

RWTH AACHEN UNIVERSITY

Advanced Machine Learning Lecture 5

Gaussian Processes 2

09.11.2015

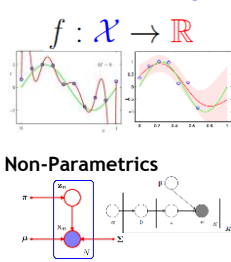
Bastian Leibe
RWTH Aachen
<http://www.vision.rwth-aachen.de/>
leibe@vision.rwth-aachen.de

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

This Lecture: Advanced Machine Learning

- Regression Approaches
 - Linear Regression
 - Regularization (Ridge, Lasso)
 - Kernels (Kernel Ridge Regression)
 - Gaussian Processes
- Bayesian Estimation & Bayesian Non-Parametrics
 - Mixture Models & EM
 - Dirichlet Processes
 - Latent Factor Models
 - Beta Processes
- SVMs and Structured Output Learning
 - SV Regression, SVDD
 - Large-margin Learning



$f : \mathcal{X} \rightarrow \mathbb{R}$

$f : \mathcal{X} \rightarrow \mathcal{Y}$

B. Leibe

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Kernels
 - Recap: Kernel trick
 - Constructing kernels
- Gaussian Processes
 - Recap: Definition
 - Prediction with noise-free observations
 - Prediction with noisy observations
 - GP Regression
 - Influence of hyperparameters
- Learning Gaussian Processes
 - Bayesian Model Selection
 - Model selection for Gaussian Processes
- Applications

3

B. Leibe

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

Recap: Kernel Ridge Regression

- Dual definition
 - Instead of working with w , substitute $w = \Phi^T a$ into $J(w)$ and write the result using the **kernel matrix** $K = \Phi \Phi^T$:

$$J(a) = \frac{1}{2} a^T K K a - a^T K t + \frac{1}{2} t^T t + \frac{\lambda}{2} a^T K a$$
 - Solving for a , we obtain

$$a = (K + \lambda I_N)^{-1} t$$
- Prediction for a new input x :
 - Writing $k(x)$ for the vector with elements $k_n(x) = k(x_n, x)$

$$y(x) = w^T \phi(x) = a^T \Phi \phi(x) = k(x)^T (K + \lambda I_N)^{-1} t$$

\Rightarrow The dual formulation allows the solution to be entirely expressed in terms of the kernel function $k(x, x')$.

4

B. Leibe

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

Recap: Properties of Kernels

- Theorem
 - Let $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a **positive definite kernel function**. Then there exists a **Hilbert Space** \mathcal{H} and a mapping $\phi: \mathcal{X} \rightarrow \mathcal{H}$ such that

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$$
 - where $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is the inner product in \mathcal{H} .
- Translation
 - Take **any** set \mathcal{X} and **any** function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$.
 - If k is a positive definite kernel, then we can use k to learn a classifier for the elements in \mathcal{X} !
- Note
 - \mathcal{X} can be any set, e.g. $\mathcal{X} =$ "all videos on YouTube" or $\mathcal{X} =$ "all permutations of $\{1, \dots, k\}$ ", or $\mathcal{X} =$ "the internet".

5

Slide credit: Christoph Lampert

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

Recap: The "Kernel Trick"

Any algorithm that uses data only in the form of inner products can be **kernelized**.

- How to kernelize an algorithm
 - Write the algorithm only in terms of inner products.
 - Replace all inner products by kernel function evaluations.

\Rightarrow The resulting algorithm will do the same as the linear version, but in the (hidden) feature space \mathcal{H} .

- Caveat: working in \mathcal{H} is not a guarantee for better performance. A good choice of k and model selection are important!

6

Slide credit: Christoph Lampert

B. Leibe

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

How to Check if a Function is a Kernel

- Problem:**
 - Checking if a given $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ fulfills the conditions for a kernel is difficult:
 - We need to prove or disprove

$$\sum_{i,j=1}^n t_i k(x_i, x_j) t_j \geq 0$$
 for any set $x_1, \dots, x_n \in \mathcal{X}$ and any $t \in \mathbb{R}^n$ for any $n \in \mathbb{N}$.
- Workaround:**
 - It is easy to construct functions k that are positive definite kernels.

Advanced Machine Learning Winter'12 | Slide credit: Christoph Lampert | B. Leibe | 7

RWTH AACHEN UNIVERSITY

Constructing Kernels

- We can construct kernels from scratch:
 - For any $\varphi : \mathcal{X} \rightarrow \mathbb{R}^m$, $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathbb{R}^m}$ is a kernel.
Example: $\varphi(x) = (\text{"# of red pixels in image } x", \text{green, blue})$.
 - Any norm $\|\cdot\| : V \rightarrow \mathbb{R}^m$ that fulfills the **parallelogram equation**

$$\|x + y\|^2 + \|x - y\|^2 = 2\|x\|^2 + 2\|y\|^2$$
 induces a kernel by **polarization**:

$$k(x, y) := (\|x + y\|^2 + \|x\|^2 - \|y\|^2)$$
 Example: $\mathcal{X} = \text{time series with bounded values}$, $\|x\|^2 = \sum_{t=1}^{\infty} \frac{1}{2^t} x_t$

Advanced Machine Learning Winter'12 | Slide credit: Christoph Lampert | B. Leibe | 8

RWTH AACHEN UNIVERSITY

Constructing Kernels (2)

- We can construct kernels from scratch:
 - If $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is **conditionally positive definite**, i.e.

$$\sum_{i,j=1}^n t_i d(x_i, x_j) t_j \geq 0 \text{ for any } t \in \mathbb{R}^n \text{ with } \sum_i t_i = 0,$$
 for $x_1, \dots, x_n \in \mathcal{X}$ for any $n \in \mathbb{N}$, then

$$k(x, x') := \exp(-d(x, x'))$$
 is a positive kernel.
Example: $d(x, x') = \|x - x'\|^2$.

$$k(x, x') = \exp\{-\|x - x'\|_{L_2}^2\}$$

Advanced Machine Learning Winter'12 | Slide credit: Christoph Lampert | B. Leibe | 9

RWTH AACHEN UNIVERSITY

Constructing Kernels (3)

- We can construct kernels from other kernels:
 - If k is a kernel and $\alpha > 0$, then αk and $k + \alpha$ are kernels.
 - if k_1, k_2 are kernels, then $k_1 + k_2$ and $k_1 \cdot k_2$ are kernels.
 - if k is a kernel, then $\exp(k)$ is a kernel.
- Examples for kernels for $\mathcal{X} = \mathbb{R}^d$:
 - Any linear combination $\sum_j \alpha_j k_j$ with $\alpha_j \geq 0$,
 - Polynomial kernels** $k(x, x') = (1 + \langle x, x' \rangle)^m$, $m > 0$
 - Gaussian** a.k.a. RBF

$$k(x, x') = \exp\left\{-\frac{\|x - x'\|^2}{2\sigma^2}\right\}$$
 with $\sigma > 0$.

Advanced Machine Learning Winter'12 | Slide credit: Christoph Lampert | B. Leibe | 10

RWTH AACHEN UNIVERSITY

Constructing Kernels (4)

- We can construct kernels from other kernels:
 - If k is a kernel and $\alpha > 0$, then αk and $k + \alpha$ are kernels.
 - if k_1, k_2 are kernels, then $k_1 + k_2$ and $k_1 \cdot k_2$ are kernels.
 - if k is a kernel, then $\exp(k)$ is a kernel.
- Examples for kernels for other \mathcal{X} :
 - $k(h, h') = \sum_{i=1}^n \min(h_i, h'_i)$ for n -bin histograms h, h' .
 - $k(p, p') = \exp(-\text{KL}(p, p'))$ with **KL-divergence** between positive probability distributions.
 - $k(s, s') = \exp(-D(s, s'))$ for strings s, s' and $D = \text{edit distance}$
- Not an example:** $\tanh(a \langle x, x' \rangle + b)$ is **not positive definite!**

Advanced Machine Learning Winter'12 | Slide credit: Christoph Lampert | B. Leibe | 11

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Kernels**
 - Recap: Kernel trick
 - Constructing kernels
- Gaussian Processes**
 - Recap: Definition
 - Prediction with noise-free observations
 - Prediction with noisy observations
 - GP Regression
 - Influence of hyperparameters
- Learning Gaussian Processes**
 - Bayesian Model Selection
 - Model selection for Gaussian Processes
- Applications**

Advanced Machine Learning Winter'12 | B. Leibe | 13

RWTH AACHEN UNIVERSITY

Recap: Gaussian Process

- Gaussian distribution
 - Probability distribution over scalars / vectors.
- Gaussian Process (generalization of Gaussian distrib.)
 - Describes properties of functions.
 - Function: Think of a function as a long vector where each entry specifies the function value $f(x_i)$ at a particular point x_i .
 - Issue: How to deal with infinite number of points?
 - If you ask only for properties of the function at a finite number of points...
 - Then inference in Gaussian Process gives you the same answer if you ignore the infinitely many other points.
- Definition
 - A Gaussian Process (GP) is a collection of random variables any finite number of which has a joint Gaussian distribution.

Slide credit: Bernt Schiele B. Leibe 14

RWTH AACHEN UNIVERSITY

Recap: Gaussian Process

- A Gaussian Process is completely defined by
 - Mean function $m(x)$ and

$$m(x) = \mathbb{E}[f(x)]$$
 - Covariance function $k(x, x')$

$$k(x, x') = \mathbb{E}[(f(x) - m(x))(f(x') - m(x')))]$$
 - We write the Gaussian Process (GP)

$$f(x) \sim \mathcal{GP}(m(x), k(x, x'))$$

Slide adapted from Bernt Schiele B. Leibe 15

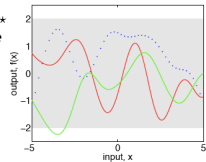
RWTH AACHEN UNIVERSITY

Recap: GPs Define Prior over Functions

- Distribution over functions:
 - Specification of covariance function implies distribution over functions.
 - I.e. we can draw samples from the distribution of functions evaluated at a (finite) number of points.
 - Procedure
 - We choose a number of input points X_*
 - We write the corresponding covariance matrix (e.g. using SE) element-wise:

$$K(X_*, X_*)$$
 - Then we generate a random Gaussian vector with this covariance matrix:

$$f_* \sim \mathcal{N}(0, K(X_*, X_*))$$



Slide credit: Bernt Schiele B. Leibe Image source: Rasmussen & Williams, 2006 16

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Kernels
 - Recap: Kernel trick
 - Constructing kernels
- Gaussian Processes
 - Recap: Definition
 - Prediction with noise-free observations
 - Prediction with noisy observations
 - GP Regression
 - Influence of hyperparameters
- Learning Gaussian Processes
 - Bayesian Model Selection
 - Model selection for Gaussian Processes
- Applications

Slide credit: Bernt Schiele B. Leibe 17

RWTH AACHEN UNIVERSITY

Prediction with Noise-free Observations

- Assume our observations are noise-free:

$$\{(x_n, f_n) \mid n = 1, \dots, N\}$$
- Joint distribution of the training outputs f and test outputs f_* according to the prior:

$$\begin{bmatrix} f \\ f_* \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right)$$
 - $K(X, X_*)$ contains covariances for all pairs of training and test points.
- To get the posterior (after including the observations)
 - We need to restrict the above prior to contain only those functions which agree with the observed values.
 - Think of generating functions from the prior and rejecting those that disagree with the observations (obviously prohibitive).

Slide credit: Bernt Schiele B. Leibe 19

RWTH AACHEN UNIVERSITY

Prediction with Noise-free Observations

- Calculation of posterior: simple in GP framework
 - Corresponds to conditioning the joint Gaussian prior distribution on the observations:

$$f_* | X_*, X, f \sim \mathcal{N}(\bar{f}_*, \text{cov}[f_*]) \quad \bar{f}_* = \mathbb{E}[f_* | X, X_*, f]$$
 - with:

$$\bar{f}_* = K(X_*, X)K(X, X)^{-1}f$$

$$\text{cov}[f_*] = K(X_*, X_*) - K(X_*, X)K(X, X)^{-1}K(X, X_*)$$
 - This uses the general property of Gaussians that

$$\mu = \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \Sigma = \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix} \Rightarrow \begin{matrix} \mu_{a|b} = \mu_a + \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b) \\ \Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba} \end{matrix}$$

Slide credit: Bernt Schiele B. Leibe 20

RWTH AACHEN UNIVERSITY

Prediction with Noise-free Observations

- Example:

Prior

Posterior using 5 noise-free observations

Slide credit: Bernt Schiele B. Leibe Image source: Rasmussen & Williams, 2006 21

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- **Kernels**
 - Recap: Kernel trick
 - Constructing kernels
- **Gaussian Processes**
 - Recap: Definition
 - Prediction with noise-free observations
 - Prediction with noisy observations
 - GP Regression
 - Influence of hyperparameters
- **Learning Gaussian Processes**
 - Bayesian Model Selection
 - Model selection for Gaussian Processes
- **Applications**

Slide credit: B. Leibe 22

RWTH AACHEN UNIVERSITY

Prediction with Noisy Observations

- Typically, we assume noise in the observations

$$t = f(x) + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
- The prior on the noisy observations becomes

$$\text{cov}[y_p, y_q] = k(x_p, x_q) + \sigma_n^2 \delta_{pq}$$
 - Written in compact form:

$$\text{cov}[y] = K(X, X) + \sigma_n^2 I$$
- Joint distribution of the observed values and the test locations under the prior is then:

$$\begin{bmatrix} t \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right)$$

Slide credit: Bernt Schiele B. Leibe 24

RWTH AACHEN UNIVERSITY

Prediction with Noisy Observations

- **Calculation of posterior:**
 - Corresponds to **conditioning** the joint Gaussian prior distribution on the observations:

$$\mathbf{f}_* | X_*, X, t \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}[\mathbf{f}_*]) \quad \bar{\mathbf{f}}_* = \mathbb{E}[\mathbf{f}_* | X, X_*, t]$$
 - with:

$$\bar{\mathbf{f}}_* = K(X_*, X) (K(X, X) + \sigma_n^2 I)^{-1} t$$

$$\text{cov}[\mathbf{f}_*] = K(X_*, X_*) - K(X_*, X) (K(X, X) + \sigma_n^2 I)^{-1} K(X, X_*)$$
 - ⇒ **This is the key result that defines Gaussian process regression!**
 - The predictive distribution is a Gaussian whose mean and variance depend on the test points X_* and on the kernel $k(x, x')$, evaluated on the training data X .

Slide credit: Bernt Schiele B. Leibe 25

RWTH AACHEN UNIVERSITY

Gaussian Process Regression

- Example

Slide credit: Bernt Schiele B. Leibe 26

RWTH AACHEN UNIVERSITY

Gaussian Process Regression

Slide credit: Bernt Schiele B. Leibe 27

Advanced Machine Learning Winter'12

Discussion

- Key result:** $f_*|X_*, X, \mathbf{t} \sim \mathcal{N}(\bar{f}_*, \text{cov}[f_*])$ with

$$\bar{f}_* = K(X_*, X) (K(X, X) + \sigma_n^2 I)^{-1} \mathbf{t}$$

$$\text{cov}[f_*] = K(X_*, X_*) - K(X_*, X) (K(X, X) + \sigma_n^2 I)^{-1} K(X, X_*)$$
- Observations**
 - The mean can be written in linear form

$$\bar{f}(\mathbf{x}_*) = k(\mathbf{x}_*, X) \underbrace{[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{t}}_{\boldsymbol{\alpha}} = \sum_{n=1}^N \alpha_n k(\mathbf{x}_*, \mathbf{x}_n).$$
 - This form is commonly encountered in the kernel literature (\rightarrow SVM)
 - The variance is the difference between two terms

$$V(\mathbf{x}_*) = \underbrace{k(\mathbf{x}_*, \mathbf{x}_*)}_{\text{Prior variance}} - \underbrace{k(\mathbf{x}_*, X) [K(X, X) + \sigma_n^2 I]^{-1} k(X, \mathbf{x}_*)}_{\text{Explanation of data } X}$$

Slide adapted from Carl Rasmussen B. Leibe 28

Advanced Machine Learning Winter'12

Computational Complexity

- Computational complexity**
 - Central operation in using GPs involves **inverting a matrix of size $N \times N$** (the kernel matrix $K(X, X)$):

$$\bar{f}_* = K(X_*, X) \underbrace{[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{t}}_{\boldsymbol{\alpha}}$$

$$\text{cov}[f_*] = K(X_*, X_*) - K(X_*, X) \underbrace{[K(X, X) + \sigma_n^2 I]^{-1}}_{\mathbf{S}} K(X, X_*)$$
 - \Rightarrow Effort in $\mathcal{O}(N^3)$ for N data points!
 - Compare this with the basis function model (\rightarrow Lecture 3)

$$p(f_* | \mathbf{x}_*, X, \mathbf{t}) \sim \mathcal{N}\left(\frac{1}{\sigma_n^2} \phi(\mathbf{x}_*)^T \mathbf{S}^{-1} \Phi(X) \mathbf{t}, \phi(\mathbf{x}_*)^T \mathbf{S}^{-1} \phi(\mathbf{x}_*)\right)$$

$$\mathbf{S} = \frac{1}{\sigma_n^2} \Phi(X) \Phi(X)^T + \Sigma_p^{-1}$$
 - \Rightarrow Effort in $\mathcal{O}(M^3)$ for M basis functions.

B. Leibe 30

Advanced Machine Learning Winter'12

Computational Complexity

- Complexity of GP model**
 - Training effort: $\mathcal{O}(N^3)$ through matrix inversion
 - Test effort: $\mathcal{O}(N^2)$ through vector-matrix multiplication
- Complexity of basis function model**
 - Training effort: $\mathcal{O}(M^3)$
 - Test effort: $\mathcal{O}(M^2)$
- Discussion**
 - If the number of basis functions M is smaller than the number of data points N , then the basis function model is more efficient.
 - However, advantage of GP viewpoint is that we can consider covariance functions that can only be expressed by an **infinite number of basis functions**.
 - Still, exact GP methods become infeasible for large training sets.

B. Leibe 31

Advanced Machine Learning Winter'12

GP Regression Algorithm

- Very simple algorithm!**

```

input: X (inputs), y (targets), k (covariance function),  $\sigma_n^2$  (noise level),
       $\mathbf{x}_*$  (test input)
2: L := cholesky(K +  $\sigma_n^2 I$ )
    $\boldsymbol{\alpha} := L^{-1}(L \backslash \mathbf{y})$ 
4:  $\bar{f}_* := \mathbf{k}_*^T \boldsymbol{\alpha}$ 
    $\mathbf{v} := L \backslash \mathbf{k}_*$ 
6:  $\mathbb{V}[f_*] := k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^T \mathbf{v}$ 
8: return:  $\bar{f}_*$  (mean),  $\mathbb{V}[f_*]$  (variance),  $\log p(\mathbf{y}|X)$  (log marginal likelihood)

```

 - Based on the following equations (Matrix inv. \leftrightarrow Cholesky fact.)

$$\bar{f}_* = \mathbf{k}_*^T (K + \sigma_n^2 I)^{-1} \mathbf{t}$$

$$\text{cov}[f_*] = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (K + \sigma_n^2 I)^{-1} \mathbf{k}_*$$

$$\log p(\mathbf{t}|X) = -\frac{1}{2} \mathbf{t}^T (K + \sigma_n^2 I)^{-1} \mathbf{t} - \frac{1}{2} \log |K + \sigma_n^2 I| - \frac{N}{2} \log 2\pi$$

B. Leibe Image source: Rasmussen & Williams, 2006 32

Advanced Machine Learning Winter'12

Influence of Hyperparameters

- Most covariance functions have some free parameters.**
 - Example:

$$k_{ij}(\mathbf{x}_p, \mathbf{x}_q) = \sigma_f^2 \exp\left\{-\frac{(\mathbf{x}_p - \mathbf{x}_q)^2}{2 \cdot l^2}\right\} + \sigma_n^2 \delta_{pq}$$
 - Parameters: (l, σ_f, σ_n)
 - Signal variance: σ_f^2
 - Range of neighbor influence (called "length scale"): l
 - Observation noise: σ_n^2

Slide credit: Bernt Schiele B. Leibe 33

Advanced Machine Learning Winter'12

Influence of Hyperparameters

$$k_{ij}(\mathbf{x}_p, \mathbf{x}_q) = \sigma_f^2 \exp\left\{-\frac{(\mathbf{x}_p - \mathbf{x}_q)^2}{2 \cdot l^2}\right\} + \sigma_n^2 \delta_{pq}$$

- Examples for different settings of the length scale**

$(l, \sigma_f, \sigma_n) =$ (σ parameters set by optimizing the marginal likelihood)

$= (0.3, 1.08, 0.00005)$ $= (1, 1, 0.1)$ $= (3.0, 1.16, 0.89)$

Slide credit: Bernt Schiele B. Leibe Image source: Rasmussen & Williams, 2006 34

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- **Kernels**
 - Recap: Kernel trick
 - Constructing kernels
- **Gaussian Processes**
 - Recap: Definition
 - Prediction with noise-free observations
 - Prediction with noisy observations
 - GP Regression
 - Influence of hyperparameters
- **Learning Gaussian Processes**
 - Bayesian Model Selection
 - Model selection for Gaussian Processes
- **Applications**

35

RWTH AACHEN UNIVERSITY

Learning Kernel Parameters

- Can we determine the length scale and noise levels from training data?

36

RWTH AACHEN UNIVERSITY

Bayesian Model Selection

- **Goal**
 - Determine/learn different parameters of Gaussian Processes
- **Hierarchy of parameters**
 - **Lowest level**
 - w - e.g. parameters of a linear model.
 - **Mid-level (hyperparameters)**
 - θ - e.g. controlling prior distribution of w .
 - **Top level**
 - Typically discrete set of model structures \mathcal{H}_i .
- **Approach**
 - Inference takes place one level at a time.

37

RWTH AACHEN UNIVERSITY

Model Selection at Lowest Level

- **Posterior of the parameters w is given by Bayes' rule**

$$p(w|t, X, \theta, \mathcal{H}_i) = \frac{p(t|X, w, \theta, \mathcal{H}_i)p(w|\theta, X, \mathcal{H}_i)}{p(t|X, \theta, \mathcal{H}_i)}$$

$$= \frac{p(t|X, w, \mathcal{H}_i)p(w|\theta, \mathcal{H}_i)}{p(t|X, \theta, \mathcal{H}_i)}$$

- **with**
 - $p(t|X, w, \mathcal{H}_i)$ likelihood and
 - $p(w|\theta, \mathcal{H}_i)$ prior parameters w ,
 - Denominator (normalizing constant) is independent of the parameters and is called **marginal likelihood**.

$$p(t|X, \theta, \mathcal{H}_i) = \int p(t|X, w, \mathcal{H}_i)p(w|\theta, \mathcal{H}_i)dw$$

38

RWTH AACHEN UNIVERSITY

Model Selection at Mid Level

- **Posterior of parameters θ is again given by Bayes' rule**

$$p(\theta|t, X, \mathcal{H}_i) = \frac{p(t|X, \theta, \mathcal{H}_i)p(\theta|X, \mathcal{H}_i)}{p(t|X, \mathcal{H}_i)}$$

$$= \frac{p(t|X, \theta, \mathcal{H}_i)p(\theta|\mathcal{H}_i)}{p(t|X, \mathcal{H}_i)}$$

- **where**
 - The marginal likelihood of the previous level $p(t|X, \theta, \mathcal{H}_i)$ plays the role of the likelihood of this level.
 - $p(\theta|\mathcal{H}_i)$ is the **hyperprior** (prior of the hyperparameters)
 - Denominator (normalizing constant) is given by:

$$p(t|X, \mathcal{H}_i) = \int p(t|X, \theta, \mathcal{H}_i)p(\theta|\mathcal{H}_i)d\theta$$

which is again a **marginal likelihood** (at the mid level).

39

RWTH AACHEN UNIVERSITY

Model Selection at Top Level

- **At the top level, we calculate the posterior of the model**

$$p(\mathcal{H}_i|t, X) = \frac{p(t|X, \mathcal{H}_i)p(\mathcal{H}_i)}{p(t|X)}$$

- **where**
 - Again, the denominator of the previous level $p(t|X, \mathcal{H}_i)$ plays the role of the likelihood.
 - $p(\mathcal{H}_i)$ is the prior of the model structure.
 - Denominator (normalizing constant) is given by:

$$p(t|X) = \sum_i p(t|X, \mathcal{H}_i)p(\mathcal{H}_i)$$

40

Advanced Machine Learning Winter'12

Bayesian Model Selection

RWTH AACHEN UNIVERSITY

- Discussion
 - Marginal likelihood is main difference to non-Bayesian methods
 - It automatically incorporates a trade-off between the model fit and the model complexity:
 - A simple model can only account for a limited range of possible sets of target values - if a simple model fits well, it obtains a high posterior.
 - A complex model can account for a large range of possible sets of target values - therefore, it can never attain a very high posterior.

Slide credit: Bernt Schiele B. Leibe Image source: Rasmussen & Williams, 2006

41

Advanced Machine Learning Winter'12

Bayesian Model Selection

RWTH AACHEN UNIVERSITY

- Computational issues
 - Requires the evaluation of several integrals, which may or may not be analytically tractable, depending on details of the models.
 - In general, one may have to resort to analytic approximations or MCMC methods. (→Lecture 7)
- Model selection for GP regression
 - GP regression models with Gaussian noise are an (important) exception:
 - Integrals over the parameters are analytically tractable and
 - At the same time, the models are flexible.

Slide credit: Bernt Schiele B. Leibe

42

Advanced Machine Learning Winter'12

Example

RWTH AACHEN UNIVERSITY

Slide credit: Bernt Schiele B. Leibe

43

Advanced Machine Learning Winter'12

Example

RWTH AACHEN UNIVERSITY

Slide credit: Bernt Schiele B. Leibe

44

Advanced Machine Learning Winter'12

Example

RWTH AACHEN UNIVERSITY

Slide credit: Bernt Schiele B. Leibe

45

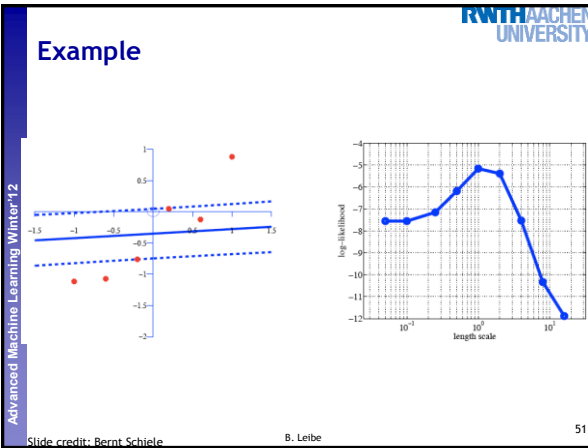
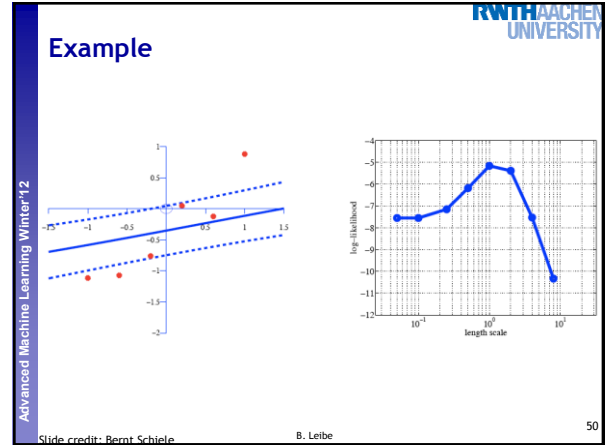
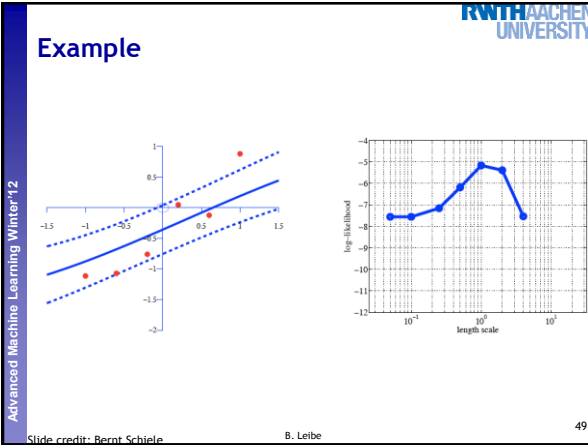
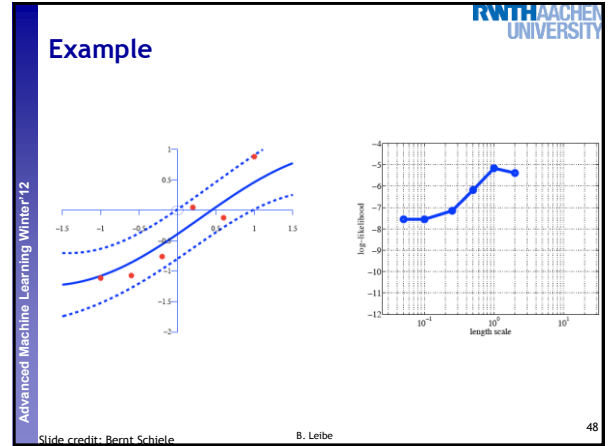
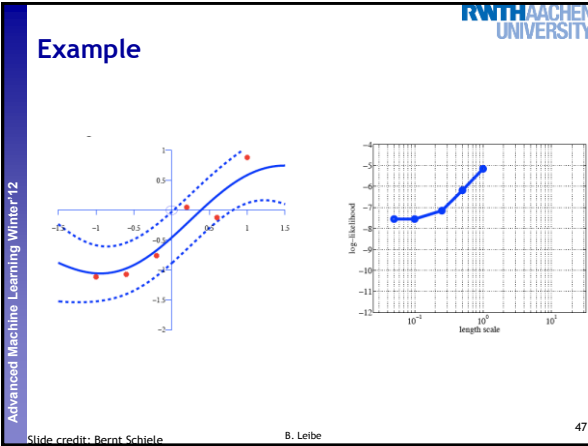
Advanced Machine Learning Winter'12

Example

RWTH AACHEN UNIVERSITY

Slide credit: Bernt Schiele B. Leibe

46



Advanced Machine Learning Winter'12

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- **Kernels**
 - Recap: Kernel trick
 - Constructing kernels
- **Gaussian Processes**
 - Recap: Definition
 - Prediction with noise-free observations
 - Prediction with noisy observations
 - GP Regression
 - Influence of hyperparameters
- **Learning Gaussian Processes**
 - Bayesian Model Selection
 - Model selection for Gaussian Processes
- **Applications**

Slide credit: Bernt Schiele

B. Leibe

65

RWTH AACHEN UNIVERSITY

Application: Non-Linear Dimensionality Reduction

2D manifold in 3D space ↔ 2D space ↔ 2D latent space

30D articulated body space ↔ 2D latent space

Slide credit: Andreas Geiger B. Leibe 66

RWTH AACHEN UNIVERSITY

Gaussian Process Latent Variable Model

- At each time step t , we express our observations y as a combination of basis functions ψ of latent variables x .

$$y_t = \sum_j b_j \psi_j(x_t) + \delta_t$$

- This is modeled as a Gaussian process...

Slide credit: Andreas Geiger B. Leibe 67

RWTH AACHEN UNIVERSITY

Example: Style-based Inverse Kinematics

Learned GPLVMs using a walk, a jump shot and a baseball pitch

Slide credit: Andreas Geiger B. Leibe 68

RWTH AACHEN UNIVERSITY

Application: Modeling Body Dynamics

- Task: estimate full body pose in m video frames.
 - High-dimensional \mathbf{Y} .
 - Model body dynamics using hierarchical Gaussian process latent variable model (hGPLVM) [Lawrence & Moore, ICML 2007].

Time (frame #) $\mathbf{T} = [t_i \in \mathbb{R}]$

Latent space $\mathbf{Z} = [z_i \in \mathbb{R}^q]$

Configuration $\mathbf{Y} = [y_i \in \mathbb{R}^D]$

Training

$$p(\mathbf{Z}|\mathbf{T}, \theta) = \prod_{i=1}^q \mathcal{N}(z_{:,i} | 0, \mathbf{K}_{\mathbf{T}})$$

$$p(\mathbf{Y}|\mathbf{Z}, \theta) = \prod_{i=1}^D \mathcal{N}(y_{:,i} | 0, \mathbf{K}_{\mathbf{z}})$$

Slide credit: Bernd Schiele B. Leibe [Andriuka, Roth, Schiele, CVPR'08] 69

RWTH AACHEN UNIVERSITY

Application: Mapping b/w Pose and Appearance

- Appearance prediction
 - Regression problem
 - High-dimensional data on both sides
 - Low-dim. representation needed for learning!
- Training with Motion-capture data possible
 - Synthesized silhouettes for training
 - Background subtraction for test

- 3D joint locations • 60-dim.
- segm. image • ~2500-dim.

Slide credit: Jaeselli, Koller-Meier, Van Gool, ACCV'07 70

RWTH AACHEN UNIVERSITY

Learning a Generative Mapping

Body Pose

X : Body Pose (high dim.)

x : Body Pose (low dim.)

reconstruct pose

Learn LLE dim. red.

dynamic prior

generative mapping

likelihood

Appearance

Y : Image (high dim.)

y : Appearance Descriptor: (low dim.)

PCA projection

Slide credit: B. Leibe [Jaeselli, Koller-Meier, Van Gool, ACCV'07] 71

Advanced Machine Learning Winter'12

Experimental Results

- Difficulties
 - Changing viewpoints
 - Low resolution (50 px)
 - Compression artifacts
 - Disturbing objects

Original video

J. Jaaeeli, Koller, Heiser, Van Gool, ACCV'07

RWTH AACHEN UNIVERSITY

Advanced Machine Learning Winter'12

Articulated Motion in Latent Space (different work)

- Gaussian Process regression from latent space to
 - Pose $[\rightarrow = p(\text{Pose} | z)]$ to recover original pose from latent space]
 - Silhouette $[\rightarrow = p(\text{Silhouette} | z)]$ to do inference on silhouettes]

Walking cycles have one main (periodic) DOF

Additional DOF encodes „walking style“

B. Leibe [Gammeter, Ess, Leibe, Schindler, Van Gool, ECCV'08]

RWTH AACHEN UNIVERSITY

Advanced Machine Learning Winter'12

Results

454 frames (~35 sec)
23 Pedestrians
20 detected by multi-body tracker

B. Leibe [Gammeter, Ess, Leibe, Schindler, Van Gool, ECCV'08]

RWTH AACHEN UNIVERSITY

Advanced Machine Learning Winter'15

References and Further Reading

- Kernels and Gaussian Processes are (shortly) described in Chapters 6.1 and 6.4 of Bishop's book.
 - Christopher M. Bishop
Pattern Recognition and Machine Learning
Springer, 2006
- A better introduction can be found in Chapters 3 and 5 of the book by Rasmussen & Williams (also available online: <http://www.gaussianprocess.org/gpml/>)
 - Carl E. Rasmussen, Christopher K.I. Williams
Gaussian Processes for Machine Learning
MIT Press, 2006

B. Leibe

RWTH AACHEN UNIVERSITY