

# Advanced Machine Learning Summer 2019

## Part 1 – Introduction 03.04.2019

Prof. Dr. Bastian Leibe

RWTH Aachen University, Computer Vision Group  
<http://www.vision.rwth-aachen.de>



## Organization

- Lecturer
  - Prof. Bastian Leibe ([leibe@vision.rwth-aachen.de](mailto:leibe@vision.rwth-aachen.de))
- Teaching Assistants
  - Jonathan Luiten ([luiten@vision.rwth-aachen.de](mailto:luiten@vision.rwth-aachen.de))
  - Ömer Sali ([sali@vision.rwth-aachen.de](mailto:sali@vision.rwth-aachen.de))
- Course webpage
  - <http://www.vision.rwth-aachen.de/courses/>
  - Slides will be made available on the webpage
  - There is also an electronic repository (moodle)
- Please subscribe to the lecture on RWTH Online!
  - Important to get email announcements and moodle access!



## Language

- Official course language will be English
  - If at least one English-speaking student is present.
  - If not... you can choose.
- However...
  - Please tell me when I'm talking too fast or when I should repeat something in German for better understanding!
  - You may at any time ask questions in German!
  - You may turn in your exercises in German.
  - You may take the oral exam in German.



## Relationship to Previous Courses

- Lecture *Machine Learning* (past winter semester)
  - Introduction to ML
  - Classification
  - Graphical models
- This course:
  - Natural continuation of ML course
  - Deeper look at the underlying concepts
  - But: will try to make it accessible also to newcomers
  - *Quick poll: Who hasn't heard the ML lecture?*
- This year: changed lecture content (compared to WS'16)
  - Large lecture block on Probabilistic Graphical Models
  - Updated with some exciting new topics (GANs, VAEs, Deep RL)



## Organization

- Structure: 3V (lecture) + 1Ü (exercises)
  - 6 EECS credits
  - Part of the area "Applied Computer Science"
- Place & Time
  - Lecture/Exercises: Wed 10:30 – 12:00 room H06
  - Lecture/Exercises: Thu 10:30 – 12:00 room H04
- Exam
  - Oral or written exam, depending on number of participants



## Course Webpage

Course Schedule			
Date	Title	Content	Material
Wed, 2019-04-03	Introduction	Introduction, Polynomial Fitting, Least-Squares Regression, Overfitting, Regularization, Ridge Regression	
Thu, 2019-04-04	Linear Regression I	Probabilistic View of Regression, Maximum Likelihood, MAP, Bayesian Curve Fitting	
Wed, 2019-04-10	Linear Regression II	Basis Functions, Sequential Learning, Multiple Outputs, Regularization, Lasso, Bias-Variance Decomposition	
Thu, 2019-04-11	Linear Regression III	Kernels, Kernel Ridge Regression	
Wed, 2019-04-17	Deep Reinforcement Learning I	Reinforcement Learning, TD Learning, Q-Learning, SARSA, Deep RL	
Thu, 2019-04-18	Deep Reinforcement Learning II	Deep RL, Deep Q-Learning, Deep Policy Gradients, Case studies	
Wed, 2019-04-24	Exercise 1	Regression, Least-Squares, Ridge, Kernel	

<http://www.vision.rwth-aachen.de/courses/>



## Exercises and Supplementary Material

### • Exercises

- Typically 1 exercise sheet every 2 weeks.
- Pen & paper and programming exercises
  - Matlab / numpy for early topics
  - Theano for Deep Learning topics
- Hands-on experience with the algorithms from the lecture.
- Send your solutions the night before the exercise class.

### • Supplementary material

- Research papers and book chapters
- Will be provided on the webpage.

7

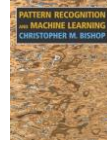
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

## Textbooks

- Many lecture topics will be covered in Bishop's book.
- Some additional topics can be found in Murphy's book



Christopher M. Bishop  
Pattern Recognition and Machine Learning  
Springer, 2006

(available in the library's "Handapparat")

Kevin P. Murphy  
Machine Learning – A Probabilistic Perspective  
MIT Press, 2012



- Research papers will be given out for some topics.
  - Tutorials and deeper introductions.
  - Application papers

8

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

## How to Find Us

### • Office:

- UMIC Research Centre
- Mies-van-der-Rohe-Strasse 15, room 124



### • Office hours

- If you have questions to the lecture, come see us.
- My regular office hours will be announced.
- Send us an email before to confirm a time slot.

*Questions are welcome!*

9

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

## Machine Learning

### • Statistical Machine Learning

- Principles, methods, and algorithms for learning and prediction on the basis of past evidence

### • Already everywhere

- Speech recognition (e.g. speed-dialing)
- Computer vision (e.g. face detection)
- Hand-written character recognition (e.g. letter delivery)
- Information retrieval (e.g. image & video indexing)
- Operation systems (e.g. caching)
- Fraud detection (e.g. credit cards)
- Text filtering (e.g. email spam filters)
- Game playing (e.g. strategy prediction)
- Robotics

10

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

Slide credit: Bert Scheile

## What Is Machine Learning Useful For?



Siri. beta  
Your wish is  
its command.



Automatic Speech Recognition

11

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

Slide adapted from Zoubin Ghahramani

## What Is Machine Learning Useful For?



Computer Vision  
(Object Recognition, Segmentation, Scene Understanding)

12

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

Slide adapted from Zoubin Ghahramani

## What Is Machine Learning Useful For?



Information Retrieval  
(Retrieval, Categorization, Clustering, ...)

13

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part I - Introduction

Slide adapted from Zoubin Ghahramani



## What Is Machine Learning Useful For?



Financial Prediction  
(Time series analysis, ...)

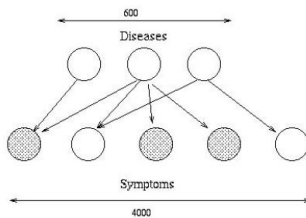
14

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part I - Introduction

Slide adapted from Zoubin Ghahramani



## What Is Machine Learning Useful For?



Medical Diagnosis  
(Inference from partial observations)

15

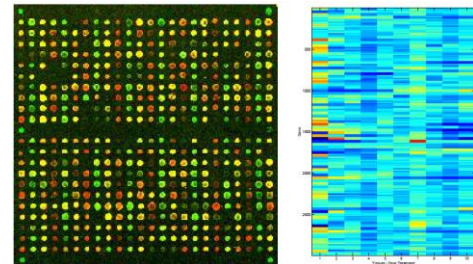
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part I - Introduction

Slide adapted from Zoubin Ghahramani



Image from Kevin Murphy

## What Is Machine Learning Useful For?



Bioinformatics  
(Modelling gene microarray data,...)

16

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part I - Introduction

Slide adapted from Zoubin Ghahramani



## What Is Machine Learning Useful For?



Robotics & Autonomous Driving

17

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part I - Introduction

Slide adapted from Zoubin Ghahramani



## Machine Learning: Core Questions

- *Learning to perform a task from experience*
- Task
  - Can often be expressed through a mathematical function
$$y = f(x; w)$$
  - $x$ : Input
  - $y$ : Output
  - $w$ : Parameters (this is what is "learned")
- Classification vs. Regression
  - Regression: continuous  $y$
  - Classification: discrete  $y$ 
    - E.g. class membership, sometimes also posterior probability

18

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part I - Introduction

Slide credit: Bernt Schiele

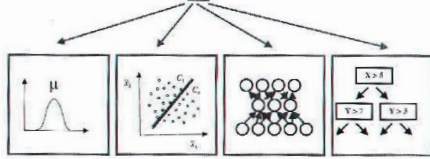


## Machine Learning: Core Questions

•  $y = f(x; w)$

- $w$ : characterizes the family of functions
- $w$ : indexes the space of hypotheses
- $w$ : vector, connection matrix, graph, ...

• Look inside box:



19

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction

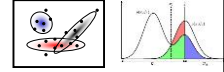
Slide credit: Bernt Schiele



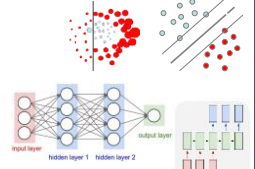
RWTH AACHEN  
UNIVERSITY

## A Look Back: Lecture *Machine Learning*

- Fundamentals
  - Bayes Decision Theory
  - Probability Density Estimation
  - Mixture Models and EM
- Classification Approaches
  - Linear Discriminant Functions
  - Support Vector Machines
  - Ensemble Methods & Boosting
- Deep Learning
  - Foundations
  - Convolutional Neural Networks
  - Recurrent Neural Networks



$f: \mathcal{X} \rightarrow \{0, 1\}$



20

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction

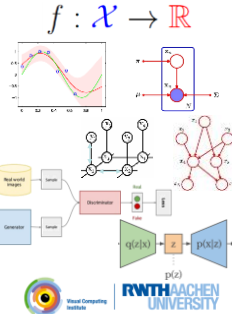


RWTH AACHEN  
UNIVERSITY

## This Lecture: *Advanced Machine Learning*

Extending lecture *Machine Learning* from last semester...

- Regression Techniques
  - Regularization (Ridge, Lasso)
  - Bayesian Regression
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders
- Deep Reinforcement Learning



21

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction



RWTH AACHEN  
UNIVERSITY

## Let's Get Started...

- Some of you already have basic ML background
  - *Who hasn't?*
- We'll start with a gentle introduction
  - I'll try to make the lecture accessible to everyone.
  - We'll review the main concepts before applying them.
  - I'll point out chapters to review from the ML lecture whenever knowledge from there is needed/helpful.
- But...
  - This *is* an advanced topics class.
  - There *will* be math involved.
  - We will take a deeper look into the theory than in the ML lecture.

22

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction



RWTH AACHEN  
UNIVERSITY

## Topics of This Lecture

- **Regression: Motivation**
  - Polynomial fitting
  - General Least-Squares Regression
  - Overfitting problem
  - Regularization
  - Ridge Regression
- **Recap: Important Concepts from ML Lecture**
  - Probability Theory
  - Bayes Decision Theory
  - Maximum Likelihood Estimation
  - Bayesian Estimation
- **A Probabilistic View on Regression**
  - Least-Squares Estimation as Maximum Likelihood

23

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction



RWTH AACHEN  
UNIVERSITY

## Regression

- Learning to predict a continuous function value
  - Given: training set  $\mathbf{X} = \{x_1, \dots, x_N\}$  with target values  $\mathbf{T} = \{t_1, \dots, t_N\}$ .
  - $\Rightarrow$  Learn a continuous function  $y(x)$  to predict the function value for a new input  $x$ .
- Steps towards a solution
  - Choose a form of the function  $y(x, w)$  with parameters  $w$ .
  - Define an error function  $E(w)$  to optimize.
  - Optimize  $E(w)$  for  $w$  to find a good solution. (This may involve math).
  - Derive the properties of this solution and think about its limitations.

24

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction



RWTH AACHEN  
UNIVERSITY

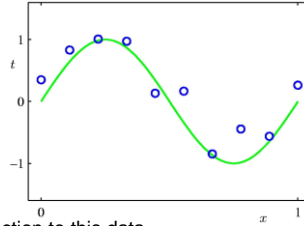
## Example: Polynomial Curve Fitting

- Toy dataset

- Generated by function

$$f(x) = \sin(2\pi x) + \epsilon$$

- Small level of random noise with Gaussian distribution added (blue dots)



- Goal: fit a polynomial function to this data

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

- Note: Nonlinear function of  $x$ , but linear function of the  $w_j$ .

25

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction



Image source: C.M. Bishop, 2006

## Error Function

- How to determine the values of the coefficients  $\mathbf{w}$ ?

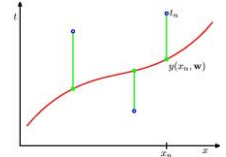
- We need to define an **error function** to be minimized.
- This function specifies how a deviation from the target value should be weighted.

- Popular choice: sum-of-squares error

- Definition

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

- We'll discuss the motivation for this particular function later...



26

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction



Image source: C.M. Bishop, 2006

## Minimizing the Error

- How do we minimize the error?

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

- Solution (Always!)

- Compute the derivative and set it to zero.

$$\frac{\partial E(\mathbf{w})}{\partial w_j} = \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\} \frac{\partial y(x_n, \mathbf{w})}{\partial w_j} \stackrel{!}{=} 0$$

- Since the error is a quadratic function of  $\mathbf{w}$ , its derivative will be linear in  $\mathbf{w}$ .

⇒ Minimization has a unique solution.

27

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction



## Least-Squares Regression

- We have given

- Training data points:  $X = \{\mathbf{x}_1 \in \mathbb{R}^d, \dots, \mathbf{x}_n\}$
- Associated function values:  $T = \{t_1 \in \mathbb{R}, \dots, t_n\}$

- Start with **linear regressor**:

- Try to enforce  $\mathbf{x}_i^T \mathbf{w} + w_0 = t_i, \quad \forall i = 1, \dots, n$
- One linear equation for each training data point / label pair.

- This is the same basic setup used for least-squares classification!
  - Only the values are now continuous.

28

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction



Slide credit: Bernt Schiele

## Least-Squares Regression

$$\mathbf{x}_i^T \mathbf{w} + w_0 = t_i, \quad \forall i = 1, \dots, n$$

- Setup

- Step 1: Define  $\tilde{\mathbf{x}}_i = \begin{pmatrix} \mathbf{x}_i \\ 1 \end{pmatrix}, \quad \tilde{\mathbf{w}} = \begin{pmatrix} \mathbf{w} \\ w_0 \end{pmatrix}$

- Step 2: Rewrite  $\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} = t_i, \quad \forall i = 1, \dots, n$

- Step 3: Matrix-vector notation

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} = \mathbf{t} \quad \text{with} \quad \tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n] \\ \mathbf{t} = [t_1, \dots, t_n]^T$$

- Step 4: Find least-squares solution

$$\|\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{t}\|^2 \rightarrow \min$$

- Solution:  $\tilde{\mathbf{w}} = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1} \tilde{\mathbf{X}}\mathbf{t}$

29

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction



Slide credit: Bernt Schiele

## Regression with Polynomials

- How can we fit arbitrary polynomials using least-squares regression?

- We introduce a feature transformation (as before in ML).

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) \\ = \sum_{i=0}^M w_i \phi_i(\mathbf{x})$$

assume  $\phi_0(\mathbf{x}) = 1$

basis functions

- E.g.:  $\phi(\mathbf{x}) = (1, x, x^2, x^3)^T$
- Fitting a cubic polynomial.

30

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction



Slide credit: Bernt Schiele

### Varying the Order of the Polynomial.

Which one should we pick?

Massive overfitting!

31 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
RWTH AACHEN UNIVERSITY  
Image source: C.M. Bishop, 2006

### Analysis of the Results

- Results for different values of  $M$ 
  - Best representation of the original function  $\sin(2\pi x)$  with  $M = 3$ .
  - Perfect fit to the training data with  $M = 9$ , but poor representation of the original function.
- Why is that???
  - After all,  $M = 9$  contains  $M = 3$  as a special case!

32 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
RWTH AACHEN UNIVERSITY  
Image source: C.M. Bishop, 2006

### Overfitting

- Problem
  - Training data contains some noise
$$f(x) = \sin(2\pi x) + \epsilon$$
  - Higher-order polynomial fitted perfectly to the noise.
  - We say it was **overfitting to the training data**.
- Goal is a good prediction of future data
  - Our target function should fit well to the training data, but also generalize.
  - Measure generalization performance on independent test set.

33 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
RWTH AACHEN UNIVERSITY  
Image source: C.M. Bishop, 2006

### Measuring Generalization

Overfitting!

- E.g., Root Mean Square Error (RMS):  $E_{RMS} = \sqrt{2E(\mathbf{w}^*)/N}$
- Motivation
  - Division by  $N$  lets us compare different data set sizes.
  - Square root ensures  $E_{RMS}$  is measured on the same scale (and in the same units) as the target variable  $t$ .

34 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
RWTH AACHEN UNIVERSITY  
Image source: C.M. Bishop, 2006

### Analyzing Overfitting

- Example: Polynomial of degree 9
  - Relatively little data: Overfitting typical (N=15)
  - Enough data: Good estimate (N=100)

⇒ Overfitting becomes less of a problem with more data.

35 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
Slide adapted from Bernt Schiele  
RWTH AACHEN UNIVERSITY  
Image source: C.M. Bishop, 2006

### What Is Happening Here?

- The coefficients get very large:
  - Fitting the data from before with various polynomials.
- Coefficients:
 

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	7.99	232.37
$w_2^*$			-25.43	-5321.83
$w_3^*$			17.37	48568.31
$w_4^*$				-231639.30
$w_5^*$				640042.26
$w_6^*$				-1061800.52
$w_7^*$				1042400.18
$w_8^*$				-557682.99
$w_9^*$				125201.43

36 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
Slide credit: Bernt Schiele  
RWTH AACHEN UNIVERSITY  
Image source: C.M. Bishop, 2006

## Regularization

- What can we do then?
  - How can we apply the approach to data sets of limited size?
  - We still want to use relatively complex and flexible models.

### Workaround: Regularization

- Penalize large coefficient values

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- Here we've simply added a **quadratic regularizer**, which is simple to optimize

$$\|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + \dots + w_M^2$$

- The resulting form of the problem is called **Ridge Regression**.
- (Note:  $w_0$  is often omitted from the regularizer.)

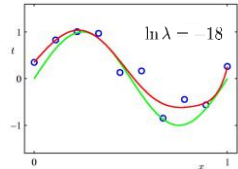
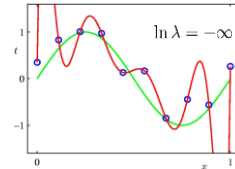
37

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction

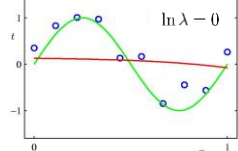


RWTH AACHEN  
UNIVERSITY

## Results with Regularization (M=9)



	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^2$	0.35	0.35	0.13
$w_1^2$	232.37	4.74	-0.05
$w_2^2$	-5321.83	-0.77	-0.06
$w_3^2$	48568.31	-31.97	-0.05
$w_4^2$	-231639.30	-3.89	-0.03
$w_5^2$	640642.26	55.28	-0.02
$w_6^2$	-1061800.52	41.32	-0.01
$w_7^2$	1042400.18	-45.95	-0.00
$w_8^2$	-557682.99	-91.53	0.00
$w_9^2$	125201.43	72.68	0.01



38

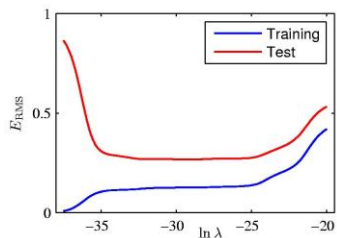
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

Image source: C.M. Bishop, 2006

## RMS Error for Regularized Case



- Effect of regularization
  - The trade-off parameter  $\lambda$  now controls the effective model complexity and thus the degree of overfitting.

39

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

Image source: C.M. Bishop, 2006

## Summary

- We've seen several important concepts
  - Linear regression
  - Overfitting
  - Role of the amount of data
  - Role of model complexity
  - Regularization
- How can we approach this more systematically?
  - Would like to work with complex models.
  - How can we prevent overfitting systematically?
  - How can we avoid the need for validation on separate test data?
  - What does it *mean* to do linear regression?
  - What does it *mean* to do regularization?

40

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

## Topics of This Lecture

- Regression: Motivation
  - Polynomial fitting
  - General Least-Squares Regression
  - Overfitting problem
  - Regularization
  - Ridge Regression
- Recap: Important Concepts from ML Lecture
  - Probability Theory
  - Bayes Decision Theory
  - Maximum Likelihood Estimation
  - New: Bayesian Estimation
- A Probabilistic View on Regression
  - Least-Squares Estimation as Maximum Likelihood

41

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

## Recap: The Rules of Probability

- Basic rules

$$\text{Sum Rule} \quad p(X) = \sum_Y p(X, Y)$$

$$\text{Product Rule} \quad p(X, Y) = p(Y|X)p(X)$$

- From those, we can derive

$$\text{Bayes' Theorem} \quad p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

$$\text{where} \quad p(X) = \sum_Y p(X|Y)p(Y)$$

42

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

### Recap: Bayes Decision Theory

$p(x|a)$   $p(x|b)$  Likelihood

$p(x|a)p(a)$   $p(x|b)p(b)$  Likelihood  $\times$  Prior

Decision boundary

$p(a|x)$   $p(b|x)$  Posterior =  $\frac{\text{Likelihood} \times \text{Prior}}{\text{NormalizationFactor}}$

43 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
Slide credit: Bernt Schiele

### Recap: Gaussian (or Normal) Distribution

- One-dimensional case
  - Mean  $\mu$
  - Variance  $\sigma^2$
- Multi-dimensional case
  - Mean  $\mu$
  - Covariance  $\Sigma$

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2}|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$

44 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
Image source: C.M. Bishop, 2006

### Side Note

- Notation
  - In many situations, it will be necessary to work with the inverse of the covariance matrix  $\Sigma$ :
 
$$\Lambda = \Sigma^{-1}$$
  - We call  $\Lambda$  the precision matrix.
  - We can therefore also write the Gaussian as
 
$$\mathcal{N}(x|\mu, \Lambda^{-1}) = \frac{1}{\sqrt{2\pi}\lambda^{-1/2}} \exp\left\{-\frac{\lambda}{2}(x-\mu)^2\right\}$$
  - $$\mathcal{N}(x|\mu, \Lambda^{-1}) = \frac{1}{(2\pi)^{D/2}|\Lambda|^{-1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Lambda (x-\mu)\right\}$$

45 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction

### Recap: Parametric Methods

- Given
  - Data  $X = \{x_1, x_2, \dots, x_N\}$
  - Parametric form of the distribution with parameters  $\theta$
  - E.g. for Gaussian distrib.:  $\theta = (\mu, \sigma)$
- Learning
  - Estimation of the parameters  $\theta$
- Likelihood of  $\theta$ 
  - Probability that the data  $X$  have indeed been generated from a probability density with parameters  $\theta$
$$L(\theta) = p(X|\theta)$$

46 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
Slide adapted from Bernt Schiele

### Recap: Maximum Likelihood Approach

- Computation of the likelihood
  - Single data point:  $p(x_n|\theta)$
  - Assumption: all data points  $X = \{x_1, \dots, x_n\}$  are independent
 
$$L(\theta) = p(X|\theta) = \prod_{n=1}^N p(x_n|\theta)$$
  - Log-likelihood\*
 
$$E(\theta) = -\ln L(\theta) = -\sum_{n=1}^N \ln p(x_n|\theta)$$
- Estimation of the parameters  $\theta$  (Learning)
  - Maximize the likelihood (=minimize the negative log-likelihood)
  - Take the derivative and set it to zero.
 
$$\frac{\partial}{\partial \theta} E(\theta) = -\sum_{n=1}^N \frac{\partial}{\partial \theta} \ln p(x_n|\theta) \stackrel{!}{=} 0$$

47 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
Slide credit: Bernt Schiele

### Recap: Maximum Likelihood Approach

- Maximum Likelihood has several significant limitations
  - It systematically underestimates the variance of the distribution!
  - E.g. consider the case
 
$$N = 1, X = \{x_1\}$$

$\Rightarrow$  Maximum-likelihood estimate:

- We say ML *overfits to the observed data*.
- We will still often use ML, but it is important to know about this effect.

48 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 - Introduction  
Slide adapted from Bernt Schiele



## Recap: Deeper Reason

- Maximum Likelihood is a **Frequentist** concept
  - In the **Frequentist view**, probabilities are the frequencies of random, repeatable events.
  - These frequencies are fixed, but can be estimated more precisely when more data is available.
- This is in contrast to the **Bayesian** interpretation
  - In the **Bayesian view**, probabilities quantify the uncertainty about certain states or events.
  - This uncertainty can be revised in the light of new evidence.
- Bayesians and Frequentists do not like each other too well...

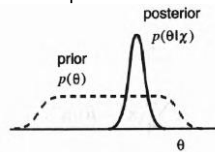


## Bayesian vs. Frequentist View

- To see the difference...
    - Suppose we want to estimate the uncertainty whether the Arctic ice cap will have disappeared by the end of the century.
    - This question makes no sense in a Frequentist view, since the event cannot be repeated numerous times.
    - In the Bayesian view, we generally have a prior, e.g. from calculations how fast the polar ice is melting.
    - If we now get fresh evidence, e.g. from a new satellite, we may revise our opinion and update the uncertainty from the prior.
- $$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$
- This generally allows to get better uncertainty estimates for many situations.
  - Main Frequentist criticism
    - The prior has to come from somewhere and if it is wrong, the result will be worse.

## Bayesian Approach to Parameter Learning

- Conceptual shift
  - Maximum Likelihood views the true parameter vector  $\theta$  to be unknown, but fixed.
  - In Bayesian learning, we consider  $\theta$  to be a random variable.
- This allows us to use knowledge about the parameters  $\theta$ 
  - i.e. to use a prior for  $\theta$
  - Training data then converts this prior distribution on  $\theta$  into a posterior probability density.
- The prior thus encodes knowledge we have about the type of distribution we expect to see for  $\theta$ .



## Bayesian Learning Approach

- Bayesian view:
  - Consider the parameter vector  $\theta$  as a random variable.
  - When estimating the parameters, what we compute is

$$p(x|X) = \int p(x, \theta|X) d\theta$$

Assumption: given  $\theta$ , this doesn't depend on  $X$  anymore

$$p(x, \theta|X) = p(x|\theta, X)p(\theta|X)$$

$$p(x|X) = \int p(x|\theta)p(\theta|X) d\theta$$

This is entirely determined by the parameter  $\theta$  (i.e. by the parametric form of the pdf).

## Bayesian Learning Approach

$$p(x|X) = \int p(x|\theta)p(\theta|X) d\theta$$

$$p(\theta|X) = \frac{p(X|\theta)p(\theta)}{p(X)} = \frac{p(\theta)}{p(X)} L(\theta)$$

$$p(X) = \int p(X|\theta)p(\theta) d\theta = \int L(\theta)p(\theta) d\theta$$

- Inserting this above, we obtain

$$p(x|X) = \int \frac{p(x|\theta)L(\theta)p(\theta)}{\int L(\theta)p(\theta) d\theta} d\theta$$

## Bayesian Learning Approach

- Discussion
  - Likelihood of the parametric form  $\theta$  given the data set  $X$ .

Estimate for  $x$  based on parametric form  $\theta$       Prior for the parameters  $\theta$

$$p(x|X) = \int \frac{p(x|\theta)L(\theta)p(\theta)}{\int L(\theta)p(\theta) d\theta} d\theta$$

Normalization: integrate over all possible values of  $\theta$

- The more uncertain we are about  $\theta$ , the more we average over all possible parameter values.

## Bayesian Density Estimation

- Discussion

$$p(x|X) = \int p(x|\theta)p(\theta|X)d\theta = \int \frac{p(x|\theta)L(\theta)p(\theta)}{\int L(\theta)p(\theta)d\theta}d\theta$$

- The probability  $p(\theta|X)$  makes the dependency of the estimate on the data explicit.

- If  $p(\theta|X)$  is very small everywhere, but is large for one  $\hat{\theta}$ , then

$$p(x|X) \approx p(x|\hat{\theta})$$

- ⇒ The more uncertain we are about  $\theta$ , the more we average over all parameter values.

- Problem

- In the general case, the integration over  $\theta$  is not possible (or only possible stochastically).

55

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

Slide credit: Bernt Schiele

## Topics of This Lecture

- Regression: Motivation

- Polynomial fitting
- General Least-Squares Regression
- Overfitting problem
- Regularization
- Ridge Regression

- Recap: Important Concepts from ML Lecture

- Probability Theory
- Bayes Decision Theory
- Maximum Likelihood Estimation
- Bayesian Estimation

- **A Probabilistic View on Regression**

- Least-Squares Estimation as Maximum Likelihood

56

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

Next lecture...

57

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY

## References and Further Reading

- More information, including a short review of Probability theory and a good introduction in Bayes Decision Theory can be found in Chapters 1.1, 1.2 and 1.5 of

Christopher M. Bishop  
Pattern Recognition and Machine Learning  
Springer, 2006



66

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 1 – Introduction



RWTH AACHEN  
UNIVERSITY