

Enabling scalable stochastic gradient-based inference for Gaussian processes by employing the Unbiased Linear System SolvEr (ULISSE)

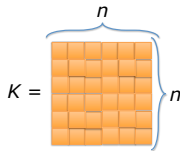
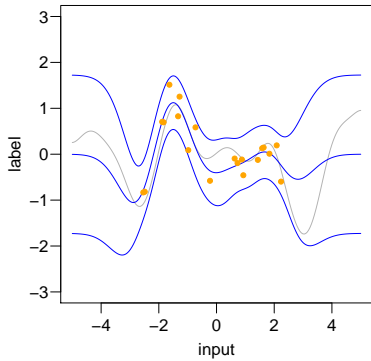
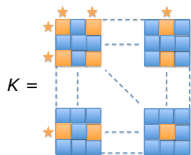
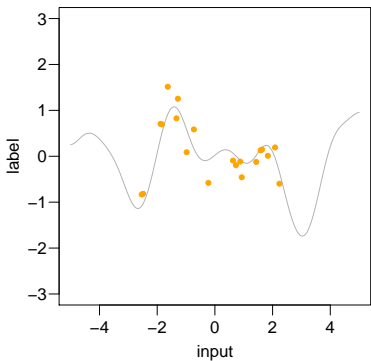
Maurizio Filippone

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

`Maurizio.Filippone@eurecom.fr`

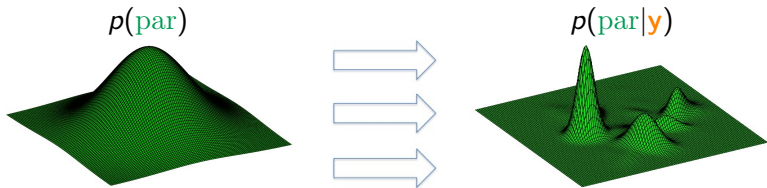
ICML 2015, Lille
July 9th, 2015

Gaussian Processes



Bayesian Inference

- Inputs = X Labels = y
- $K = K(X, \text{par})$



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

Markov chain Monte Carlo - Random walk example

$$\text{Acceptance probability : } \min \left(1, \frac{p(\mathbf{y}|\text{par}')p(\text{par}')}{p(\mathbf{y}|\text{par})p(\text{par})} \right)$$

Metropolis et al., *JoCP*, 1953 - Hastings, *Biometrika*, 1970

- Gaussian likelihood case

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

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

Should we care about inference of covariance parameters?

MCMC for Variationally Sparse Gaussian Processes

James Hensman
Department of Computer Science
University of Sheffield
Sheffield, UK

Alexander G. de G. Matthews
Department of Engineering
University of Cambridge
Cambridge, UK

Maurizio Filippone
School of Computing Science
University of Glasgow
Glasgow, UK

Zoubin Ghahramani
Department of Engineering
University of Cambridge
Cambridge, UK

Gaussian process (GP) models for
sideable research effort has been
to compute efficiently when the
when the likelihood is not Gaussian
posterior. This paper simultaneously
to the posterior which is sparse in
result is a Hybrid Monte-Carlo sam-
imation over the function values at
computations based on inductive
this paper will be available shortly.

1. Introduction

Gaussian process models are attract
parametric nature. By combining a
machine learning tasks can be tackle
to consider when using a GP model
the likelihood is non-Gaussian), con-
trix, which scales poorly in the num-
covariance function parameters. A
for efficient computation when the
retaining a sub-set of the data [2, 3]
point approach, where the model is

2214

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, VOL. 36, NO. 11, NOVEMBER 2014

Pseudo-Marginal Bayesian Inference for Gaussian Processes

Maurizio Filippone and Mark Girolami

Abstract—The main challenges that arise when approximating Gaussian process priors in probabilistic models are how to carry out exact Bayesian inference and how to account for uncertainty. Using pseudo regression as an illustrative work, pseudo-marginal approach to Markov chain Monte Carlo paper show improvements over existing sampling covariance function of the Gaussian Process prior. Inference of Gaussian Process based hierarchies at integration of all model parameters is actually less predictions. Extensive comparisons with respect to

Index Terms—Hierarchic Bayesian models, Gaussian processes, approximate Bayesian inference

1 INTRODUCTION

NON-PARAMETRIC or kernel based models
successful class of statistical modelling
methods. To focus ideas throughout the paper
the working example of predictive classification
the methodology presented, however, is appli-
hierarchical Bayesian models in general and thus
Gaussian process (GP) priors in particular. Im-
ples of kernel-based classifiers are the sup-
machine (SVM) [1], [2], the relevance vector ma-
[3], [4], and the Gaussian process classifier [5].
These classifiers are based on different model-
tions and paradigms of statistical inference, the
terized by a kernel function or covariance of
allows one to build nonlinear classifiers able to
lengthy problems [6], [7], [8], [9], [10], [11].

In order to allow these classifiers to be flexi-
sary to parameterize the kernel (or covariance) f-
set of so called hyper-parameters. After observ-
training data, the aim is to estimate or infer
parameters. In the case of SVMs point estimat-
parameters are obtained by optimizing a cost
error. This makes optimization viable only in
very few hyper-parameters, as grid search
employed, and is limited by the available sam-
ple classification instead the probabilistic

The Annals of Applied Statistics
2012, Vol. 6, No. 4, 1883–1905
DOI: 10.1214/12-AAS023
© Institute of Mathematical Statistics, 2012

PROBABILISTIC PREDICTION OF NEUROLOGICAL DISORDERS WITH A STATISTICAL ASSESSMENT OF NEUROIMAGING DATA MODALITIES

BY M. FILIPPONE, A. F. MARQUAND¹, C. R. V. BLAIN, S. C. R. WILLIAMS,
J. MOURÃO-MIRANDA² AND M. GIROLAMI³

University of Glasgow, King's College London, King's College London, King's
College London, King's College London and University College London

For many neurological disorders, prediction of disease state is an im-
portant clinical aim. Neuroimaging provides detailed information about brain
structure and function from which such predictions may be statistically de-
rived. A multinomial logit model with Gaussian process priors is proposed to:
(i) predict disease state based on whole-brain neuroimaging data and
(ii) analyze the relative informativeness of different image modalities and
brain regions. Advanced Markov chain Monte Carlo methods are employed
to perform posterior inference over the model. This paper reports a statistical
assessment of multiple neuroimaging modalities applied to the discrimination
of three Parkinsonian neurological disorders from one another and healthy
controls, showing promising predictive performance of disease states when
compared to nonprobabilistic classifiers based on multiple modalities. The
statistical analysis also quantifies the relative importance of different neu-
roimaging measures and brain regions in discriminating between these dis-
eases and suggests that for prediction there is little benefit in acquiring mul-
tiple neuroimaging sequences. Finally, the predictive capability of different

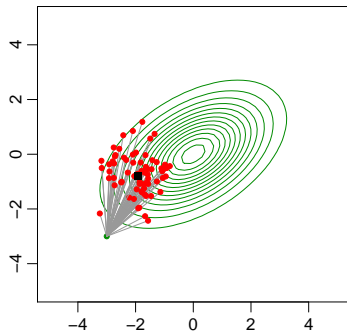
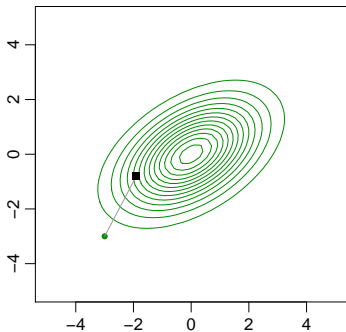
arXiv:1506.04000v1 [stat.ML] 12 Jun 2015



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

Stochastic Gradient ascent

$$\mathbb{E} \left\{ \widetilde{\nabla}_{\text{par}} \log[p(\mathbf{y}|\text{par})] \right\} = \nabla_{\text{par}} \log[p(\mathbf{y}|\text{par})]$$



Robbins and Monro, *AoMS*, 1951

Stochastic Gradient ascent

$$\text{par}' = \text{par} + \frac{\alpha_t}{2} \widetilde{\nabla}_{\text{par}} \log[p(\mathbf{y}|\text{par})p(\text{par})] \quad \alpha_t \rightarrow 0$$

Robbins and Monro, *AoMS*, 1951

Stochastic Gradient Langevin Dynamics (SGLD) algorithm

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

- Traditionally, in SGLD stochastic gradients

$$\widetilde{\nabla}_{\text{par}} \log[p(\mathbf{y}|\text{par})p(\text{par})]$$

are computed based on mini-batches of data

- In GPs the likelihood DOES NOT factorize
- What can we do?

- Marginal likelihood

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

- Derivatives wrt par

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

- Stochastic estimate of the trace

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

with $\mathbb{E}[\mathbf{r}\mathbf{r}^T] = I$

- For example r_j drawn from $\{-1, 1\}$ with $p = 1/2$

- Stochastic estimate of the trace

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

with $\mathbb{E}[\mathbf{r}\mathbf{r}^T] = I$

- For example r_j drawn from $\{-1, 1\}$ with $p = 1/2$
- Stochastic gradient

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

- **Linear systems only!**

- Linear systems:

$$K\mathbf{s} = \mathbf{b}$$

- Can be solved using conjugate gradient:

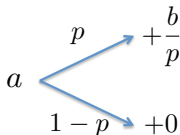
$$\mathbf{s} = \arg \min_{\mathbf{x}} \left(\frac{1}{2} \mathbf{x}^T K \mathbf{x} - \mathbf{x}^T \mathbf{b} \right)$$

- Iterative update $\mathbf{s} = \mathbf{s}_0 + \boldsymbol{\delta}_1 + \dots + \boldsymbol{\delta}_T$
- Requires only $K\mathbf{v}$ multiplications! $O(n^2)$ time
- No need to store K ! $O(n)$ space

The Unbiased Linear System SolvEr - ULISSE

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

- Accelerate the solution of dense linear systems
- ... returning an unbiased estimate of the solution
- Basic idea - unbiased estimator for generic sums $a + b$:

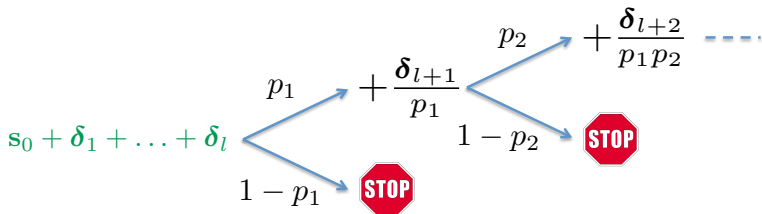


The Unbiased Linear System SolvEr - ULISSE

- Full CG solution:

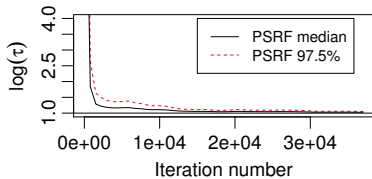
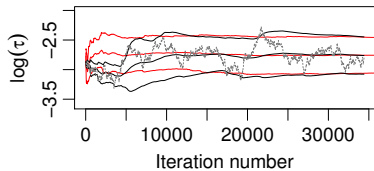
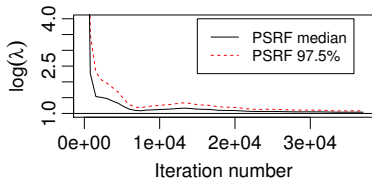
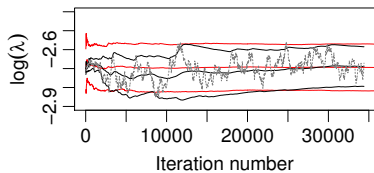
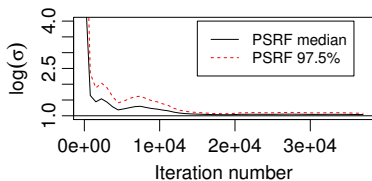
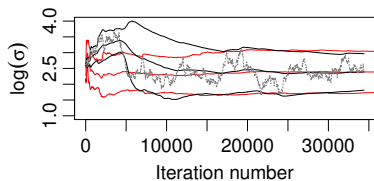
$$\mathbf{s} = \mathbf{s}_0 + \delta_1 + \dots + \delta_l + \delta_{l+1} \dots + \delta_T$$

- ULISSE:

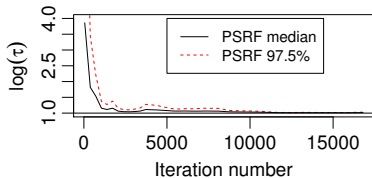
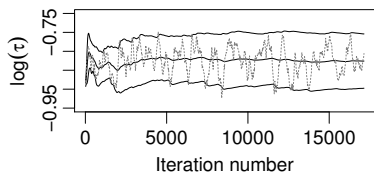
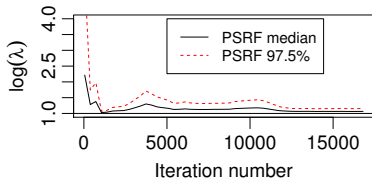
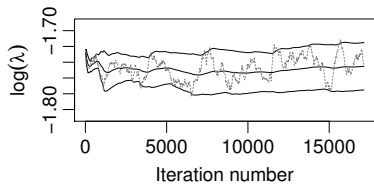
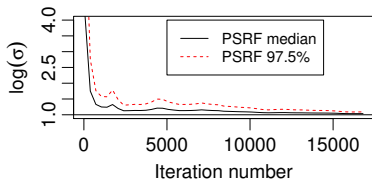
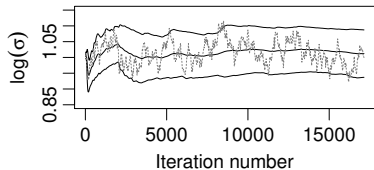


- Final solution is an unbiased estimate of \mathbf{s} !

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



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



- “Noisy” MCMC offers a practical and scalable way to carry out “exact” Bayesian computations for GPs

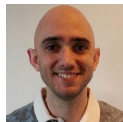
Conclusions and ongoing work

- “Noisy” MCMC offers a practical and scalable way to carry out “exact” Bayesian computations for GPs
- Novel adaptation of SGLD when the likelihood does not factorize
- Novel linear solver ULISSE to speed up computations of stochastic gradients

Conclusions and ongoing work

- “Noisy” MCMC offers a practical and scalable way to carry out “exact” Bayesian computations for GPs
- Novel adaptation of SGLD when the likelihood does not factorize
- Novel linear solver ULISSE to speed up computations of stochastic gradients
- General likelihoods?
- Preconditioners?

- Joint work with Raphael Engler



Andre Marquand
Radboud



Guido Sanguinetti
Edinburgh



James Hensman
Sheffield



Mark Girolami
Warwick



Alessandro Vinciarelli
Glasgow



Dirk Husmeier
Glasgow

- [1] M. Filippone, A. F. Marquand, C. R. V. Blain, S. C. R. Williams, J. Mourão-Miranda, and M. Girolami. Probabilistic prediction of neurological disorders with a statistical assessment of neuroimaging data modalities. *Annals of Applied Statistics*, 6(4):1883-1905, 2012.
- [2] M. Filippone, M. Zhong, and M. Girolami. A comparative evaluation of stochastic-based inference methods for Gaussian process models. *Machine Learning*, 93(1):93-114, 2013.
- [3] 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.
- [4] M. Filippone. Bayesian inference for Gaussian process classifiers with annealing and pseudo-marginal MCMC, In *ICPR 2014*, pages 614-619. *IEEE*, 2014.
- [5] J. Hensman, A. G. de G. Matthews, M. Filippone, and Z. Ghahramani. MCMC for variationally sparse Gaussian processes, arXiv:1506.04000, 2015.