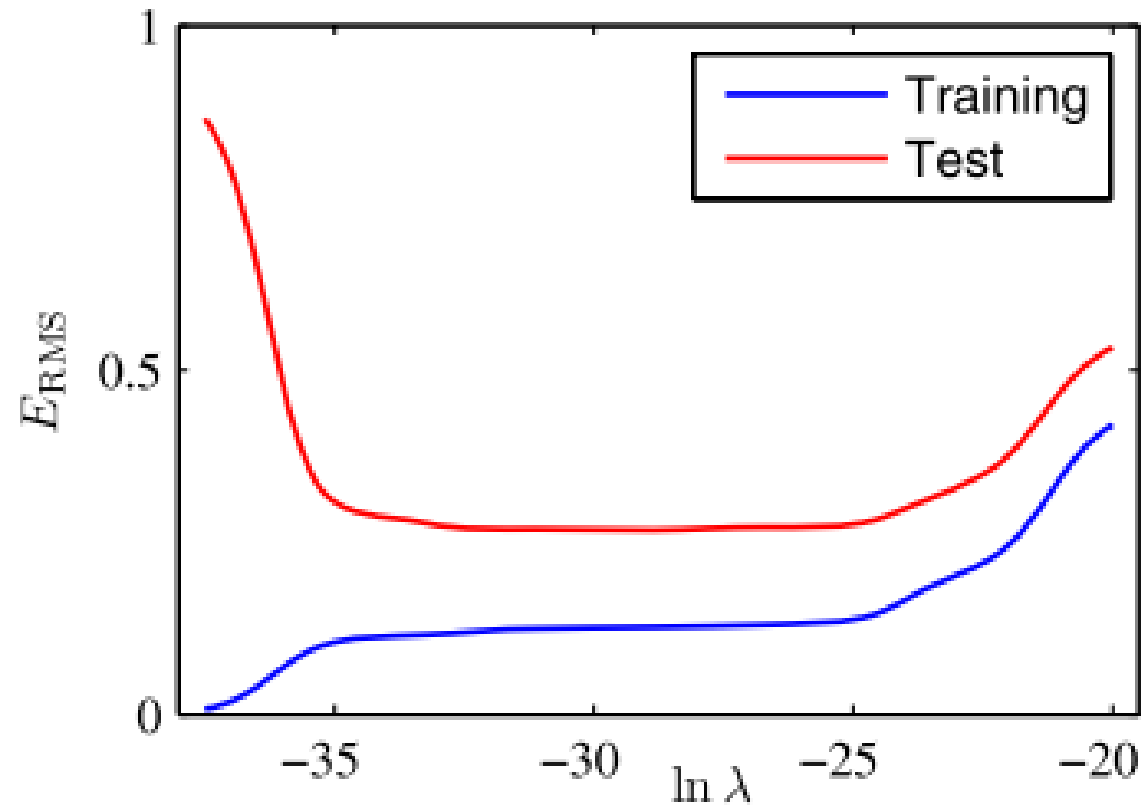
- <https://www.analyticsvidhya.com/blog/2016/01/complete-tutorial-ridge-lasso-regression-python/>
- <https://codingstartups.com/practical-machine-learning-ridge-regression-vs-lasso/>

# Application of Regularization

- Use high dimensional input space for test model
- Ridge regression is applied
- Regularized solutions perform far better at interpolation
- Note: You must evaluate at points BETWEEN your sample points



# Tuning Regularization with Validation Data



←  
Increased model complexity

Figure source: Bishop –  
Pattern Recognition and  
Machine Learning

## Bias Variance Trade Off

- The issue of over- and underfitting can be formalized
- Study the predictor performance on a previously unseen dataset
- Assume that the observations  $y = f(x) + \epsilon$  are corrupted by noise
- The variance of the test data is then written as

$$\mathbb{E} \left( \left( y - \hat{f}(x) \right)^2 \right) = \underbrace{\mathbb{E} \left( f(x) - \hat{f}(x) \right)^2}_{\text{Bias}} + \underbrace{\mathbb{E} \left( \hat{f}(x)^2 \right) - \left( \mathbb{E} \left( \hat{f}(x) \right) \right)^2}_{\text{Variance}} + \underbrace{\sigma^2}_{\text{Noise}}$$

## Additional Information

Derivation of the bias-variance decomposition following

<https://towardsdatascience.com/mse-and-bias-variance-decomposition-77449dd2ff55>

$$\mathbb{E} \left( \left( y - \hat{f}(x) \right)^2 \right) \\ \mathbb{E} (y^2) + \mathbb{E} \left( \hat{f}(x)^2 \right) - \mathbb{E} \left( 2y\hat{f} \right)$$

Using the fact that  $\mathbb{E} (a^2) + \text{Var}(x) + \mathbb{E}(a)^2$  holds for the first and the second term we obtain

$$\text{Var}(y) + \mathbb{E}(y)^2 + \text{Var} \left( \hat{f}(x) \right) + \mathbb{E} \left( \hat{f}(x) \right)^2 - \mathbb{E}(2y\hat{f}(x))$$

$$\text{Var}(y) + \mathbb{E}(y)^2 + \text{Var} \left( \hat{f}(x) \right) + \mathbb{E} \left( \hat{f}(x) \right)^2 - 2\mathbb{E} \left( f(x)\hat{f}(x) \right) - 2\mathbb{E} \left( \epsilon\hat{f}(x) \right)$$

Since we evaluate on previously unseen data, the last term vanishes (the noise of new measurements is independent of the fitted model and zero mean).

$$\text{Var} (f(x)) + \text{Var}(\epsilon) + \mathbb{E}(y)^2 + \text{Var} \left( \hat{f}(x) \right) + \mathbb{E} \left( \hat{f}(x) \right)^2 - 2\mathbb{E} \left( f(x)\hat{f}(x) \right)$$

$$\text{Var} (f(x)) + \text{Var}(\epsilon) + \mathbb{E}(y)^2 + \text{Var} \left( \hat{f}(x) \right) + \mathbb{E} \left( \hat{f}(x) \right)^2 - 2\mathbb{E} (f(x)) \mathbb{E} \left( \hat{f}(x) \right) - 2\text{cov} \left( f(x), \hat{f}(x) \right)$$

$$\text{Var} \left( f(x) - \hat{f}(x) \right) + \text{Var}(\epsilon) + \mathbb{E} (y)^2 + \mathbb{E} \left( \hat{f}(x) \right)^2 - 2\mathbb{E} (f(x)) \mathbb{E} \left( \hat{f}(x) \right)$$

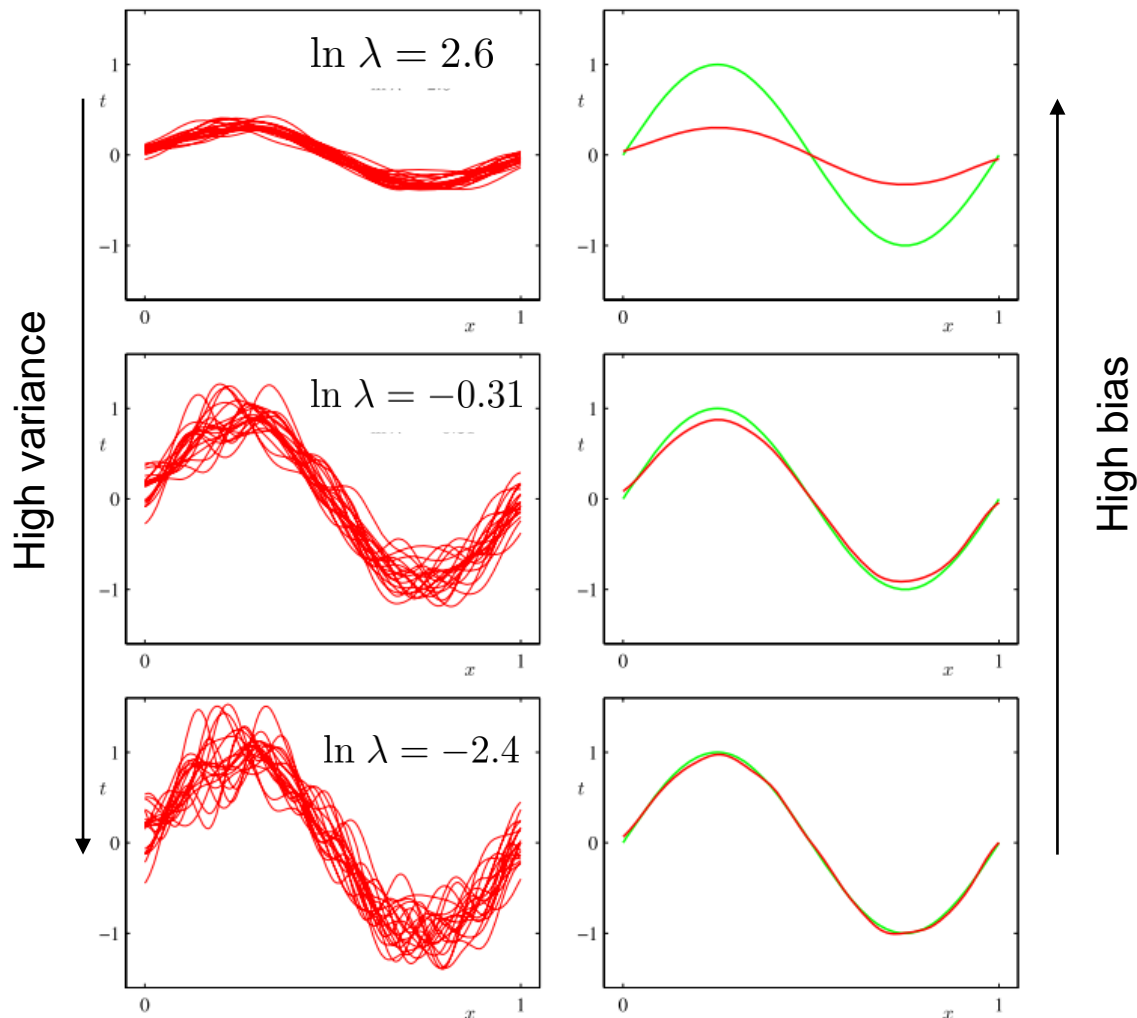
$$\text{Var} \left( f(x) - \hat{f}(x) \right) + \text{Var}(\epsilon) + \mathbb{E} (f(x))^2 + \mathbb{E} (\epsilon)^2 + 2\mathbb{E} (f(x)) \mathbb{E} (\epsilon) + \mathbb{E} \left( \hat{f}(x) \right)^2 - 2\mathbb{E} (f(x)) \mathbb{E} \left( \hat{f}(x) \right)$$

$$\text{Var} \left( f(x) - \hat{f}(x) \right) + \text{Var}(\epsilon) + \left( \mathbb{E} (f(x)) - \mathbb{E} \left( \hat{f}(x) \right) \right)^2 + \mathbb{E} (\epsilon)^2 + 2\mathbb{E} (f(x)) \mathbb{E} (\epsilon)$$

$$\text{Var} \left( f(x) - \hat{f}(x) \right) + \left( \mathbb{E} (f(x)) - \mathbb{E} \left( \hat{f}(x) \right) \right)^2 + \text{Var}(\epsilon)$$

# Regularization in Context of Bias-Variance Trade-Off

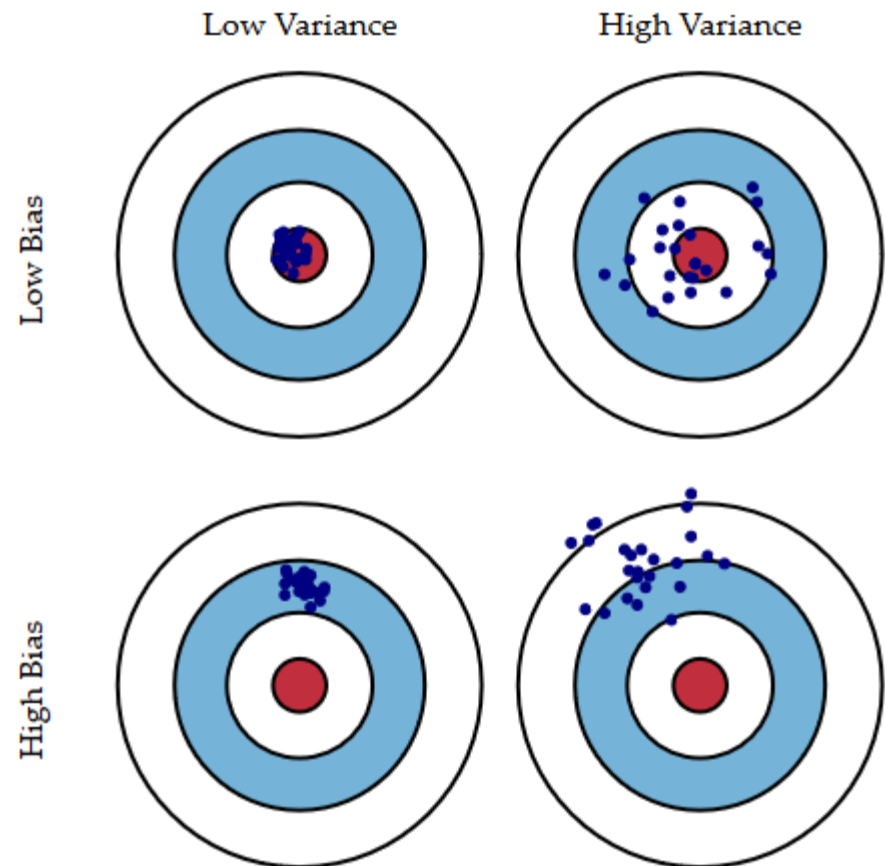
- Imagine you can repeat the fitting process lots of times
- This renders the fitting process of a random experiment
- Model bias: The mean error over all possible fits of the model
- Model variance: The variance of all possible fits of the model





# Regularization in Context of Bias-Variance Trade-Off

- Bulls-eye: True model
- Each hit is a realization of prediction model
- In general, we cannot ensure to reach low bias and low variance at the same time
- We have to balance bias and variance according to our objective
- Low bias and low variance require large, high quality datasets



Source: Scott Fortmann-Roe

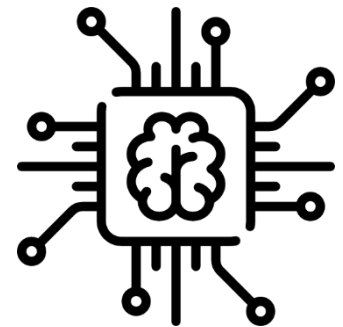
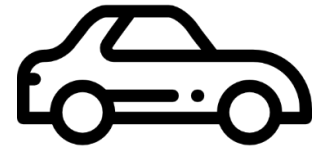
## Supervised Learning: Regression

Johannes Betz / Prof. Dr. Markus Lienkamp / Prof. Dr. Boris Lohmann  
(Alexander Wischnewski, M. Sc.)

### Agenda

---

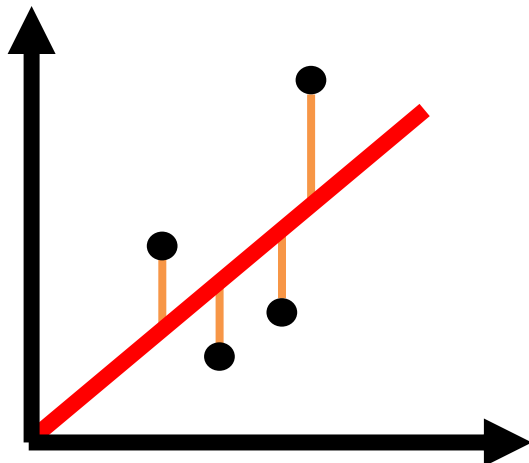
1. Chapter: Motivation
2. Chapter: Linear Models
3. Chapter: Loss Functions
4. Chapter: Regularization & Validation
- 5. Chapter: Practical Considerations**
6. Chapter: Summary



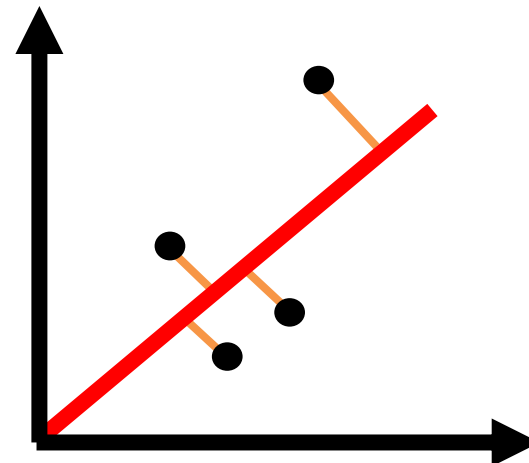
# Comparison Regression and Intuitive Solution

- Regression: Minimizing error in output domain
- Independent Variables are assumed as noise free
- Often not true in engineering applications

Standard linear regression



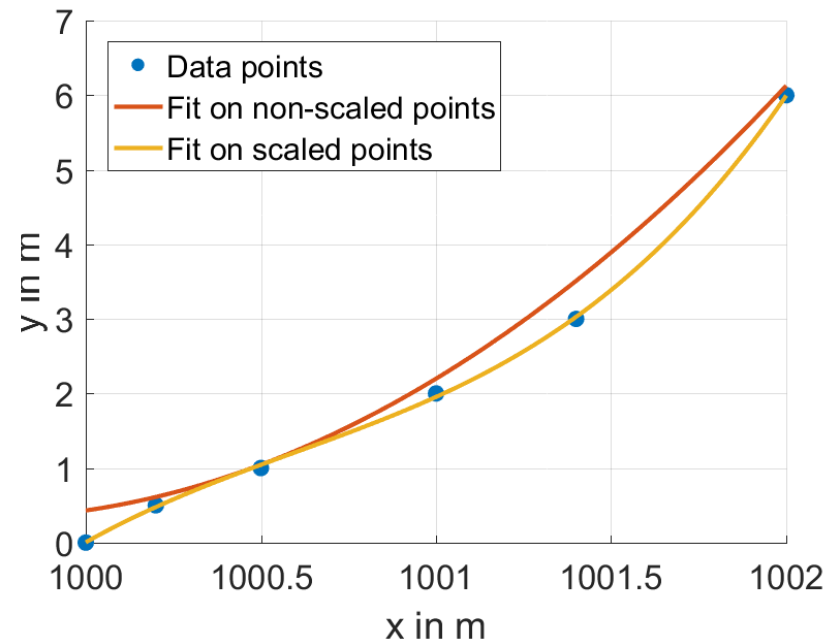
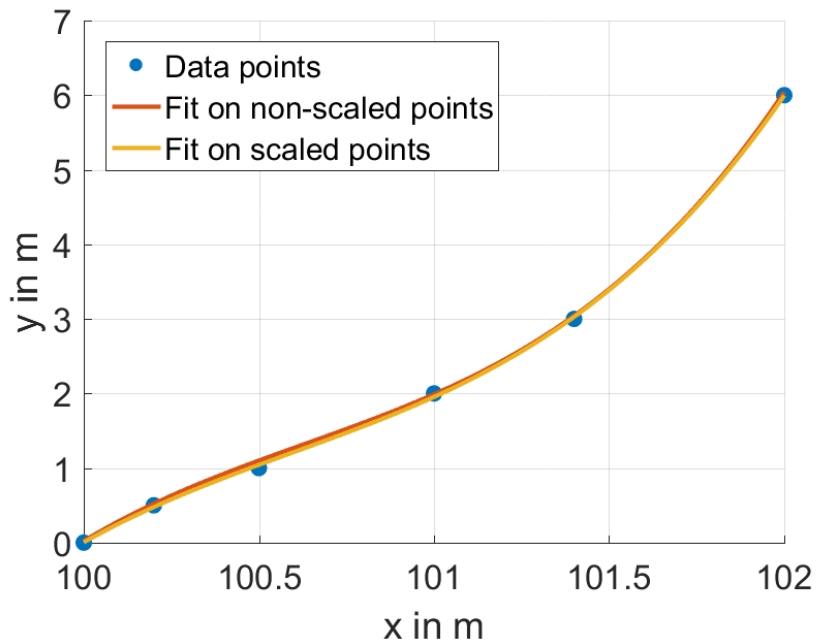
Total least squares  
Principal component analysis



## Normalization – Input Data

- Regression implies inversion of design matrix
- Numerically expensive and can get unstable
- In general, all matrix entries variables should lie within the same order of magnitude!
- Examples for scaling and mean value normalization
  - Multiple length measures should be transformed to the same unit
  - Variables which do not have any physical relation should be transformed based on their min/max values
  - Temperatures are likely to be far away from zero, they should be normalized by a mean temperature

# Normalization – Input Data



## Additional Information

The plots in the slide before were obtained by fitting a regression based on a simple polynomial basis to the data points. The basis functions in this example are

$$1 \quad x \quad x^2$$

This leads to the fact that the design matrix becomes ill-conditioned very fast for large input data. It results from the polynomials which lead to a design matrix (depicted for  $x = 1000$  and  $1001$ ) of

$$\begin{bmatrix} 1 & 1000 & 1000000 & 1000000000 \\ 1 & 1001 & 1002001 & 1003003001 \\ \dots & \dots & \dots & \dots \end{bmatrix}$$

Instead of (depicted for  $x = 0$  and  $x = 1$ )

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

The latter matrix is far easier to solve by numeric computations. The intuition behind this is that the "spread" between matrix entries is much smaller. This helps the computer to obtain a numerically stable result since precision is limited even in floating point operations. As a result from this, the maximum "spread" leading to reliable results is depending on the available computation precision. This is much more critical on embedded hardware.

## Normalization – Output Data

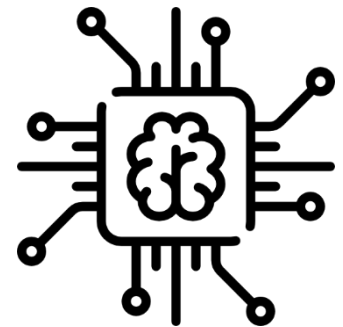
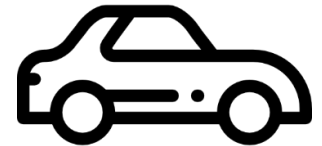
- In engineering applications we often face the task to fit parameters such that multiple objectives are met
- These may not lie within the same units or scale
- Example:
  - Estimate the driven curvature  $\kappa$  based on speed  $v$ , yaw rate  $\dot{\psi}$  and lateral acceleration  $a_y$
  - Model equations:  $\dot{\psi} = v\kappa$  and  $a_y = v^2\kappa$
  - Scale varies heavily with speed and is inconsistent between the two equations
  - Better use:  $\frac{\dot{\psi}}{v} = \kappa$  and  $\frac{a_y}{v^2} = \kappa$
- Scaling the output equations can also be used to weight equations which are more important than others

## Supervised Learning: Regression

Johannes Betz / Prof. Dr. Markus Lienkamp / Prof. Dr. Boris Lohmann  
(Alexander Wischnewski, M. Sc.)

### Agenda

1. Chapter: Motivation
2. Chapter: Linear Models
3. Chapter: Loss Functions
4. Chapter: Regularization & Validation
5. Chapter: Practical Considerations
6. **Chapter: Summary**





# Summary

- We learned about:
  - The characteristics of a **regression** problem and learned about the **differences** of **clustering** and **classification**
  - **Use cases** and **applications** for regression in automotive technology
  - Different applications for **local and global basis functions**
  - Different **loss functions**
  - **Iterative** and **analytic** solution methods for model training
  - **Regularization** techniques and **validation** techniques
  - Partition of datasets into training and validation sets

## Related Literature

- Christopher M. Bishop, *Machine Learning and Pattern Recognition*, Springer, 2006
- Trevor Hastie, Robert Tibshirani, Jerome Friedman, *The Elements of Statistical Learning*, Springer, 2009
- Kevin P. Murphy, *Machine Learning – A Probabilistic Perspective*, MIT Press 2012