

Scaling Bayesian inference for Gaussian processes using the Stochastic Gradient Langevin Dynamics algorithm

Maurizio Filippone

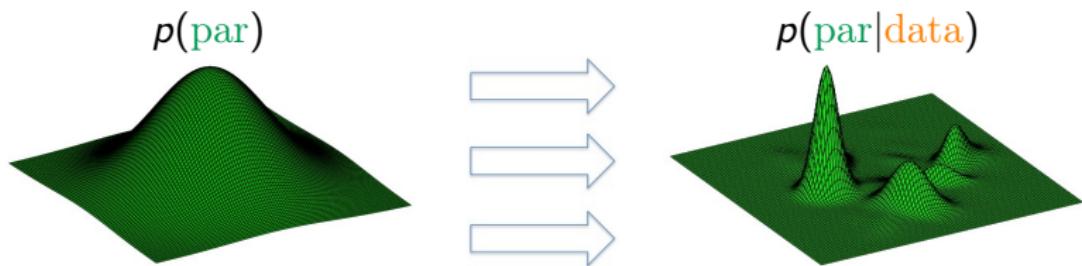
EURECOM, Sophia Antipolis, France
&
University of Glasgow, Glasgow, UK

maurizio.filippone@glasgow.ac.uk

June 15th, 2015

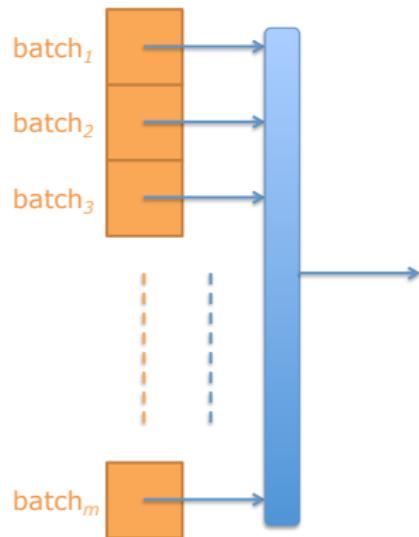
Bayesian Inference

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



Scaling Inference using Mini-Batches

$$p(\text{data} | \text{par}) = \prod_{i=1}^n p(\text{data}_i | \text{par})$$

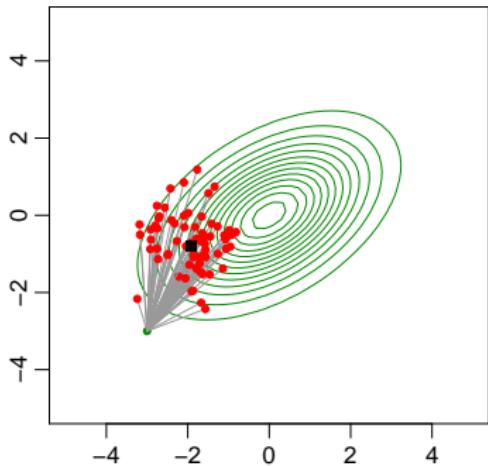
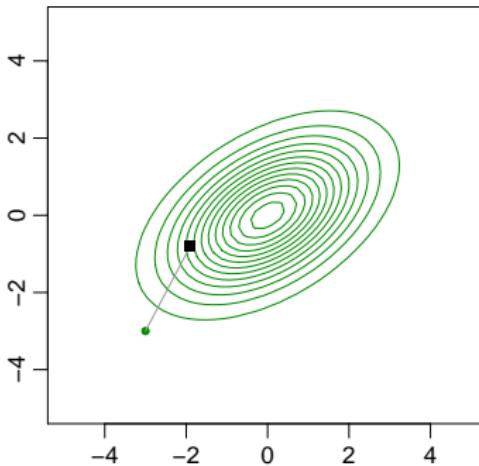


Gradient ascent

$$\text{par}' = \text{par} + \frac{\alpha}{2} \nabla_{\text{par}} \log[p(\text{data}|\text{par})p(\text{par})]$$

Stochastic Gradient ascent

$$E \left\{ \widetilde{\nabla_{\text{par}}} \log[p(\text{data}|\text{par})] \right\} = \nabla_{\text{par}} \log[p(\text{data}|\text{par})]$$



Robbins and Monro, AoMS, 1951

Stochastic Gradient ascent

$$\text{par}' = \text{par} + \frac{\alpha_t}{2} \widetilde{\nabla_{\text{par}}} \log[p(\text{data}|\text{par})p(\text{par})] \quad \alpha_t \rightarrow 0$$

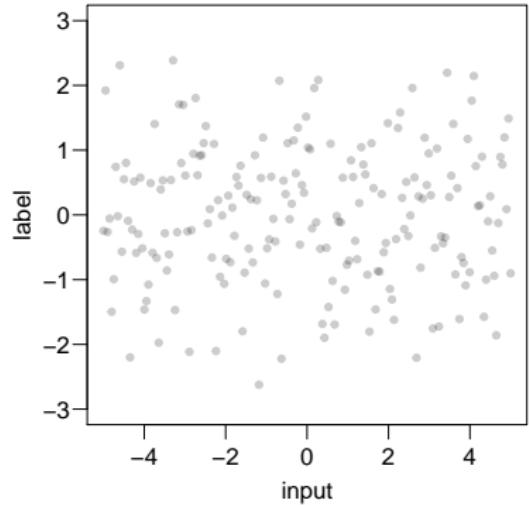
Robbins and Monro, AoMS, 1951

Stochastic Gradient Langevin Dynamics algorithm

$$\text{par}' = \text{par} + \frac{\alpha_t}{2} \widetilde{\nabla_{\text{par}}} \log[p(\text{data}|\text{par})p(\text{par})] + \eta_t \quad \eta_t \sim \mathcal{N}(0, \alpha_t)$$

What if $p(\text{data}|\text{par})$ does NOT factorize into $\prod_{i=1}^n p(\text{data}_i|\text{par})$?

Gaussian Processes

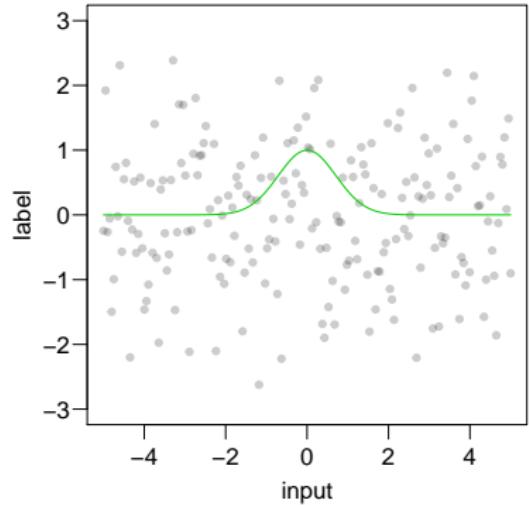


$$K = \begin{matrix} & \text{---} \\ \text{---} & \end{matrix}$$

A diagram illustrating a covariance matrix K . It is represented as a square grid of blue squares, with dashed lines indicating the boundaries. A dashed diagonal line runs from the top-left corner to the bottom-right corner, representing the main diagonal of the matrix.

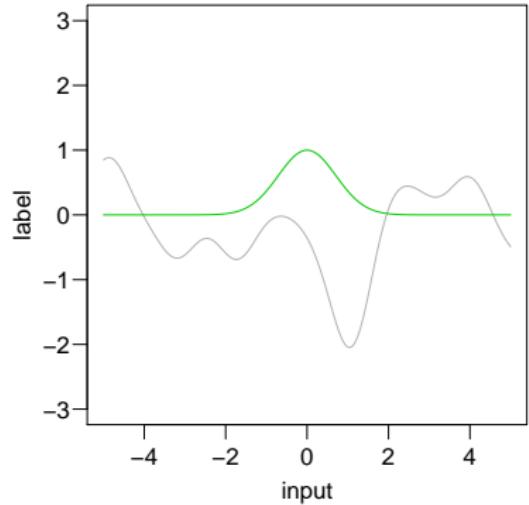
$K =$

Gaussian Processes



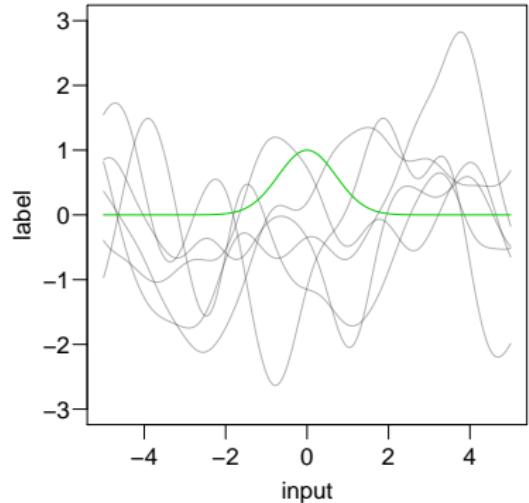
$$K = \begin{matrix} & \text{---} \\ \text{---} & \begin{matrix} \text{---} & \text{---} \\ \text{---} & \text{---} \end{matrix} \end{matrix}$$

Gaussian Processes



$$K = \begin{matrix} & \begin{matrix} \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \end{matrix} \\ \begin{matrix} \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \end{matrix} & \begin{matrix} \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \end{matrix} \end{matrix}$$

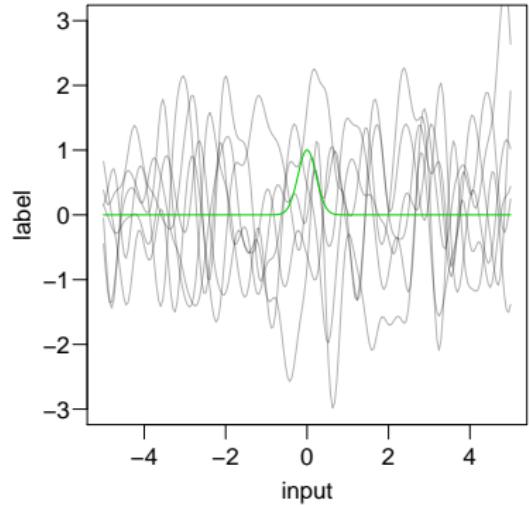
Gaussian Processes



$$K = \begin{matrix} & \begin{matrix} \text{---} & \end{matrix} \\ \begin{matrix} \text{---} & \end{matrix} & \begin{matrix} \text{---} & \end{matrix} \end{matrix}$$

A diagram illustrating the covariance matrix K of a Gaussian process. It shows a 4x4 grid of blue squares representing the matrix elements. Dashed lines indicate the connections between corresponding entries in the same row and column, representing the covariance between different input points.

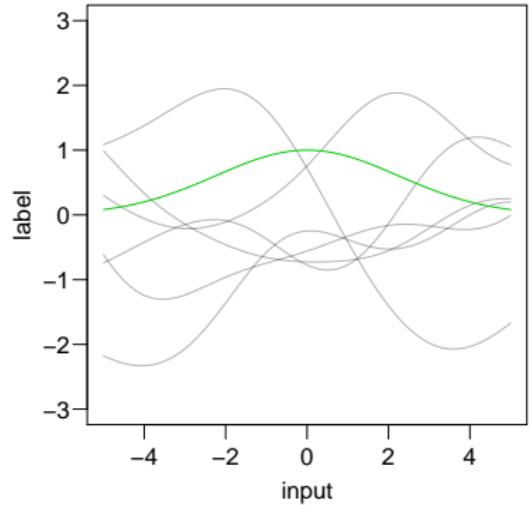
Gaussian Processes



$$K = \begin{matrix} \text{---} & \text{---} \\ \text{---} & \text{---} \end{matrix}$$

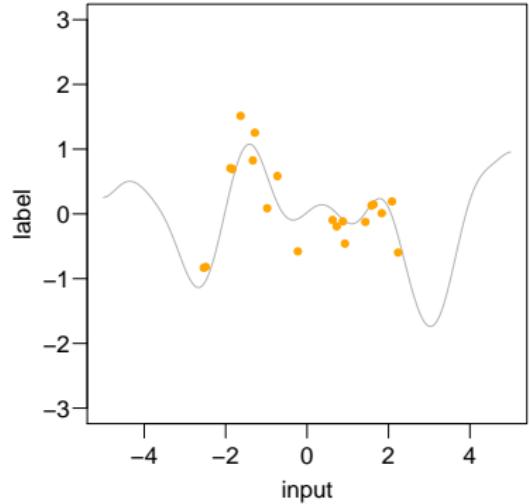
A diagram illustrating a covariance matrix K . It consists of four 3x3 grids of blue squares arranged in a 2x2 pattern. Dashed lines connect the centers of the top-left and bottom-right grids, and also connect the center of the top-left grid to the center of the bottom-right grid, forming a cross-like shape.

Gaussian Processes



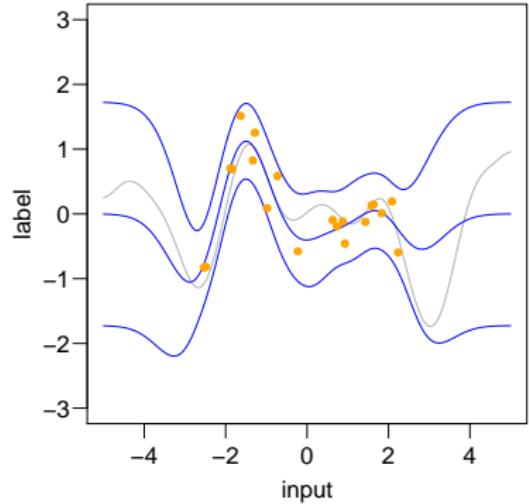
$$K = \begin{matrix} & \begin{matrix} \text{---} & \end{matrix} \\ \begin{matrix} \text{---} & \end{matrix} & \begin{matrix} \text{---} & \end{matrix} \end{matrix}$$

Gaussian Processes



$$K = \begin{matrix} & \star & \star \\ \star & \begin{matrix} \textcolor{blue}{\square} & \textcolor{orange}{\square} \\ \textcolor{orange}{\square} & \textcolor{blue}{\square} \end{matrix} & \begin{matrix} \textcolor{blue}{\square} & \textcolor{orange}{\square} \\ \textcolor{orange}{\square} & \textcolor{blue}{\square} \end{matrix} \\ \star & \begin{matrix} \textcolor{blue}{\square} & \textcolor{orange}{\square} \\ \textcolor{orange}{\square} & \textcolor{blue}{\square} \end{matrix} & \begin{matrix} \textcolor{blue}{\square} & \textcolor{orange}{\square} \\ \textcolor{orange}{\square} & \textcolor{blue}{\square} \end{matrix} \\ & \star & \star \end{matrix}$$

Gaussian Processes

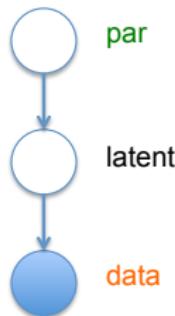


$$K = \underbrace{\begin{matrix} & & \\ & & \\ \text{---} & \text{---} & \text{---} \\ & & \\ & & \end{matrix}}_n \quad \underbrace{\begin{matrix} & & \\ & & \\ \text{---} & \text{---} & \text{---} \\ & & \\ & & \end{matrix}}_n$$

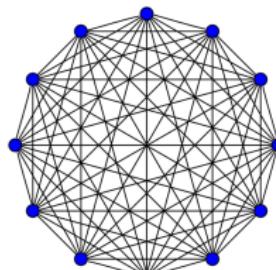
Gaussian Process models

- Gaussian Process models

$$p(\text{data} | \text{latent}) \quad p(\text{latent} | \text{par}) \quad p(\text{par})$$



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



Non-factorizing likelihood

- Marginal likelihood

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

can only be computed if $p(\text{data}|\text{latent})$ is Gaussian

- ... even then

$$\log[p(\text{data}|\text{par})] = -\frac{1}{2} \log |K| - \frac{1}{2} \mathbf{y}^T K^{-1} \mathbf{y} + \text{const.}$$

where $K = K(\text{par})$ is an $n \times n$ dense matrix!

Stochastic Gradients in GP regression

- Marginal likelihood

$$\log[p(\text{data}|\text{par})] = -\frac{1}{2} \log |K| - \frac{1}{2} \mathbf{y}^T K^{-1} \mathbf{y} + \text{const.}$$

- Derivatives wrt par

$$\frac{\partial \log[p(\text{data}|\text{par})]}{\partial \text{par}_i} = -\frac{1}{2} \text{Tr} \left(K^{-1} \frac{\partial K}{\partial \text{par}_i} \right) + \frac{1}{2} \mathbf{y}^T K^{-1} \frac{\partial K}{\partial \text{par}_i} K^{-1} \mathbf{y}$$

Stochastic Gradients in GP regression

- Stochastic estimate of the trace

$$\text{Tr} \left(K^{-1} \frac{\partial K}{\partial \text{par}_i} \right) = \text{Tr} \left(K^{-1} \frac{\partial K}{\partial \text{par}_i} E[\mathbf{r}\mathbf{r}^T] \right) = E \left[\mathbf{r}^T K^{-1} \frac{\partial K}{\partial \text{par}_i} \mathbf{r} \right]$$

with $E[\mathbf{r}\mathbf{r}^T] = I$ - e.g., r_j drawn from $\{-1, 1\}$ with $p = 1/2$

Stochastic Gradients in GP regression

- Stochastic estimate of the trace

$$\text{Tr} \left(K^{-1} \frac{\partial K}{\partial \text{par}_i} \right) = \text{Tr} \left(K^{-1} \frac{\partial K}{\partial \text{par}_i} E[\mathbf{r}\mathbf{r}^T] \right) = E \left[\mathbf{r}^T K^{-1} \frac{\partial K}{\partial \text{par}_i} \mathbf{r} \right]$$

with $E[\mathbf{r}\mathbf{r}^T] = I$ - e.g., r_j drawn from $\{-1, 1\}$ with $p = 1/2$

- Stochastic gradient

$$-\frac{1}{2N_r} \sum_{i=1}^{N_r} \mathbf{r}^{(i)T} K^{-1} \frac{\partial K}{\partial \theta_i} \mathbf{r}^{(i)} + \mathbf{y}^T K^{-1} \frac{\partial K}{\partial \theta_i} K^{-1} \mathbf{y}$$

- Only linear systems!

Solving linear systems

- Linear systems:

$$Ks = b$$

- Can be solved using conjugate gradient:

$$s = \arg \min_x \left(\frac{1}{2} x^T K x - x^T b \right)$$

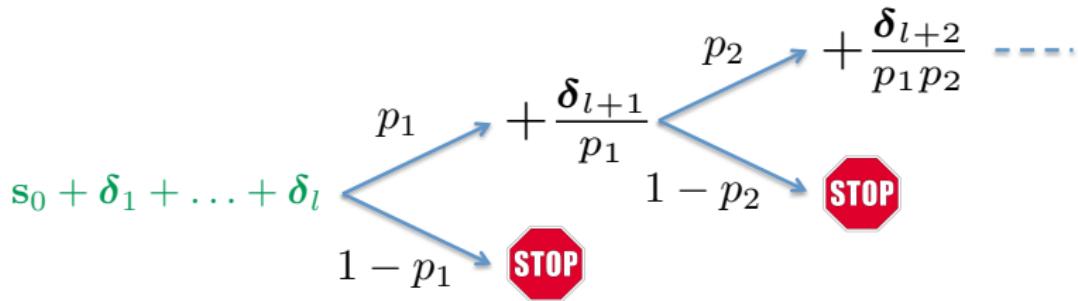
- Iterative update $s = s_0 + \delta_1 + \dots + \delta_T$
- Requires only Kv multiplications! $O(n^2)$ time
- No need to store K ! $O(n)$ space

- Accelerate the solution of dense linear systems
- ... returning an unbiased estimate of the solution

- Full CG solution:

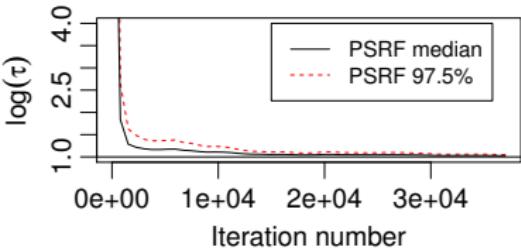
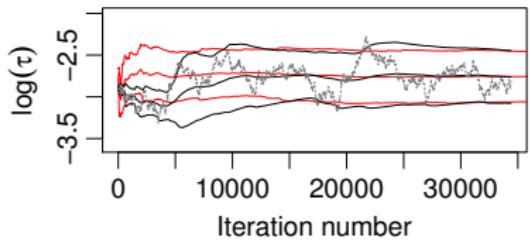
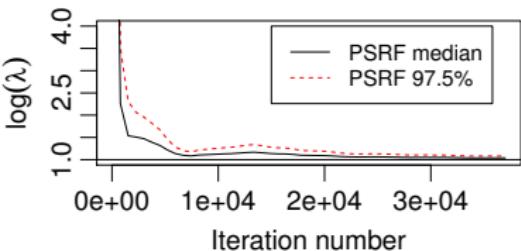
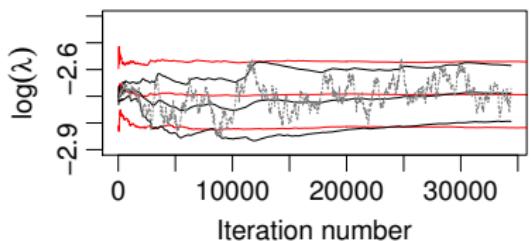
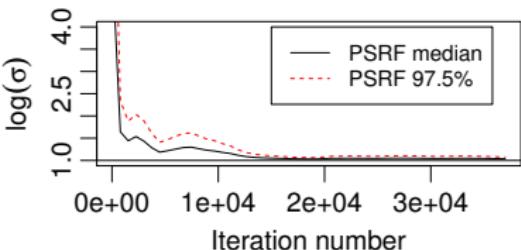
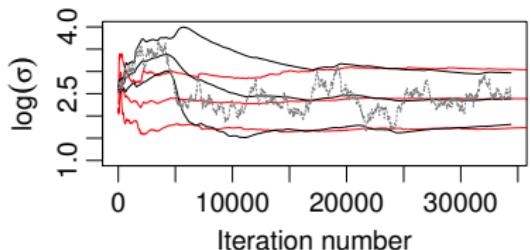
$$\mathbf{s} = \mathbf{s}_0 + \boldsymbol{\delta}_1 + \dots + \boldsymbol{\delta}_l + \boldsymbol{\delta}_{l+1} \dots + \boldsymbol{\delta}_T$$

- ULISSE:

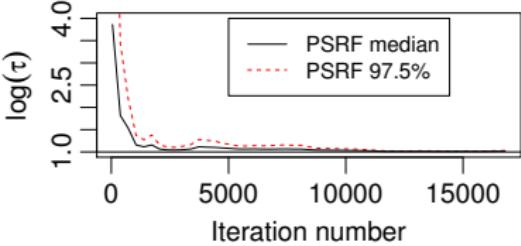
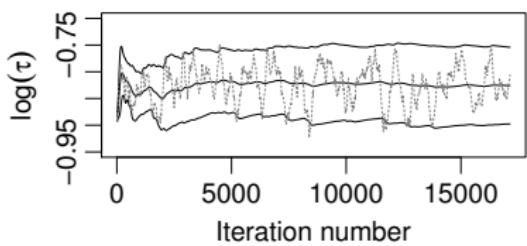
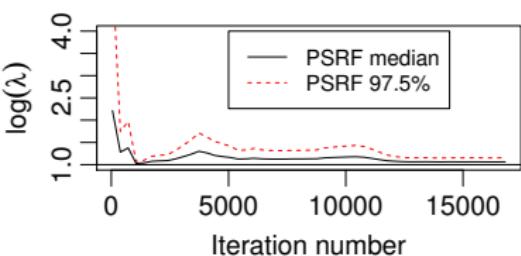
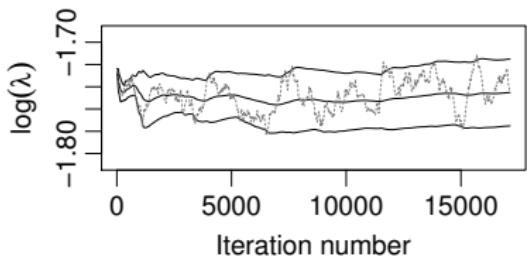
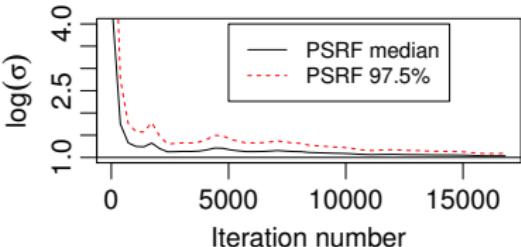
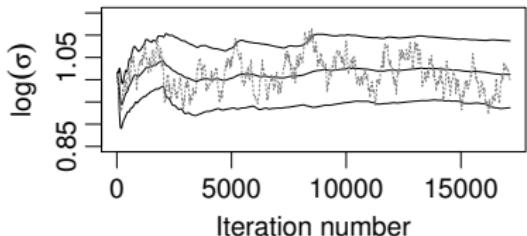


- Final solution is an unbiased estimate of \mathbf{s} !

Comparison with MCMC - Concrete dataset - $n \approx 1K$



Larger n - Census dataset - $n \approx 23K$



Conclusions and ongoing work

- Gaussian Processes yield flexible and interpretable nonparametric models
 - Bayesian inference to accurately quantifying uncertainty in GP models
 - “Noisy” MCMC offers a practical and scalable way to carry out “exact” Bayesian computations for GPs

Acknowledgements & References



Andre Marquand
Radboud



Guido Sanguinetti
Edinburgh



James Hensman
Sheffield



Mark Girolami
Warwick



Alessandro Vinciarelli
Glasgow



Dirk Husmeier
Glasgow

- [1] M. Filippone and R. Engler. Enabling scalable stochastic gradient-based inference for Gaussian processes by employing the Unbiased Linear System SolvEr (ULISSE), In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, July 6-11, 2015.* 2015.
- [2] M. Filippone and M. Girolami. Pseudo-Marginal Bayesian inference for Gaussian processes, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(11):2214-2226, 2014.
- [3] M. Filippone. Bayesian inference for Gaussian process classifiers with annealing and pseudo-marginal MCMC, In *Proceedings of the 22nd International Conference on Pattern Recognition, ICPR 2014, Stockholm, Sweden, August 24-28, 2014, pages 614-619.* IEEE, 2014.
- [4] 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.
- [5] A. F. Marquand et al. Automated, high accuracy classification of Parkinsonian disorders: a pattern recognition approach. *PLoS ONE*, 8(7):e69237+, 2013.
- [6] M. Filippone et al. A comparative evaluation of stochastic-based inference methods for Gaussian process models. *Machine Learning*, 93(1):93-114, 2013.