

Pattern Recognition and Machine Learning

Chapter 6: Kernel Methods

Vasil Khalidov Alex Kläser

December 13, 2007

Training Data: Keep or Discard?

- ▶ **Parametric** methods (linear/nonlinear) so far:
 - ▶ learn parameter vector \mathbf{w} or posterior distribution $p(\mathbf{w}|\mathbf{t})$
 - ▶ discard training data \mathbf{t}
- ▶ **Non-parametric** methods:
 - ▶ **Parzen** probability density model: set of kernel functions centered on training data points
 - ▶ **Nearest neighbours** technique: closest example from the training set
 - ▶ **Memory-based** methods: similar examples from the training set
- ▶ **Kernel** methods:
 - ▶ prediction is based on linear combinations of a *kernel function* evaluated at the training data points

Kernel Function

- ▶ for models based on feature space mapping $\phi(\mathbf{x})$:

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') \quad (6.1)$$

- ▶ symmetric function: $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$
- ▶ simple example - *linear* kernel: $\phi(\mathbf{x}) = \mathbf{x}$
- ▶ *stationary* kernel: $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$
- ▶ *homogeneous* kernel: $k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$

Algorithm, expressed in terms of scalar products can be reformulated using *kernel substitution*: PCA, nearest-neighbour classifiers, Fisher discriminant

Dual Representation

Many linear models for regression and classification can be reformulated in terms of a dual representation in which the kernel function arises naturally.

- ▶ linear regression model ($\lambda \geq 0$):

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \quad (6.2)$$

- ▶ set the gradient to zero:

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \} \phi(\mathbf{x}_n) = \Phi^T \mathbf{a} \quad (6.3)$$

- ▶ substitute \mathbf{w} and define the *Gram* matrix $\mathbf{K} = \Phi \Phi^T$:

$$K_{nm} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) \quad (6.6)$$

Solution for Dual Problem

- ▶ in terms of new parameter vector \mathbf{a} :

$$\mathbf{w} = \frac{1}{2}\mathbf{a}^T\mathbf{K}\mathbf{K}\mathbf{a} - \mathbf{a}^T\mathbf{K}\mathbf{t} + \frac{1}{2}\mathbf{t}^T\mathbf{t} + \frac{\lambda}{2}\mathbf{a}^T\mathbf{K}\mathbf{a} \quad (6.7)$$

- ▶ set the gradient to zero:

$$\mathbf{a} = (\mathbf{K} + \lambda\mathbf{I}_N)^{-1}\mathbf{t} \quad (6.8)$$

- ▶ prediction for a new input \mathbf{x}

$$y(\mathbf{x}) = \mathbf{w}^T\phi(\mathbf{x}) = \mathbf{a}^T\Phi\phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T(\mathbf{K} + \lambda\mathbf{I}_N)^{-1}\mathbf{t} \quad (6.9)$$

where

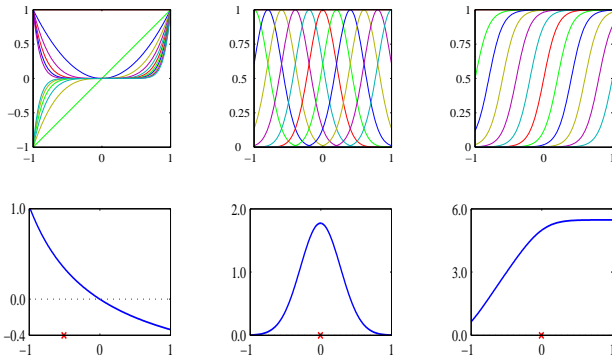
$$\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}_1, \mathbf{x}), \dots, k(\mathbf{x}_N, \mathbf{x}))$$

- ▶ inverting $N \times N$ matrix instead of $M \times M$

Constructing Kernels - First Approach

- ▶ choose feature space mapping $\phi(x)$

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x') \quad (6.10)$$



Constructing Kernels - Second Approach

- ▶ construct kernel function directly and verify its validity
- ▶ simple example

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2 \quad (6.11)$$

in 2-D case corresponds to

$$k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z}) \quad (6.12)$$

with $\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)^T$

To test validity without having to construct the function $\phi(\mathbf{x})$ explicitly, one can use the condition:

Function $k(\mathbf{x}, \mathbf{x}')$ is a valid kernel $\iff \mathbf{K} \geq 0 \quad \forall \{x_n\}$

Combining Kernels

Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$ the following kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

with corresponding conditions on $c, f, q, \phi, k_3, \mathbf{A}, \mathbf{x}_a, \mathbf{x}_b, k_a, k_b$

Examples

- ▶ Polynomial kernels:

$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$ contains only terms of degree 2

$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^2, c > 0$

$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^M$ contains all monomials of order M

$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^M, c > 0$

- ▶ *Gaussian* kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2/2\sigma^2) \quad (6.23)$$

Note: can substitute $\mathbf{x}^T \mathbf{x}'$ with a nonlinear kernel $\kappa(\mathbf{x}, \mathbf{x}')$

- ▶ Kernel on nonvectorial space:

$$k(A_1, A_2) = 2^{|A_1 \cap A_2|} \quad (6.27)$$

- ▶ Sigmoidal kernel:

$$k(\mathbf{x}, \mathbf{x}') = \tanh(a\mathbf{x}^T \mathbf{x}' + b) \quad (6.37)$$

Probabilistic generative models

- ▶ Kernel for generative model $p(\mathbf{x})$:

$$k(\mathbf{x}, \mathbf{x}') = p(\mathbf{x})p(\mathbf{x}') \quad (6.28)$$

$$k(\mathbf{x}, \mathbf{x}') = \sum_i p(\mathbf{x}|i)p(\mathbf{x}'|i)p(i) \quad (6.29)$$

$$k(\mathbf{x}, \mathbf{x}') = \int p(\mathbf{x}|z)p(\mathbf{x}'|z)p(z)dz \quad (6.30)$$

- ▶ Kernel for HMM:

$$k(\mathbf{X}, \mathbf{X}') = \sum_{\mathbf{Z}} p(\mathbf{X}|\mathbf{Z})p(\mathbf{X}'|\mathbf{Z})p(\mathbf{Z}) \quad (6.31)$$

$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_L\}$ - observations

$\mathbf{Z} = \{z_1, \dots, z_L\}$ - hidden states

Fisher kernel

- ▶ parametric generative model $p(\mathbf{x}|\boldsymbol{\theta})$
- ▶ *Fisher score*:

$$\mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) = \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x}|\boldsymbol{\theta}) \quad (6.32)$$

- ▶ *Fisher kernel* and information matrix:

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^T \mathbf{F}^{-1} \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}') \quad (6.33)$$

$$\mathbf{F} = \mathbb{E}_{\mathbf{x}} [\mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) \mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^T | \boldsymbol{\theta}] \quad (6.34)$$

- ▶ *Note*: the kernel is invariant under $\boldsymbol{\theta} \rightarrow \boldsymbol{\psi}(\boldsymbol{\theta})$
- ▶ Simplify matrix calculation:

$$\mathbf{F} \simeq \frac{1}{N} \sum_{n=1}^N \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}_n) \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}_n)^T \quad (6.35)$$

Or simply omit the Fisher information matrix

Radial Basis Functions

- ▶ By definition $\phi_j(\mathbf{x}) = h(\|\mathbf{x} - \boldsymbol{\mu}_j\|)$
- ▶ Originally were introduced for the problem of exact interpolation $f(\mathbf{x}_n) = t_n$:

$$f(\mathbf{x}) = \sum_{n=1}^N w_n h(\|\mathbf{x} - \mathbf{x}_n\|) \quad (6.38)$$

- ▶ Green's functions for an isotropic differential operator in regularizer
- ▶ Interpolation problem with noisy inputs

$$E = \frac{1}{2} \sum_{n=1}^N \int \{y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n\}^2 \nu(\boldsymbol{\xi}) d\boldsymbol{\xi} \quad (6.39)$$

Interpolation With Noisy Inputs

- ▶ Sum of squares function

$$E = \frac{1}{2} \sum_{n=1}^N \int \{y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n\}^2 \nu(\boldsymbol{\xi}) d\boldsymbol{\xi} \quad (6.39)$$

- ▶ Optimal value

$$y(\mathbf{x}_n) = \sum_{n=1}^N t_n h(\mathbf{x} - \mathbf{x}_n) \quad (6.40)$$

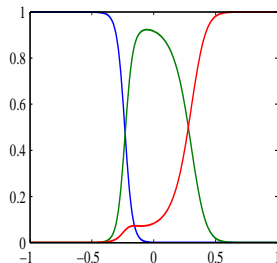
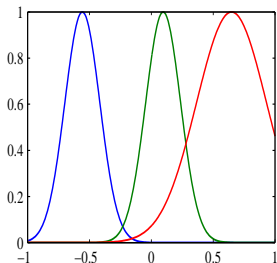
with basis functions given by normalized functions

$$h(\mathbf{x} - \mathbf{x}_n) = \frac{\nu(\mathbf{x} - \mathbf{x}_n)}{\sum_{n=1}^N \nu(\mathbf{x} - \mathbf{x}_n)} \quad (6.41)$$

- ▶ If noise distribution $\nu(\boldsymbol{\xi})$ is isotropic, basis functions would be radial

Normalization effect

Avoids regions in an input space where all of the basis functions take small values



Reducing Size of the Basis

- ▶ Keep number of basis functions M smaller than input data size N
- ▶ Centers locations μ_i are determined based on the input data $\{\mathbf{x}_n\}$ alone
- ▶ Coefficients $\{w_i\}$ are determined by least squares
- ▶ Choice of centers:
 - ▶ random
 - ▶ orthogonal least squares - greatest error reduction
 - ▶ using clustering algorithms

Parzen Density Estimator

- ▶ Prediction of linear regression model - linear combination of t_n with 'equivalent kernel' values
- ▶ Same result starting from Parzen density estimator

$$p(\mathbf{x}, t) = \frac{1}{N} \sum_{n=1}^N f(\mathbf{x} - \mathbf{x}_n, t - t_n) \quad (6.42)$$

- ▶ Regression function

$$\begin{aligned} y(\mathbf{x}) &= \mathbb{E}[t|\mathbf{x}] = \int_{-\infty}^{\infty} tp(t|\mathbf{x})dt = \\ &= \frac{\int tp(\mathbf{x}, t)dt}{\int p(\mathbf{x}, t)dt} = \\ &= \frac{\sum_n \int tf(\mathbf{x} - \mathbf{x}_n, t - t_n)dt}{\sum_m \int f(\mathbf{x} - \mathbf{x}_m, t - t_m)dt} \end{aligned} \quad (6.43)$$

Nadaraya-Watson Model

Assume that the component density functions have zero mean so that

$$\int_{-\infty}^{\infty} t f(\mathbf{x}, t) dt = 0 \quad (6.44)$$

for all values of \mathbf{x} . Then by variable change

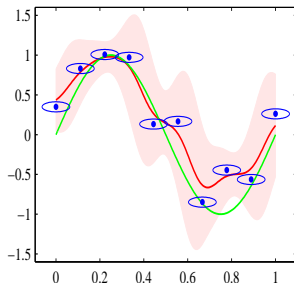
$$\begin{aligned} y(\mathbf{x}) &= \frac{\sum_n g(\mathbf{x} - \mathbf{x}_n) t_n}{\sum_m g(\mathbf{x} - \mathbf{x}_m)} = \\ &= \sum_n k(\mathbf{x}, \mathbf{x}_n) t_n \end{aligned} \quad (6.45)$$

with $g(\mathbf{x}) = \int_{-\infty}^{\infty} f(\mathbf{x}, t) dt$ and $k(\mathbf{x}, \mathbf{x}_n)$ given by

$$k(\mathbf{x}, \mathbf{x}_n) = \frac{g(\mathbf{x} - \mathbf{x}_n)}{\sum_m g(\mathbf{x} - \mathbf{x}_m)} \quad (6.46)$$

Illustration of the Nadaraya-Watson Model

single input variable x
 $f(x, t)$ is a zero-mean
isotropic Gaussian with
variance σ^2



Conditional distribution

$$p(t|\mathbf{x}) = \frac{p(t, \mathbf{x})}{\int p(t, \mathbf{x}) dt} = \frac{\sum_n f(\mathbf{x} - \mathbf{x}_n, t - t_n)}{\sum_m \int f(\mathbf{x} - \mathbf{x}_m, t - t_m) dt} \quad (6.47)$$

is given by a mixture of Gaussians

Gaussian processes: Key Idea

- ▶ The idea is similar to linear regression with a fixed set of basis functions (Chapter 3):

$$y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad (6.49)$$

- ▶ However, we forget about the parametric model
- ▶ Instead we take a infinite number of basis functions given by a probability distribution over functions
- ▶ Might look difficult to handle, but it is not . . . we only have to consider the values of the functions at training/test data points

Linear regression revisited

Model defined by linear combination of M fixed basis functions:

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) \quad (6.49)$$

A Gaussian prior distribution over the weight vector \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}) \quad (6.50)$$

governed by the hyperparameter α (precision) induces then a probability distribution over functions $y(\mathbf{x})$.

Linear regression revisited (2)

Evaluating $y(\mathbf{x})$ for a set of training data points $\mathbf{x}_1, \dots, \mathbf{x}_N$, we have a joint distribution

$$\mathbf{y} = \Phi \mathbf{w}, \quad \text{with elements} \quad y_n = y(\mathbf{x}_n) = \mathbf{w}^T \phi(\mathbf{x}_n) \propto \mathcal{N}, \quad (6.51)$$

where Φ is the design matrix with elements $\Phi_{nk} = \phi_k(\mathbf{x}_n)$.

$p(\mathbf{y})$ is Gaussian, and its mean and covariance can be shown to be

$$\mathbb{E}[\mathbf{y}] = \mathbf{0} \quad (6.52)$$

$$\text{cov}[\mathbf{y}] = \frac{1}{\alpha} \Phi \Phi^T = \mathbf{K} \quad (6.53)$$

where \mathbf{K} is the Gram matrix with elements

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \frac{1}{\alpha} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m). \quad (6.54)$$

Up to now we *only* took data points + prior, but no target values.

Linear regression revisited: summary

- ▶ The model we have seen is one example for a Gaussian process
- ▶ “Gaussian process is defined as a probability distribution over functions $y(\mathbf{x})$ such that the set of values of $y(\mathbf{x})$ evaluated at an arbitrary set of points $\mathbf{x}_1, \dots, \mathbf{x}_N$ jointly have a Gaussian distribution”
- ▶ Key point: the joint distribution is defined completely by second-order statistics (mean, covariance)
- ▶ Note, usually the mean is taken to be zero, then we only need the covariance, i.e., the kernel-function:

$$\mathbb{E}[y(\mathbf{x}_n)y(\mathbf{x}_m)] = k(\mathbf{x}_n, \mathbf{x}_m) \quad (6.55)$$

- ▶ Actually, instead of choosing (a limited number of) basis functions, we can directly choose a kernel function (which may result in an infinite number of basis functions)

Gaussian process regression

To use Gaussian processes for regression, we need to model noise:

$$t_n = y_n + \epsilon_n \quad \text{with} \quad y_n = y(\mathbf{x}_n) \quad (6.57)$$

For noise processes with a Gaussian distribution we obtain

$$p(t_n|y_n) = \mathcal{N}(t_n|y_n, \beta^{-1}). \quad (6.58)$$

The joint distribution for $\mathbf{t} = (t_1, \dots, t_n)^T$ and $\mathbf{y} = (y_1, \dots, y_n)^T$ is then (since the noise is assumed to be independent)

$$p(\mathbf{t}|\mathbf{y}) = \mathcal{N}(\mathbf{t}|\mathbf{y}, \beta^{-1}\mathbf{I}_N). \quad (6.59)$$

Gaussian process regression (2)

And from the definition of a Gaussian process we have

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}), \quad (6.60)$$

i.e., points that are more similar (given the kernel function k) will be stronger correlated. For the marginal distribution $p(\mathbf{t})$, we need to integrate over \mathbf{y} (see Section 2.3.3):

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y})d\mathbf{y} = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) \quad \text{with} \quad \mathbf{C} = \mathbf{K} + \beta^{-1}\mathbf{I}. \quad (6.61)$$

Making predictions

- ▶ So far we have a model for the joint probability distribution over sets of data points
- ▶ For predictions of a new input variable x_{N+1} , we need to evaluate the predictive distribution $p(t_{N+1}|\mathbf{t})$
- ▶ By partitioning (see Section 2.3.1) the joint Gaussian distribution over $\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}_{N+1}$, we obtain $p(t_{N+1}|\mathbf{t})$ given by its mean and covariance:

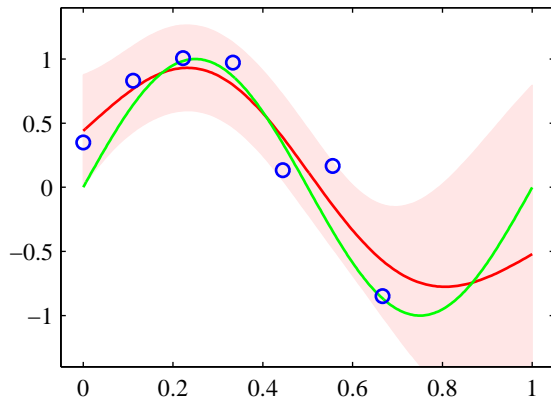
$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}^{-1} \mathbf{t} \quad \text{with (6.66)}$$

$$\mathbf{k} = (k(\mathbf{x}_1, \mathbf{x}_{N+1}), \dots, k(\mathbf{x}_N, \mathbf{x}_{N+1}))^T$$

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}^{-1} \mathbf{k} \quad \text{with (6.67)}$$

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

Making predictions (2)

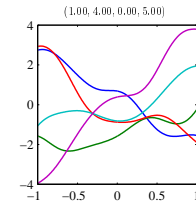
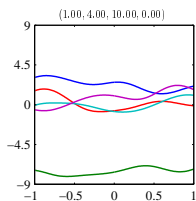
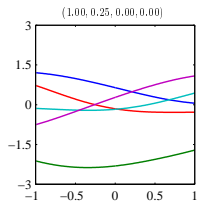
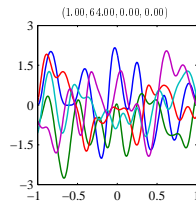
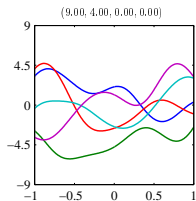
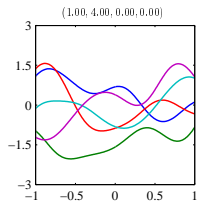


Green curve: original sinusoidal function; blue points: sampled training data points with additional noise; red line: mean estimate; shaded regions: $\pm 2\sigma$

Distribution over functions ... hm?

Sample functions from prior $p(\mathbf{y})$ with common kernel function

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m \quad (6.63)$$



Learning hyperparameters θ

- ▶ In practice, it can be preferable not to fix parameters, but to infer them from the data
- ▶ Parameters θ are, e.g.: length scale of correlations, precision of noise (β)
- ▶ Simplest approach:
 - ▶ Maximizing θ for the log-likelihood (maximum likelihood): $\ln p(\mathbf{t}|\theta)$
 - ▶ Problem: $\ln p(\mathbf{t}|\theta)$ is in general non-convex and can have multiple maxima
- ▶ Introduce a prior $p(\theta)$ and maximize the log-posterior (maximum a posteriori): $\ln p(\mathbf{t}|\theta) + \ln p(\theta)$
- ▶ To be Bayesian, we need the actual distribution and have to marginalize; this is not tractable \Rightarrow approximations
- ▶ The noise might not be additive but dependent on \mathbf{x}
 - ▶ A second Gaussian process can be introduced to represent the dependency of β on \mathbf{x}

Automatic relevance detection (ARD)

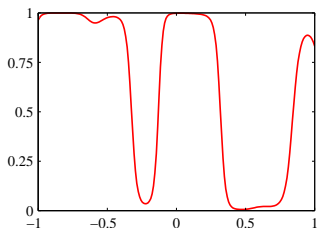
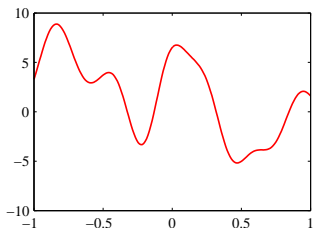
- ▶ An additional hyperparameter can be introduced for each input dimension, e.g.:

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left(-\frac{1}{2} \sum_{i=1}^D \eta_i (x_{ni} - x_{mi})^2 \right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m \quad (6.72)$$

- ▶ Hyperparameter optimization by maximum likelihood allows then a different weighting of each dimension
- ▶ Unrelevant dimensions (with small weights) can be detected and discarded

Gaussian process classification

- ▶ **Objective:** model posterior probabilities of the target variable for a new input
- ▶ **Problem:** we need to map values to interval $(0; 1)$
- ▶ **Solution:** use a Gaussian process together with a non-linear activation function (e.g., sigmoid)



Gaussian process classification (2)

Consider two-class problem with target values $t \in \{0, 1\}$. Define a Gaussian process over a function $a(\mathbf{x})$ and transform a using the logistic sigmoid to

$$y = \sigma(a(\mathbf{x})).$$

Similar to before, we need to predict the conditional distribution:

$$\begin{aligned} p(t_{N+1} = 1 | \mathbf{t}) &= \int p(t_{N+1} = 1 | a_{N+1}) p(a_{N+1} | \mathbf{t}) da_{N+1} \\ &= \int \sigma(a_{N+1}) p(a_{N+1} | \mathbf{t}) da_{N+1}. \end{aligned} \quad (6.76)$$

However, the integral is analytically intractable. Approximations can be of numerical or analytical nature.

Gaussian process classification (3)

- ▶ Problem 1: we need to compute a weird integral
- ▶ Solution 1: we know how to compute the convolution of an Gaussian and a sigmoid function (Eq. (4.153)) \Rightarrow approximate the posterior distribution $p(a_{N+1}|\mathbf{t})$ as Gaussian
- ▶ Problem 2: Hm . . . but how do we approximate the posterior?
- ▶ Solution 2: the Laplace approximation (among others)

Laplace approximation

We can rewrite the posterior over a_{N+1} using Bayes's theorem:

$$\begin{aligned} p(a_{N+1}|\mathbf{t}) &= \int p(a_{N+1}, \mathbf{a}|\mathbf{t})d\mathbf{a} \\ &= \dots = \int p(a_{N+1}|\mathbf{a})p(\mathbf{a}|\mathbf{t})d\mathbf{a}. \end{aligned} \quad (6.77)$$

We know how to compute mean and covariance for $p(a_{N+1}|\mathbf{a})$.
Now we have to find a Gaussian approximation only for $p(\mathbf{a}|\mathbf{t})$.
This is done noting that $p(\mathbf{a}|\mathbf{t}) \propto p(\mathbf{t}|\mathbf{a})p(\mathbf{a})$, thus

$$\Psi(\mathbf{a}) = \ln p(\mathbf{a}|\mathbf{t}) = \ln p(\mathbf{t}|\mathbf{a}) + \ln p(\mathbf{a}) + \text{const.} \quad (6.80)$$

Laplace approximation (2)

- ▶ We can use the iterative reweighted least squares (IRLS) algorithm (Sec. 4.3.3) to find the mode of $\Psi(\mathbf{a})$ (first and second derivative have to be evaluated)
- ▶ It can be shown that $\Psi(\mathbf{a})$ is convex and thus has only one mode :-)
- ▶ The mode position \mathbf{a}^* and the Hessian matrix \mathbf{H} at this position define our Gaussian approximation

$$q(\mathbf{a}) = \mathcal{N}(\mathbf{a}|\mathbf{a}^*, \mathbf{H}^{-1}) \quad (6.86)$$

- ▶ Now we can go back to the formulas and compute the integrals and finally also $p(t_{N+1}|\mathbf{t})$ from Eq. (6.76)

Learning hyperparameters

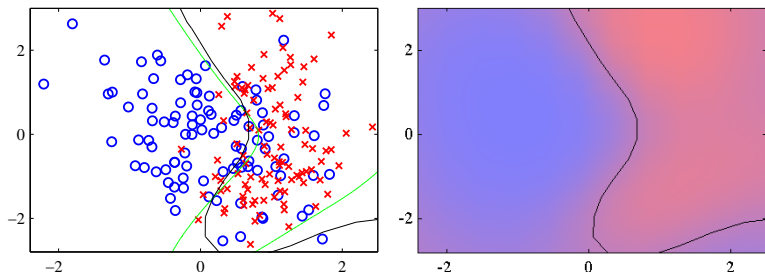
- ▶ To determine the parameters θ , we can maximize the likelihood function:

$$p(\mathbf{t}|\theta) = \int p(\mathbf{t}|\mathbf{a})p(\mathbf{a}|\theta)d\mathbf{a} \quad (6.89)$$

- ▶ Again, the integral is analytically intractable, so the Laplace approximation can be applied again
- ▶ We need an expression for the gradient of the logarithm of $p(\mathbf{t}|\theta)$
- ▶ Hm ... this is a bit trickier, but nevertheless doable
- ▶ ...by now we might be happy to skip the exact details :-)

Gaussian process classification (finishing up)

Just to finish up, and now we have a binary classifier based on a Gaussian process



Connection to neural networks

▶ Neural Networks

- ▶ The range of representable functions is governed by the number M of hidden units
- ▶ Within the maximum likelihood framework, they overfit as M comes close to the number of training samples

▶ Bayesian Neural Networks

- ▶ The prior over \mathbf{w} in conjunction with the network function $f(\mathbf{x}, \mathbf{w})$ produces a prior distribution over functions from $y(\mathbf{x})$
- ▶ The distribution of functions will tend to a Gaussian process in the limit $M \rightarrow \infty$
- ▶ The property that the outputs share hidden units (and thus 'borrow statistical strength' from each other) is lost in this limit
- ▶ ... and some other details