

# Statistical Machine Learning

## Lecture 13: Kernel Regression and Gaussian Processes

**Kristian Kersting**

**TU Darmstadt**

Summer Term 2020

# Today's Objectives

- Make you understand how to use kernels for regression both from a frequentist and Bayesian point of view
- Covered Topics
  - Why kernel methods?
  - Radial basis function networks
  - What is a kernel?
  - Dual representation
  - Gaussian Process Regression

# Outline

- 1. Kernel Methods for Regression**
- 2. Gaussian Processes Regression**
- 3. Bayesian Learning and Hyperparameters**
- 4. Wrap-Up**

# Outline

## 1. Kernel Methods for Regression

## 2. Gaussian Processes Regression

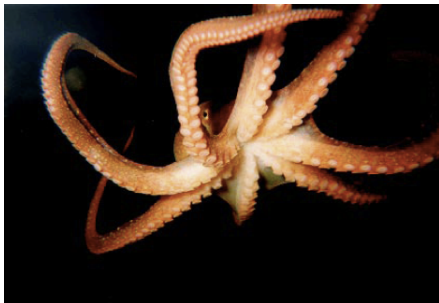
## 3. Bayesian Learning and Hyperparameters

## 4. Wrap-Up

# Why Kernels and not Neural Networks?

- Multi-Layer Perceptrons use univariate projections to “span” the space of the data (like an “octopus”)

$$y = g(\mathbf{w}^T \mathbf{x})$$



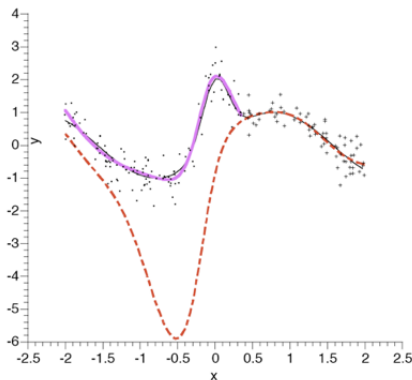
# Why Kernels and not Neural Networks?

## ■ Pros

- Universal function approximation
- Large range generalization (extrapolation)
- Good for high dimensional data

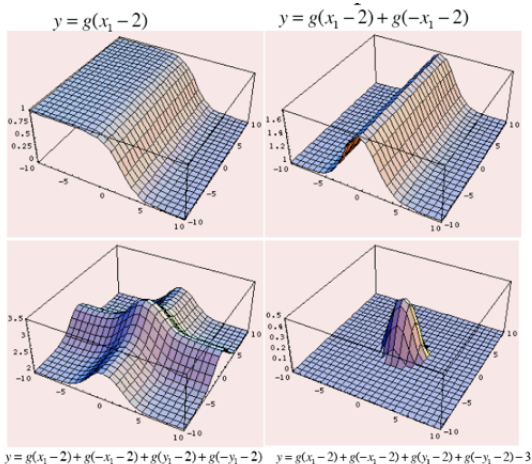
## ■ Cons

- Hard to train
- Danger of interference



# Radial Basis Function Networks

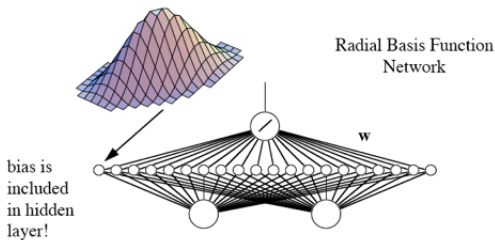
- Use spatially localized kernels for learning
- Note: there are other basis functions that are not spatially localized



# Radial Basis Function Networks

- For instance with Gaussian kernels

$$y = \sum_{i=1}^K w_i \phi(\|\mathbf{x} - \mathbf{x}_i\|) = \mathbf{w}^T \Phi(\|\mathbf{x} - \mathbf{x}_i\|)$$



$$\phi(\mathbf{x}, \mathbf{c}_k) = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{c}_k)^T \mathbf{D}(\mathbf{x} - \mathbf{c}_k)\right)$$

with  $\mathbf{D}$  positive definite



# Radial Basis Function Networks

- The “output layer” is just a linear regression
- Often needs regularization (e.g., ridge regression)

$$J = \frac{1}{2} (\mathbf{t} - \mathbf{y})^\top (\mathbf{t} - \mathbf{y}) = \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^\top (\mathbf{t} - \Phi \mathbf{w})$$

$$\mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}, \quad \Phi = \begin{bmatrix} \phi_{11} & \phi_{12} & \dots & \phi_{1m} \\ \phi_{21} & \phi_{22} & \dots & \phi_{2m} \\ \dots & \dots & \dots & \dots \\ \phi_{n1} & \phi_{n2} & \dots & \phi_{nm} \end{bmatrix}$$

$$\mathbf{w} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{t}$$

# Radial Basis Function Networks

- The “input layer” can be optimized by gradient descent with respect to distance metric and centers of RBFs

$$\frac{\partial J}{\partial \mathbf{c}_k} = (t - y) \left( -\frac{\partial y}{\partial \mathbf{c}_k} \right) = -(t - y) w_k \left( \frac{\partial \Phi}{\partial \mathbf{c}_k} \right)$$

$$\begin{aligned} \frac{\partial \Phi}{\partial \mathbf{c}_k} &= \frac{\partial}{\partial \mathbf{c}_k} \exp \left( -\frac{1}{2} (\mathbf{x} - \mathbf{c}_k)^T \mathbf{D} (\mathbf{x} - \mathbf{c}_k) \right) \\ &= \exp \left( -\frac{1}{2} (\mathbf{x} - \mathbf{c}_k)^T \mathbf{D} (\mathbf{x} - \mathbf{c}_k) \right) (\mathbf{x} - \mathbf{c}_k)^T \mathbf{D} \end{aligned}$$

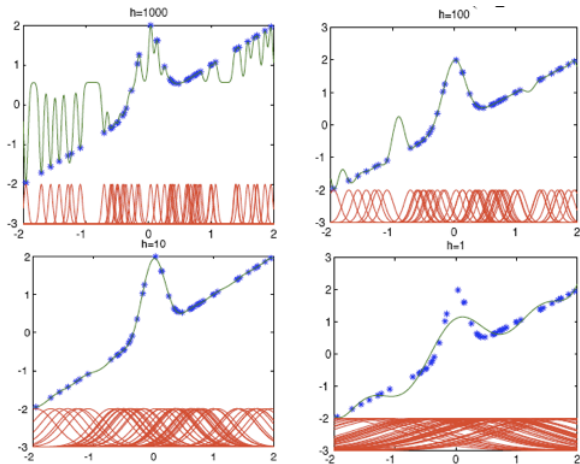
$$\frac{\partial J}{\partial \mathbf{D}_k} = (t - y) \left( -\frac{\partial y}{\partial \mathbf{D}_k} \right) = -(t - y) w_k - \frac{\partial \Phi}{\partial \mathbf{D}_k}$$

$$\begin{aligned} \frac{\partial \Phi}{\partial \mathbf{D}_k} &= \frac{\partial}{\partial \mathbf{D}_k} \exp \left( -\frac{1}{2} (\mathbf{x} - \mathbf{c}_k)^T \mathbf{D}_k (\mathbf{x} - \mathbf{c}_k) \right) \\ &= \exp \left( -\frac{1}{2} (\mathbf{x} - \mathbf{c}_k)^T \mathbf{D}_k (\mathbf{x} - \mathbf{c}_k) \right) (\mathbf{x} - \mathbf{c}_k) (\mathbf{x} - \mathbf{c}_k)^T \end{aligned}$$

- Gradient descent can make  $\mathbf{D}$  non positive definite  $\implies$  use Cholesky Decomposition
- An iterative procedure is needed to for optimization, i.e., alternate update of  $\mathbf{w}$  and update of  $\mathbf{c}_k$  and  $\mathbf{D}_k$

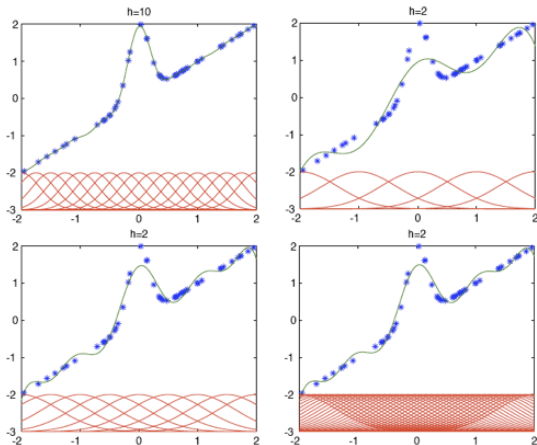
# Radial Basis Function Networks

- Sensitivity to kernel width (bandwidth, dist. metric) of  $\phi(x, c_k) = \exp\left(-\frac{1}{2}(x - c_k)^2 h\right)$



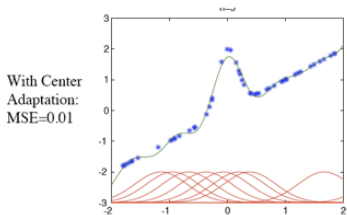
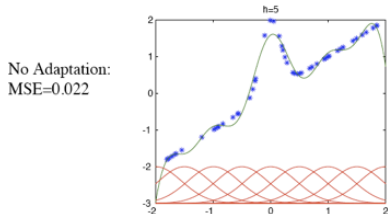
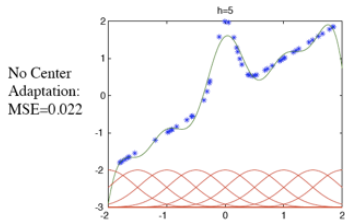
# Radial Basis Function Networks

- Sensitivity to number of kernels and metric of  $\phi(x, c_k) = \exp\left(-\frac{1}{2}(x - c_k)^2 h\right)$

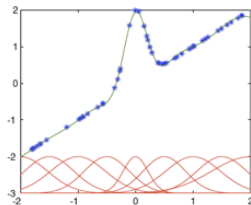


# Radial Basis Function Networks

## ■ Benefits of center and metric adaptation



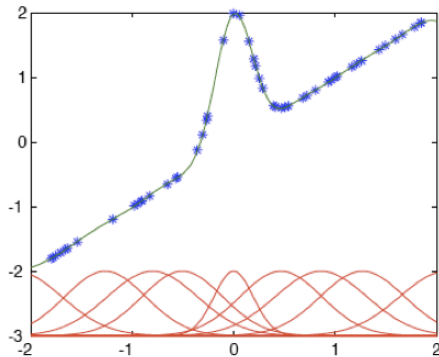
Distance Metric Adaptation:  
MSE=0.0001



# Radial Basis Function Networks

- All adaptations turned on
- Note: RBF tend to grow wider with a lot of overlap, and learning rates are sensitive

Distance Metric  
and Center  
Adaptation:  
MSE=0.00001



# Radial Basis Function Networks - Summary

- RBFs are a powerful and efficient learning tool
- Number of RBFs and hyperparameter optimization is important and a bit difficult to tune
- Theoretical remark
  - Poggio and Girosi (1990) showed that RBF networks arise naturally from minimizing the penalized cost function

$$J = \frac{1}{2} \sum_n (t^n - y(x^n))^2 + \frac{1}{2} \gamma \int |G(\mathbf{x})|^2 d\mathbf{x}$$

with, e.g.,  $G(\mathbf{x}) = \frac{\partial^2 y}{\partial \mathbf{x}^2}$ , a smoothless prior

# Kernel Methods in General

- What is a kernel?
  - Most intuitive approach for a fixed nonlinear feature space: an inner product of feature vectors

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

- A kernel is symmetric

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$$

- Examples

- Stationary kernels:  $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$
- Linear kernel:  $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$
- Homogeneous kernels:  $k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$



# Dual Representation of Linear Regression

- The dual representation gives natural rise to the kernel functions

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}, \quad \text{where } \lambda \geq 0$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n) \phi(\mathbf{x}_n) + \lambda \mathbf{w} = 0$$

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n) \phi(\mathbf{x}_n) = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a}$$

where  $\Phi = [\phi(\mathbf{x}_1)^T \dots \phi(\mathbf{x}_N)^T] \in \mathbb{R}^{N \times D}$

- Thus,  $\mathbf{w}$  is a linear combination of  $\phi(\mathbf{x}_n)$
- The dual representation focuses on solving for  $\mathbf{a}$ , and not  $\mathbf{w}$

# Dual Representation of Linear Regression

- Insert the dual representation into the cost function

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

$$\begin{aligned} J(\mathbf{a}) &= \frac{1}{2} \sum_{n=1}^N (\mathbf{a}^T \Phi \phi(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \\ &= \frac{1}{2} \sum_{n=1}^N \mathbf{a}^T \Phi \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \mathbf{a} + \frac{1}{2} \sum_{n=1}^N t_n^2 - \sum_{n=1}^N \mathbf{a}^T \Phi \phi(\mathbf{x}_n) t_n + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \\ &= \frac{1}{2} \mathbf{a}^T \underbrace{\Phi \Phi^T}_{\mathbf{K}} \Phi \Phi^T \mathbf{a} + \frac{1}{2} \mathbf{t}^T \mathbf{t} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \\ &= \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} + \frac{1}{2} \mathbf{t}^T \mathbf{t} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a} \end{aligned}$$

- $\mathbf{K} = \Phi \Phi^T$  is the **Gram Matrix**, and  $K_{ij} = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j)$

# Dual Representation of Linear Regression

- Solve the dual problem for  $\mathbf{a}$

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^\top \mathbf{K} \mathbf{K} \mathbf{a} + \frac{1}{2} \mathbf{t}^\top \mathbf{t} - \mathbf{a}^\top \mathbf{K} \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^\top \mathbf{K} \mathbf{a}$$
$$\frac{\partial J(\mathbf{a})}{\partial \mathbf{a}} = \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{K} \mathbf{t} + \lambda \mathbf{K} \mathbf{a} = \mathbf{K} (\mathbf{K} \mathbf{a} - \mathbf{t} + \lambda \mathbf{a}) = 0$$
$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t}$$

- Side note: since by definition of a kernel matrix,  $\mathbf{K}$  is Positive Semi-Definite,  $\mathbf{K}^{-1}$  exists

# Dual Representation of Linear Regression

- Compute the prediction as

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

where  $\mathbf{k}(\mathbf{x}) = [k(\mathbf{x}, \mathbf{x}_1) \dots k(\mathbf{x}, \mathbf{x}_N)]^\top$

- All computations can be expressed in terms of the kernel function  $k$

# Pros and Cons of the Dual Representation

## ■ Cons

- Need to invert a  $N \times N$  matrix

## ■ Pros

- Can work entirely in feature space with the help of kernels
- Can even consider infinite feature spaces, as the kernel function does only have the inner product of feature vectors, which is a scalar, even for infinite feature spaces
- Many novel algorithms can be derived from the dual representation
- Many old problems of RBFs (how many kernels, which metric, which centers) can be solved in a principled way

## Some Useful Kernels

- Polynomial kernels
  - E.g., 2nd order:  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z})^2$
  - N-th order with offset:  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + c)^N$
- Gaussian Kernel (also called Radial Basis Function - RBF)

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|^2\right)$$

- Arises from a feature space with an **INFINITE** number of radial basis functions

$$\int_{-\infty}^{+\infty} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|^2\right) \exp\left(-\frac{1}{2\sigma^2} \|\tilde{\mathbf{x}} - \mathbf{x}'\|^2\right) d\tilde{\mathbf{x}}$$

$$\propto \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|^2\right)$$

# Outline

1. Kernel Methods for Regression

**2. Gaussian Processes Regression**

3. Bayesian Learning and Hyperparameters

4. Wrap-Up

# Dual Representation of Linear Regression

- Classical linear (ridge/regularized) regression

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$$

$$\mathbf{w} = (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top \mathbf{t}$$

$$\mathbf{t} = \mathbf{y} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \lambda)$$

- Dual representation of linear regression

$$y(\mathbf{x}) = \mathbf{a}^\top \mathbf{k}(\mathbf{x})$$

$$\mathbf{a} = (\mathbf{K} + \lambda I)^{-1} \mathbf{t}$$

$$\mathbf{k}(\mathbf{x}) = [k(\mathbf{x}, \mathbf{x}_1) \quad \dots \quad k(\mathbf{x}, \mathbf{x}_n)]^\top$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

$$\mathbf{t} = \mathbf{y} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \lambda)$$



# Bayesian Linear Regression Revisited

- Regression model

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$$

- Parameter Distribution

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha^{-1}I)$$

- Thus, for any  $\mathbf{w}$ , one particular function of  $\mathbf{x}$  is defined
- The distribution over  $\mathbf{w}$  thus induces a distribution over functions
- Goal: evaluate the function at some values of  $\mathbf{x}$ , e.g., the training set  $\mathbf{x}_1, \dots, \mathbf{x}_n$

$$\mathbf{y} = \Phi \mathbf{w}$$

and predict the joint probability  $p(\mathbf{y}_1, \dots, \mathbf{y}_n)$

# Bayesian Linear Regression Revisited

$$\mathbf{y} = \Phi \mathbf{w} \quad p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I})$$

- $\mathbf{y}$  is a linear combination of Gaussian random variables, and thus Gaussian itself
- To obtain the joint distribution of all  $\mathbf{y}$ , we only need the mean and covariance

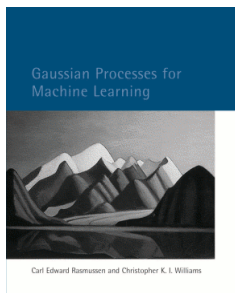
$$\mathbb{E}\{\mathbf{y}\} = \mathbb{E}\{\Phi \mathbf{w}\} = \Phi \mathbb{E}\{\mathbf{w}\} = \mathbf{0}$$

$$\text{cov}\{\mathbf{y}\} = \mathbb{E}\{\mathbf{y}\mathbf{y}^T\} = \Phi \mathbb{E}\{\mathbf{w}\mathbf{w}^T\} \Phi^T = \frac{1}{\alpha} \Phi \Phi^T = \mathbf{K}$$

$$\text{where } K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{\alpha} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

# Gaussian Processes

- A Gaussian Process (GP) is a probability distribution over functions  $y(\mathbf{x})$ , such that any finite set of function values  $y(\mathbf{x})$  evaluated at inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  is jointly Gaussian distributed
- A Gaussian Process over  $n$  variables  $y_1, \dots, y_n$  is completely specified by the 2nd order statistics, i.e., mean and covariance
- Rasmussen and Williams, 2006, *Gaussian Processes for Machine Learning* (<http://www.gaussianprocess.org/gpml/>)
- Good introduction to GPs by Carl Rasmussen: [http://videlectures.net/mlss09uk\\_rasmussen\\_gp/](http://videlectures.net/mlss09uk_rasmussen_gp/)

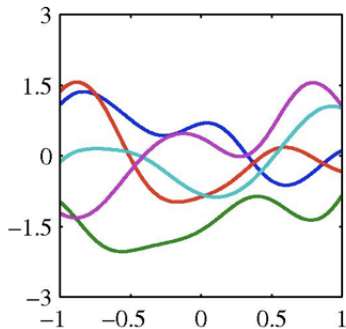


# Gaussian Processes

- A GP is fully specified by a mean function and a covariance function (kernel)
  - Prior mean function: expected function before observing any data
  - Covariance function: encodes some structural assumptions (e.g., smoothness) (e.g., multivariate Gaussian kernel)
- Most applications assume the prior mean of  $\mathbf{y}$  to be zero
  - Corresponds to a mean-zero prior of  $\mathbf{w}$
- Thus, a GP is completely defined by

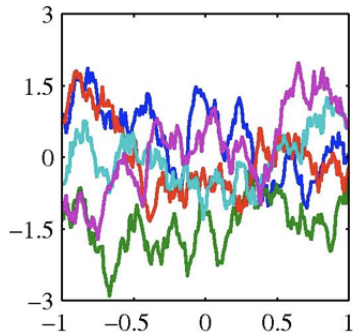
$$\begin{aligned}\mathbb{E}\{\mathbf{y}\} &= \mathbb{E}\{\Phi\mathbf{w}\} = \Phi\mathbb{E}\{\mathbf{w}\} = \mathbf{0} \\ \mathbb{E}\{y(\mathbf{x}_i)y(\mathbf{x}_j)\} &= k(\mathbf{x}_i, \mathbf{x}_j)\end{aligned}$$

# GPs - Different Covariance Functions



$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right)$$

**Gaussian/RBF Kernel**



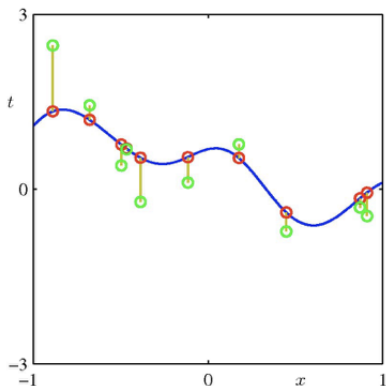
$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\theta \|\mathbf{x}_i - \mathbf{x}_j\|)$$

**Ornstein-Uhlenbeck Process  
(Brownian Motion)**

# GPs for Regression

- Generative model:  $t_n = y(x_n) + \epsilon$
- Noise model:  $p(t_n | y_n) = \mathcal{N}(t_n | y_n, \beta^{-1})$
- Prior distribution over function values:  $p(y) = \mathcal{N}(y | 0, K)$
- The kernel function that determines  $K$  is typically chosen to express the property that, for similar points  $x_n$  and  $x_m$  the corresponding values  $y(x_n)$  and  $y(x_m)$  will be more strongly correlated than for dissimilar points. The definition of similarity depends on the application

# GPs for Regression - Sampling Example



- Illustration of the sampling of data points  $t_n$  from a GP. The blue curve shows a sample function  $y$  from the GP posterior over functions. The red points show the values of  $y_n$  obtained by evaluating the function at a set of input values  $x_n$ . The corresponding values of  $t_n$ , shown in green, are obtained by adding independent Gaussian noise to each of the  $y_n$

# Inferring Functions with GPs

- Prior over functions (GP):  $p(y)$
- Likelihood (measurement/noise model):  $p(t | y)$
- Posterior over functions via Bayes theorem

$$p(y|t) = \frac{p(t|y)p(y)}{p(t)}$$



# GPs Regression - Prediction for New Data Points

- Training set

$$\mathbf{t}_n = (t_1, \dots, t_n)^\top \text{ with corresponding } \mathbf{x}_1, \dots, \mathbf{x}_n$$

- **Predict**  $t_{n+1}$  for  $x_{n+1}$
- Approach: evaluate the predictive distribution

$$p(t_{n+1} \mid x_{n+1}, t_{1:n}, x_{1:n})$$

- For the derivation, remember that **GPs assume that**  
 $p(t_1, t_2, \dots, t_n, t_{n+1})$  is jointly Gaussian
- Therefore, **the conditional distribution**  $p(t_{n+1} \mid x_{n+1}, t_{1:n}, x_{1:n})$  is also Gaussian distributed

# Gaussian Conditioning

- Assume  $\mathbf{x}$  is Gaussian distributed and it can be partitioned in two disjoint subsets  $\mathbf{x}_a$  and  $\mathbf{x}_b$ . We can rewrite the distribution in terms of the mean and covariance matrices of  $\mathbf{x}_a$  and  $\mathbf{x}_b$

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{bmatrix}$$

- The conditional distribution is also Gaussian

$$p(\mathbf{x}_a \mid \mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a \mid \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$$

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba}$$

# Gaussian Conditioning and Marginalization

- With Gaussian distributions the following holds

$$p(\mathbf{x}, \mathbf{y}) \sim \mathcal{N}$$

$$\Downarrow$$

$$p(\mathbf{x} | \mathbf{y}) \sim \mathcal{N}$$

$$p(\mathbf{y} | \mathbf{x}) \sim \mathcal{N}$$

$$p(\mathbf{x}) \sim \mathcal{N}$$

$$p(\mathbf{y}) \sim \mathcal{N}$$

# GPs Regression - Prediction for New Data Points

$$p(\mathbf{t}_{n+1}) = \mathcal{N}(\mathbf{t}_{n+1} \mid \mathbf{0}, \mathbf{C}_{n+1})$$

$$\mathbf{C}_{n+1} = \begin{pmatrix} \mathbf{C}_n & \mathbf{k} \\ \mathbf{k} & c \end{pmatrix}$$

where

$$\mathbf{k} = [k(\mathbf{x}_1, \mathbf{x}_{n+1}) \quad \dots \quad k(\mathbf{x}_n, \mathbf{x}_{n+1})]^\top$$

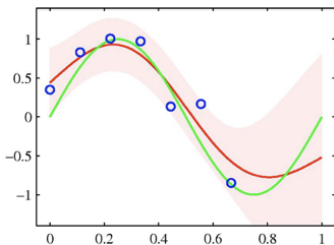
$$c = k(\mathbf{x}_{n+1}, \mathbf{x}_{n+1}) + \beta^{-1}$$

# GPs Regression - Prediction for New Data Points

## ■ Prediction Equations

$$m(\mathbf{x}_{n+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}$$

$$\sigma^2(\mathbf{x}_{n+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}$$



- Example of Sinusoidal Data Set (green: true function; blue: noisy data; red: GPR predictive mean; shaded:  $\pm 2\sigma$ )

# GPs Regression - Notes

- Interpretation as RBFs

$$m(\mathbf{x}_{n+1}) = \mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{t} = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x}_{n+1})$$

- Computational Complexity
  - For building the model:  $O(N^3)$
  - For prediction of one function value:  $O(N^2)$  (for the variance)
- Key advantage of GPR: non-parametric and probabilistic

# GPs Regression - Notes

- Naive methods can deal with  $\sim 10.000 - 20.000$  data points
- Advanced methods (e.g., Sparse GPs) for more than 50.000 data points
- **IMPORTANT:** Hyperparameter optimization (parameters of the kernel / covariance function). E.g., for squared-exponential kernel

$$k(x_i, x_j) = \sigma_f^2 \exp\left(-\frac{1}{2l^2} (x_i - x_j)^2\right) + \sigma_n^2 \delta_{ij}$$

where  $\sigma_f^2$  is the signal variance,  $l$  is the length-scale and  $\sigma_n^2$  is the noise variance

## GPs - Summary

- GPs are a Bayesian approach to regression with possibly infinite feature spaces
- Resulting prediction equations are very straightforward and obtained in closed-form because of the Gaussian properties
- Hyperparameter optimization more complex and expensive
- While GP for Regression is computationally very expensive, it is one of the most principled approaches to statistical learning for regression



# Outline

1. Kernel Methods for Regression

2. Gaussian Processes Regression

**3. Bayesian Learning and Hyperparameters**

4. Wrap-Up

## Bayesian Learning - Pros

- Bayesian methods are a superset of many learning methods
- Regularization is a natural consequence
- No need for splitting into training and test sets
- Confidence intervals and error bars can be obtained
- Regularization can be obtained automatically
- Model comparison
- Active learning (determine where to sample next)
- Automatic relevance detection (which inputs are important)
- Black-box learning approaches
- Theoretically among the most powerful method

# Bayesian Learning - Cons

- Requires to choose prior distributions, mostly based on analytical convenience rather than real knowledge about the problem
- Computationally intractable
  - Posterior probabilities involve the computation of an integral

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

- On the contrary, in non-Bayesian statistics we estimate parameters with maximum likelihood estimation. MLE is usually easier because it involves finding the maximum of the likelihood function, for which you can still use gradient descent, if there is no analytical solution

# Bayesian Learning - Key Issues

- All parameters are treated probabilistically, i.e., avoid “point estimates”
- The probabilistic treatment allows integrating out all unknown parameters
- The problem of infinite regress?

# Bayesian Learning - Key Issues

- Quantities of most interest in Bayesian approaches
  - Model Evidence

$$p(D) = \int p(D, \theta) d\theta = \int p(D | \theta) p(\theta) d\theta$$

- Posterior of parameters

$$p(\theta | D) = \frac{p(D | \theta) p(\theta)}{p(D)}$$

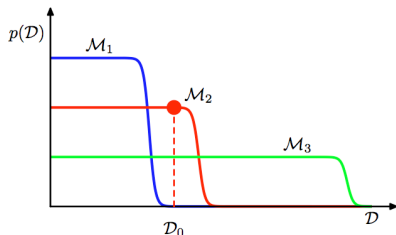
- Predictive distribution

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

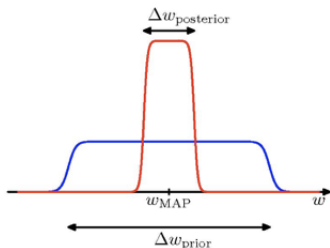
# The Philosophy of Bayesian Model Selection

- Due to probability measure, models can be compared using the evidence
  - Complex models have lower probability density over large range of data sets
  - Simple models have high probability density over a small range of data sets
  - Thus, there should be a compromise in terms of complexity and confidence in the model

$$p(M_i | D) = \frac{p(D | M_i) p(M_i)}{\sum_j p(D | M_j) p(M_j)}$$



# Why the Evidence Achieves Regularization



- Approximate evidence in a one parameter scenario

$$p(D) = \int p(D | w) p(w) dw \approx p(D | w_{\text{MAP}}) \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}}$$

$$\ln p(D) \approx \ln p(D | w_{\text{MAP}}) + \ln \left( \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right)$$

- Note that the 2nd term penalizes the model complexity according to how finely the posterior is tuned

# Why the Evidence Achieves Regularization

- For  $M$  parameters

$$\ln p(D) \approx \ln p(D | w_{\text{MAP}}) + M \ln \left( \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right)$$

- The penalty increases with the number of parameters



# Bayesian Learning

- Parameters are modeled by probability distributions
- The conditional distribution of a new data point  $x$  given the training data  $D$  can be written as the marginalized joint distribution

$$\begin{aligned} p(x | D) &= \int p(x, \mathbf{w} | D) d\mathbf{w} = \int p(x | D, \mathbf{w}) p(\mathbf{w} | D) d\mathbf{w} \\ &= \int p(x | \mathbf{w}) p(\mathbf{w} | D) d\mathbf{w} \end{aligned}$$

- Hence the Bayesian approach performs a weighted average over all values of  $\mathbf{w}$

# Bayesian Learning

- Connection to maximum likelihood estimation

$$\begin{aligned} p(x | D) &= \int p(x | \mathbf{w}, D) p(\mathbf{w} | D) d\mathbf{w} \\ &\approx p(x | \hat{\mathbf{w}}, D) \underbrace{\int p(\mathbf{w} | D) d\mathbf{w}}_{=1} \\ &= p(x | \hat{\mathbf{w}}, D) \end{aligned}$$

- The approximation usually holds for sufficiently many training data points

# Bayesian Learning

- How to perform Bayesian Updates

$$p(\mathbf{w} | D) = \frac{p(D | \mathbf{w}) p(\mathbf{w})}{p(D)} = \frac{p(\mathbf{w})}{p(D)} \prod_{n=1}^N p(\mathbf{x}^n | D)$$

$$p(D) = \int p(\mathbf{w}') \prod_{n=1}^N p(\mathbf{x}^n | \mathbf{w}') d\mathbf{w}'$$

- To obtain predictions, one has to evaluate the integrals

$$p(D) = \int p(\mathbf{w}') \prod_{n=1}^N p(\mathbf{x}^n | \mathbf{w}') d\mathbf{w}'$$

$$p(x | D) = \int p(x | \mathbf{w}) p(\mathbf{w} | D) d\mathbf{w}$$

- Generally this is a very complex computation
- Analytical solutions exist only if the posterior has the same parametric form as the prior (**conjugate priors**, reproducing densities)

# Example - Bayesian Density Estimation with a Gaussian

- Determine the mean of the Gaussian by Bayesian Learning

$$\text{Prior: } p_0(\mu) = \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left(-\frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right)$$

$$\text{Gaussian model: } p(x | \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

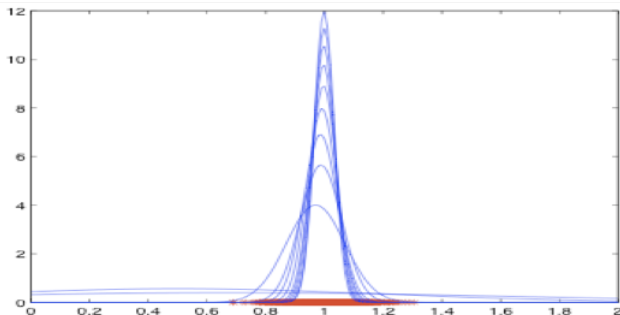
$$p_N(\mu | \mathbf{X}) = \frac{p_0(\mu)}{p(\mathbf{X})} \prod_{n=1}^N p(x^n | \mu)$$

- We get a Gaussian posterior with parameters

$$\mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \bar{x} + \frac{\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_0, \quad \frac{1}{\sigma_N^2} = \frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}, \quad \bar{x} = \frac{1}{N} \sum_{i=1}^N x^i$$

# Example - Bayesian Density Estimation with a Gaussian

- Evolution of the posterior probability of the mean (blue) with increasing number of data points (red)



# Bayesian Learning

- Assume probability distribution over network weights and some prior
- Need to interpret outputs probabilistically
- **Bayesian Learning Procedure**
  - Start with prior distribution  $p(w)$  and choose appropriate parameters (usually broad distribution to reflect uncertainty)
  - Observe data, compute posterior of parameters with Bayes rule
  - Continue updating if more data comes in, replacing the prior with the posterior
  - In order to make a prediction, the expectation given the posterior distribution has to be found (might be very complex)

## Gaussian Priors

- As usual, Gaussian priors are most convenient to deal with (for real numbers), although they may be unrealistic

$$p_0(w) = \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left(-\frac{(w-w_0)^2}{2\sigma_0^2}\right)$$

- More generally ...

$$p(\mathbf{w}) = \frac{1}{Z_w(\alpha)} \exp(-\alpha E_w)$$

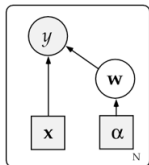
where

$$Z_w(\alpha) = \int \exp(-\alpha E_w) d\mathbf{w} = \left(\frac{2\pi}{\alpha}\right)^{W/2}$$

$$E_w = \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|^2 = \frac{1}{2} \sum_{i=1}^W (w_i - w_{0,i})^2, \quad E_w = \frac{1}{2} \|\mathbf{w}\|^2 = \frac{1}{2} \sum_{i=1}^W w_i^2$$

$\alpha$  is a hyperparameter

## Example - Logistic Regression



$$p_N(\mathbf{w} | D) = \frac{p_0(\mathbf{w})}{p(D)} \prod_{n=1}^N p(t^n | \mathbf{w})$$

$$\text{Prior: } p(\mathbf{w}) = \exp\left(-\frac{\alpha}{2} \|\mathbf{w}\|^2\right) / \left(\frac{2\pi}{\alpha}\right)^{W/2}$$

$$\text{Likelihood: } p(t^n | \mathbf{w}) = y(\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

$$p(D) = \int y(\mathbf{x}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

- Note:  $p(D)$  is difficult to estimate, but we do not need it as long as we do not attempt model comparison, since it is only a scaling factor

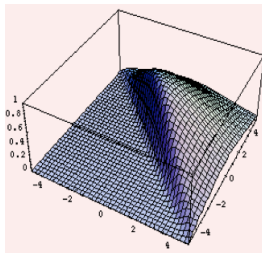


## Example - Logistic Regression

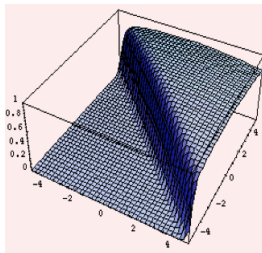
- Consider the data set

$$X = \begin{bmatrix} 5 & 5 \\ -5 & -5 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad T = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

- Use only two data points



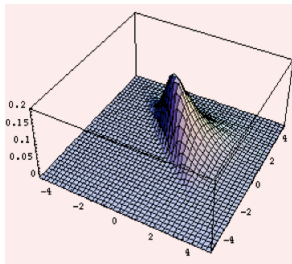
$\alpha = 0.1$



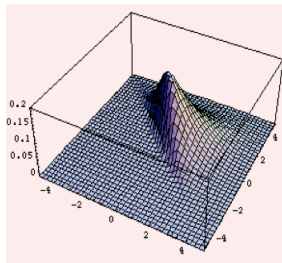
$\alpha = 0.01$

# Example - Logistic Regression

- Use all data points



$\alpha = 0.1$



$\alpha = 0.01$

# Gaussian Noise Models

- Make Gaussian Assumption for Likelihood

$$p(D | \mathbf{w}) = \frac{1}{Z_D(\beta)} \exp(-\beta E_D)$$

$$\text{where } Z_D(\beta) = \int \exp(-\beta E_D)$$

- For instance, for regression assume

$$p(t | \mathbf{x}, \mathbf{w}) \propto \exp\left(-\frac{1}{2}\beta (y(\mathbf{x}, \mathbf{w}) - t)^2\right)$$

# Gaussian Noise Models

- Posterior Distributions of Weights
  - Since all distributions are Gaussian, the posterior must be Gaussian and can be written as

$$p(\mathbf{w} | D) = \frac{1}{Z_S} \exp(-\beta E_D - \alpha E_W) = \frac{1}{Z_S} \exp(-S(\mathbf{w}))$$

$$\text{where } Z_D(\alpha, \beta) = \int \exp(-\beta E_D - \alpha E_W) d\mathbf{w}$$

- What is the parameter vector maximizing the posterior? This can be achieved by minimizing

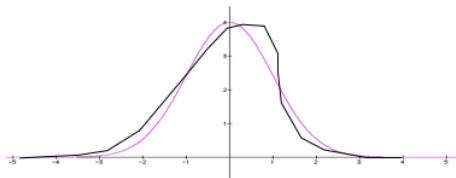
$$S(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N (y(\mathbf{x}^n; \mathbf{w}) - t^n)^2 + \frac{\alpha}{2} \sum_{i=1}^W w_i^2$$

# Gaussian Approximation of Posterior Distribution

- Assume that there is an **analytically intractable** distribution
  - E.g., after obtaining the posterior of the parameters, the likelihood of the model may be desirable

$$p(y | D) = \int p(y | \mathbf{w}) p(\mathbf{w} | D) d\mathbf{w}$$

- E.g., the posterior is required in Gaussian form
- Way out: e.g., approximate the intractable distribution with a Gaussian



# Laplace Approximation

- Assume the generic probability distribution

$$p(z) = \frac{1}{Z} f(z), \text{ with } Z = \int f(z) dz$$

- Goal: Approximate  $p(z)$  with a Gaussian distribution, centered at the mode  $z_0$

$$\left. \frac{df(z)}{dz} \right|_{z=z_0}$$

- We get a 2nd order Taylor series expansion

$$\ln f(z) \approx \ln f(z_0) - \frac{1}{2} \left( - \frac{d^2}{dz^2} \ln f(z) \Big|_{z=z_0} \right) (z - z_0)^2$$

# Laplace Approximation

- Taking the exp we get

$$f(\mathbf{z}) \approx f(\mathbf{z}_0) \exp\left(-\frac{1}{2}\mathbf{A}(\mathbf{z} - \mathbf{z}_0)^2\right)$$

$$q(\mathbf{z}) \approx \left(\frac{\mathbf{A}}{2\pi}\right)^{1/2} \exp\left(-\frac{1}{2}\mathbf{A}(\mathbf{z} - \mathbf{z}_0)^2\right)$$

- For the multivariate case

$$\ln f(\mathbf{z}) \approx \ln f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0)$$

$$\mathbf{A} = -\nabla\nabla \ln f(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0}$$

$$f(\mathbf{z}) \approx f(\mathbf{z}_0) \exp\left(-\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0)\right)$$

$$q(\mathbf{z}) \approx \left(\frac{|\mathbf{A}|}{(2\pi)^M}\right)^{1/2} \exp\left(-\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0)\right) = \mathcal{N}\left(\mathbf{z} \mid \mathbf{z}_0, \mathbf{A}^{-1}\right)$$

# Laplace Approximation

- Illustration for approximation of logistic regression



## Statistical Machine Learning

### Lecture 13: Kernel Regression and Gaussian Processes

Jan Peters  
TU Darmstadt

Summer Semester 2019



# Dealing with Hyperparameters

- Augment Framework to also model the hyperparameters probabilistically

$$\begin{aligned} p(\mathbf{w} | D) &= \int \int p(\mathbf{w}, \alpha, \beta | D) d\alpha d\beta \\ &= \int \int p(\mathbf{w} | \alpha, \beta, D) p(\alpha, \beta | D) d\alpha d\beta \end{aligned}$$

- Assuming a sharp peak of the distributions of the hyperparameters

$$p(\mathbf{w} | D) \approx p(\mathbf{w} | D, \alpha_{\text{MP}}, \beta_{\text{MP}}) \int \int p(\alpha, \beta | D) d\alpha d\beta = p(\mathbf{w} | D, \alpha_{\text{MP}}, \beta_{\text{MP}})$$

- These assumptions offer the possibility to first find the hyperparameters that maximize the posterior, and then perform the remaining calculations with these optimized hyperparameters
- Note that there are also other methods to obtain the hyperparameters

# Hyperparameters in Gaussian Processes

- What are the hyperparameters in GPs?
  - E.g, exponential-quadratic kernel

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^\top \mathbf{x}_m$$

- Approach: optimize the evidence w.r.t. the hyperparameters

$$p(\mathbf{t}) = \int p(\mathbf{t} | \mathbf{y}) p(\mathbf{y}) d\mathbf{y} = \mathcal{N}(\mathbf{t} | \mathbf{0}, \mathbf{C})$$

$$\text{with } \mathbf{C}(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1} \delta_{nm}$$

- E.g., by gradient descent

$$\frac{\partial}{\partial \theta_i} \log p(\mathbf{t} | \theta) = -\frac{1}{2} \text{Tr}\left(\mathbf{C}_n^{-1} \frac{\partial \mathbf{C}_n}{\partial \theta_i}\right) + \frac{1}{2} \mathbf{t}^\top \mathbf{C}_n^{-1} \frac{\partial \mathbf{C}_n}{\partial \theta_i} \mathbf{C}_n^{-1} \mathbf{t}$$

# Hyperparameters - Summary

- Bayesian learning offers an automatic way of regularization and meta-parameter tuning
- The evidence framework for model selection offers a principled tool to compare different learning systems
- Most of the time, Bayesian learning is analytically intractable
- Approximation methods exist to deal with the intractable components (Bayesian “hacking”)

# Outline

1. Kernel Methods for Regression

2. Gaussian Processes Regression

3. Bayesian Learning and Hyperparameters

**4. Wrap-Up**

## 4. Wrap-Up

You know now:

- What RBF Networks are
- What Kernels are, how to construct them and why they are beneficial
- How to derive the dual formulation of linear regression, and what are its pros and cons
- What GPs are, and the assumptions behind them
- With GPs we can predict the value for a new point in closed form, because of the Gaussian conditionals
- Doing regression with GPs we get a mean value and a variance (uncertainty) of the estimate
- Generally methods with kernels do not scale well with data
- The ideas behind Bayesian Learning, its pros and cons

# Self-Test Questions

- Why kernel methods for regression?
- How do you get from radial basis functions to kernels?
- What is the role of the two pseudo-inverses in kernel regression?
- Why are kernel regression methods very computationally expensive?
- Why is kernel regression the dual to linear regression?
- What is the major advantage of GPs over Kernel Ridge Regression?
- Why are GPs a Bayesian approach?
- What principle allowed deriving GPs from a Bayesian regression point of view?
- How to get the hyperparameters in a Bayesian setup?

# Extra Material & Homework

- Extra material
  - Goertler, et al., "A Visual Exploration of Gaussian Processes", Distill, 2019 (<https://distill.pub/2019/visual-exploration-gaussian-processes/>)
- Reading Assignment for next lecture
  - Bishop 8