# Enabling scalable stochastic gradient-based inference for Gaussian processes by employing the Unbiased LInear System SolvEr (ULISSE)
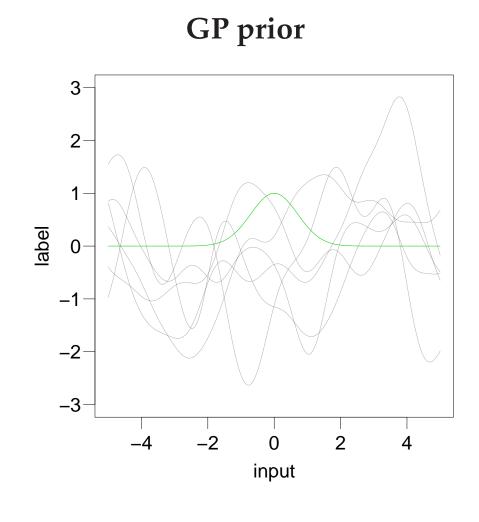
M. Filippone[1,2], R. Engler[2]

1 - EURECOM, Sophia Antipolis, France. email: maurizio.filippone@eurecom.fr
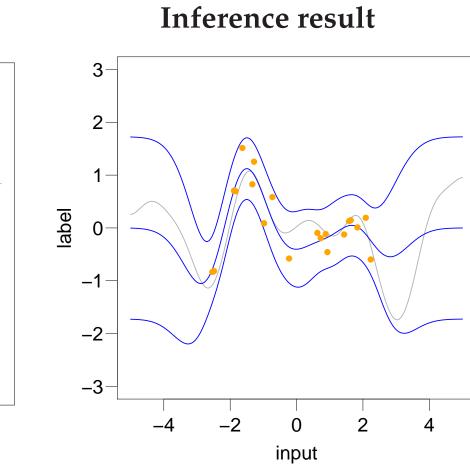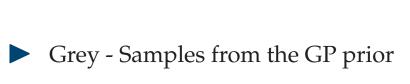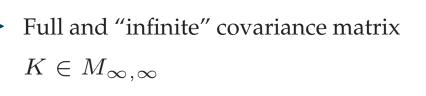
2 - University of Glasgow, Glasgow, UK.

## Gaussian Process (GP) Regression - Illustration

**GP prior**    **GP regression example**    **Inference result**



- ▶ Grey - Samples from the GP prior
- ▶ Green - Radial Basis Function covariance
- ▶ Full and "infinite" covariance matrix $K \in M_{\infty,\infty}$

- ▶ Data generated from the model
- ▶ Marginal distributions of multivariate Gaussian are Gaussian
- ▶ It is sufficient to look at the values of the covariance at the input locations

- ▶ Blue - Mean prediction $\pm$ 2 Std devs
- ▶ Conditional distributions of multivariate Gaussian are Gaussian
- ▶ Predictions can be made by calculating conditional distributions

## Bayesian Inference for GPs

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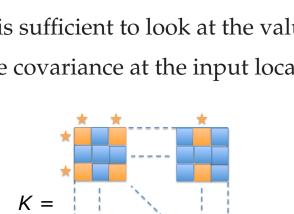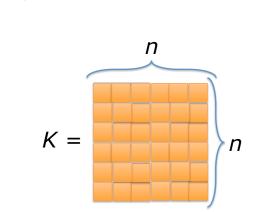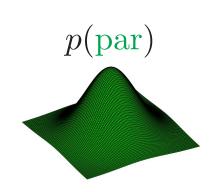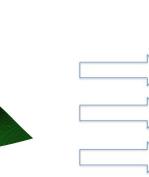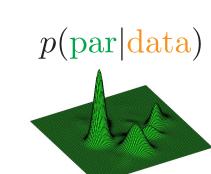
$p(\mathrm{par})$    $p(\mathrm{par}|\mathrm{data})$

## Marginal likelihood

- ▶ Marginal likelihood

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

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

- ▶ ... even then

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

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

## Stochastic Gradient Langevin Dynamics (SGLD) algorithm

- ▶ Stochastic gradient ascent optimization with injected noise $\eta_t$

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

- ▶ First phase – $\alpha_t$ large – Optimization phase
  - – Injected noise $\eta_t$ is smaller than the gradient-based update
  - – Behavior similar to stochastic gradient ascent
- ▶ Second phase – $\alpha_t$ small – Langevin dynamics phase
  - – Injected noise $\eta_t$ dominates gradient-based update
  - ✓ Acceptance rate reaches one so no need to accept/reject
  - ✓ No need to evaluate $p(\mathrm{data}|\mathrm{par})$
  - ✓ We only need stochastic gradients to obtain samples from $p(\mathrm{par}|\mathrm{data})$

## Stochastic gradients for GPs

- ▶ Marginal likelihood

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

- ▶ Derivatives wrt par

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

- ▶ Stochastic estimate of the trace

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

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

- ▶ Stochastic gradient

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

- ▶ Linear systems only!

## Solving linear systems

- ▶ Linear systems:

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

- ▶ Can be solved using the Conjugate Gradient algorithm:

$$\mathbf{s} = \arg\min_{\mathbf{x}}\left(\frac{1}{2}\mathbf{x}^{\mathrm{T}}K\mathbf{x} - \mathbf{x}^{\mathrm{T}}\mathbf{b}\right)$$
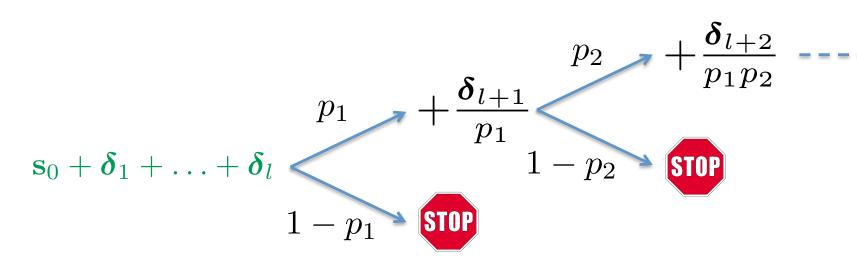
- ▶ Iterative update $\mathbf{s} = \mathbf{s}_0 + \boldsymbol{\delta}_1 + \ldots + \boldsymbol{\delta}_T$
- ▶ Requires only Covariance Matrix Vector Products (CMVPs)! $O(n^2)$ time
- ▶ No need to store $K$! $O(n)$ space

## ULISSE - the Unbiased LInear System SolvEr

- ▶ Accelerate the solution of dense linear systems
- ▶ ... returning an unbiased estimate of the solution
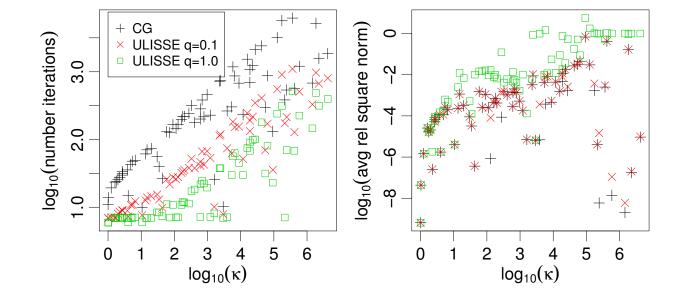- ▶ Full CG solution:

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

- ▶ ULISSE:



In this work:
$$p_i = \exp(-\beta i)$$

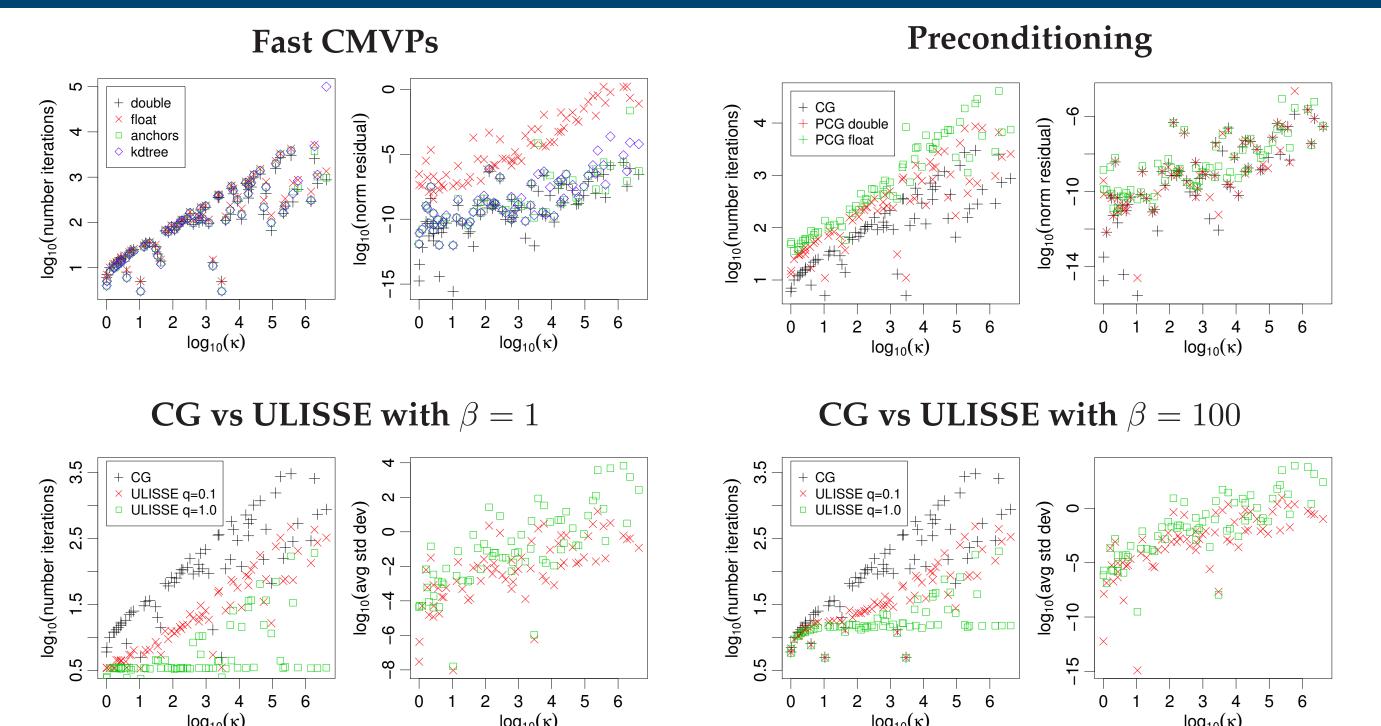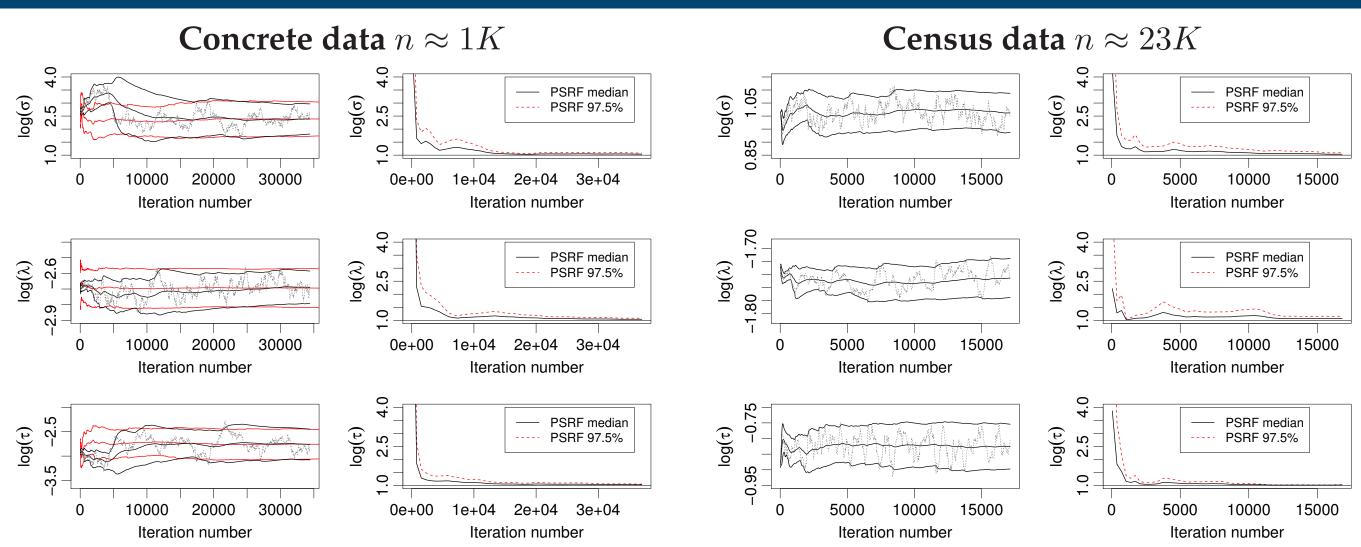- ▶ Final solution is an unbiased estimate of s!



- ✓ Fast computation of stochastic gradients
- ✓ Small relative error wrt exact gradients

$$\text{rel square norm} = \frac{\|\mathbf{g}(\boldsymbol{\theta}) - \tilde{\mathbf{g}}(\boldsymbol{\theta})\|^2}{\|\mathbf{g}(\boldsymbol{\theta})\|^2}$$

## Traditional solvers vs ULISSE

**Fast CMVPs**    **Preconditioning**

**CG vs ULISSE with $\beta = 1$**    **CG vs ULISSE with $\beta = 100$**



## Inference Results

**Concrete data $n \approx 1K$**    **Census data $n \approx 23K$**



## Conclusions

- ▶ Novel adaptation of SGLD to infer covariance parameters in Gaussian processes
  - ✓ Accurate in characterizing the posterior distribution over covariance parameters
  - ✓ Scales with $O(n)$ in space and with $O(n^2)$ in time
  - ✓ Massively parallelizable
  - ✓ Without assuming factorization of the likelihood (mini-batches)
  - ✓ Without considering subsets of the data or inducing points
  - ✓ Without considering subsets of the spectrum of the covariance
  - ✓ Without imposing sparsity on the covariance or its inverse
- ▶ Novel linear solver - ULISSE
  - ✓ Early stop of iterative linear solver that yields an unbiased solution
  - ✓ Can be adopted to accelerate **any** iterative solver
- ▶ Ongoing work
  - – How to extend this work to other likelihoods
  - – Tuning of a preconditioner in SGLD
  - – Mixed precision calculations within the Conjugate Gradient algorithm

## References

[1] 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.

[2] 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.

[3] M. N. Gibbs, *Bayesian Gaussian processes for regression and classification*. PhD thesis, University of Cambridge, 1997.

[4] M. Welling and Y. W. Teh, Bayesian Learning via Stochastic Gradient Langevin Dynamics. *ICML 2011*, pp. 681–688. Omnipress, 2011.