CS 540: Machine Learning Lecture 8: Kernel Methods

AD

February 2008

AD () February 2008 1 / 34

Feature space

- 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 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

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

• This is a symetric function of its arguments

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

February 2008

Kernels

- 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. $k(\mathbf{x}, \mathbf{x}') = \Phi^{\mathsf{T}}(\mathbf{x}) \Phi(\mathbf{x}')$.
- For example assume $\mathbf{x} = (x_1, x_2)$ and

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^{\mathsf{T}} \mathbf{x}')^{2}$$

$$= (x_{1}x'_{1} + x_{2}x'_{2})^{2}$$

$$= x_{1}^{2} (x'_{1})^{2} + x_{2}^{2} (x'_{2})^{2} + 2x_{1}x'_{1}x_{2}x'_{2}$$

$$= (x_{1}^{2}, \sqrt{2}x_{1}x_{2}, x_{2}^{2}) ((x'_{2})^{2}, \sqrt{2}x'_{1}x'_{2}, (x'_{2})^{2})$$

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

where

$$\Phi\left(\mathbf{x}\right) = \left(x_1^2, \sqrt{2}x_1x_2, x_2^2\right).$$

AD () February 2008

Positive Semi-definite Kernels

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

$$\sum_{i}\sum_{j}\alpha_{i}\alpha_{j}k\left(\mathbf{x}_{i},\mathbf{x}_{i}\right)\geq0$$

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

$$k\left(\mathbf{x},\mathbf{x}'\right) = \sum_{i=1}^{\infty} \lambda_{i} \Phi_{i}\left(\mathbf{x}\right) \Phi_{i}\left(\mathbf{x}'\right).$$

More later...

AD () February 2008

Kernel trick

- In many situations, as mentioned earlier, we actually only use $\Phi(\mathbf{x})$ through $\Phi^{\mathsf{T}}(\mathbf{x}) \Phi(\mathbf{x}')$.
- ullet Moreover it is often very difficult to design good features $\Phi\left(\mathbf{x}\right)$.
- 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 $\Phi(\mathbf{x})$.

AD () February 2008 5 / 3-

Dual Representation of Linear Regression

Consider

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

where $\lambda > 0$.

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

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

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

$$\boldsymbol{\Phi} = \left(\begin{array}{c} \boldsymbol{\Phi}^\mathsf{T} \left(\boldsymbol{\mathsf{x}}_1 \right) \\ \vdots \\ \boldsymbol{\Phi}^\mathsf{T} \left(\boldsymbol{\mathsf{x}}_{\textit{N}} \right) \end{array} \right)$$

AD () February 2008

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

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

where $K = \Phi \Phi^{\mathsf{T}}$.

K is the Gram matrix

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

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

AD () February 2008 7 / 34

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

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

It follows that

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

where

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

- 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!

AD () February 2008

Constructing kernels

- 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. $\int k(\mathbf{x}, \mathbf{x}') \Phi_i(\mathbf{x}) d\mathbf{x} = \lambda_i \Phi_i(\mathbf{x}).$

AD () February 2008 9 / 34

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

AD () February 2008 10 / 34

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.

AD () February 2008 11 / 34

Gaussian kernel

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

$$||\mathbf{x} - \mathbf{x}'||^2 = (\mathbf{x} - \mathbf{x}')^{\mathsf{T}} (\mathbf{x} - \mathbf{x}')$$
$$= \mathbf{x}^{\mathsf{T}} \mathbf{x} + \mathbf{x}'^{\mathsf{T}} \mathbf{x}' - 2\mathbf{x}^{\mathsf{T}} \mathbf{x}'$$

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

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

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

AD () February 2008

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.

AD () February 2008 13 / 34

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

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

February 2008

Fisher Kernel

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

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

where

$$g(\theta, \mathbf{x}) = \nabla_{\theta} \log p(\mathbf{x}|\theta)$$

$$F = \mathbb{E}_{\mathbf{x}}[g(\theta, \mathbf{x})^{\mathsf{T}} g(\theta, \mathbf{x}')]$$

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

AD () February 2008 15 / 34

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 $(y(\mathbf{x}_1), \ldots, y(\mathbf{x}_n))$ 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.

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

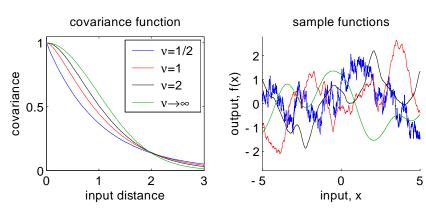
• For any $\{x_n\}$, we have

$$y(\mathbf{x}_{1:n}) \sim \mathcal{N}\left(m(\mathbf{x}_{1:n}), K(\mathbf{x}_{1:n})\right)$$
 where $y(\mathbf{x}_{1:n}) = (y(\mathbf{x}_1), \dots, y(\mathbf{x}_n))^\mathsf{T}$,
$$m(\mathbf{x}_{1:n}) = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_n))^\mathsf{T}$$
,
$$\left[K(\mathbf{x}_{1:n})\right]_{i,i} = k(\mathbf{x}_i, \mathbf{x}_i).$$

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

AD () February 2008

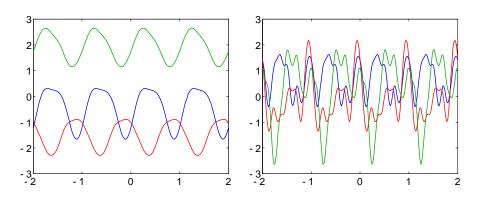
Samples from the prior for Matern covariance



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

AD () February 2008 17 / 34

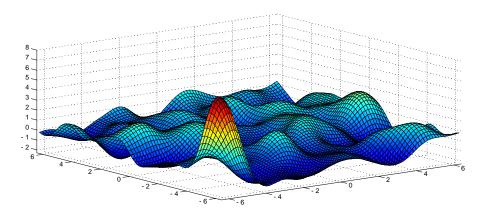
Samples from the prior for a periodic covariance



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

AD () February 2008 18 / 34

Samples from the prior with a Gaussian covariance



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

AD () February 2008 19 /

Bayesian linear regression & Gaussian Processes

• Consider the linear regression model where

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

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

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

and

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

• 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 $y(\mathbf{x})$.

AD () February 2008

Bayesian regression with Gaussian Processes

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

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

and

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

We have

$$y\left(\mathbf{x}\right)|D\sim\textit{GP}\left(m_{\mathsf{post}}\left(\mathbf{x}\right),k_{\mathsf{post}}\left(\mathbf{x},\mathbf{x}'\right)\right)$$

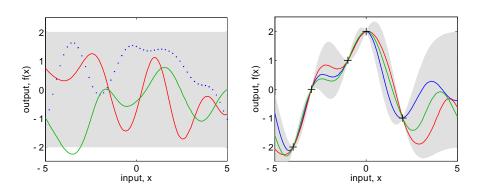
where

$$m_{\text{post}}(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}_{1:N}) \left[K(\mathbf{x}_{1:N}, \mathbf{x}_{1:N}) + \sigma^2 I \right]^{-1} \mathbf{t}_{1:N},$$

$$k_{\text{post}}(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}') - k(\mathbf{x}, \mathbf{x}_{1:N}) \left[K(\mathbf{x}_{1:N}, \mathbf{x}_{1:N}) + \sigma^2 I \right]^{-1} k(\mathbf{x}_{1:N}, \mathbf{x}')$$

AD () February 2008

From the prior to the posterior



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

AD () February 2008 22 / 34

Predictive distribution and Interpretation

Given x*, we have

$$ho\left(\left.t^{*}\right|\mathbf{x}_{1:N},\mathbf{t}_{1:N},\mathbf{x}^{*}
ight)=\mathcal{N}\left(t^{*};\mu\left(\mathbf{x}^{*}
ight),\sigma^{2}\left(\mathbf{x}^{*}
ight)
ight)$$

where

$$\mu\left(\mathbf{x}^{*}\right) = k\left(\mathbf{x}^{*}, \mathbf{x}_{1:N}\right) \left[K\left(\mathbf{x}_{1:N}, \mathbf{x}_{1:N}\right) + \sigma^{2}I\right]^{-1} \mathbf{t}_{1:N},$$

$$\sigma^{2}\left(\mathbf{x}^{*}\right) = k\left(\mathbf{x}^{*}, \mathbf{x}^{*}\right) + \sigma^{2}$$

$$-k\left(\mathbf{x}^{*}, \mathbf{x}_{1:N}\right) \left[K\left(\mathbf{x}_{1:N}, \mathbf{x}_{1:N}\right) + \sigma^{2}I\right]^{-1} k\left(\mathbf{x}_{1:N}, \mathbf{x}^{*}\right)$$

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

$$\mu\left(\mathbf{x}^{*}\right) = \sum_{i=1}^{n} a_{i} t_{i} = \sum_{i=1}^{n} b_{i} K\left(\mathbf{x}^{*}, \mathbf{x}_{i}\right)$$

The variance is of the form

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

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

AD () February 2008

Computational Complexity

- 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.

AD () February 2008 24 / 34

Learning the hyperparameters

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

$$\log p\left(\mathbf{t}_{1:N} \middle| \theta, \mathbf{x}_{1:N}\right) = -\frac{1}{2} \log \left| \mathcal{K}_{N}^{\theta} \middle| -\frac{1}{2} \mathbf{t}_{1:N}^{\mathsf{T}} \left[\mathcal{K}_{N}^{\theta} \right]^{-1} \mathbf{t}_{1:N} - \frac{N}{2} \log 2\pi$$

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

The gradient of the log-likelihood is given by

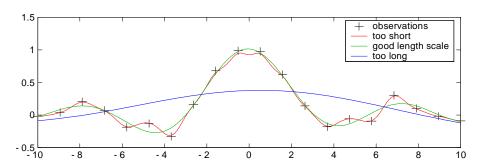
$$\frac{\partial \log p\left(\mathbf{t}_{1:N} \middle| \theta, \mathbf{x}_{1:N}\right)}{\partial \theta_{i}} = -\frac{1}{2} \operatorname{Tr} \left(\left[K_{N}^{\theta} \right]^{-1} \frac{\partial K_{N}^{\theta}}{\partial \theta_{i}} \right) \\
+ \frac{1}{2} \mathbf{t}_{1:N}^{\mathsf{T}} \left[K_{N}^{\theta} \right]^{-1} \frac{\partial K_{N}^{\theta}}{\partial \theta_{i}} \left[K_{N}^{\theta} \right]^{-1} \mathbf{t}_{1:N}$$

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

AD () February 2008

Example: Fitting the length scale parameter

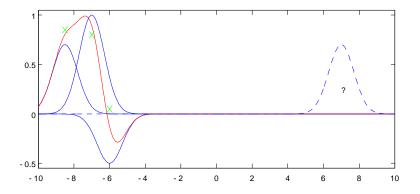
• Parameterized covariance function: $k(\mathbf{x}, \mathbf{x}') = \nu \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{I}\right)$.



• 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.

AD () February 2008

Using a finite number of basis functions can be dangerous



AD () February 2008 27 / 34

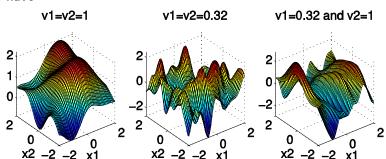
Automatic Relevance Determination

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

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

where
$$heta=\left(
u_{0}^{2},
u_{1}^{2},...,
u_{D}^{2}
ight)$$
 .

We have



AD () February 2008

Gaussian Processes for Binary Classification

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

$$\Pr\left(\left.t=1\right|\mathbf{x}\right)=\sigma\left(a\left(\mathbf{x}\right)\right)$$
.

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

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

We are interested in computing

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

February 2008

Laplace Approximation

We have

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

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

We have

$$\begin{array}{lcl} p\left(\left. a\left(\mathbf{x}^{*} \right) \right| \left. a\left(\mathbf{x}_{1:N} \right) \right) &=& \mathcal{N}(a\left(\mathbf{x}^{*} \right); k^{\mathsf{T}}\left(\mathbf{x}^{*}, \mathbf{x}_{1:N} \right) K_{N}^{-1} a\left(\mathbf{x}_{1:N} \right), \\ && k\left(\mathbf{x}, \mathbf{x} \right) - k^{\mathsf{T}}\left(\mathbf{x}^{*}, \mathbf{x}_{1:N} \right) K_{N}^{-1} k\left(\mathbf{x}^{*}, \mathbf{x}_{1:N} \right) \end{array}$$

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

AD () February 2008

• The unnormalized posterior is given by

$$\begin{split} &\log p\left(a\left(\mathbf{x}_{1:N}\right),\mathbf{t}_{1:N}\right) \\ &= &\log p\left(a\left(\mathbf{x}_{1:N}\right),\mathbf{t}_{1:N}\right) + \log p\left(\mathbf{t}_{1:N} \middle| a\left(\mathbf{x}_{1:N}\right)\right) \\ &= & -\frac{1}{2}a^{\mathsf{T}}\left(\mathbf{x}_{1:N}\right)K_{N}^{-1}a\left(\mathbf{x}_{1:N}\right) - \frac{N}{2}\log\left(2\pi\right) - \frac{1}{2}\log\left|K_{N}\right| \\ &+ \mathbf{t}_{1:N}^{\mathsf{T}}a\left(\mathbf{x}_{1:N}\right) - \sum_{n=1}^{N}\log\left(1 + \exp a\left(\mathbf{x}_{N}\right)\right) + cst \end{split}$$

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

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

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

and

$$abla
abla \log p\left({a\left({{f x}_{1:N}}
ight),{f t}_{1:N}}
ight) = - {\it W_N} - {\it K_N^{ - 1}}$$

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

AD () February 2008

• The Newton-Raphson formula takes the form

$$\mathbf{a}^{(k+1)}\left(\mathbf{x}_{1:N}\right) = K_{N}\left(I + W_{N}K_{N}\right)^{-1}\left\{\mathbf{t}_{1:N} - \sigma_{1:N} + W_{N}\mathbf{a}\left(\mathbf{x}_{1:N}\right)\right\}$$

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

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

• The Gaussian approximation is given by

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

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

$$\begin{split} \mathbb{E}\left(\left. a\left(\mathbf{x}^{*}\right) \right| \mathbf{t}_{1:N}\right) &= k\left(\mathbf{x}^{*}, \mathbf{x}_{1:N}\right) \left(\mathbf{t}_{1:N} - \sigma_{1:N}\right), \\ \mathbb{V}\left(\left. a\left(\mathbf{x}^{*}\right) \right| \mathbf{t}_{1:N}\right) &= k\left(\mathbf{x}^{*}, \mathbf{x}^{*}\right) \\ &- k^{\mathsf{T}}\left(\mathbf{x}^{*}, \mathbf{x}_{1:N}\right) \left(W_{N}^{-1} + K_{N}\right)^{-1} k\left(\mathbf{x}^{*}, \mathbf{x}_{1:N}\right) \end{split}$$

AD () February 2008 32 / 34

Now we finally use the approximation combining logistic and Gaussian

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

which states that

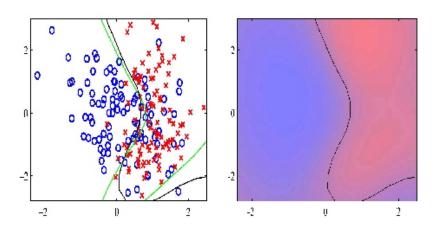
$$\int \sigma\left(\mathbf{a}\right) \mathcal{N}\left(\mathbf{a}; \mu, \sigma^{2}\right) d\mathbf{a} \simeq \sigma\left(\frac{\mu}{\sqrt{1 + \pi\sigma^{2}/8}}\right)$$

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

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

AD () February 2008 33 / 34

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

AD () February 2008