



Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

Sparse GP inference

Gaussian Process Inference using sparse approximations

Lehel Csató

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

9 October 2018



Table of Contents

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- 1 Modelling using latent variables
- 2 Gaussian Processes
- 3 Sparse Representation
- 4 Applications of GP inference
- 5 Research questions



Contents

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- 1 Modelling using latent variables
- 2 Gaussian Processes
- 3 Sparse Representation
- 4 Applications of GP inference
- 5 Research questions



Modelling Data

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

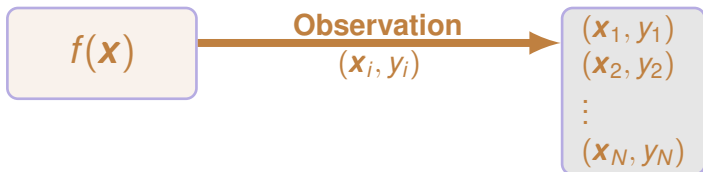
Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions



- 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

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

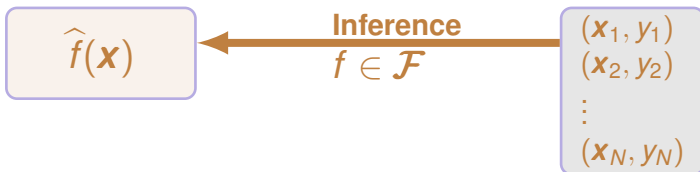
Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions



- **Data set** – collected.
- **We have to assume:**
 - a function class (polynomial, Fourier, Wavelet);
 - Observation process that **encodes** the noise;
- **Goal:** Find the optimal function from the class.



Latent variable models II

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- We have the **data set** $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$.
- We **assume** a function class:
 - 1 Linear or generalised linear models,
 - 2 Fourier expansions up to a given order.
- **Assume** an observation (or generation) process, eg:

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

$$y_n = \text{sign}(f(\mathbf{x}|\theta) + \epsilon)$$

Latent variables

Since $f(\mathbf{x}|\theta)$ is never seen, they are latent functions.



Contents

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- 1 Modelling using latent variables
- 2 Gaussian Processes
 - Inference and Prediction
 - Posterior Approximations
 - Prediction with Gaussian Processes
 - Optimising hyper-parameters
- 3 Sparse Representation
- 4 Applications of GP inference
- 5 Research questions



Sparse GPs

Léhel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

≡ functional latent variable models.

- GP models

≡ “simple” random functions.

- Appear in the likelihood:

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

Local dependencies only: $\mathbf{f} \longrightarrow \mathbf{f}_{\mathbf{x}}$



Gaussian processes I

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

Gaussian process: a “generalised” Gaussian distribution.

- \mathbf{f} Gaussian random function.

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

- \mathcal{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.:

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



Gaussian processes II

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

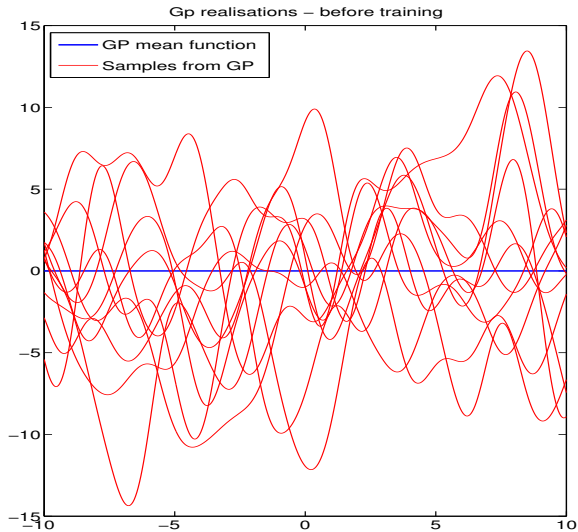
Inverse problems

Questions

Gaussian process: random function

“Parameters”

- mean function
- covariance kernel





Gaussian process parameters:

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- mean function – usually is 0.
- the **kernel** function, usually containing other parameters,

- Example:

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

- Always generate valid covariance matrices.

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

$$\mathbf{K}_{\mathcal{X}\mathcal{X}} = [K(\mathbf{x}_i, \mathbf{x}_j)]_{i,j=1}^T$$

is a positive definite matrix.



Gaussian Process Inference I

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- \mathcal{GP} inference is the application of Bayes' rule.

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

- For **any** set \mathcal{X} of inputs, the posterior distribution is:

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

where

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

is the data probability.

Problem: the posterior process cannot be analytically represented – not a \mathcal{GP} .



Motivating the approximations to the posterior

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- **If likelihood non-Gaussian** \Rightarrow posterior is not analytic.
(No “summarising” statistics)
- **Non-Gaussian** posterior process: we have to build approximations to them.
- **Representation** – How to represent the posterior?
Must be finite;



A Lemma for posterior processes

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

Property of Gaussian averages:

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

where $\alpha(i) = \partial_{\langle f_i \rangle_0} \ln P(\mathcal{D})$

Parametrisation lemma

We have a **parametrisation** similar to the Kimeldorf-Wahba parametrisation lemma, augmented with the kernel:



A Lemma for posterior processes

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

Property of Gaussian averages:

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

where $\alpha(i) = \partial_{\langle f_i \rangle_0} \ln P(\mathcal{D})$

Parametrisation lemma

We have a **parametrisation** similar to the Kimeldorf-Wahba parametrisation lemma, augmented with the kernel:

$$K_{\text{post}}(\mathbf{x}, \mathbf{x}') = K_0(\mathbf{x}, \mathbf{x}') + \sum_{i,j} K_0(\mathbf{x}, \mathbf{x}_i) C(ij) K_0(\mathbf{x}_j, \mathbf{x}') \\ \text{with } C(ij) = \partial_{\langle f_i \rangle_0, \langle f_j \rangle_0} \ln P(\mathcal{D})$$



Methods to approximate the posterior

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

Obtaining the posterior approximation:

- Sampling;
- Using modes – Laplace approximation;
- **KL-based approximation** – Expectation-propagation.

▶ skip others



Sampling from the posterior I

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

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:

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

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



Sampling from the posterior II

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

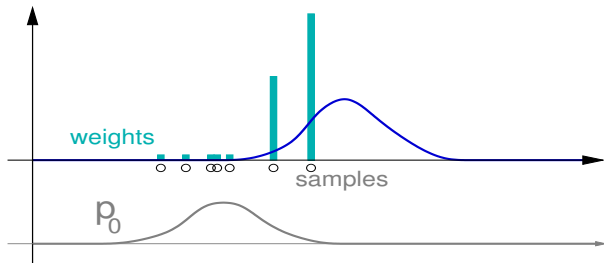
Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions



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 – I

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D.

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

Log-Posterior:

$$\log p_{\text{post}}(\mathbf{f}_{\mathcal{X}}, \mathbf{f}_{\mathcal{D}}) = K + \log P(\mathcal{D}|\mathbf{f}_{\mathcal{D}}) + \log p_0(\mathbf{f}_{\mathcal{D}}, \mathbf{f}_{\mathcal{X}})$$

Finding maximum of:

$$\hat{\mathbf{f}}_{\mathcal{D}} = \arg \max g(\mathbf{f}_{\mathcal{D}})$$

$$g(\mathbf{f}_{\mathcal{D}}) = \log P(\mathcal{D}|\mathbf{f}_{\mathcal{D}}) + \log p_0(\mathbf{f}_{\mathcal{D}})$$

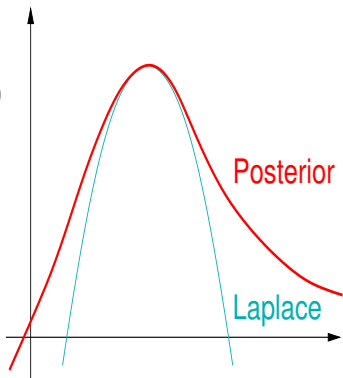
Taylor expansion around $\hat{\mathbf{f}}_{\mathcal{D}}$:

$$\hat{p}_{\text{post}}(\mathbf{f}_{\mathcal{X}}, \mathbf{f}_{\mathcal{D}}) \propto$$

$$\left(\hat{\mathbf{f}}_{\mathcal{D}} - \mathbf{f}_{\mathcal{D}}\right)^T \left[H_g(\hat{\mathbf{f}}_{\mathcal{D}})\right] \left(\hat{\mathbf{f}}_{\mathcal{D}} - \mathbf{f}_{\mathcal{D}}\right)$$

$$+ \log \log p_0(\mathbf{f}_{\mathcal{D}}, \mathbf{f}_{\mathcal{X}}) - \log p_0(\mathbf{f}_{\mathcal{D}})$$

⇒ Gaussian





Laplace approximation – II

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

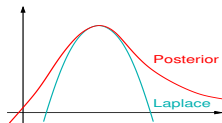
Classification

Multi-class

Inverse problems

Questions

Log-quadratic approximation:



$$\hat{\rho}_{\text{post}}(\mathbf{f}_X, \mathbf{f}_D) \propto \left(\hat{\mathbf{f}}_D - \mathbf{f}_D\right)^T \left[H_g(\hat{\mathbf{f}}_D)\right] \left(\hat{\mathbf{f}}_D - \mathbf{f}_D\right) + \log \rho_0(\mathbf{f}_D, \mathbf{f}_X) - \log \rho_0(\mathbf{f}_D)$$

Defines an approximation to the **likelihood**:

$$\hat{P}(D|\mathbf{f}_D) \propto \left(\hat{\mathbf{f}}_D - \mathbf{f}_D\right)^T \left[H_g(\hat{\mathbf{f}}_D)\right] \left(\hat{\mathbf{f}}_D - \mathbf{f}_D\right) - \log \rho_0(\mathbf{f}_D)$$

The approximation is **Gaussian**.



Expectation propagation – I

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

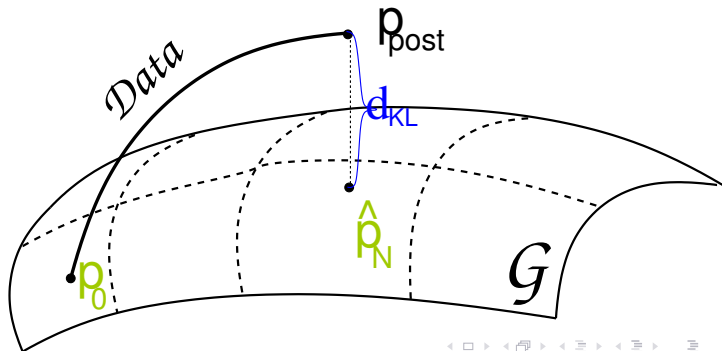
Multi-class

Inverse problems

Questions

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

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





Expectation propagation – II

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

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:

$$\begin{aligned} \langle \mathbf{f}_x \rangle_{\mathcal{GP}^*} &\stackrel{\text{def}}{=} \langle \mathbf{f}_x \rangle_{\text{post}} \\ K_{\mathcal{GP}^*}(\mathbf{x}, \mathbf{x}') &\stackrel{\text{def}}{=} K_{\text{post}}(\mathbf{x}, \mathbf{x}') \end{aligned}$$

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



Expectation propagation – III

Sparse GPs

Léhel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

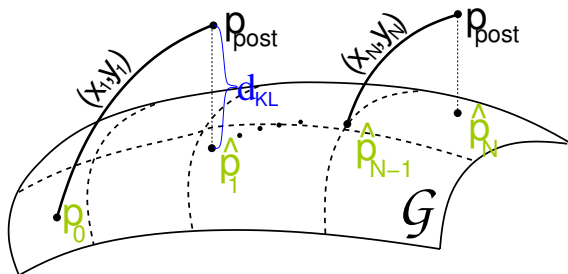
Inverse problems

Questions

Bayesian Online Learning (recursive)

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

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





Expectation propagation – IV

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

Inference is propagating the mean and the kernel:

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

with q (and r also) function of the **single** likelihood:

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

Expectation propagation:

Instead of \mathcal{GP}_{t+1} , at each step, an “actual” $\mathcal{GP}_{\setminus \ell}$, with $\ell = t + 1$, is computed.

$$\mathcal{GP}_{\setminus \ell} \stackrel{\text{def}}{=} \mathcal{GP}_{\mathcal{D}} / t_{\ell}(q_{\ell}, r_{\ell})$$

It is used as an actual prior in computing the updates.

Consequence: potential to several iterations.



Expectation propagation – V

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

Expectation propagation:

- Keeps the update parameters q_ℓ and r_ℓ for each data – called “sites”.
- Updates the “site parameters” using the steps:
 - 1 Removing the actual site from the posterior,
 - 2 Computing the new approximation and inferring the updated site parameters.
- The algorithm runs until convergence.
- It is stable for log-concave likelihood functions – proof?

► Unification Perspective



Prediction with Gaussian processes

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

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

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

where $\mathbf{f}^* = \mathbf{f}_{\mathbf{x}^*}$ – random variable “at” \mathbf{x}^* .



Optimising hyper-parameters I

Sparse GPs

Léhel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

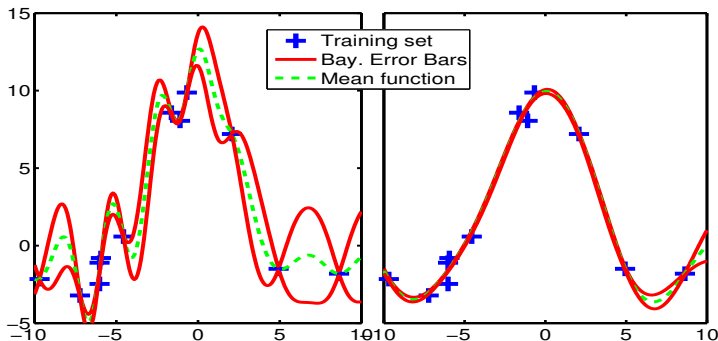
Inverse problems

Questions

GP kernel parameters \Leftrightarrow model choice.

Exemplu:

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



$\beta = 0.1$

$\beta = 1.$



Optimising hyper-parameters II

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

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



Contents

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- 1 Modelling using latent variables
- 2 Gaussian Processes
- 3 Sparse Representation**
 - Motivating Sparsity
 - Linear Algebra
- 4 Applications of GP inference
- 5 Research questions



Sparse representations – Motivation

Sparse GPs

Léhel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

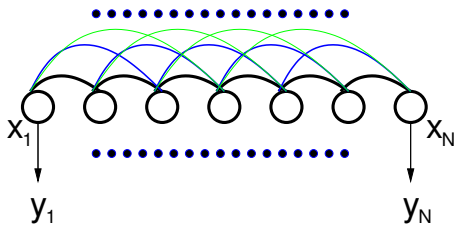
Classification

Multi-class

Inverse problems

Questions

Gaussian Processes are fully connected graphical models.



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

$$\langle \mathbf{f}_x \rangle_{\text{post}} = \mathbf{y}^T \left(\mathbf{K}_N + \sigma_o^2 \mathbf{I}_N \right)^{-1} \mathbf{k}_x$$

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



Motivating sparsity using Linear Algebra

1

Sparse GPs

Lehel Csató

Modelling using
latent variablesGaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

 \mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

$K_0(\mathbf{x}, \mathbf{x}')$ defines a *feature space* with
 $\phi_{\mathbf{x}}, \phi_{\mathbf{x}'} \in \mathcal{F}$ and $K_0(\mathbf{x}, \mathbf{x}') = \phi_{\mathbf{x}}^T \phi_{\mathbf{x}'}$.

Using \mathcal{F} and the scalar product:

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

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

**By identification, this means that \mathcal{GP} inference is
the estimation of a *Gaussian distribution in \mathcal{F} .***



Sparse GPs

Lehel Csató

Modelling using
latent variablesGaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

 \mathcal{GP} applications

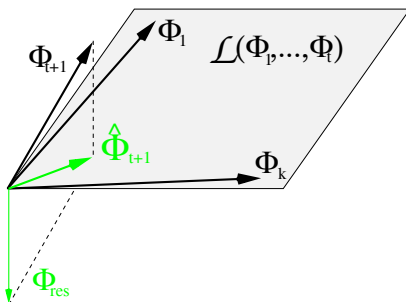
Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions



Projection

- is on the space of **inducing variables**;
- loss is measured using the KL-divergence between two posteriors:

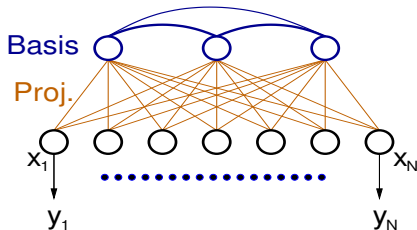
$$\text{KL}(\mathcal{GP}(\mathbf{f}_x) \parallel \mathcal{GP}(\mathbf{f}_x | \mathbf{f}_{\mathcal{B}V})) \quad (1)$$



Sparse representations – a solution

Condition **all training locations** on a **set of basis locations** or **inducing variables**¹.

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



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

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

where \mathbf{P} is the projection matrix:

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

¹Term used by Rasmussen, Candela, Snelson, etc.

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions



Sparse algorithm

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

The similarity with Gaussian distributions in an Euclidean space allowed to:

- Optimally remove any single component from the basis of expansion;
- Compute a “loss” associated to removing a component;
- Doing it whilst performing the other sparsification steps.

Results



Contents

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- 1 Modelling using latent variables
- 2 Gaussian Processes
- 3 Sparse Representation
- 4 Applications of GP inference
 - Regression with various noise models
 - Classification
 - Multi-class classification
 - Inverse problems
- 5 Research questions



Gaussian Regression

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

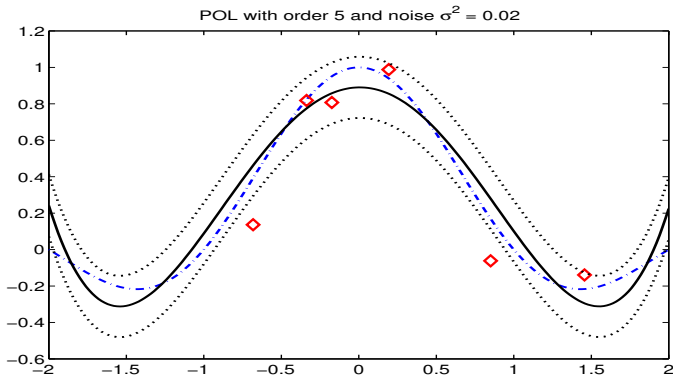
Multi-class

Inverse problems

Questions

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

Sparse GPs

Léhel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

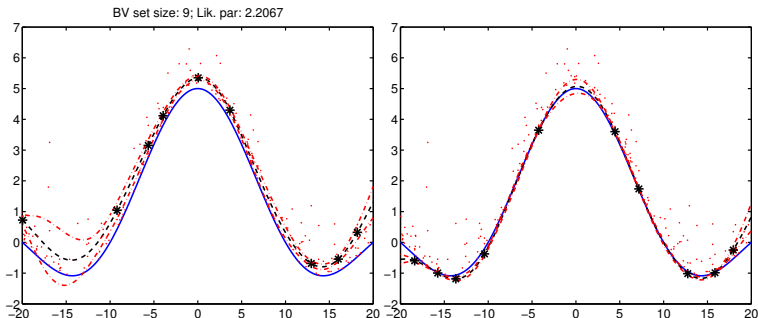
Multi-class

Inverse problems

Questions

Exponential, one-sided, additive noise.

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





Classification

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

For each input \mathbf{x} we have ± 1 and we use the **probit** likelihood:

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

Erf is the incomplete cumulative 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.



Toy Classification

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

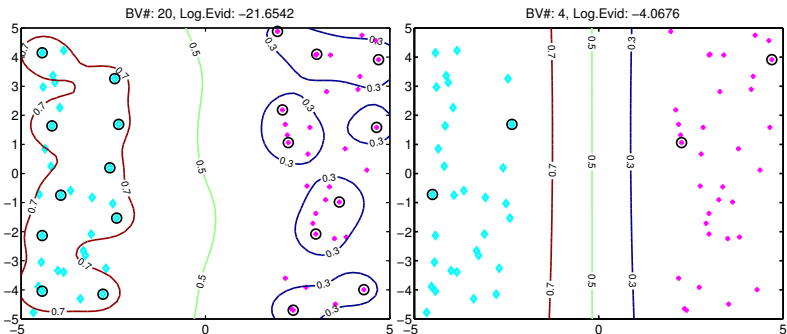
Multi-class

Inverse problems

Questions

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

behaviour of the ARD parameters β_i





Classification

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

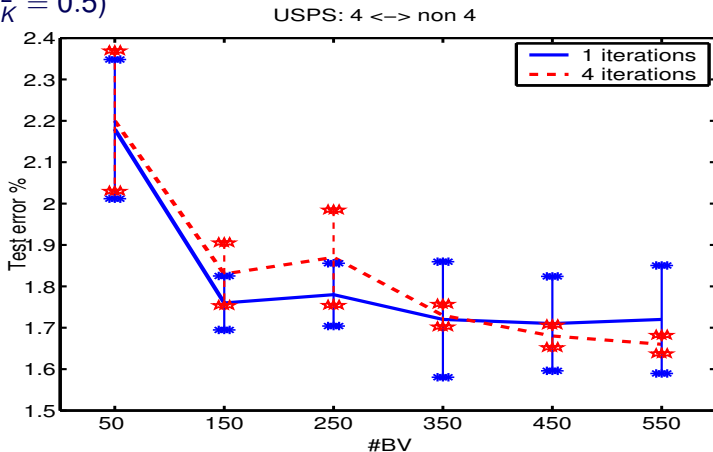
Multi-class

Inverse problems

Questions

USPS data-set

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





Crab data-set

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

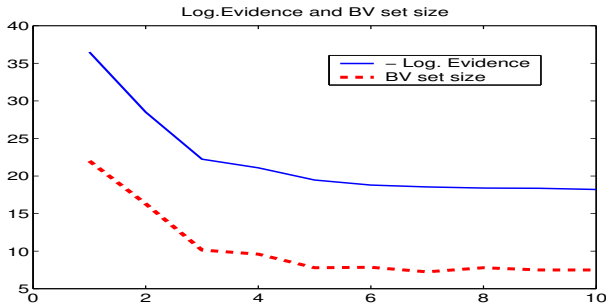
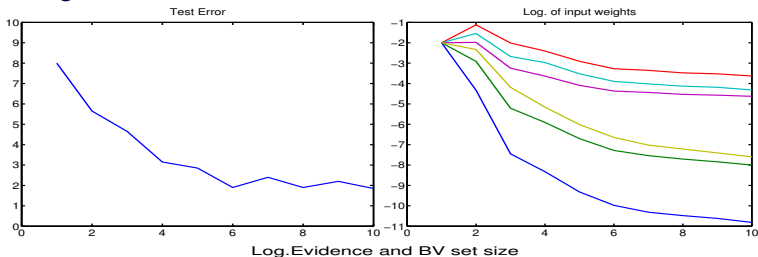
Classification

Multi-class

Inverse problems

Questions

Using RBF kernels





Multiclass Classification

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

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:

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

- The posterior processes are not independent.



Multiclass Classification

Sparse GPs

Léhel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

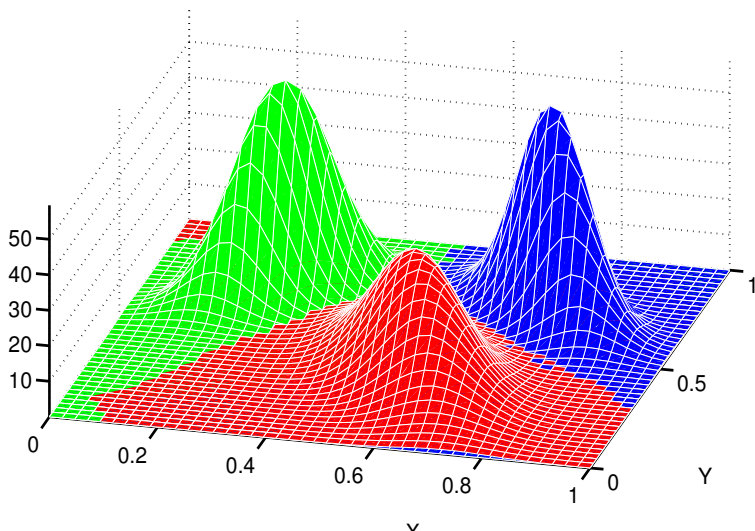
Classification

Multi-class

Inverse problems

Questions

Two-dimensional demo: class-conditional distributions





Multiclass Classification

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

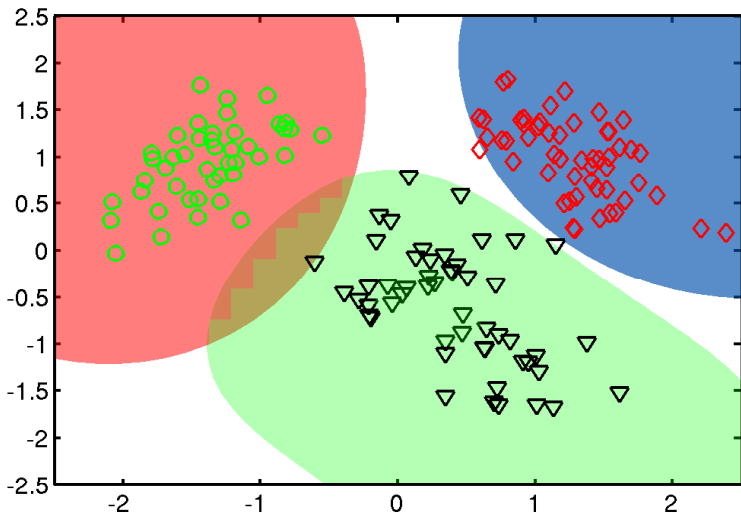
Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions





Sparse GPs

Léhel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

Likelihood:

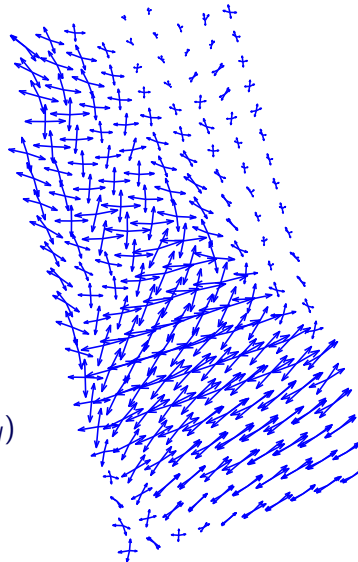
local observations of global wind-fields (u_i, v_i) .

Probabilistic framework preferred due to lack of direct observations.

Uncertainty captured in **Mixture density networks:**

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

$\beta_{ik}, \mu_{ik}, \mathbf{A}_{ik}$ local parameters.





Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

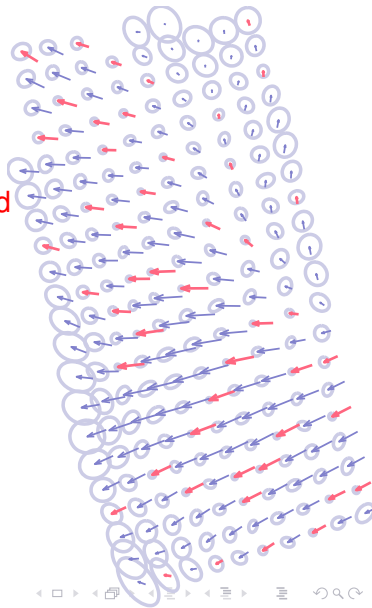
Classification

Multi-class

Inverse problems

Questions

- **Few Inducing Vectors retained**
- Approximation preserves information about local uncertainty
- The inference process is fast





Contents

Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- 1 Modelling using latent variables
- 2 Gaussian Processes
- 3 Sparse Representation
- 4 Applications of GP inference
- 5 Research questions**



Research questions

Sparse GPs

Lehel Csató

Modelling using latent variables

Gaussian Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

GP applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

- Does the EP algorithm converge?
 - An empirical observation was that for *log-concave* likelihood functions the EP produces the same results.
 - Performance did not depend on the data ordering.
 - *But EP* is a variational approximation.
- Is there a better way to optimise \mathcal{GP} parameters?
 - *Hyper-parameters* govern the quality of sparse \mathcal{GP} 's.
 - When \mathcal{BV} set size is restricted, the *locations* also become parameters to the MLII approximation.
 - Is the gradient descent still correct?

It is, but are the results correct?

Depends on the kernels \Rightarrow topology between inputs



Sparse GPs

Lehel Csató

Modelling using
latent variables

Gaussian
Processes

Inference

Approximations

Prediction

A.R.D

Sparsity

Motivating Sparsity

Linear Algebra

\mathcal{GP} applications

Non-Gaussian Noise

Classification

Multi-class

Inverse problems

Questions

Thank you!