

Machine Learning: Waseda University

Kernel Methods

AD

June 2011

- Most/all of the algorithms we have discussed rely on a finite dimensional vector of *features* $\Phi(\mathbf{x})$.
- In this way, a model that is linear in \mathbf{x} may be made nonlinear by using a nonlinear mapping $\Phi(\mathbf{x})$.
- In many situations, we only rely on $\Phi(\mathbf{x})$ through the scalar product

$$k(\mathbf{x}, \mathbf{x}') = \Phi^T(\mathbf{x}) \Phi(\mathbf{x}')$$

- This is a symmetric function of its arguments

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

- A valid kernel is a function $k(\mathbf{x}, \mathbf{x}')$ that corresponds to a scalar (inner) product in some (perhaps infinite dimensional) feature space, i.e. $k(\mathbf{x}, \mathbf{x}') = \Phi^T(\mathbf{x}) \Phi(\mathbf{x}')$.
- For example assume $\mathbf{x} = (x_1, x_2)$ and

$$\begin{aligned}k(\mathbf{x}, \mathbf{x}') &= (\mathbf{x}^T \mathbf{x}')^2 \\&= (x_1 x_1' + x_2 x_2')^2 \\&= x_1^2 (x_1')^2 + x_2^2 (x_2')^2 + 2x_1 x_1' x_2 x_2' \\&= (x_1^2, \sqrt{2}x_1 x_2, x_2^2) \left((x_1')^2, \sqrt{2}x_1' x_2', (x_2')^2 \right) \\&= \Phi^T(\mathbf{x}) \Phi(\mathbf{x}')\end{aligned}$$

where

$$\Phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2).$$

Positive Semi-definite Kernels

- Loosely speaking, a kernel $k(\mathbf{x}, \mathbf{x}')$ can be written as a scalar product possibly in an infinite-dimensional space if it is positive semidefinite; that is for any n , $(\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathcal{X}^n$ and $(\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$ then

$$\sum_i \sum_j \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$$

- Indeed for continuous symmetric positive semidefinite kernel, we have Mercer's theorem. There exists a positive sequence $\{\lambda_i\}$ and functions $\Phi_i(\mathbf{x})$ such that

$$k(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \lambda_i \Phi_i(\mathbf{x}) \Phi_i(\mathbf{x}').$$

- More later...

- In many situations, as mentioned earlier, we actually only use $\Phi(\mathbf{x})$ through $\Phi^T(\mathbf{x}) \Phi(\mathbf{x}')$.
- Moreover it is often very difficult to design good features $\Phi(\mathbf{x})$.
- Wherever we have $\Phi^T(\mathbf{x}) \Phi(\mathbf{x}')$, we can 'kernelize' the algorithm and replace it by $k(\mathbf{x}, \mathbf{x}')$ where $k(\mathbf{x}, \mathbf{x}')$ is a p.s.d. kernel.
- So we can use infinite number of features.
- We can think of $k(\mathbf{x}, \mathbf{x}')$ as a similarity measure: it can be easier to design $k(\mathbf{x}, \mathbf{x}')$ than $\Phi(\mathbf{x})$.

Dual Representation of Linear Regression

- Consider

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

where $\lambda > 0$.

- By setting $\frac{\partial J}{\partial \mathbf{w}} = \mathbf{0}$ we obtain

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

where $a_n = -\frac{1}{\lambda} (\mathbf{w}^\top \Phi(\mathbf{x}_n) - t_n)$ and Φ is the design matrix

$$\Phi = \begin{pmatrix} \Phi^\top(\mathbf{x}_1) \\ \vdots \\ \Phi^\top(\mathbf{x}_N) \end{pmatrix}$$

- We now write $\mathbf{w} = \Phi^T \mathbf{a}$ and plug this expression in $J(\mathbf{w})$ so

$$\begin{aligned} J(\mathbf{a}) &= \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} - \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \\ &= \frac{1}{2} \mathbf{a}^T K K \mathbf{a} - \mathbf{a}^T K \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T K \mathbf{a} \end{aligned}$$

where $K = \Phi \Phi^T$.

- K is the Gram matrix

$$[K]_{i,j} = \Phi^T(\mathbf{x}_i) \Phi(\mathbf{x}_j)$$

- Note that by construction, K is a p.s.d. matrix; that is $\alpha^T K \alpha \geq 0$ for all α .

- Solving $\frac{\partial J}{\partial \mathbf{a}} = 0$ yields

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

- It follows that

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

where

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

- We now have to invert an $N \times N$ matrix instead of an $M \times M$ matrix (where $\Phi(\mathbf{x}) \in \mathbb{R}^M$).
- Now if we let $k(\mathbf{x}, \mathbf{x}')$ be a p.s.d. then you can still define $y(\mathbf{x}, \mathbf{w})$ whereas M is infinite!

- Mercer's theorem reformulated: $k(\mathbf{x}, \mathbf{x}')$ is a valid kernel iff the Gram matrix $K = [k(\mathbf{x}_n, \mathbf{x}_m)]$ is positive semi definite for all possible $\{\mathbf{x}_n\}$.
- A matrix A is psd iff $\alpha^T A \alpha \geq 0$ for all α .
- The corresponding features $\Phi(\cdot)$ are eigenfunctions of k , i.e.
$$\int k(\mathbf{x}, \mathbf{x}') \Phi_i(\mathbf{x}) d\mathbf{x} = \lambda_i \Phi_i(\mathbf{x}).$$

Example Kernels

- Stationary: $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$.
- Isotropic: $k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$.
- Monomials of order M : $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}')^M$.
- Monomials of order up to M : $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^M$
- “Gaussian” $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2)$.
- Sigmoid “kernel” (does not satisfy Mercer’s theorem!):
 $k(\mathbf{x}, \mathbf{x}') = \tanh(a\mathbf{x}^\top \mathbf{x}' + b)$.

Combining Kernels

- Assume $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$ are p.s.d. kernels then we can combine them in multiple ways to obtain new kernels.
- For any $\alpha, \beta > 0$ $k(\mathbf{x}, \mathbf{x}') = \alpha k_1(\mathbf{x}, \mathbf{x}') + \beta k_2(\mathbf{x}, \mathbf{x}')$ is p.s.d.
- $k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x}) k_1(\mathbf{x}, \mathbf{x}') f(\mathbf{x}')$ is p.s.d.
- $k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$ is p.s.d.
- $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') k_2(\mathbf{x}, \mathbf{x}')$ is p.s.d.
- $k(\mathbf{x}, \mathbf{x}') = k_1(\Phi(\mathbf{x}), \Phi(\mathbf{x}'))$ is p.s.d.

Gaussian kernel

- The Gaussian kernel $\exp(-||\mathbf{x} - \mathbf{x}'||^2/2\sigma^2)$ might be the most used kernel in practice.
- It is not limited to Euclidean space. Consider that

$$\begin{aligned} ||\mathbf{x} - \mathbf{x}'||^2 &= (\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}') \\ &= \mathbf{x}^T \mathbf{x} + \mathbf{x}'^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}' \end{aligned}$$

then we can consider a nonlinear kernel where

$$||\mathbf{x} - \mathbf{x}'||^2 \longleftrightarrow k_1(\mathbf{x}, \mathbf{x}) + k_1(\mathbf{x}, \mathbf{x}') - 2k_1(\mathbf{x}, \mathbf{x}')$$

- We then consider the kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\sigma^2} (k_1(\mathbf{x}, \mathbf{x}) + k_1(\mathbf{x}, \mathbf{x}') - 2k_1(\mathbf{x}, \mathbf{x}'))\right)$$

- Any algorithm where a distance appears can be kernelized...

Kernels on graphs, sets, strings etc

- Over the past few years, there has been a lot of work on defining kernels between non-Euclidean objects.
- The aim is to come up with a p.s.d. kernel.
- It is not though because a kernel is p.s.d. that it is a 'good' measure of similarity.

Kernels derived from probabilistic models

- Generative models (eg HMMs) provide a way to deal with variable-dimension objects (eg strings of different lengths).
- We can then use these for discriminative learning by defining kernels.
- For example for a generative model $p(\mathbf{x})$, we could define

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

or

$$k(\mathbf{x}, \mathbf{x}') = \int p(\mathbf{x}|\theta) p(\mathbf{x}'|\theta) p(\theta) d\theta$$

- Consider a parametric generative model $p(\mathbf{x}|\theta)$.
- We introduce the kernel which uses a feature vector of size $|\theta|$

$$k(\mathbf{x}, \mathbf{x}') = g(\theta, \mathbf{x}) F^{-1} g(\theta, \mathbf{x}')$$

where

$$\begin{aligned} g(\theta, \mathbf{x}) &= \nabla_{\theta} \log p(\mathbf{x}|\theta) \\ F &= \mathbb{E}_{\mathbf{x}}[g(\theta, \mathbf{x})^T g(\theta, \mathbf{x}')] \end{aligned}$$

- F is the Fisher information matrix, the kernel is invariant to the parametrization of θ .

Gaussian Processes

- A stochastic process is a collection of RVs indexed by the input vector \mathbf{x} . A Gaussian Process is a stochastic process for which $(y(\mathbf{x}_1), \dots, y(\mathbf{x}_n))$ is jointly Gaussian for any $\{\mathbf{x}_n\}$.
- A GP can be characterized by its mean function $m(\mathbf{x})$ (often assumed 0) and its covariance function $k(\mathbf{x}, \mathbf{x}')$; i.e.

$$\mathbb{E}[y(\mathbf{x})] = m(\mathbf{x}), \quad \text{cov}[y(\mathbf{x}), y(\mathbf{x}')] = k(\mathbf{x}, \mathbf{x}')$$

- For any $\{\mathbf{x}_n\}$, we have

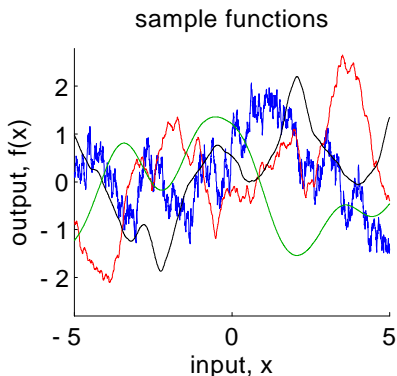
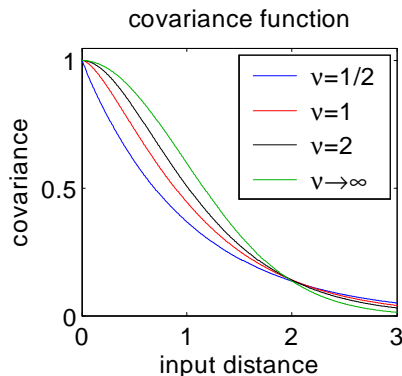
$$y(\mathbf{x}_{1:n}) \sim \mathcal{N}(m(\mathbf{x}_{1:n}), K(\mathbf{x}_{1:n}))$$

$$\text{where } y(\mathbf{x}_{1:n}) = (y(\mathbf{x}_1), \dots, y(\mathbf{x}_n))^T, \\ m(\mathbf{x}_{1:n}) = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_n))^T,$$

$$[K(\mathbf{x}_{1:n})]_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j).$$

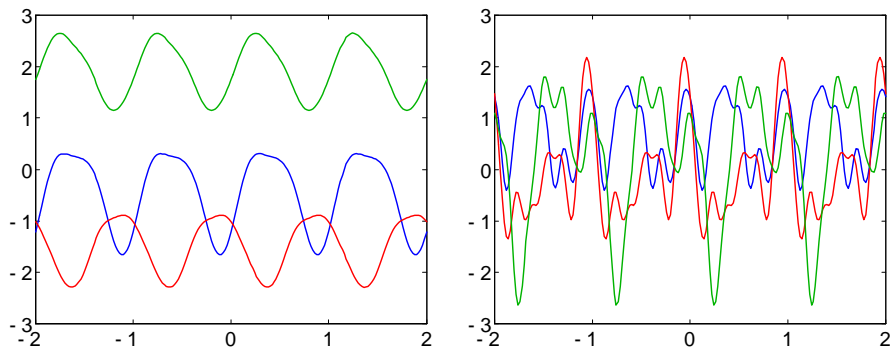
- A GP gives a prior on the space of functions.

Samples from the prior for Matern covariance



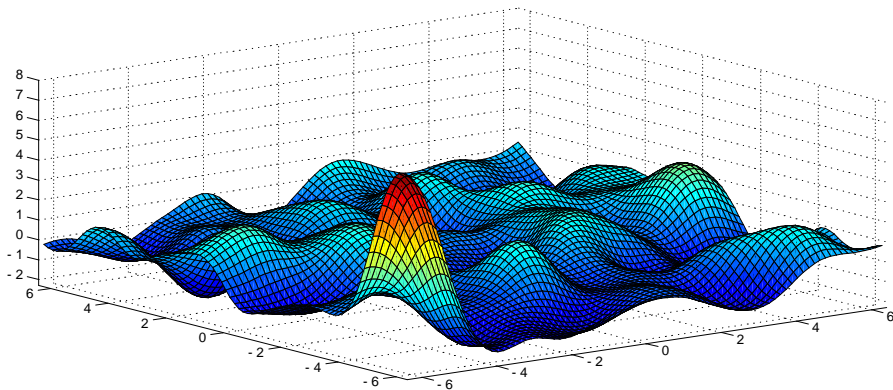
Covariance function $k_\nu(\mathbf{x}, \mathbf{x}') = k_\nu(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{l} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}r}{l} \right)$ for $\|\mathbf{x} - \mathbf{x}'\| = r$ (left) and sample paths (right)

Samples from the prior for a periodic covariance



Sample paths from the prior for $l > 1$ (left) and $l < 1$ (right) where
$$k_v(\mathbf{x}, \mathbf{x}') = \exp \left(-2 \sin^2 (\pi (\mathbf{x} - \mathbf{x}')) / l^2 \right)$$

Samples from the prior with a Gaussian covariance



Sample surface for $k(\mathbf{x}, \mathbf{x}') = \exp \left(-\nu^2 \|\mathbf{x} - \mathbf{x}'\|^2 \right)$

Bayesian linear regression & Gaussian Processes

- Consider the linear regression model where

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

and we set $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha^{-1}I)$.

- $y(\mathbf{x}, \mathbf{w})$ is a linear combination of Gaussians rvs so it is a GP with

$$\mathbb{E}[y(\mathbf{x}, \mathbf{w})] = \mathbb{E}[\mathbf{w}^T] \Phi(\mathbf{x}) = 0$$

and

$$\begin{aligned} \text{cov}[y(\mathbf{x}, \mathbf{w}), y(\mathbf{x}', \mathbf{w})] &= \Phi^T(\mathbf{x}) \mathbb{E}[\mathbf{w} \mathbf{w}^T] \Phi(\mathbf{x}) \\ &= \alpha^{-1} \Phi^T(\mathbf{x}) \Phi(\mathbf{x}'). \end{aligned}$$

- Instead of introducing a prior on $y(\mathbf{x})$ by defining a prior on \mathbf{w} and introducing a finite dimensional vector of features, we can directly introduce a GP prior on $y(\mathbf{x})$.

Bayesian regression with Gaussian Processes

- Consider the data $D = \{\mathbf{x}_n, t_n\}_{n=1}^N$ where

$$t_n = t(\mathbf{x}_n) = y(\mathbf{x}_n) + \varepsilon_n \text{ where } \varepsilon_n \sim \mathcal{N}(0, \sigma^2)$$

and

$$y(\mathbf{x}) \sim GP(m(\mathbf{x}) = 0, k(\mathbf{x}, \mathbf{x}'))$$

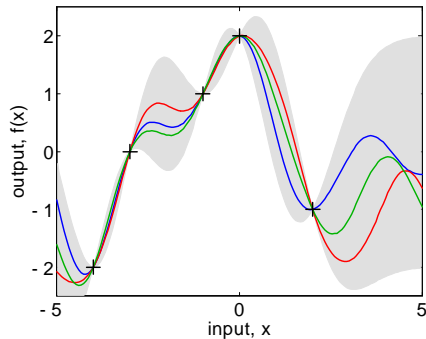
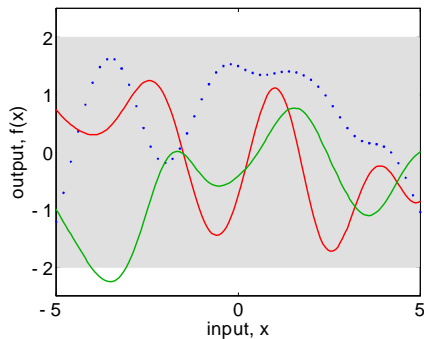
- We have

$$y(\mathbf{x}) | D \sim GP(m_{\text{post}}(\mathbf{x}), k_{\text{post}}(\mathbf{x}, \mathbf{x}'))$$

where

$$\begin{aligned} m_{\text{post}}(\mathbf{x}) &= k(\mathbf{x}, \mathbf{x}_{1:N}) [K(\mathbf{x}_{1:N}, \mathbf{x}_{1:N}) + \sigma^2 I]^{-1} \mathbf{t}_{1:N}, \\ k_{\text{post}}(\mathbf{x}, \mathbf{x}') &= k(\mathbf{x}, \mathbf{x}') - \\ &\quad k(\mathbf{x}, \mathbf{x}_{1:N}) [K(\mathbf{x}_{1:N}, \mathbf{x}_{1:N}) + \sigma^2 I]^{-1} k(\mathbf{x}_{1:N}, \mathbf{x}') \end{aligned}$$

From the prior to the posterior



Random draws from the prior (left) and the posterior (right): The shaded area represents the pointwise mean \pm twice the standard deviation.

Predictive distribution and Interpretation

- Given \mathbf{x}^* , we have

$$p(t^* | \mathbf{x}_{1:N}, \mathbf{t}_{1:N}, \mathbf{x}^*) = \mathcal{N}(t^*; \mu(\mathbf{x}^*), \sigma^2(\mathbf{x}^*))$$

where

$$\begin{aligned}\mu(\mathbf{x}^*) &= k(\mathbf{x}^*, \mathbf{x}_{1:N}) [K(\mathbf{x}_{1:N}, \mathbf{x}_{1:N}) + \sigma^2 I]^{-1} \mathbf{t}_{1:N}, \\ \sigma^2(\mathbf{x}^*) &= k(\mathbf{x}^*, \mathbf{x}^*) + \sigma^2 \\ &\quad - k(\mathbf{x}^*, \mathbf{x}_{1:N}) [K(\mathbf{x}_{1:N}, \mathbf{x}_{1:N}) + \sigma^2 I]^{-1} k(\mathbf{x}_{1:N}, \mathbf{x}^*)\end{aligned}$$

- The mean $\mu(\mathbf{x}^*)$ is linear in two ways

$$\mu(\mathbf{x}^*) = \sum_{i=1}^n a_i t_i = \sum_{i=1}^n b_i K(\mathbf{x}^*, \mathbf{x}_i)$$

- The variance is of the form

$$\sigma^2(\mathbf{x}^*) = \text{prior variance} - \text{positive terms dependent on } \mathbf{x}_{1:N}$$

- Remark:* the variance is independent of the observations $\mathbf{t}_{1:N}$.

- The central computation operation in using GP involves inverting a $N \times N$ matrix. Standard methods requires $O(N^3)$ operations.
- In the finite basis function model with M basis, we have to invert a $M \times M$ matrix.
- So if the number M of basis functions is smaller than N then we are better off with the standard method.
- If the kernel considered corresponds to an infinite M , we do not have the choice!
- Several techniques have been developed to perform approximate inference.

Learning the hyperparameters

- In practice, we often parametrize the kernel by some parameters θ .
- To estimate θ , we can maximize the marginal log-likelihood

$$\log p(\mathbf{t}_{1:N} | \theta, \mathbf{x}_{1:N}) = -\frac{1}{2} \log |K_N^\theta| - \frac{1}{2} \mathbf{t}_{1:N}^\top [K_N^\theta]^{-1} \mathbf{t}_{1:N} - \frac{N}{2} \log 2\pi$$

using $[K_N^\theta]_{i,j} = K^\theta(\mathbf{x}_i, \mathbf{x}_j)$.

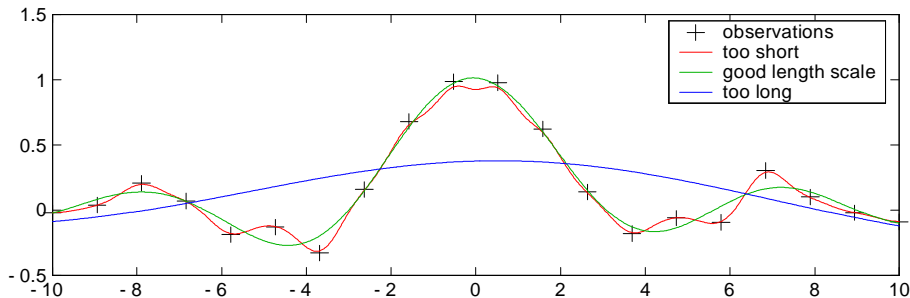
- The gradient of the log-likelihood is given by

$$\begin{aligned} \frac{\partial \log p(\mathbf{t}_{1:N} | \theta, \mathbf{x}_{1:N})}{\partial \theta_i} &= -\frac{1}{2} \text{Tr} \left([K_N^\theta]^{-1} \frac{\partial K_N^\theta}{\partial \theta_i} \right) \\ &\quad + \frac{1}{2} \mathbf{t}_{1:N}^\top [K_N^\theta]^{-1} \frac{\partial K_N^\theta}{\partial \theta_i} [K_N^\theta]^{-1} \mathbf{t}_{1:N} \end{aligned}$$

- The log-likelihood is typically not concave in θ .

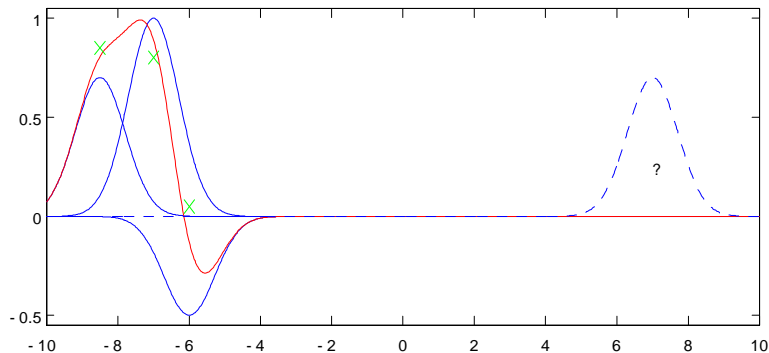
Example: Fitting the length scale parameter

- Parameterized covariance function: $k(\mathbf{x}, \mathbf{x}') = \nu \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{l}\right)$.



- The mean posterior predictive distribution is plotted for 3 different length scales (the green curve corresponds to optimizing the likelihood). Note that we can get an almost perfect fit for a small length scale but the marginal likelihood does not favour it.

Using a finite number of basis functions can be dangerous



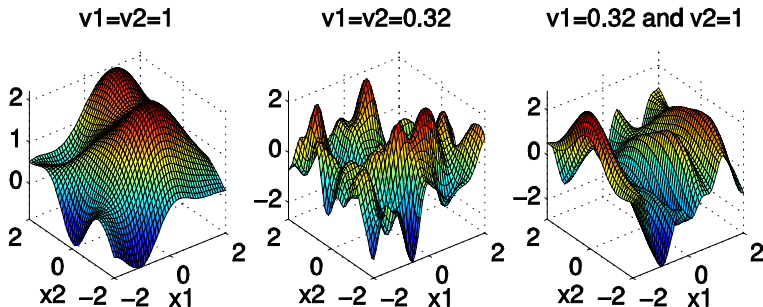
Automatic Relevance Determination

- We can extend the technique described before to select automatically the relevant input variables; i.e. say

$$k(\mathbf{x}, \mathbf{x}') = v_0^2 \exp \left(- \frac{\sum_{i=1}^D (x_i - x'_i)^2}{2v_i^2} \right)$$

where $\theta = (v_0^2, v_1^2, \dots, v_D^2)$.

- We have



Gaussian Processes for Binary Classification

- The input is given by \mathbf{x} and the output $t \in \{0, 1\}$ with

$$\Pr(t = 1 | \mathbf{x}) = g(a(\mathbf{x})).$$

- We model $a(\mathbf{x})$ through a Gaussian process define by

$$\mathbb{E}[a(\mathbf{x})] = 0 \text{ and } \text{cov}[a(\mathbf{x}) a(\mathbf{x}')] = k(\mathbf{x}, \mathbf{x}') = m(\mathbf{x}, \mathbf{x}') + \nu \delta(\mathbf{x} - \mathbf{x}')$$

- We are interested in computing

$$\begin{aligned} p(t^* | \mathbf{x}_{1:N}, \mathbf{t}_{1:N}, \mathbf{x}^*) &= \int p(t^* | a(\mathbf{x}^*)) p(a(\mathbf{x}^*) | \mathbf{t}_{1:N}) da(\mathbf{x}^*) \\ &= \int g(a(\mathbf{x}^*)) p(a(\mathbf{x}^*) | \mathbf{t}_{1:N}) da(\mathbf{x}^*) \end{aligned}$$

Laplace Approximation

- We have

$$\begin{aligned} p(a(\mathbf{x}^*) | \mathbf{t}_{1:N}) &= \int p(a(\mathbf{x}^*), a(\mathbf{x}_{1:N}) | \mathbf{t}_{1:N}) da(\mathbf{x}_{1:N}) \\ &= \int p(a(\mathbf{x}^*) | a(\mathbf{x}_{1:N})) p(a(\mathbf{x}_{1:N}) | \mathbf{t}_{1:N}) da(\mathbf{x}_{1:N}) \end{aligned}$$

- We have

$$\begin{aligned} p(a(\mathbf{x}^*) | a(\mathbf{x}_{1:N})) &= \mathcal{N}(a(\mathbf{x}^*); k^\top(\mathbf{x}^*, \mathbf{x}_{1:N}) K_N^{-1} a(\mathbf{x}_{1:N}), \\ &\quad k(\mathbf{x}, \mathbf{x}) - k^\top(\mathbf{x}^*, \mathbf{x}_{1:N}) K_N^{-1} k(\mathbf{x}^*, \mathbf{x}_{1:N})) \end{aligned}$$

- We make a Gaussian approximation of $p(a(\mathbf{x}_{1:N}) | \mathbf{t}_{1:N})$ using Laplace.

- The unnormalized posterior is given by

$$\begin{aligned}
 & \log p(a(\mathbf{x}_{1:N}), \mathbf{t}_{1:N}) \\
 = & \log p(a(\mathbf{x}_{1:N}), \mathbf{t}_{1:N}) + \log p(\mathbf{t}_{1:N} | a(\mathbf{x}_{1:N})) \\
 = & -\frac{1}{2} \mathbf{a}^\top(\mathbf{x}_{1:N}) K_N^{-1} \mathbf{a}(\mathbf{x}_{1:N}) - \frac{N}{2} \log(2\pi) - \frac{1}{2} \log |K_N| \\
 & + \mathbf{t}_{1:N}^\top \mathbf{a}(\mathbf{x}_{1:N}) - \sum_{n=1}^N \log(1 + \exp a(\mathbf{x}_N)) + cst
 \end{aligned}$$

as $g(a)^t (1 - g(a))^{1-t} = \exp(at) g(-a)$

- We perform a Taylor expansion of the $\log p(a(\mathbf{x}_{1:N}), \mathbf{t}_{1:N})$ around its mode which can be computed using a Newton-Raphson method where

$$\nabla \log p(a(\mathbf{x}_{1:N}), \mathbf{t}_{1:N}) = \mathbf{t}_{1:N} - \sigma_{1:N} - K_N^{-1} \mathbf{a}(\mathbf{x}_{1:N})$$

and

$$\nabla \nabla \log p(a(\mathbf{x}_{1:N}), \mathbf{t}_{1:N}) = -W_N - K_N^{-1}$$

where $W_N = \text{diag}(g(a(\mathbf{x}_N))(1 - g(a(\mathbf{x}_N))))$.

- The Newton-Raphson formula takes the form

$$a^{(k+1)}(\mathbf{x}_{1:N}) = K_N (I + W_N K_N)^{-1} \{\mathbf{t}_{1:N} - \sigma_{1:N} + W_N a(\mathbf{x}_{1:N})\}$$

- Once the mode $a^*(\mathbf{x}_{1:N})$ has been found, we compute the associated

$$H = -\nabla \nabla \log p(a(\mathbf{x}_{1:N}), \mathbf{t}_{1:N}) = W_N + K_N^{-1}$$

- The Gaussian approximation is given by

$$q(a(\mathbf{x}_{1:N})) = \mathcal{N}(a(\mathbf{x}_{1:N}); a^*(\mathbf{x}_{1:N}), H)$$

- It follows that we obtain a Gaussian approximation of $p(a(\mathbf{x}^*) | \mathbf{t}_{1:N})$ with

$$\begin{aligned} \mathbb{E}(a(\mathbf{x}^*) | \mathbf{t}_{1:N}) &= k(\mathbf{x}^*, \mathbf{x}_{1:N}) (\mathbf{t}_{1:N} - \boldsymbol{\mu}_{1:N}), \\ \mathbb{V}(a(\mathbf{x}^*) | \mathbf{t}_{1:N}) &= k(\mathbf{x}^*, \mathbf{x}^*) \\ &\quad - k^\top(\mathbf{x}^*, \mathbf{x}_{1:N}) (W_N^{-1} + K_N)^{-1} k(\mathbf{x}^*, \mathbf{x}_{1:N}) \end{aligned}$$

- Now we finally use the approximation combining logistic and Gaussian

$$p(t^* | \mathbf{x}_{1:N}, \mathbf{t}_{1:N}, \mathbf{x}^*) = \int g(a(\mathbf{x}^*)) p(a(\mathbf{x}^*) | \mathbf{t}_{1:N}) da(\mathbf{x}^*)$$

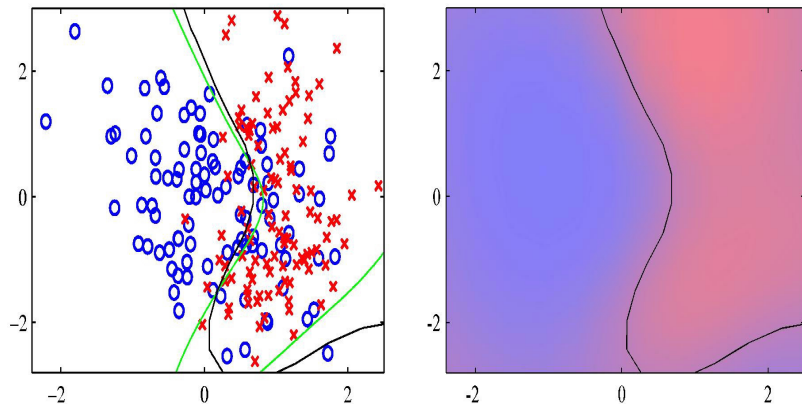
which states that

$$\int g(a) \mathcal{N}(a; \mu, \sigma^2) da \simeq g\left(\frac{\mu}{\sqrt{1 + \pi\sigma^2/8}}\right)$$

- The Laplace approximation also yields an approximation of the log-marginal likelihood

$$\log p(\mathbf{t}_{1:N}) \simeq \log p(a^*(\mathbf{x}_{1:N}), \mathbf{t}_{1:N}) - \frac{1}{2} |H| + \frac{N}{2} \log(2\pi)$$

Example of Binary Classification using GP



Left: Optimal decision boundary (green) and GP classifier (black).
Right: predicted posterior proba for the blue and red classes