# Pseudo-Marginal Bayesian Inference for Gaussian Processes
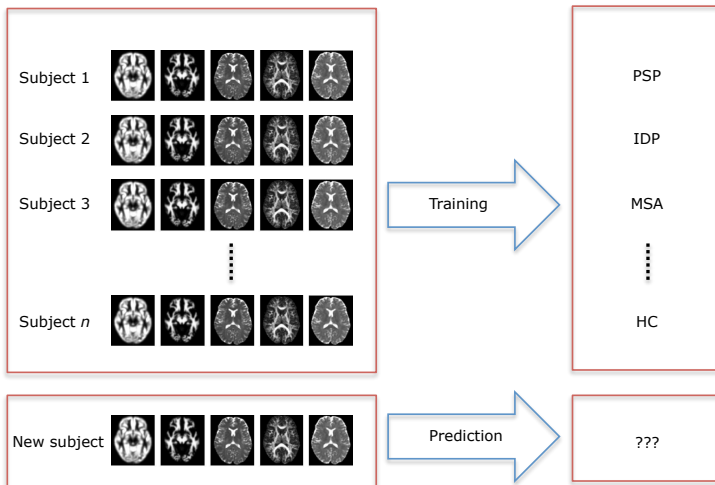
Maurizio Filippone

School of Computing Science
University of Glasgow
maurizio.filippone@glasgow.ac.uk
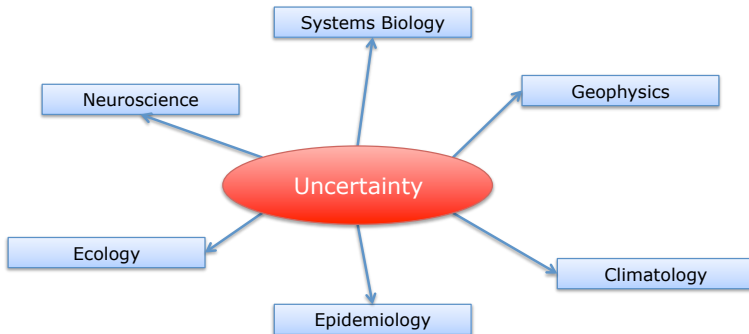
October 16th, 2014

# Motivating Application



Subject 1

Subject 2

Subject 3

Subject $n$

Training → PSP

IDP

MSA

HC

New subject

Prediction → ???

HC - Healthy control

MSA - Multiple system atrophy

PSP - Progressive Supranuclear Palsy

IDP - Idiopathic Parkinson's disease

- How do we estimate model parameters?
- How do we assess that a model is preferable over another?
- How do we incorporate knowledge by experts?
- How can we attach confidence intervals to our predictions and parameter estimates?

Probabilistic modeling offers an answer to these questions

- Data viewed as random variables



- Probabilities as degrees of belief

- Mapping input to labels

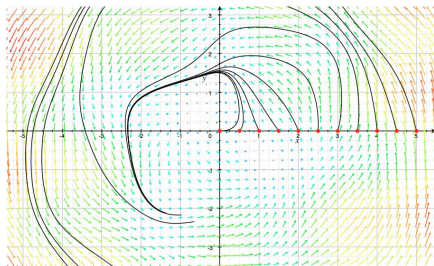$$\text{label} = \text{function}(\text{input}, \text{par}, \text{noise})$$

- Mapping input to labels

$$\text{label} = \text{function}(\text{input}, \text{par}, \text{noise})$$
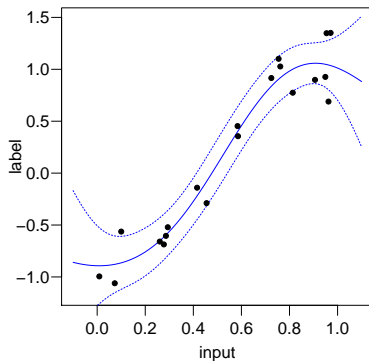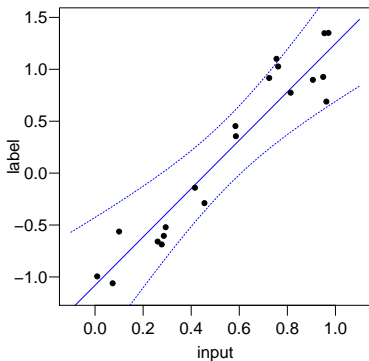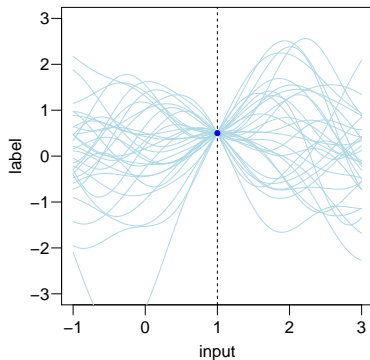
- function can describe a physical system

$$\frac{d}{dt} \left[ \begin{array}{c} x \\ y \end{array} \right] = \left[ \begin{array}{c} -y \\ x \end{array} \right] + \left[ \begin{array}{c} \cos(\text{par}_1 x) \\ \sin(\text{par}_2 y) \end{array} \right] - \text{par}_3 \left[ \begin{array}{c} x \\ y \end{array} \right]$$
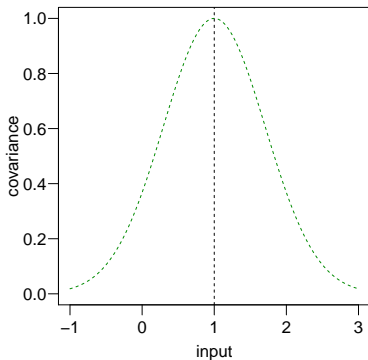
- Mapping input to labels

$$\text{label} = \text{function}(\text{input}, \text{par}, \text{noise})$$
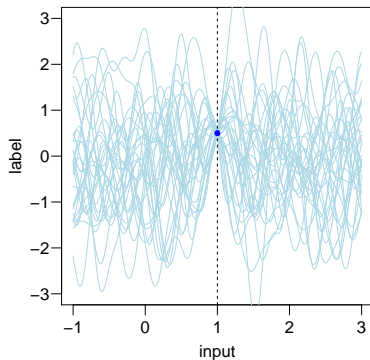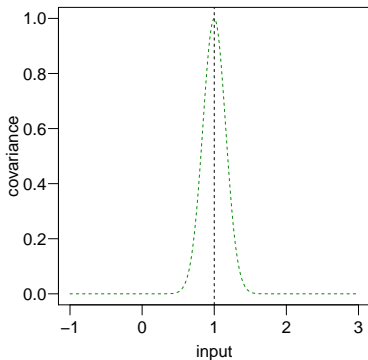
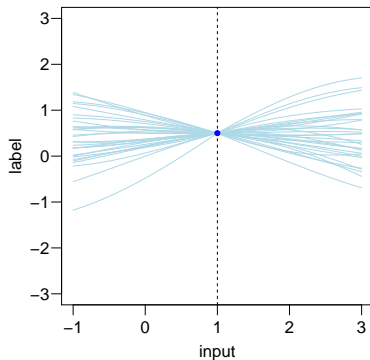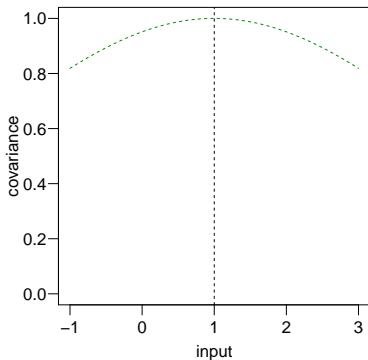- We may have no clue about function - we need assumptions

# Gaussian Processes

- Gaussians with distance dependent covariance

# Gaussian Processes

- Gaussians with distance dependent covariance

# Gaussian Processes

- Gaussians with distance dependent covariance

# Gaussian Process Models

- Class of hierarchical models
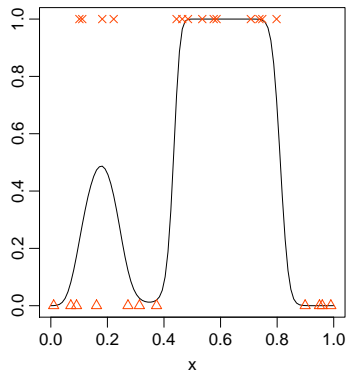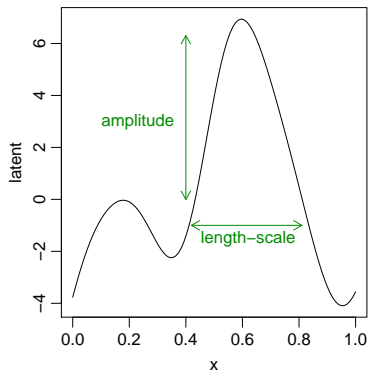
$$p(\mathrm{data}|\mathrm{latent}) \qquad p(\mathrm{latent}|\mathrm{par}) \qquad p(\mathrm{par})$$

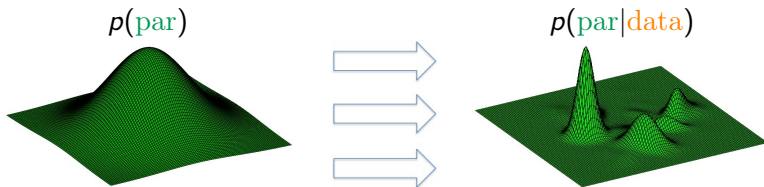- $p(\mathrm{latent}|\mathrm{par}) = \text{Gaussian Process}$

# Gaussian Process Models - Classification example

- Inference using Bayes theorem:

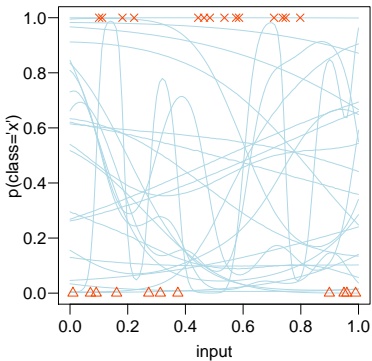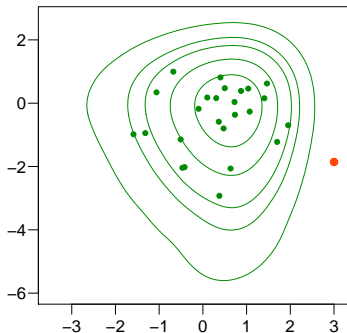$$p(\mathrm{par}|\mathrm{data}) = \frac{p(\mathrm{data}|\mathrm{par})p(\mathrm{par})}{\displaystyle\int p(\mathrm{data}|\mathrm{par})p(\mathrm{par})d\mathrm{par}}$$

$p(\mathrm{par})$

$p(\mathrm{par}|\mathrm{data})$

# Bayesian Inference and Predictions

- Predictions for new data

$$p(\text{label}_*|\text{label}) = \int p(\text{label}_*|\text{par})p(\text{par}|\text{label}) \, d\text{par}$$

- Monte Carlo integration:

$$\int p(\text{label}_*|\text{par})p(\text{par}|\text{label}) \, d\text{par} \simeq \frac{1}{N} \sum_{i=1}^{N} p(\text{label}_*|\text{par}^{(i)})$$

with $\text{par}^{(i)}$ drawn from $p(\text{par}|\text{data})$

- Draw samples according to the posterior density

- Bayesian inference

$$p(\mathrm{par}|\mathrm{data}) = \frac{p(\mathrm{data}|\mathrm{par})p(\mathrm{par})}{\displaystyle\int p(\mathrm{data}|\mathrm{par})p(\mathrm{par})d\mathrm{par}}$$

- Random walk sampler - accept a proposal with probability

$$\min\left(1, \frac{p(\mathrm{par}'|\mathrm{data})}{p(\mathrm{par}|\mathrm{data})}\right)$$

- Explore the parameter space according to the density

- Explore the parameter space according to the density

- Explore the parameter space according to the density

- Explore the parameter space according to the density

- Explore the parameter space according to the density

- Explore the parameter space according to the density

- Explore the parameter space according to the density

- Explore the parameter space according to the density

- Explore the parameter space according to the density

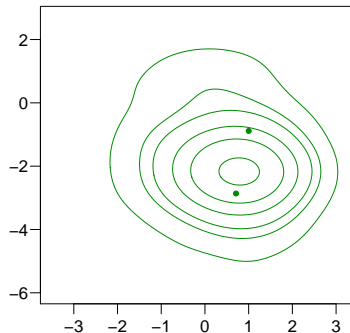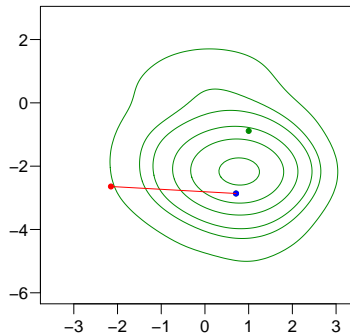- Explore the parameter space according to the density

Marginal likelihood

$$p(\mathrm{data}|\mathrm{par}) = \int p(\mathrm{data}|\mathrm{latent})p(\mathrm{latent}|\mathrm{par})d\mathrm{latent}$$

is unavailable analytically. Options:

- Approximate $p(\mathrm{data}|\mathrm{par})$ within MCMC
- Sample from $p(\mathrm{par}, \mathrm{latent}|\mathrm{data})$
- Pseudo-Marginal MCMC

# Gaussian Approximations for marginal likelihood

Gaussian approximation to $p(\text{latent}|\text{data}, \text{par})$

- Laplace Approximation
- Expectation Propagation
- Variational Bayes
- ...

- No exact Gibbs steps - need to employ Metropolis within Gibbs steps - waste of computations when rejecting
- Updates of $\mathrm{par}$ cost $O(n^3)$
- $\mathrm{par}$ can be large dimensional (e.g., Automatic Relevance Determination (ARD) covariance function)
- There are $n$ latent variables (as many as the number of observations)

# Challenges in MCMC for GPMs - Structure

Obvious iterative scheme (aka Sufficient Augmentation (SA) scheme). Alternate between:

- Drawing from $p(\text{latent}|\text{par}, \text{data})$
- Drawing from $p(\text{par}|\text{latent})$ (**bad idea** - see figure)

Ancillary Augmentation (AA) scheme - reparametrization:

$$K = LL^{\mathrm{T}} \qquad \text{ancillary} = L^{-1}\,\text{latent}$$

- Replace sampling of par with $p(\text{par}|\text{ancillary}, \text{data})$

# Mitigating coupling effect through reparameterization

Surrogate data model (SURR):

- Introduce set of auxiliary variables informed by the posterior over $\mathrm{latent}$



$$\mathrm{surrogate} = f(\mathrm{latent}, \mathrm{par})$$

- Replace posterior by unbiased estimate

$$\min\left(1, \frac{\tilde{p}(\mathrm{par}'|\mathrm{data})}{\tilde{p}(\mathrm{par}|\mathrm{data})}\right)$$

- Achieved by using an unbiased estimate of $\tilde{p}(\mathrm{data}|\mathrm{par})$

# Importance Sampling estimator

- Approximate posterior over latent variables using $q(\text{latent})$
- Then

$$\tilde{p}(\text{data}|\text{par}) = \frac{1}{N} \sum_{i=1}^{N} \frac{p(\text{data}|\text{latent}^{(i)})p(\text{latent}^{(i)}|\text{par})}{q(\text{latent}^{(i)})}$$

Abalone data set (two classes) $n = 2835$ - inference of length-scale

# Annealed Importance Sampling estimator

- Annealing from an approximation

- Annealing from an approximation



- Annealing from the prior

Analysis of the variance of the AIS and IS estimators

- $r$ is the variance of the $\log_{10}$ marginal likelihood

# Motivating Application



HC - Healthy control

MSA - Multiple system atrophy

PSP - Progressive Supranuclear Palsy

IDP - Idiopathic Parkinson's disease

- Multiclass classification based on GPs

$$p(\text{disease} = c|\text{sources}) = \text{unknown function}$$

- unknown function modeled using GPs

- Covariance based on source-dependent covariances $S_k$

$$\sum_{k=1}^{K} w_{ck} S_k(subject_i, subject_j)$$

# Parkinson syndromes data - multi source

| Method | Accuracy |
|--------|----------|
| GP classifier | 0.598 |
| SimpleMKL | 0.418 |

## Multiclass classification with multiple regions

Analysis of brain regions

1. brainstem
2. bilateral cerebellum
3. bilateral caudate
4. bilateral middle occipital gyrus
5. bilateral putamen
6. all other regions

# Multiclass classification with multiple regions

| Method | Accuracy |
|---|---|
| GP classifier | 0.614 |
| SimpleMKL | 0.229 |

# Brain Maps

- Gaussian Processes yield flexible and interpretable nonparametric models

- Gaussian Processes yield flexible and interpretable nonparametric models
- Bayesian inference to accurately quantifying uncertainty in such models

- Gaussian Processes yield flexible and interpretable nonparametric models
- Bayesian inference to accurately quantifying uncertainty in such models
- Pseudo-Marginal MCMC offers a practical way to carry out exact Bayesian computations

## Conclusions and ongoing work

- Gaussian Processes yield flexible and interpretable nonparametric models
- Bayesian inference to accurately quantifying uncertainty in such models
- Pseudo-Marginal MCMC offers a practical way to carry out exact Bayesian computations
- How to make exact Bayesian computations for Gaussian Processes scalable?

# Acknowledgements

- Dr Andre F. Marquand (Radboud University)
- Prof Mark Girolami (University of Warwick)
- Dr Guido Sanguinetti (University of Edinburgh)
- Dr Alessandro Vinciarelli (University of Glasgow)

[1] M. Filippone and M. Girolami. Pseudo-Marginal Bayesian inference for Gaussian processes, *IEEE Transactions on Pattern Analysis and Machine Intelligence, in press.*

[2] M. Filippone. Bayesian inference for Gaussian process classifiers with annealing and pseudo-marginal MCMC, In *ICPR, 2014.*

[3] M. Filippone et al. Probabilistic prediction of neurological disorders with a statistical assessment of neuroimaging data modalities. *Annals of Applied Statistics, 6(4):1883-1905, 2012.*

[4] A. F. Marquand et al. Automated, high accuracy classification of Parkinsonian disorders: a pattern recognition approach. *PLoS ONE, 2013.*

[5] M. Filippone et al. A comparative evaluation of stochastic-based inference methods for Gaussian process models. *Machine Learning, 93(1):93-114, 2013.*