



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

GP applications

Ref

Functional modelling with Gaussian Processes

Lehel Csató

Faculty of Mathematics and Informatics
Babeş-Bolyai University, Cluj-Napoca,

Matematika és Informatika Tanszék
Babeş-Bolyai Tudományegyetem, Kolozsvár

September 2008



Table of Contents

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

GP applications

Ref

- 1 Modelling
- 2 Non-parametric models
- 3 Support Vectors Machines
- 4 Gaussian Processes
- 5 Models using Gaussian Processes
- 6 References



Modelling

Functional Modelling

Lehel Csató

Modelling

Machine Learning

Latent variable models

Estimation methods

Summary

Nonparametrics

SVM

GP's

GP applications

Ref

- 1 Modelling
 - Machine Learning
 - Latent variable models
 - Estimation methods
 - Drawbacks of parametric models
- 2 Non-parametric models
- 3 Support Vectors Machines
- 4 Gaussian Processes
- 5 Models using Gaussian Processes
- 6 References



Motivation

Functional Modelling

Lehel Csató

Modelling

Machine Learning

Latent variable models

Estimation methods

Summary

Nonparametrics

SVM

GP's

GP applications

Ref

D. Donoho

Data, Data, Data!

Challenges and opportunities of the coming **Data Deluge**

Several data types:

- classification problem – needed in decision systems: frequently there are data of very high dimension and several hundred classes to take into account;
- regression/prediction problems.



Machine learning

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

Historical background / Motivation:

- Huge amount of **data**, that should **automatically** be processed,
- Mathematics provides general solutions, solutions are i.e. **not for a given problem**,
- Need for “science”, that uses mathematics machinery for solving **practical** problems.



Definition of Machine Learning

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

Machine learning

Collection of methods (from statistics, probability theory) to solve problems **met in practice**.

- noise filtering for
 - non-linear regression and/or
 - non-Gaussian noise
- Classification:
 - binary,
 - multiclass,
 - partially labeled
- Clustering,
- Inversion problems,
- density estimation, novelty detection.

Generally, we need to **model the data**,



Modelling Data

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

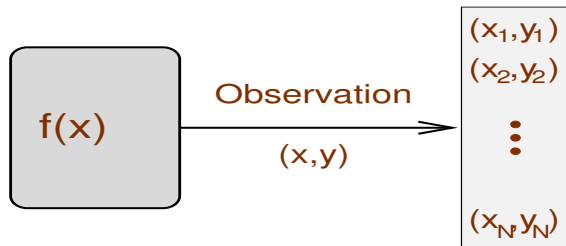
Nonparametrics

SVM

GP's

GP applications

Ref



- Real world: there “**is**” a function $y = f(x)$
- Observation process: a **corrupted** datum is collected for a sample x_n :

$$t_n = y_n + \epsilon \quad \text{additive noise}$$

$$t_n = h(y_n, \epsilon) \quad h \text{ distortion function}$$

- **Problem: find function** $y = f(x)$



Latent variable models

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

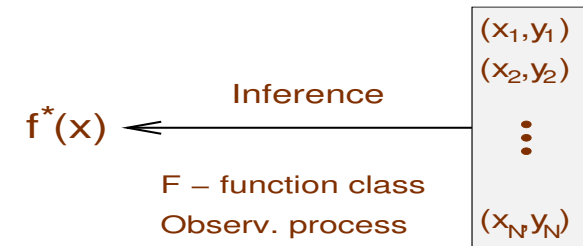
Nonparametrics

SVM

GP's

GP applications

Ref



- **Data set** – collected.
- Assume a function class.
 - polynomial,
 - Fourier expansion,
 - Wavelet;
- Observation process – **encodes** the noise;
- Find the optimal function from the class.



Latent variable models II

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

- We have the **data set** $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$.

- Consider a function class:

$$(1) \quad \mathcal{F} = \{\mathbf{w}^T \mathbf{x} + b \mid \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}\}$$

$$(2) \quad \mathcal{F} = \left\{ a_0 + \sum_{k=1}^K a_k \sin(2\pi kx) + \sum_{k=1}^K b_k \cos(2\pi kx) \mid \mathbf{a}, \mathbf{b} \in \mathbb{R}^K, a_0 \in \mathbb{R} \right\}$$

- Assume an observation process:

$$y_n = f(\mathbf{x}_n) + \epsilon \quad \text{with } \epsilon \sim N(0, \sigma^2).$$



Latent variable models III

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

- 1 The **data set**: $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$.

- 2 Assume a function class:

$$\mathcal{F} = \{f(\mathbf{x}, \boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \mathbb{R}^p\}$$

\mathcal{F} – polynomial, etc.

- 3 Assume an observation process. Define a **loss function**:

$$L(y_n, f(\mathbf{x}_n, \boldsymbol{\theta}))$$

For the Gaussian noise:

$$L(y_n, f(\mathbf{x}_n, \boldsymbol{\theta})) = (y_n - f(\mathbf{x}_n, \boldsymbol{\theta}))^2.$$



Parameter estimation

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

Estimating parameters:

Finding the **optimal value to $\boldsymbol{\theta}$** :

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \Omega} L(\mathcal{D}, \boldsymbol{\theta})$$

where

- Ω is the domain of the parameters.
- $L(\mathcal{D}, \boldsymbol{\theta})$ is a “loss function” for the data set.

Example:

$$L(\mathcal{D}, \boldsymbol{\theta}) = \sum_{n=1}^N L(y_n, f(\mathbf{x}_n, \boldsymbol{\theta}))$$



Parameter estimation – M.L.

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

$L(\mathcal{D}, \boldsymbol{\theta})$ – (log)likelihood function.

Maximum likelihood estimation of the model:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \Omega} L(\mathcal{D}, \boldsymbol{\theta})$$

Example – quadratic regression:

$$L(\mathcal{D}, \boldsymbol{\theta}) = \sum_{n=1}^N (y_n - f(\mathbf{x}_n, \boldsymbol{\theta}))^2 \quad \text{– factorisation}$$

Drawback: produces **perfect** fit to the data – **over-fitting**.



Maximum Likelihood – Over-fitting

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

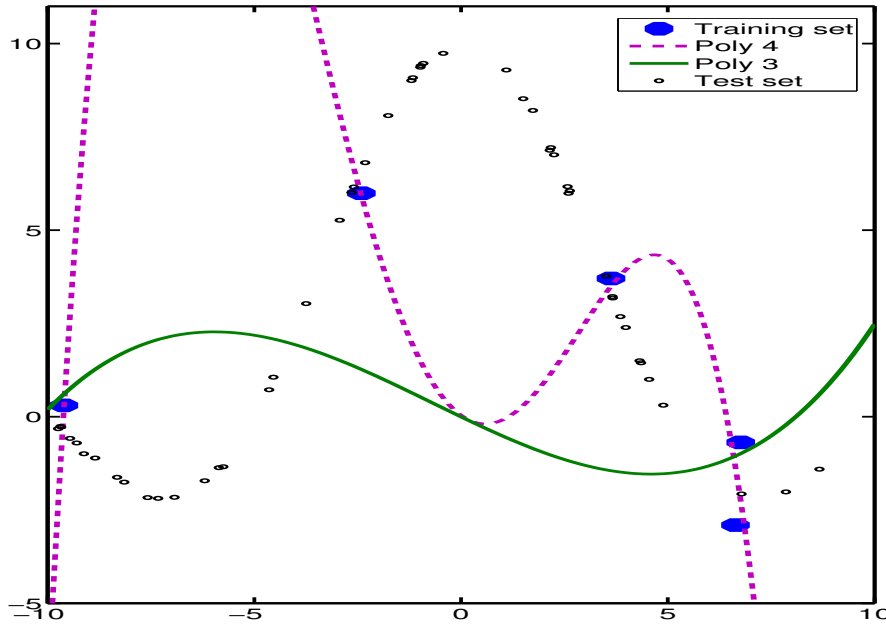
Nonparametrics

SVM

GP's

GP applications

Ref



Maximum a-posteriori estimation

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

M.A.P. – assigning **probabilities** to the

- Data \mathcal{D} : log-likelihood function: the probability of the data drawn using θ

$$P(y_n | \mathbf{x}_n, \theta, \mathcal{F}) \propto \exp[-L(y_n, f(\mathbf{x}_n, \theta))]$$

\propto – a normalisation constant missing.

- Parameters θ : probability that θ could have had a given value

$$p_0(\theta) \propto \exp\left[-\frac{\|\theta\|^2}{2}\right]$$

prior probability.



M.A.P. estimation II

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

Combining prior with observation – likelihood – using Bayes' rule:

A-posteriori probability of the parameters:

$$p(\theta | \mathcal{D}, \mathcal{F}) = \frac{P(\mathcal{D} | \theta) p_0(\theta)}{p(\mathcal{D} | \mathcal{F})}$$

$$p(\mathcal{D} | \mathcal{F}) = \int_{\Omega} d\theta P(\mathcal{D} | \theta) p_0(\theta)$$

$p(\mathcal{D} | \mathcal{F})$ – probability of the **data set** under a given family of models.



M.A.P. estimation III

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

M.A.P. estimation – aims at finding θ with the largest probability:

$$\theta_{MAP}^* = \arg \max_{\theta \in \Omega} p(\theta | \mathcal{D}, \mathcal{F})$$

Example:

Using $L(y_n, f(\mathbf{x}_n, \theta))$ in defining the likelihood and Gaussian prior:

$$\theta_{MAP}^* = \arg \max_{\theta \in \Omega} K - \sum_n L(y_n, f(\mathbf{x}_n, \theta)) - \frac{\|\theta\|^2}{2\sigma_0^2}$$

For $\sigma_0^2 \rightarrow \infty$ we have **maximum likelihood**

after a change of sign and $\max \rightarrow \min$



M.A.P. – Example

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

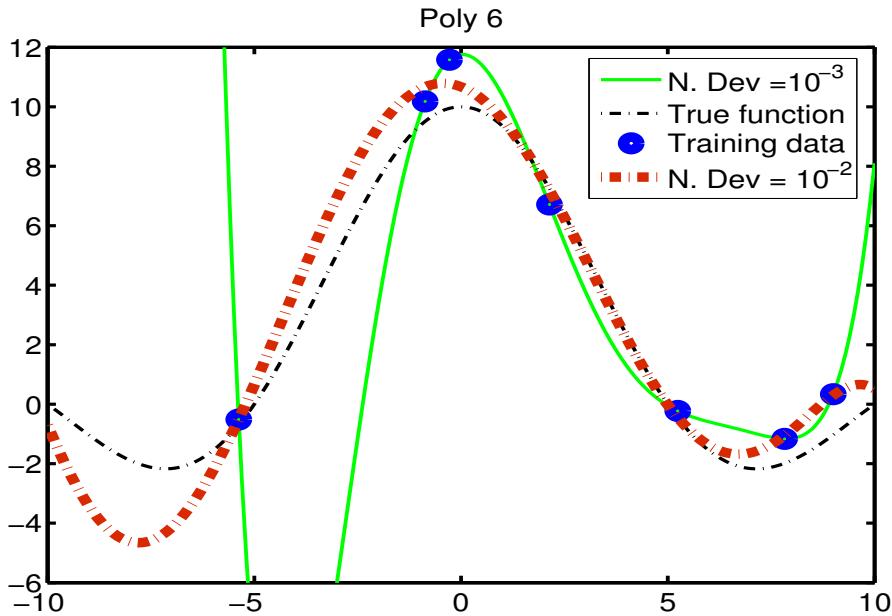
Nonparametrics

SVM

GP's

GP applications

Ref



Parameter estimation – Bayes

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

We use Bayes' rule:

$$p(\theta|\mathcal{D}, \mathcal{F}) = \frac{P(\mathcal{D}|\theta)p_0(\theta)}{p(\mathcal{D}|\mathcal{F})}$$

$$p(\mathcal{D}|\mathcal{F}) = \int_{\Omega} d\theta P(\mathcal{D}|\theta)p_0(\theta)$$

and **try to store** the whole distribution of the possible values.

We operate therefore with

$$p_{\text{post}}(\theta|\mathcal{D}, \mathcal{F})$$



Bayes estimation II

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

When computing $p_{\text{post}}(\theta|\mathcal{D}, \mathcal{F})$ we assumed that the posterior **can be represented** analytically.

This is not the case.

Approximations are needed for the

- posterior distribution
- predictive distribution

In Bayesian modelling an important issue is **how** we approximate the posterior distribution.



Bayes – Example

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

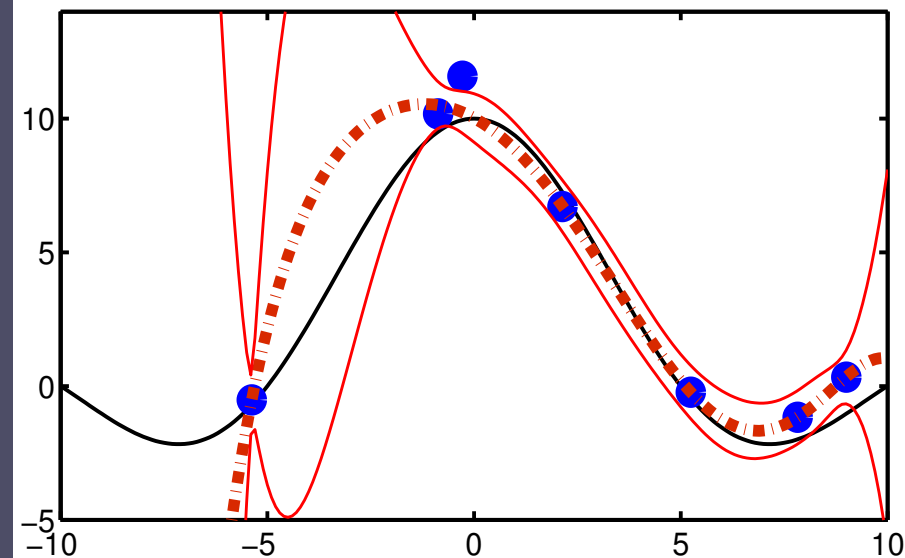
SVM

GP's

GP applications

Ref

Pol. 6 – N.var $\sigma^2 = 1$





Graphical Models

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

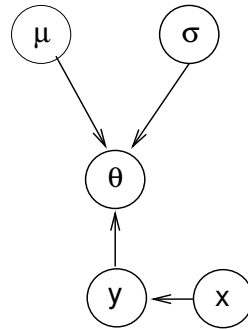
SVM

GP's

GP applications

Ref

- M.L. estimation: no prior
- M.A.P. estimation: no distribution
- Bayes est.



Bayes' models – if approximations used – we have **Level II Maximum Likelihood**.



Drawbacks of Parametric Models

Functional Modelling

Lehel Csató

Modelling

Machine Learning
Latent variable models
Estimation methods
Summary

Nonparametrics

SVM

GP's

GP applications

Ref

Model complexity – is an important “parameter”:

- choosing \mathcal{F} is an important decision in modelling the data:
- Example – for polynomial functions:
 - linear – too simple
 - quadratic – “good” for medium-sized data
 - ...
- The model complexity **should be** changed if we have more data available.



Non-parametric models

Functional Modelling

Lehel Csató

Modelling

Nonparametrics
Density Estimation
Regr./Class.

SVM

GP's

GP applications

Ref

- 1 Modelling
- 2 Non-parametric models
 - Density Estimation
 - Regression and classification
- 3 Support Vectors Machines
- 4 Gaussian Processes
- 5 Models using Gaussian Processes
- 6 References



Non-parametric Models

Functional Modelling

Lehel Csató

Modelling

Nonparametrics
Density Estimation
Regr./Class.

SVM

GP's

GP applications

Ref

Non-parametric models: use the available data for prediction.

- 1 Density estimation:
 - Histogram method;
 - Parzen window;
- 2 Regression/Classification:
 - K-Nearest Neighbour Rule;
 - Support Vector Machine;
 - **Gaussian Processes;**



Density Estimation

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

Density Estimation
Regr./Class.

SVM

GP's

GP applications

Ref

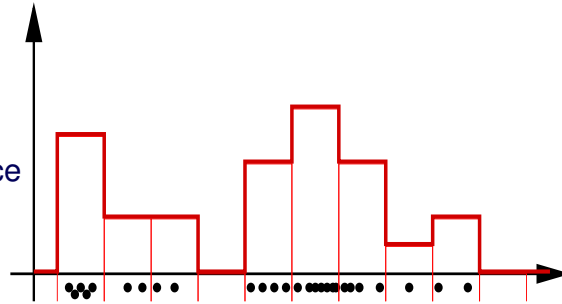
Histogram method:

“Parameters”

- bin width;
- locations;

Sensitive to the choice of parameters.

Not usable for higher dimensions



Density Estimation II

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

Density Estimation
Regr./Class.

SVM

GP's

GP applications

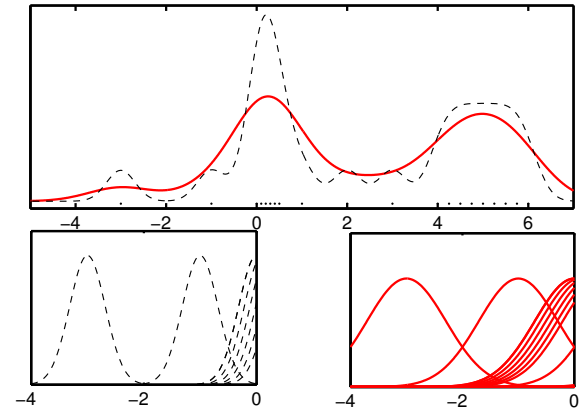
Ref

Parzen window:

“Parameters”

- function width;
- data locations

Not usable for large data-sets.



$$h(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_n, \mathbf{x})$$

Nonparametric: summation scales with N .



Non-parametric Classification

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

Density Estimation
Regr./Class.

SVM

GP's

GP applications

Ref

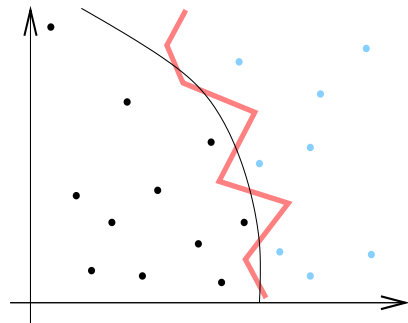
K-Nearest Neighbour: (Knn)

“Parameters”

- # of neighbours

Slow for large data.

Nonparametric: all data taken into account.



Knn for high-dimensions

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

Density Estimation
Regr./Class.

SVM

GP's

GP applications

Ref

We generally:

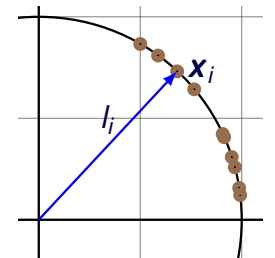
- normalise and center the data:

$$\mathbf{x}_i \leftarrow \mathbf{x}_i - \bar{\mathbf{x}} \quad \bar{\mathbf{x}} - \text{mean}$$

$$\mathbf{x}_{ij} \leftarrow \mathbf{x}_{ij} / \sigma_j \quad \sigma_j^2 \text{ is var. for } j\text{-th comp.}$$

- Each x_{ij} has zero mean and unit variance.
- the length of the random vector \mathbf{x}_i is

$$l_i^2 = \sum_j x_{ij}^2$$

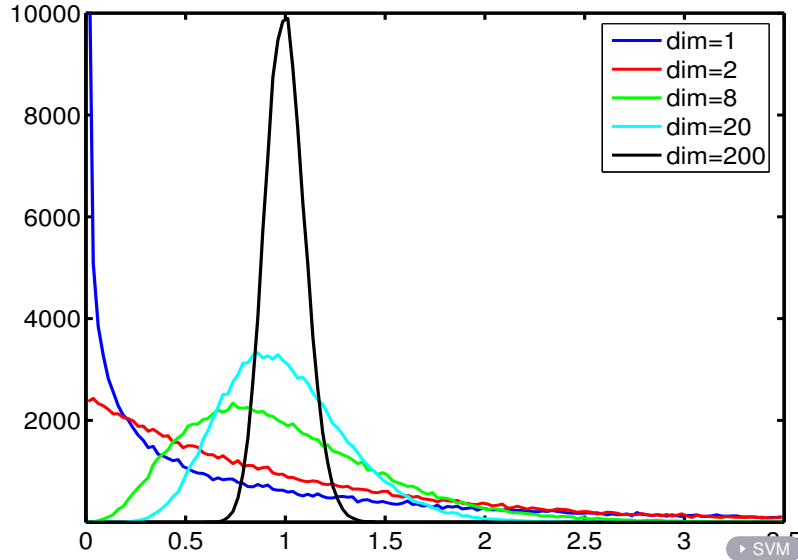


cf. central limit theorem

the larger the dimension, the more concentrated the average length is around the mean.



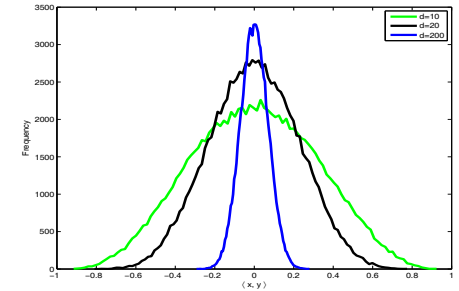
Example: `dim=20;x=randn(100000,dim);s2=sum(x.^2,2)/dim;g=histc(s2,bin)`



The angle:

$$\langle \mathbf{x}_i, \mathbf{x}_j \rangle = \sum_{\ell=1}^d x_{i\ell} x_{j\ell}$$

with average value 0.



Central limit theorem

The larger the dimension

- the more orthogonal random vectors are;
- the more difficult is to select a representative.



- 1 Modelling
- 2 Non-parametric models
- 3 Support Vectors Machines
 - Loss functions
 - Function class
 - Kernels
- 4 Gaussian Processes
- 5 Models using Gaussian Processes
- 6 References



Motivation:

Knn

- is too demanding for large datasets;
- works only in low-dimensions;

▶ HD

Algorithm that works “well” in high-dimensions

$$f(\mathbf{x}) = \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \sum_n L(y_n, f(\mathbf{x}_n)) + \|\mathbf{P}f(\cdot)\|^2 \right\}$$

Learning Algorithm

Design of a method that simultaneously minimises the empirical error (**first term**) and selects “cleverly” (**second term**) from a large family \mathcal{F} of available functions.

Details of the elements: ...

▶ Skip HighDim



Algorithm that works "well" in high-dimensions

$$f(\mathbf{x}) = \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \sum_n L(y_n, f(\mathbf{x}_n)) + \|\mathbf{P}f(\cdot)\|^2 \right\}$$

Learning Algorithm

Design of a method that simultaneously minimises the empirical error (**first term**) and selects "cleverly" (**second term**) from a large family \mathcal{F} of available functions.

Within the **SVM framework** we look for candidates

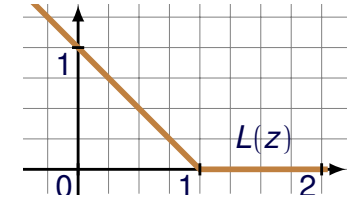
- in a **large** family of functions; and
- we penalise the **complexity** of the functions;

Explain loss function and \mathcal{F} .



Hinge Loss:

$$L(z) = \begin{cases} 0 & \text{if } z \geq 1 \\ 1 - z & \text{if } z < 1 \end{cases}$$



For **classification** with labels $y_n = \pm 1$ the loss function for the **data-set** \mathcal{D} :

$$L(\mathcal{D}) \stackrel{\text{def}}{=} \sum_n L(y_n f(\mathbf{x}_n))$$

Minimality

$L(\mathcal{D})$ is minimal if all elements separated **with margin** ≥ 1 .



Loss functions:

Hinge loss – in classification – penalises data away from boundary;

- assuming class labels $y_i = \pm 1$, we have

$$L(y_n, f(\mathbf{x}_n)) = H_+(1 - y_n f(\mathbf{x}_n))$$

Logit loss – in binary classification, returns the log-probabilities:

$$L_{\text{logit}}(y_n, f(\mathbf{x}_n)) = -\log(1 + \exp(-y_n f(\mathbf{x}_n)))$$

- when $y_n f(\mathbf{x}_n) \rightarrow \infty$, $L_{\text{logit}} \rightarrow 0$
- in **probit model** $(1 + \exp(\cdot))^{-1}$ replaced with $\Phi(\cdot)$.

Quadratic loss – quadratic error:

$$L_{\text{mse}}(y_n, f(\mathbf{x}_n)) = (y_n - f(\mathbf{x}_n))^2 \equiv (1 - y_n f(\mathbf{x}_n))^2$$

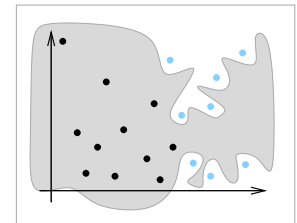
- can also be applied to regression problems.



Which family of functions? – we have **only** data \mathcal{D} .

Possibilities:

- **linear** – might be too simple;
- **complex** – might be too complex;



In general we

- want a flexible function class, but
- do **not** want a too complex solution.

Solution:

use a large function class \mathcal{F} **and** define penalties on the complexity of the solution.



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

Loss functions

Function class

Kernels

GP's

GP applications

Ref

- Loss functions help us find the best function
- **from a family of functions.**

The family of functions is important

It should be flexible enough

- for large data-sets
- and to allow for different length-scales.

Solution is of the form (**P** penalty):

$$f(\mathbf{x}) = \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \sum_n L(y_n, f(\mathbf{x}_n)) + \|\mathbf{P}f(\cdot)\|^2 \right\}$$

$\|\mathbf{P}f(\cdot)\|^2 = 0 \Rightarrow$ null-space, gives possible solutions.



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

Loss functions

Function class

Kernels

GP's

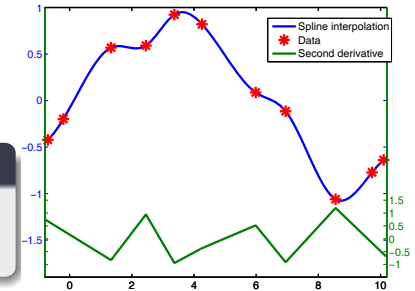
GP applications

Ref

Example:

\mathcal{F} – **twice** diff. functions with cont. second derivatives;
 $\mathbf{P} \stackrel{\text{def}}{=} (\partial_x^2 f(\cdot))^2$ – the sum of the second derivatives.

Result: interpolating splines with knots at the data points.



Message

Functional optimisation plausible with reliable results.

Presented in Wahba'90: Splines Models for Observational Data



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

Loss functions

Function class

Kernels

GP's

GP applications

Ref

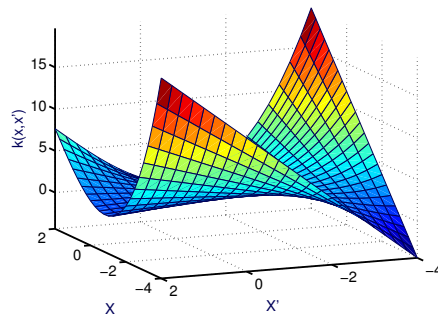
Result of the optimisation problem:

$$K_s(\mathbf{x}, \mathbf{x}') = \frac{1}{3} \min(\mathbf{x}, \mathbf{x}')^3 + \frac{1}{2} |\mathbf{x} - \mathbf{x}'| \min(\mathbf{x}, \mathbf{x}')^2 + \mathbf{x}\mathbf{x}' + 1$$

Solution of the argmin:

$$\hat{f}(\mathbf{x}) = \sum_n \alpha_n K_s(\mathbf{x}, \mathbf{x}_n)$$

Linear combination of polynomials.



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

Loss functions

Function class

Kernels

GP's

GP applications

Ref

Reproducing Kernel Hilbert space:

Assume \mathcal{X} an index set, $\mathcal{H} = \{f \mid f : \mathcal{X} \rightarrow \mathbb{R}\}$ with $\langle \cdot, \cdot \rangle$ the dot product s.t. $\|f\|^2 = \langle f, f \rangle$.

If there exists $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ such that

- 1 k has the reproducing property:

$$\forall f \in \mathcal{H} \quad \langle f, k(x, \cdot) \rangle = f(x);$$

- 2 and k spans \mathcal{H} .

Not a feature map

RKHS is **not equivalent** with a feature map. More than a **single** feature map for the same $k(\cdot, \cdot)$.

$k(\cdot, \cdot)$ spans $\mathcal{H} \Rightarrow$ every function $f \in \mathcal{H}$ is a linear combination of $\{k(\mathbf{x}_n, \cdot)\}_{n=1}^{N_{\mathcal{H}}}$:

$$f(\mathbf{x}) = \sum_{n=1}^{N_{\mathcal{H}}} w_n k(\mathbf{x}, \mathbf{x}_n) \quad (N_{\mathcal{H}} = \infty \text{ allowed})$$

Note:

- we choose the *support set*:
 - the location of the points and
 - the set **size**;
- we choose the *weights*.

Empirical kernel map

The support set \mathcal{X} is the training data set:

$$\mathcal{H}_e = \left\{ \sum_n \alpha_n k(\mathbf{x}, \mathbf{x}_n) \mid \alpha_n \in \mathbb{R} \right\}$$

Kernel functions: Let $\Phi : \mathbb{R} \rightarrow \mathbb{R}^3$ be given by:

$$\Phi(\mathbf{x}) = \begin{bmatrix} 1 \\ \sqrt{2}\mathbf{x} \\ \mathbf{x}^2 \end{bmatrix}$$

and compute

$$\begin{aligned} \Phi(\mathbf{x})^T \Phi(\mathbf{x}') &= 1 \cdot 1 + \sqrt{2}\mathbf{x} \cdot \sqrt{2}\mathbf{y} + \mathbf{x}^2 \mathbf{y}^2 \\ &= (1 + \mathbf{xy})^2 \stackrel{\text{def}}{=} K(\mathbf{x}, \mathbf{y}) \end{aligned}$$

Kernel trick

We can translate **linear** algorithms into nonlinear ones using a **kernel** function – **represented** as a scalar product.

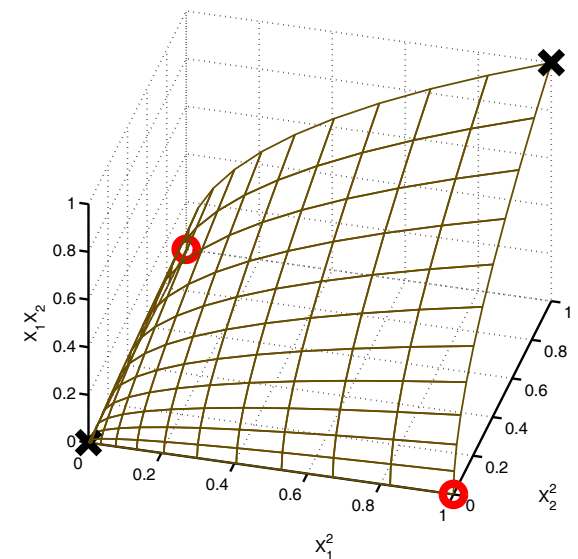
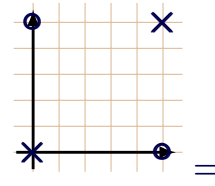
Mercer kernel

For a positive-definite $k(\cdot, \cdot)$ there exists a set $\{\phi_j\}$ of orthogonal functions and $\{\lambda_j\}$ positive constants s.t.

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

Note:

- The function $k(\cdot, \cdot)$ defines:
 - the **eigen-functions** $\{\phi_j\}$
 - the corresponding eigen-values $\{\lambda_j\}$;
- independent of the data-set we are using;
- convergence guarantee: $\sum_j \lambda_j^2 < \infty$.





Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

Loss functions

Function class

Kernels

GP's

GP applications

Ref

Kernels – two-argument functions that are the generalisation of a matrix to **non-countable** index sets.

$$K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

Valid **kernel functions** are positive definite: for any $\mathbf{a} = [a_1, a_2, \dots, a_L]$ and $\mathcal{D} = [\mathbf{x}_1, \dots, \mathbf{x}_L]^T$:

$$\sum_{k=1}^L \sum_{l=1}^L a_k K(\mathbf{x}_k, \mathbf{x}_l) a_l \geq 0$$

Proof idea: kernel function is $K(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x})^T \Phi(\mathbf{y})$

and we have

$$\sum_{k,l=1}^L a_k \Phi(\mathbf{x}_k)^T \Phi(\mathbf{x}_l) a_l = (\sum a_i \Phi(\mathbf{x}_i))^T (\sum a_i \Phi(\mathbf{x}_i)) = \mathbf{s}^T \mathbf{s} \geq 0$$

where

$$\mathbf{s} = \sum_{i=1}^L a_i \Phi(\mathbf{x}_i)$$

What is the dimensionality of $\Phi(\mathbf{x})$.



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

Loss functions

Function class

Kernels

GP's

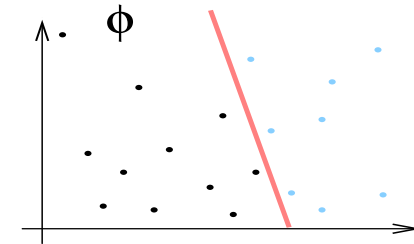
GP applications

Ref

Developed for classification.

Idea – **KERNEL TRICK**:

- use linear algorithms;
- project to Φ ;
- solution in Φ ;
- back-project;
- **NON-linear** solution;



Solution to the problem is of the form:

$$f(\mathbf{x}) = \Phi_{\mathbf{x}}^T \left(\sum_{i \in SV} \alpha_i \Phi_{\mathbf{x}_i} \right) = \sum_{i \in SV} \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

Nonparametric: number of parameters is not fixed.



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

Loss functions

Function class

Kernels

GP's

GP applications

Ref

The projection “trick” is exploited.

- Based on the success of the S.V.M.-s;
- General recipe:
 - Find/Construct a linear algorithm;
 - Re-express it in Φ – the space of “features”;
 - Write the – non-linear – solution in the space of inputs and use $K(\mathbf{x}, \mathbf{x}')$.
- Algorithms: Kernel ...
 - ... regression, ... ridge regression;
 - ... Principal/Independent Components;
 - ... Fisher Discriminants;



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

Loss functions

Function class

Kernels

GP's

GP applications

Ref

Definition: non-parametric methods \equiv number of parameters/complexity **might** change with data.

• Advantages:

- ↑ Greater complexity;
- ↑ Built-in regularisation;
- ↑ Performant algorithms.

• Disadvantages:

- ↓ Small data sets only;
- ↓ No selection of parameters;



Non-parametric models

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

- 1 Modelling
- 2 Non-parametric models
- 3 Support Vectors Machines
- 4 Gaussian Processes
 - Kernels within GP
 - Inference and Prediction
 - Gaussian Regression
 - Posterior Approximations
 - Optimising hyper-parameters
 - Sparse Representation
- 5 Models using Gaussian Processes



Bayesian Nonparametric Methods

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

- GP models \equiv functional latent variable models.
- GP models \equiv “simple” random functions.
- Appear in the likelihood:

$$P(\mathcal{D}|\mathbf{f}_{\mathcal{D}}) = \prod_i P(y_i|\mathbf{x}_i, \mathbf{f})$$

$$= \prod_i P(y_i|\mathbf{x}_i, f_{\mathbf{x}_i})$$

Local dependencies only: $\mathbf{f} \rightarrow f_{\mathbf{x}}$



Gaussian processes I

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Gaussian process: generalisation of a Gaussian.

- \mathbf{f} Gaussian random function.

$$\mathbf{f} = [f_{\mathbf{x}_1}, f_{\mathbf{x}_2}, \dots, f_{\mathbf{x}_N}, \dots]^T, \quad \mathbf{x}_n \in \text{domain.}$$

- GP prior $p_0(\mathbf{f})$ characterised with

- mean function $\langle f_{\mathbf{x}} \rangle_0$,
- covariance kernel $K_0(\mathbf{x}, \mathbf{x}')$.

Property - for **any** sample set \mathcal{D} , a **joint** Gaussian r.v.:

$$\mathbf{f}_{\mathcal{D}} = [f_{\mathbf{x}_1}, \dots, f_{\mathbf{x}_N}] \sim \mathcal{N}(\mathbf{f}_{\mathcal{D}} | \langle f_{\mathcal{D}} \rangle_0, \mathbf{K}_0)$$



Gaussian processes II

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

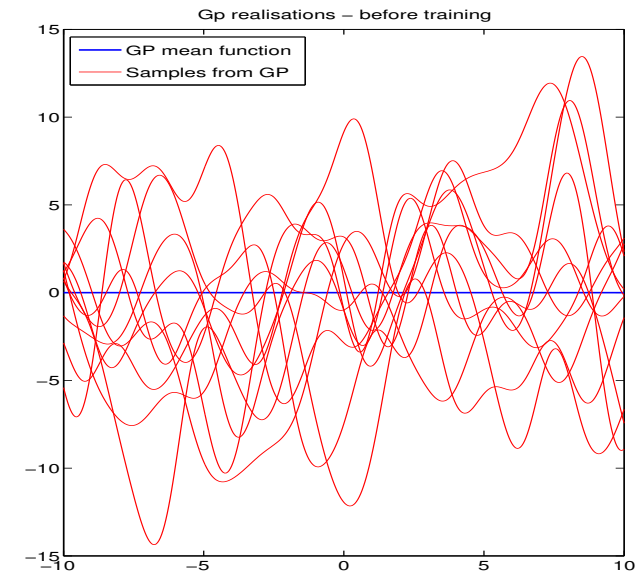
GP applications

Ref

Gaussian process: random function

“Parameters”

- mean function
- covariance kernel



DEMO



GP parameters

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Gaussian process parameters

- mean function – usually is 0.
- **parameter:** the **class** of the kernel function
- parameters hidden into the kernel function.
Example:

$$K(\mathbf{x}, \mathbf{x}') = \exp \left[\theta_0 - \frac{1}{2} \sum_{i=1}^d \theta_i (x_i - x'_i)^2 \right]$$

$\boldsymbol{\theta} = [\theta_0, \theta_1, \dots, \theta_d]^T$ – parameter vector.



Kernel Functions

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Kernel functions:

- Generate the covariance matrix.
- Need to be positive definite functions/matrices

$$\forall \mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$$

$$\mathbf{K}_{\mathcal{X}} = \begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \dots & K(\mathbf{x}_1, \mathbf{x}_T) \\ \vdots & \ddots & \vdots \\ K(\mathbf{x}_T, \mathbf{x}_1) & \dots & K(\mathbf{x}_T, \mathbf{x}_T) \end{bmatrix}$$

must be a positive definite matrix.

- A construction of kernels as covariances:



Kernel constructions

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Let $\{x_k = k\Delta_N\}_{k=1}^N$ be the **location** of Gaussian independent r.v.-s

$$v_k \sim \mathcal{N}(0, \sigma_0^2 \Delta_N)$$

where $\Delta_N = 1/N$. Let t denote the the index where $t = k_t \Delta_N$.

If $b_t \stackrel{\text{def}}{=} \sum_{i=1}^t v_i$, then $\langle b_t \rangle = 0$ and the covariance:

$$\langle b_s b_t \rangle = \left\langle \sum_{i_s=1}^s \sum_{i_t=1}^t v_{i_s} v_{i_t} \right\rangle = \sum_{i_s=1}^{\min(s,t)} \langle v_{i_s}^2 \rangle$$

$$= \sum_{i_s=1}^{\min(s,t)} \sigma_0^2 \Delta_N = \min(s, t) \sigma_0^2$$

in the following we take $\sigma_0^2 = 1$.

Brownian motion

A stochastic process with covariance kernel $K(s, t) = \min(s, t)$ is a Brownian motion.

Its derivative – v_t – is the Wiener process.



Kernel constructions

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

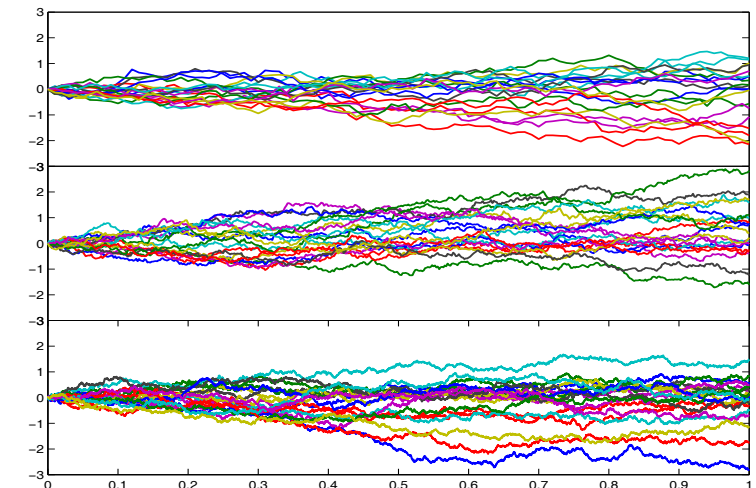
A.R.D

Sparsity

GP applications

Ref

Images of the Brownian motion at different resolutions (different N -s).





Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Let us integrate (sum) the Brownian motion.

Define

$$s_t \stackrel{\text{def}}{=} \sum_{i_t=1}^t b_{i_t}$$

leading to $\langle s_t \rangle = 0$ and

$$\begin{aligned} K(s, t) &= \langle s_s s_t \rangle = \left\langle \sum_{i_s=1}^s \sum_{i_t=1}^t s_{i_s} s_{i_t} \right\rangle \\ &= \sum_{i_s=1}^s \sum_{i_t=1}^t \langle s_{i_s} s_{i_t} \rangle = \int_0^s dz_s \int_0^t dz_t \min(z_s, z_t) \end{aligned}$$

For the integration we assume

$$s < t \implies z_s < t \implies [0, t] = [0, z_s] \cup [z_s, t].$$



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Using $[0, t] = [0, z_s] \cup [z_s, t]$

$$\begin{aligned} K(s, t) &= \int_0^s dz_s \left(\int_0^{z_s} dz_t z_t + \int_{z_s}^t dz_t z_s \right) = \int_0^s dz_s \left(\frac{z_s^2}{2} + z_s(t - z_s) \right) \\ &= \int_0^s dz_s \left(z_s t - \frac{z_s^2}{2} \right) \\ &= \frac{s^2 t}{2} - \frac{s^3}{6} \quad \text{assuming } s < t \end{aligned}$$

After symmetrization (writing the $s > t$ case and unifying)

$$\begin{aligned} K(s, t) &= \frac{st \min(s, t)}{2} - \frac{\min(s, t)^3}{6} \\ &= \frac{1}{2} \min(s, t)^2 |s - t| + \frac{1}{3} \min(s, t)^3 \end{aligned}$$



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

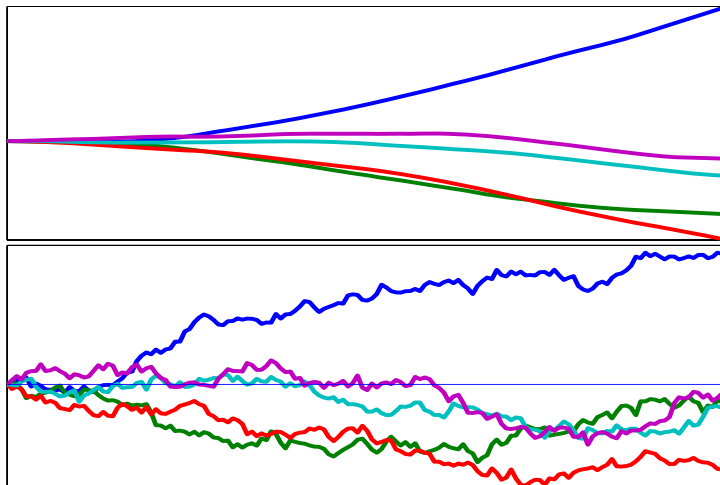
A.R.D

Sparsity

GP applications

Ref

Samples from the integrated Brownian motion
(Brownian motion on the bottom).



Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

- The integrated Brownian motion: $s(0) = s'(0) = 0$.
- For generality we add a *constant* and a *linear* term:

$$s_2(x) = w_0 + w_1 x + s(x)$$

where w_0, w_1 are i.i.d. Gaussian r.v.s.

- Means that the kernel is:

$$K_2(s, t) = \langle s_2(s) s_2(t) \rangle = 1 + st + K_s(s, t)$$

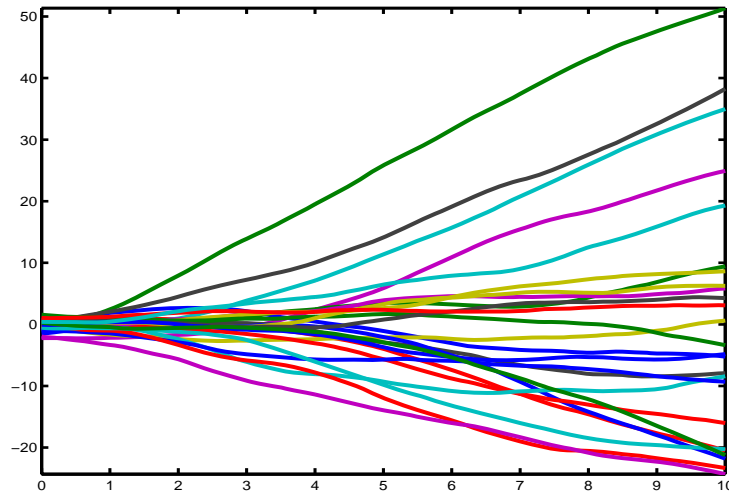
$$K_2(s, t) = 1 + st + \frac{1}{2} \min(s, t)^2 |s - t| + \frac{1}{3} \min(s, t)^3$$



Functional Modelling

Lehel Csató

Random splines:



Modelling
Nonparametrics
SVM
GP's
Kernels II
Inf./Pred.
Regression
Approximations
A.R.D
Sparsity
GP applications
Ref



Functional Modelling

Lehel Csató

- Consider the Brownian motion: $K_b(s, t) = \min(s, t)$.
- We generate random functions with $x_1 = 0$.
Called **Brownian bridge**
- To calculate the covariance, we have to condition the r.v.-s x_s and x_t on $x_1 = 0$.

Modelling
Nonparametrics
SVM
GP's
Kernels II
Inf./Pred.
Regression
Approximations
A.R.D
Sparsity
GP applications
Ref

We identify the kernel from the conditioned **joint** Gaussian distribution

$$\begin{aligned} p(x_s, x_t | x_1 = 0) &\propto \exp \left\{ -\frac{1}{2} \begin{bmatrix} x_s \\ x_t \\ 0 \end{bmatrix}^T \begin{bmatrix} K_b(s, s) & K_b(s, t) & s \\ K_b(t, s) & K_b(t, t) & t \\ 0 & 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} x_s \\ x_t \\ 0 \end{bmatrix} \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \begin{bmatrix} x_s \\ x_t \\ 0 \end{bmatrix}^T \left[\begin{pmatrix} K_b(s, s) & K_b(s, t) \\ K_b(t, s) & K_b(t, t) \end{pmatrix} - \begin{bmatrix} s \\ t \end{bmatrix} \begin{bmatrix} s \\ t \end{bmatrix}^T \right]^{-1} \begin{bmatrix} E \\ F \\ 0 \end{bmatrix} \begin{bmatrix} x_s \\ x_t \\ 0 \end{bmatrix} \right\} \end{aligned}$$

where we used the matrix inversion lemma

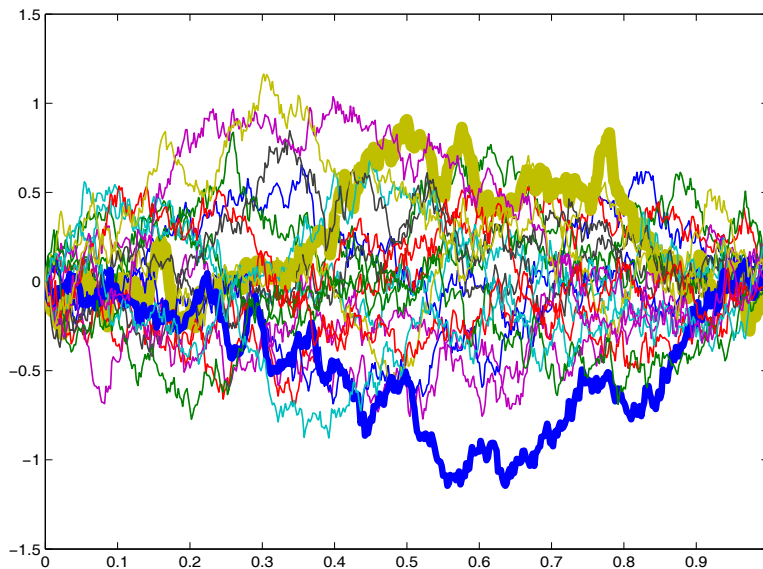
$$\Rightarrow K_0(s, t) = K_b(s, t) - st$$



Functional Modelling

Lehel Csató

Samples from a Brownian bridge:



Modelling
Nonparametrics
SVM
GP's
Kernels II
Inf./Pred.
Regression
Approximations
A.R.D
Sparsity
GP applications
Ref



Functional Modelling

Lehel Csató

Exercise:

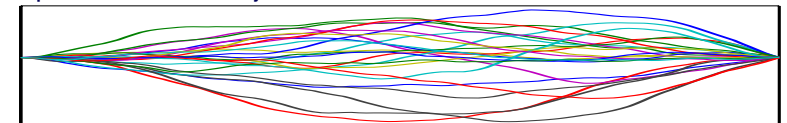
- Find the mean and kernel functions corresponding to the Brownian bridge where $x_1 = 1$.
- Consider the spline kernel

$$K_s(s, t) = \frac{1}{2} \min(s, t)^2 |s - t| + \frac{1}{3} \min(s, t)^3$$

Similarly to the Brownian bridge, find the kernel function for the splines conditioned on $x_1 = 0$.

Modelling
Nonparametrics
SVM
GP's
Kernels II
Inf./Pred.
Regression
Approximations
A.R.D
Sparsity
GP applications
Ref

Samples from the second family look like this:





Generating random samples

XI

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

```

1 clear all;
  N = 100; T=25; D = 0.99;
  t = linspace(0,D,N+1);
  t = t(2:end);

6 % put covariance function here
  k_bb= inline('min(s,t)-s*t','s','t');

  Ks = zeros(N,N);
  for ii=1:N;
11     for jj=1:N;
         Ks(ii,jj) = k_sb(t(ii),t(jj));
       end;
     end;

16 kks = chol(Ks);
   yr = randn(T,N);
   ys = zeros(N+2,T);
   ys(2:N+1,:) = (yr*kks)';

21 t0=[0,t,1];
   figure(1); cla; box on; hold on;
   plot(t0,ys);

```



Gaussian Process Inference I

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

- GP inference: application of Bayes' rule.

$$p_{\text{post}}(\mathbf{f}, \mathbf{f}_D) \propto P(\mathcal{D}|\mathbf{f}_D) p_0(\mathbf{f}_D, \mathbf{f})$$

- For any collection of indexes \mathcal{X} the posterior:

$$p_{\text{post}}(\mathbf{f}_\mathcal{X}) = \frac{1}{Z_D} \int d\mathbf{f}_D P(\mathcal{D}|\mathbf{f}_D) p_0(\mathbf{f}_D, \mathbf{f}_\mathcal{X})$$

where

$$Z_D = \int d\mathbf{f}_D P(\mathcal{D}|\mathbf{f}_D) p_0(\mathbf{f}_D)$$

Data probability – conditioned on the model.

A.R.D.

obs: NO specific $P(\mathcal{D}|\mathbf{f}_D)$



Gaussian Process Inference II

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Problems with the representation

$$p_{\text{post}}(\mathbf{f}_\mathcal{X}) \propto \int d\mathbf{f}_D P(\mathcal{D}|\mathbf{f}_D) p_0(\mathbf{f}_D, \mathbf{f}_\mathcal{X})$$

Integral evaluation necessary for posterior distribution.

- Representation – How to represent the posterior?

- Finite representation of the posterior process;
- Non-Gaussian posterior processes: approximations to them



Gaussian Process Parametrisation I

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Property of Gaussian averages:

$$\langle \mathbf{f}_\mathbf{x} \rangle_{\text{post}} = \langle \mathbf{f}_\mathbf{x} \rangle_0 + \sum_i K_0(\mathbf{x}, \mathbf{x}_i) \alpha(i)$$

Where coefficients:

$$\alpha(i) = \frac{\partial}{\partial \langle \mathbf{f}_i \rangle_0} \ln \left\langle P(\mathcal{D}|\mathbf{f}_D) \right\rangle_0$$

Provide parametrisation (see Kimeldorf-Wahba).



Gaussian Process Parametrisation II

Functional Modelling

Lehel Csató

For the posterior kernel:

$$K_{\text{post}}(\mathbf{x}, \mathbf{x}') = K_0(\mathbf{x}, \mathbf{x}') + \sum_{ij} K_0(\mathbf{x}, \mathbf{x}_i) C(ij) K_0(\mathbf{x}_j, \mathbf{x}')$$

Where coefficients:

$$C(ij) = \frac{\partial^2}{\partial \langle f_i \rangle_0 \partial \langle f_j \rangle_0} \ln \left\langle P(\mathcal{D} | \mathbf{f}_{\mathcal{D}}) \right\rangle_0$$

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref



GPs in Feature space I

Functional Modelling

Lehel Csató

$K_0(\mathbf{x}, \mathbf{x}')$ defines a *feature space* \mathcal{F} :

$$\phi_{\mathbf{x}}, \phi_{\mathbf{x}'} \in \mathcal{F} \quad \text{and} \quad K_0(\mathbf{x}, \mathbf{x}') = \phi_{\mathbf{x}}^T \phi_{\mathbf{x}'}$$

Using \mathcal{F} and the scalar product:

$$\langle f_{\mathbf{x}} \rangle_{\text{post}} = \phi_{\mathbf{x}}^T \sum_{i=1}^N \alpha(i) \phi_i = \phi_{\mathbf{x}}^T \boldsymbol{\mu}_{\text{post}}$$

$$K_{\text{post}}(\mathbf{x}, \mathbf{x}') = \phi_{\mathbf{x}}^T \left(\mathbf{I}_{\mathcal{F}} + \sum_{ij=1}^N \phi_i C(ij) \phi_j^T \right) \phi_{\mathbf{x}'} = \phi_{\mathbf{x}}^T \boldsymbol{\Sigma}_{\text{post}} \phi_{\mathbf{x}'}$$

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref



GPs in Feature space II

Functional Modelling

Lehel Csató

$K_0(\mathbf{x}, \mathbf{x}')$ defines a *feature space* \mathcal{F} :

$$\phi_{\mathbf{x}}, \phi_{\mathbf{x}'} \in \mathcal{F} \quad \text{and} \quad K_0(\mathbf{x}, \mathbf{x}') = \phi_{\mathbf{x}}^T \phi_{\mathbf{x}'}$$

$$\begin{aligned} \langle f_{\mathbf{x}} \rangle_{\text{post}} &\iff \boldsymbol{\mu}_{\text{post}} \\ K_{\text{post}}(\mathbf{x}, \mathbf{x}') &\iff \boldsymbol{\Sigma}_{\text{post}} \end{aligned}$$

GP inference:

Estimating a **Gaussian distribution** in \mathcal{F}

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref



Prediction with Gaussian processes

Functional Modelling

Lehel Csató

Given: \mathbf{x}^* - for which we **require** answer y^* .

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathcal{D}) &= \int df^* \int d\mathbf{f}_{\mathcal{D}} p(y^*, \mathbf{f}_{\mathcal{D}}, f^* | \mathbf{x}^*, \mathcal{D}) \\ &= \int df^* P(y^* | \mathbf{x}^*, f^*) \int d\mathbf{f}_{\mathcal{D}} p_{\text{post}}(\mathbf{f}_{\mathcal{D}}, f^* | \mathcal{D}) \\ &= \int df^* P(y^* | \mathbf{x}^*, f^*) p_{\text{post}}(f^* | \mathcal{D}) \end{aligned}$$

where $f^* = f_{\mathbf{x}^*}$ – random variable associated to \mathbf{x}^* .

We use posterior process:

- irrespective of the likelihood;
- if **not** Gaussian, we approximate.

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref



Regression with Gaussian noise

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

- Gaussian noise:

$$P(\mathbf{y}_D | \mathbf{f}_D) \propto \exp \left[-\frac{1}{2\sigma_0^2} \|\mathbf{y}_D - \mathbf{f}_D\|^2 \right]$$

- Gaussian latent variables:

$$P(\mathbf{f}_D | K(\cdot, \cdot | \boldsymbol{\theta})) \propto \exp \left[-\frac{1}{2} (\mathbf{f}_D - \boldsymbol{\mu}_D)^T \mathbf{K}_D^{-1} (\mathbf{f}_D - \boldsymbol{\mu}_D) \right]$$

- Combining: product quadratic \implies Gaussian



Posterior distribution – Gaussian noise

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Gaussian distribution – reading off the coefficients:

$$\boldsymbol{\mu}_{\text{post}} = \left(\mathbf{K}_D^{-1} + \frac{1}{\sigma_0^2} \mathbf{I}_N \right)^{-1} \left[\mathbf{K}_D^{-1} \boldsymbol{\mu}_D + \frac{1}{\sigma_0^2} \mathbf{y}_D \right]$$

$$\boldsymbol{\Sigma}_{\text{post}} = \left(\mathbf{K}_D^{-1} + \frac{1}{\sigma_0^2} \mathbf{I}_N \right)^{-1}$$

The **joint** distribution of all r.v.-s is a Gaussian:

$$\mathbf{f}_D \sim \mathcal{N}(\boldsymbol{\mu}_{\text{post}}, \boldsymbol{\Sigma}_{\text{post}})$$

The distribution of the r.v.-s at **training locations**.



Posterior distribution Brownian bridge

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

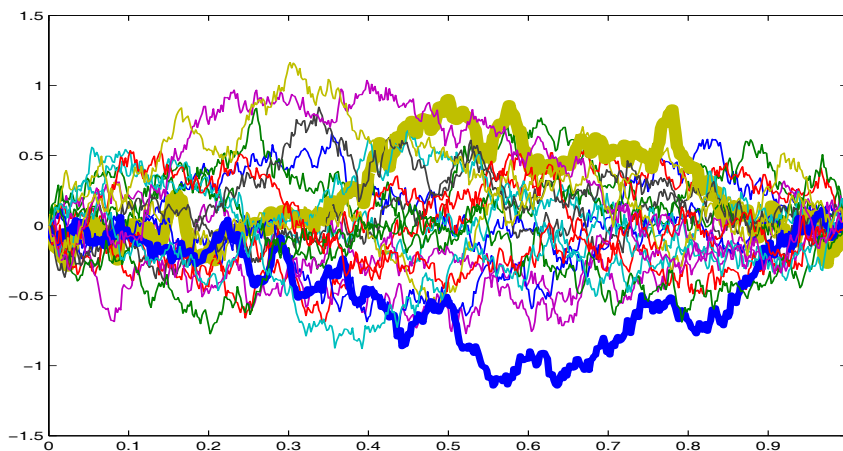
A.R.D

Sparsity

GP applications

Ref

- **Assume** that there was an observation for the Brownian motion $k(\mathbf{x}, \mathbf{x}') = \min(\mathbf{x}, \mathbf{x}')$;
- at 1 the value of the process is 0.
- It generates a new process: $k_B(\mathbf{x}, \mathbf{x}') = \min(\mathbf{x}, \mathbf{x}') - \mathbf{x}\mathbf{x}'$



Predictive distributions – Gaussian noise

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

For test location \mathbf{x}^* :

$$\begin{aligned} \boldsymbol{\mu}^* &= \boldsymbol{\alpha}_D^T \mathbf{k}_* & \boldsymbol{\alpha}_D &= \mathbf{C}_D * \mathbf{y}_D \\ \sigma^* &= K(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{k}_*^T \mathbf{C}_D \mathbf{k}_* & \text{with } \mathbf{C}_D &= \left(\mathbf{K}_D + \sigma_0^2 \mathbf{I}_N \right)^{-1} \end{aligned}$$

where $\mathbf{k}_* = [K(\mathbf{x}^*, \mathbf{x}_1), \dots, K(\mathbf{x}^*, \mathbf{x}_N)]^T$

Posterior mean and covariance functions:

$$\begin{aligned} \boldsymbol{\mu}(\mathbf{x}) &= \boldsymbol{\alpha}_D^T \mathbf{k}_x \\ K_{\text{post}}(\mathbf{x}, \mathbf{x}') &= K(\mathbf{x}, \mathbf{x}') - \mathbf{k}_x^T \mathbf{C}_D \mathbf{k}_{x'} \end{aligned}$$

where $\mathbf{k}_x = [K(\mathbf{x}, \mathbf{x}_1), \dots, K(\mathbf{x}, \mathbf{x}_N)]^T$.



Non-computable Posteriors

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

- **If likelihood non-Gaussian** \Rightarrow posterior does not have analytical form. (No “summarising” statistics)

- **Methods to obtain posterior:**
 - Sampling;
 - Analytic approximations;



Sampling from the posterior I

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Sampling

$$p_{\text{post}}(\mathbf{f}_X) = \frac{1}{Z_D} \int d\mathbf{f}_D P(\mathcal{D}|\mathbf{f}_D) p_0(\mathbf{f}_D, \mathbf{f}_X)$$

- In practise:**
- joint sampling from $p_{\text{post}}(\mathbf{f}_X, \mathbf{f}_D)$,
 - keeping only \mathbf{f}_X .

Implementation: sampling from $p_0(\mathbf{f}_D, \mathbf{f}_X)$ + weighting:

$$p_{\text{post}}(\mathbf{f}_X) \approx \frac{1}{C_T} \sum_{t=1}^T P(\mathbf{y}_N|\mathbf{f}_D^{(i)}) \delta(\mathbf{f}_X - \mathbf{f}_X^{(i)})$$

with $C_T = \sum_{t=1}^T P(\mathbf{y}_N|\mathbf{f}_D^{(i)})$



Sampling from the posterior II

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

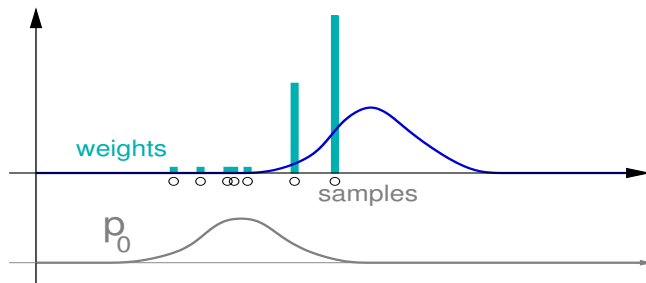
Approximations

A.R.D

Sparsity

GP applications

Ref



Sampling methods:

- **powerful** i.e. allow flexibility in modelling
- Hard to assess convergence
- Sampling algorithms suited for different models.
- Can be **incredibly slow** (tempering, MCMC)



Laplace Approximation

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Log-Posterior:

$$\begin{aligned} \log p_{\text{post}}(\mathbf{f}_X, \mathbf{f}_D) &= K + \log P(\mathcal{D}|\mathbf{f}_D) + \log p_0(\mathbf{f}_D, \mathbf{f}_X) \\ &= \underbrace{\log P(\mathcal{D}|\mathbf{f}_D) + \log p_0(\mathbf{f}_D)}_{g_D(\mathbf{f}_D)} + \underbrace{\log p_0(\mathbf{f}_X|\mathbf{f}_D)}_{g_X(\mathbf{f}_X)} \end{aligned}$$

Finding maximum of \mathbf{f}_X and \mathbf{f}_D :

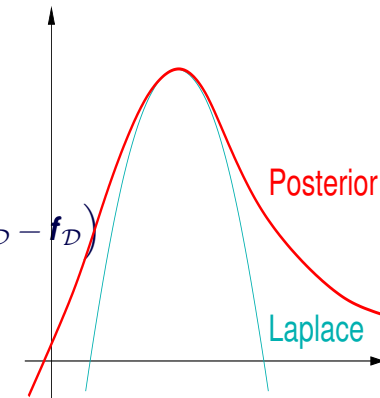
$$\hat{\mathbf{f}}_X = \mathbf{P}_{X\mathcal{D}} \mathbf{f}_D$$

$$\hat{\mathbf{f}}_D = \arg \max g_D(\mathbf{f}_D)$$

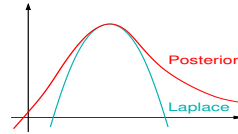
Taylor expansion around $\hat{\mathbf{f}}_D$:

$$\begin{aligned} g_D(\mathbf{f}_D) &\approx (\hat{\mathbf{f}}_D - \mathbf{f}_D)^T [H_g(\hat{\mathbf{f}}_D)] (\hat{\mathbf{f}}_D - \mathbf{f}_D) \\ &+ \mathbf{0} (\hat{\mathbf{f}}_D - \mathbf{f}_D) + g_D(\hat{\mathbf{f}}_D) \end{aligned}$$

\Rightarrow **Gaussian**



Gaussian approximation:



$$\hat{p}_{\text{post}}(\mathbf{f}_X, \mathbf{f}_D) \propto p_0(\mathbf{f}_X, \mathbf{f}_D) \underbrace{\frac{\mathcal{N}_L(\mathbf{f}_D | \hat{\mathbf{f}}_D, [H_g(\hat{\mathbf{f}}_D)]^{-1})}{p_0(\mathbf{f}_D)}}_{\hat{P}(D|\mathbf{f}_D)}$$

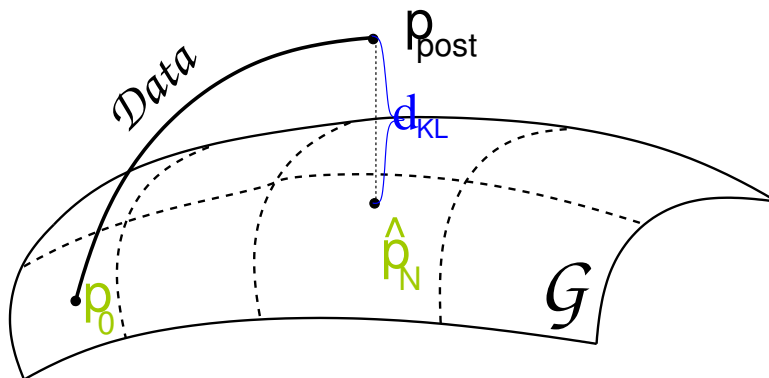
Defines an approximation to the **likelihood**:

$$\hat{P}(D|\mathbf{f}_D) \propto \frac{\mathcal{N}_L(\mathbf{f}_D | \hat{\mathbf{f}}_D, [H_g(\hat{\mathbf{f}}_D)]^{-1})}{p_0(\mathbf{f}_D)}$$

Analytic approximations – I

Aim: approximate the posterior **distribution** – or the posterior process.

GP prior \rightarrow GP **approximation to posterior**.
projection – closest GP.



The Laplace approximation:

- \oplus generates an approximation to the *likelihood*;
- \ominus applicable only for differentiable likelihood functions;
- \oplus defines an approximation to the *whole* process;
- \ominus the Hessian has to be positive definite and **“smaller” than the prior**

Analytic approximations – II

Choice of projection: Kullback-Leibler divergence

$$\text{KL}(\mathcal{GP}_{\text{post}} \parallel \mathcal{GP}) = \int d\mathcal{GP}_{\text{post}}(\mathbf{f}) \log \frac{d\mathcal{GP}_{\text{post}}(\mathbf{f})}{d\mathcal{GP}(\mathbf{f})}$$

$$\mathcal{GP}^* = \arg \min_{\mathcal{GP}} \text{KL}(\mathcal{GP}_{\text{post}} \parallel \mathcal{GP})$$

The minimiser:

$$\langle \mathbf{f}_X \rangle_{\mathcal{GP}^*} \stackrel{\text{def}}{=} \langle \mathbf{f}_X \rangle_{\text{post}}$$

$$K_{\mathcal{GP}^*}(\mathbf{x}, \mathbf{x}') \stackrel{\text{def}}{=} K_{\text{post}}(\mathbf{x}, \mathbf{x}')$$

Implies that the KL-approximation the $\mathcal{GP} \Leftrightarrow (\alpha_D, \mathbf{C}_D)$.



Computing KL-distances

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Between GPs with the same prior:

$\mathcal{GP}_1 = \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)$ and $\mathcal{GP}_2 = \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2)$, the KL-distance is:

$$2\text{KL}(\mathcal{GP}_1 \parallel \mathcal{GP}_2) = (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T \Sigma_2^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \text{tr} \left(\Sigma_1 \Sigma_2^{-1} - I_{\mathcal{F}} \right) - \ln \left| \Sigma_1 \Sigma_2^{-1} \right|$$

With parameters $(\mathbf{Q}_{B\mathcal{V}} = \mathbf{K}_{B\mathcal{V}}^{-1})$:

$$(\boldsymbol{\alpha}_2 - \boldsymbol{\alpha}_1) (\mathbf{C}_2 + \mathbf{Q}_{B\mathcal{V}})^{-1} (\boldsymbol{\alpha}_2 - \boldsymbol{\alpha}_1) + \text{tr} \left[(\mathbf{C}_1 - \mathbf{C}_2) (\mathbf{C}_2 + \mathbf{Q}_{B\mathcal{V}})^{-1} \right] - \ln \left| (\mathbf{C}_1 + \mathbf{Q}_{B\mathcal{V}}) (\mathbf{C}_2 + \mathbf{Q}_{B\mathcal{V}})^{-1} \right|$$

Assumptions: the kernel matrix on the $B\mathcal{V}$ set is non-singular $|\mathbf{K}_{B\mathcal{V}}| \neq 0$.



Bayesian Online Learning

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref

Learning: propagating the mean and the kernel:

$$\langle f_{\mathbf{x}} \rangle_{t+1} = \langle f_{\mathbf{x}} \rangle_t + q K_t(\mathbf{x}, \mathbf{x}_{t+1})$$
$$K_{t+1}(\mathbf{x}, \mathbf{x}') = K_t(\mathbf{x}, \mathbf{x}') + r K_t(\mathbf{x}, \mathbf{x}_{t+1}) K_t(\mathbf{x}_{t+1}, \mathbf{x}')$$

q, r functions of the **single** likelihood:

$$q = q^{(t+1)} = \frac{\partial}{\partial \langle f_{t+1} \rangle_t} \ln \langle P(y_{t+1} | f_{t+1}) \rangle_t$$

where $\langle \cdot \rangle_t$ average w.r.to $f_{t+1} \sim \mathcal{N}(\langle f_{t+1} \rangle_t, \sigma_{t+1}^2)$.



Analytic Approximations – III

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

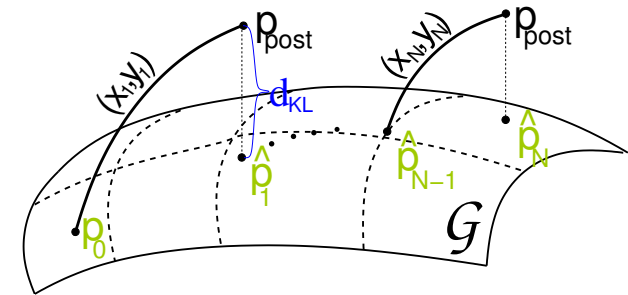
GP applications

Ref

Bayesian Online Learning (recursive)

- Instead of $|\mathcal{D}| = N$ uses $\mathcal{D} = (\mathbf{x}_{t+1}, y_{t+1})$ and;
- For prior process $\langle f_{\mathbf{x}} \rangle_t, K_t(\mathbf{x}, \mathbf{x}')$.

$$\text{KL}(\mathcal{GP}_{\text{post}}^{t+1} \parallel \mathcal{GP}^*) \quad \text{smaller approximation}$$



Optimising hyper-parameters I

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

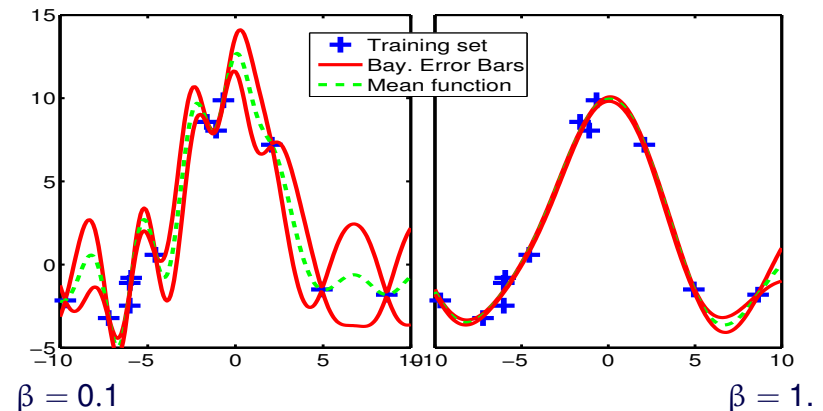
Sparsity

GP applications

Ref

GP **kernel** parameters \Leftrightarrow model choice.
Exemplu:

$$\text{RBF kernel: } K(\mathbf{x}, \mathbf{x}') = A \exp \left[- \sum (x - x')^2 \beta \right]$$





Optimising hyper-parameters II

Functional Modelling

Lehel Csató

Model evidence:

$$Z_{\mathcal{D}}(\boldsymbol{\theta}) = P(\mathcal{D}|\boldsymbol{\theta}) = \int d\mathbf{f}_{\mathcal{D}} P(\mathcal{D}|\mathbf{f}_{\mathcal{D}}) p_0(\mathbf{f}_{\mathcal{D}}|\boldsymbol{\theta})$$

Maximum Likelihood II inference

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \Omega} P(\boldsymbol{\theta}|\mathcal{D}) = \arg \min_{\boldsymbol{\theta} \in \Omega} \frac{P(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{M})}{p(\mathcal{D}|\mathcal{M})}$$

if $p(\boldsymbol{\theta}|\mathcal{M})$ “flat”

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \Omega} Z_{\mathcal{D}}(\boldsymbol{\theta})$$

Evidence maximisation.

Gradient/conj.grad. methods are used.

► MacKay ► Evidence



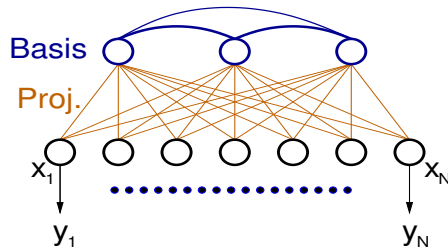
Sparse representations – a solution

Functional Modelling

Lehel Csató

Condition **all training** locations on a set of **basis** locations.

$$\mathbf{f}_{\mathcal{B}\mathcal{Y}} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathcal{B}\mathcal{Y}}, \boldsymbol{\Sigma}_{\mathcal{B}\mathcal{Y}})$$



The pseudo-latents $\mathbf{f}_{\mathcal{X}}$ are conditioned on $\mathbf{f}_{\mathcal{B}\mathcal{Y}}$:

$$\mathbf{f}_{\mathcal{X}}|\mathbf{f}_{\mathcal{B}\mathcal{Y}} \sim \mathcal{N}(\mathbf{P} \boldsymbol{\mu}_{\mathcal{B}\mathcal{Y}}, \mathbf{P} \boldsymbol{\Sigma}_{\mathcal{B}\mathcal{Y}} \mathbf{P}^T)$$

where \mathbf{P} is the projection matrix:

$$\mathbf{P} = \mathbf{P}_{\mathcal{X},\mathcal{B}\mathcal{Y}} = \mathbf{K}_{\mathcal{X},\mathcal{B}\mathcal{Y}} \mathbf{K}_{\mathcal{B}\mathcal{Y}}^{-1}$$

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

A.R.D

Sparsity

GP applications

Ref



Sparse representations – Motivation

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

Kernels II

Inf./Pred.

Regression

Approximations

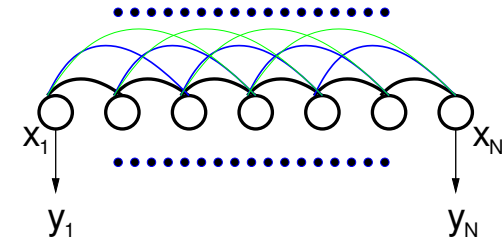
A.R.D

Sparsity

GP applications

Ref

Gaussian Processes are fully connected graphical models.



⇒ Computing estimates is difficult. E.g for the posterior mean:

$$\langle \mathbf{f}_{\mathcal{X}} \rangle_{\text{post}} = \mathbf{y}^T (\mathbf{K}_{\mathcal{N}} + \sigma_o^2 \mathbf{I}_{\mathcal{N}})^{-1} \mathbf{k}_{\mathcal{X}}$$

inversion requires $\mathcal{O}(N^3)$ time.

► Repr



Gaussian Regression

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

GP applications

Non-Gaussian Noise

Classification

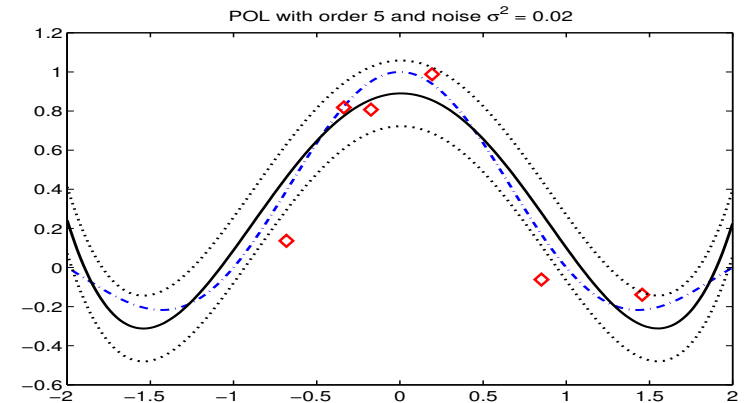
Multi-class

Inverse problems

Ref

Artificial data: $y = \sin(x)/x$ and polynomial kernel $K_0(x, x') = (1 + \mathbf{x}^T \mathbf{x}')^k$.

Number of training points: 1000 with added Gaussian noise $\sigma^2 = 0.02$





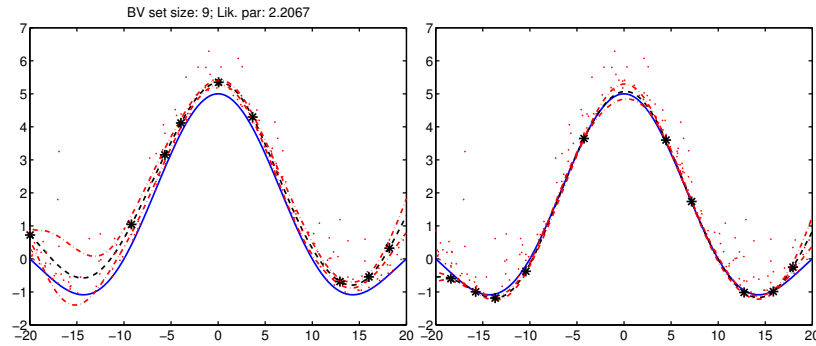
Robust one-sided regression

Functional Modelling

Lehel Csató

Exponential, one-sided, additive noise.

$$P(y|f_x) = \begin{cases} \lambda \exp[-\lambda(y - f_x)] & \text{if } y > f_x. \\ 0 & \text{otherwise.} \end{cases}$$



Modelling
Nonparametrics
SVM
GP's
GP applications
Non-Gaussian Noise
Classification
Multi-class
Inverse problems
Ref



Classification

Functional Modelling

Lehel Csató

For each location \mathbf{x} we have ± 1 . The likelihood function for this model is:

$$P(y|f(\mathbf{x})) = \text{Erf}\left(\frac{y f_x}{\sigma_0}\right)$$

Erf the incomplete Gaussian (\sim sigmoid):

$$\text{Erf}(x) = \int_{-\infty}^x dt \exp(-t^2/2)/\sqrt{2\pi}$$

- Posterior is **not** Gaussian.
- For **single** data, mean-var computable \Rightarrow iterative methods can be used.

Modelling
Nonparametrics
SVM
GP's
GP applications
Non-Gaussian Noise
Classification
Multi-class
Inverse problems
Ref



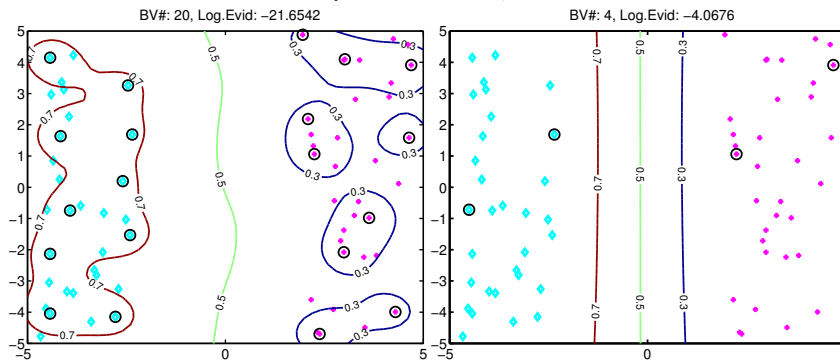
Toy Classification

Functional Modelling

Lehel Csató

$$\text{RBF kernel: } K(\mathbf{x}, \mathbf{x}') = \exp\left[-b - \sum_{i=1}^d (x_i - x'_i)^2 \beta_i\right]$$

behaviour of the ARD parameters β_i



Modelling
Nonparametrics
SVM
GP's
GP applications
Non-Gaussian Noise
Classification
Multi-class
Inverse problems
Ref



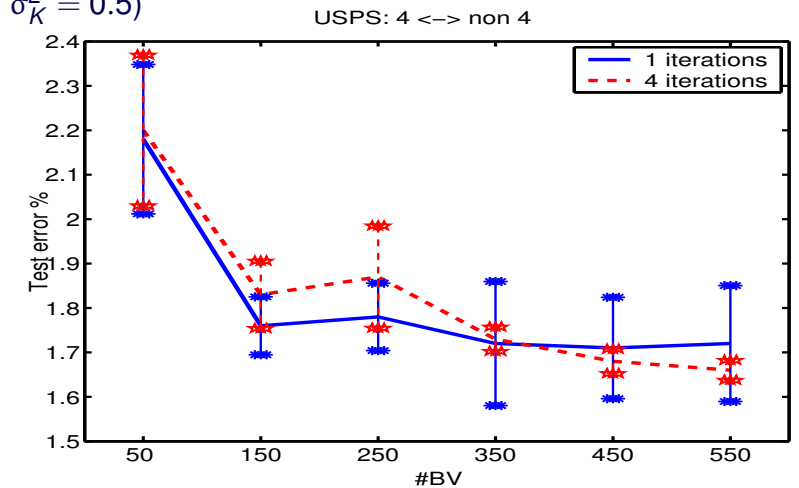
Classification

Functional Modelling

Lehel Csató

USPS data-set

Handwritten image data-set of gray-scale images with 7291 training and 2007 test patterns. (RBF kernel with $\sigma_K^2 = 0.5$)



Modelling
Nonparametrics
SVM
GP's
GP applications
Non-Gaussian Noise
Classification
Multi-class
Inverse problems
Ref

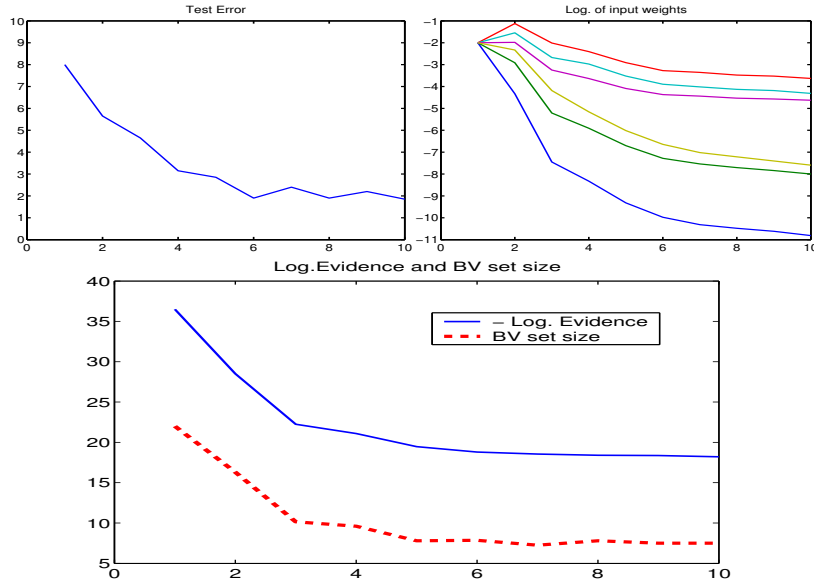


Crab data-set

Functional Modelling

Lehel Csató

Using RBF kernels



Modelling
 Nonparametrics
 SVM
 GP's
 GP applications
 Non-Gaussian Noise
 Classification
 Multi-class
 Inverse problems
 Ref



Multiclass Classification

Functional Modelling

Lehel Csató

Problem setup:

- For each location \mathbf{x} we have $y \in \{1, \dots, K\}$.
- Transforming it into $y \in \{0, 1\}^K$ Coding:

$$y = [0, \dots, 0, 1, 0, \dots]^T \quad \text{on the } k\text{-th position}$$

- K independent GP's are used. Indep. is **a-priori**.
- The **likelihood function** is:

$$P(y|f(\mathbf{x})) = \frac{y^T \mathbf{s}}{\mathbf{1}^T \mathbf{s}} \quad \text{where } \mathbf{s} = \exp([f_1(\mathbf{x}), \dots, f_K(\mathbf{x})]^T).$$

- The posterior processes are not independent.

Modelling
 Nonparametrics
 SVM
 GP's
 GP applications
 Non-Gaussian Noise
 Classification
 Multi-class
 Inverse problems
 Ref

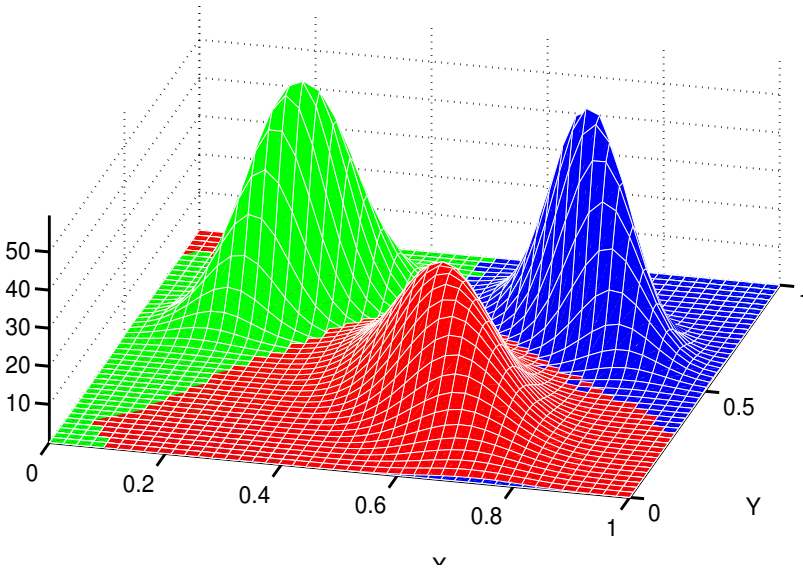


Multiclass Classification

Functional Modelling

Lehel Csató

Two-dimensional demo: class-conditional distributions



Modelling
 Nonparametrics
 SVM
 GP's
 GP applications
 Non-Gaussian Noise
 Classification
 Multi-class
 Inverse problems
 Ref

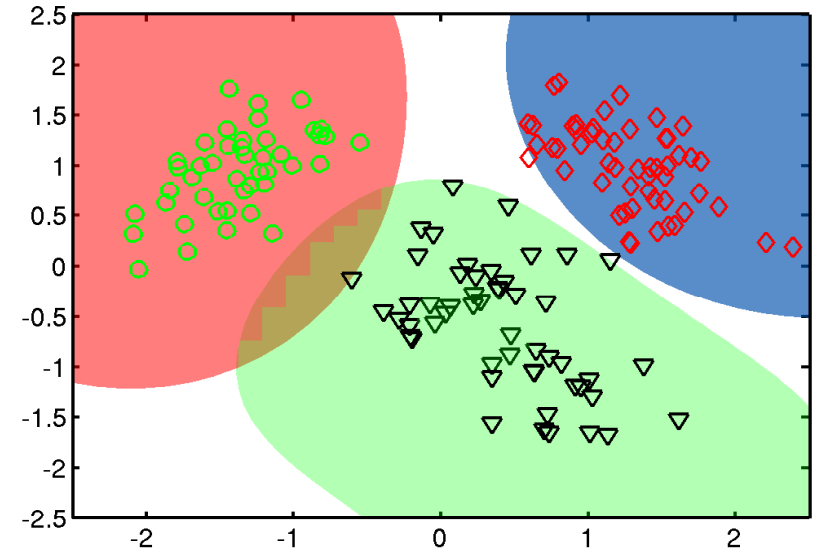


Multiclass Classification

Functional Modelling

Lehel Csató

Multiclass Classification



Modelling
 Nonparametrics
 SVM
 GP's
 GP applications
 Non-Gaussian Noise
 Classification
 Multi-class
 Inverse problems
 Ref



Modelling Inverse Problems

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Ref

Likelihood:

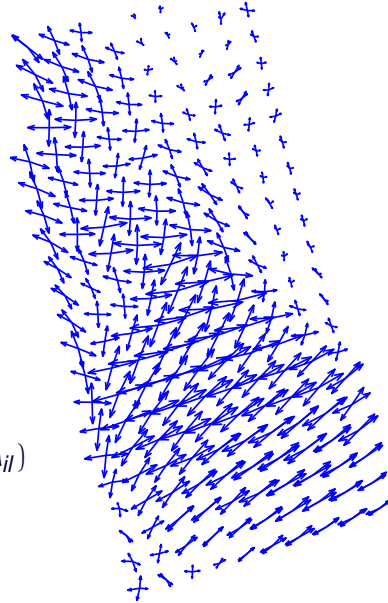
local observations of global wind-fields (u_i, v_i) .

Probabilistic framework preferred due to lack of direct observations.

Uncertainty captured in **Mixture density networks:**

$$P(u_i, v_i | obs) = \sum_{k=1}^4 \beta_{ik} \mathcal{N}(u_i, v_i | \mu_{ik}, A_{ik})$$

$\beta_{ik}, \mu_{ik}, A_{ik}$ local parameters.



Modelling Inverse Problems

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

GP applications

Non-Gaussian Noise

Classification

Multi-class

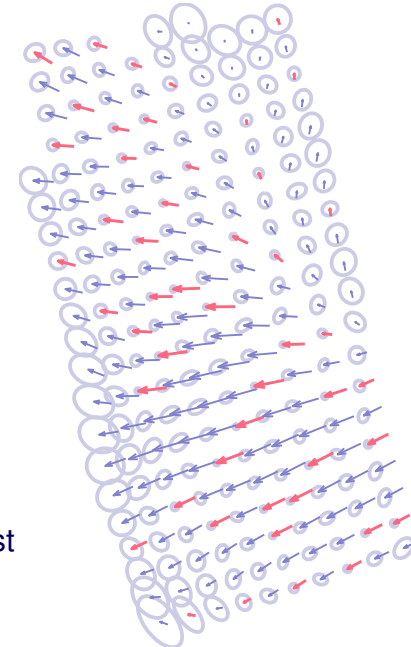
Inverse problems

Ref

- **Red Basis Vectors retained**

- Approximation preserves information about local uncertainty

- The inference process is fast



Bibliography I

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

GP applications

Ref

- **D.J.C. MacKay: *The evidence framework ...***

<http://wol.ra.phy.cam.ac.uk/mackay/class.nc.ps.gz>

H: <http://wol.ra.phy.cam.ac.uk/mackay/>

- **C. Rasmussen: *Gaussian Processes for Regression.***

<http://www.kyb.mpg.de/publications/pss/ps2468.ps>

H: <http://www.kyb.mpg.de/~carl>

- **C.K.I. Williams: *Prediction with Gaussian Processes: ...***

http://www.dai.ed.ac.uk/homes/ckiwi/postscript/NCRG_97_012.ps.gz

H: <http://www.dai.ed.ac.uk/homes/ckiwi>



Bibliography II

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

GP applications

Ref

- **R. Neal: *Priors for infinite networks***

<http://www.cs.utoronto.ca/~radford/ftp/pin.ps>

H: <http://www.cs.utoronto.ca/~radford/>

- **M. Seeger: *Bayesian Gaussian Process Models: ...***

<http://www.kyb.tuebingen.mpg.de/bs/people/seeger/papers/thesis.pdf>

H: <http://www.kyb.tuebingen.mpg.de/~seeger>

- **L. Csato: *Gaussian Processes: ... Sparse Approximations***

<http://www.ncrg.aston.ac.uk/Publications>

H: <http://www.cs.ubbcluj.ro/~csato1>



Bibliography III

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

GP applications

Ref

Gaussian Processes for Machine Learning



Carl Edward Rasmussen and Christopher K. I. Williams



NIPS'05 I

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

GP applications

Ref

NIPS'05 GP-related articles:

- J. Murillo-Fuentes, F. Perez-Cruz: Gaussian Processes for Multiuser Detection in CDMA Receivers
- Y. Shen, A. Ng, M. Seeger: Fast Gaussian Process Regression Using KD-Trees
- A. Shon, K. Grochow, A. Hertzmann, R. Rao: Gaussian Process CCA for Image Synthesis and Robotic Imitation
- R. Der, D. Lee: Beyond Gaussian Processes: On the Distributions of Infinite Networks
- D. Fleet, J. Wang, A. Hertzmann: Gaussian Process Dynamical Models



NIPS'05 II

Functional Modelling

Lehel Csató

Modelling

Nonparametrics

SVM

GP's

GP applications

Ref

NIPS'05 GP-related articles:

- Y. Engel, P. Szabo, D. Volkinshtein: Learning With Gaussian Process Temporal Difference Methods
- M. Kuss, C. Rasmussen: Assessing Approximations for Gaussian Process Classification
- E. Snelson, Z. Ghahramani: Sparse Parametric Gaussian Processes
- S. Kakade, M. Seeger, D. Foster: Worst-Case Bounds for Gaussian Process Models
- S. Keerthi, W. Chu: A Matching Pursuit Approach to Sparse Gaussian Process Regression
- E. Meeds, S. Osindero: An Alternative Infinite Mixture Of Gaussian Process Experts