

Scalable Gaussian processes with a twist of Probabilistic Numerics

Kurt Cutajar

EURECOM, Sophia Antipolis, France

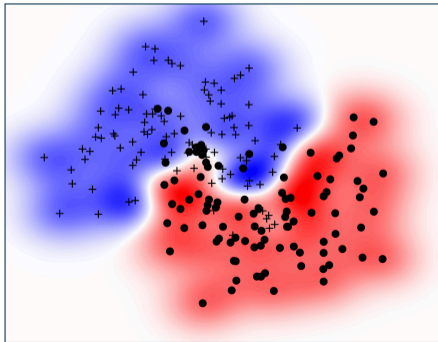
Data Science Meetup - October 30th 2017

Agenda

- Kernel Methods
- Scalable Gaussian Processes (using Preconditioning)
- Probabilistic Numerics

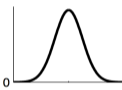


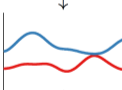
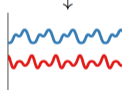
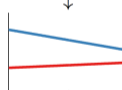
Kernel Methods

- Operate in a high-dimensional, implicit feature space
- Rely on the construction of an $n \times n$ Gram matrix K
- E.g. RBF : $k(\mathbf{x}_i, \mathbf{x}_j) = \sigma^2 \exp(-\frac{1}{2}d^2)$
where $d^2 = (\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{\Lambda} (\mathbf{x}_i - \mathbf{x}_j)$



Kernel Methods

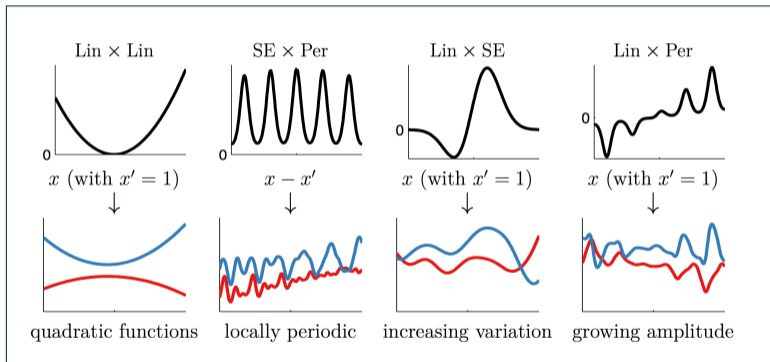
- Wide variety of kernel functions available

Kernel name:	Squared-exp (SE)	Periodic (Per)	Linear (Lin)
$k(x, x') =$	$\sigma_f^2 \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$	$\sigma_f^2 \exp\left(-\frac{2}{\ell^2} \sin^2\left(\pi \frac{x-x'}{p}\right)\right)$	$\sigma_f^2 (x - c)(x' - c)$
Plot of $k(x, x')$:			
	$x - x'$ ↓	$x - x'$ ↓	x (with $x' = 1$) ↓
Functions $f(x)$ sampled from GP prior:			
	x	x	x
Type of structure:	local variation	repeating structure	linear functions

Taken from David Duvenaud's PhD Thesis

Kernel Methods

- Choice is not always straightforward!



Taken from David Duvenaud's PhD Thesis

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}$$

$$p(\text{par}|X, \mathbf{y}) = \frac{p(\mathbf{y}|X, \text{par}) \times p(\text{par})}{p(\mathbf{y}|X)}$$

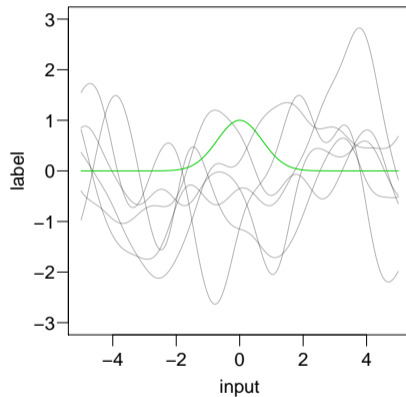
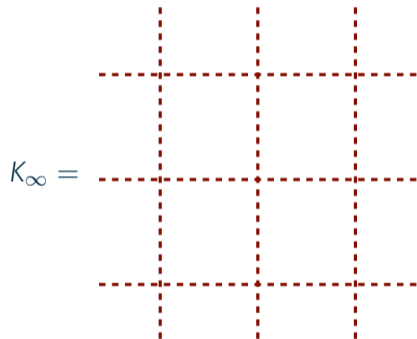
All About that ~~B~~ass Bayes - Making Predictions

- We average over all possible parameter values, weighted by their posterior probability

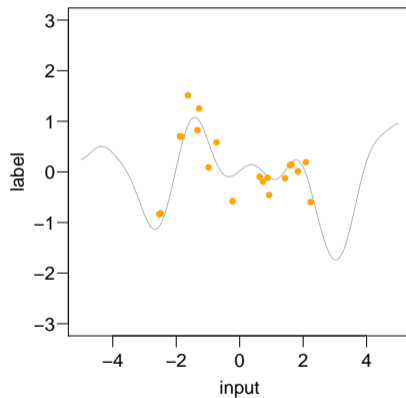
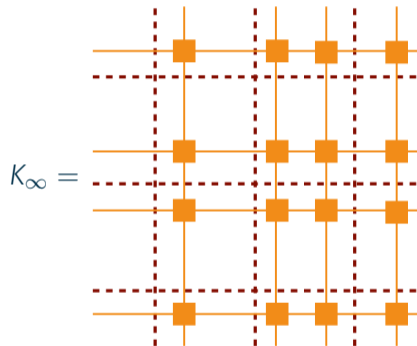
$$\begin{aligned} p(\mathbf{y}^* | \mathbf{x}^*, X, \mathbf{y}) &= \int p(\mathbf{y}^* | \mathbf{x}^*, \text{par}) p(\text{par} | X, \mathbf{y}) \, d\text{par} \\ &= \mathcal{N}(\mathbb{E}[\mathbf{y}^*], \mathbb{V}[\mathbf{y}^*]) \end{aligned}$$

Gaussian Processes

Gaussian Processes - Prior Distribution over Functions

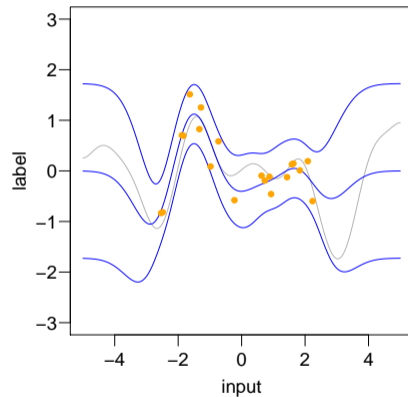


Gaussian Processes - Conditioned on Observations



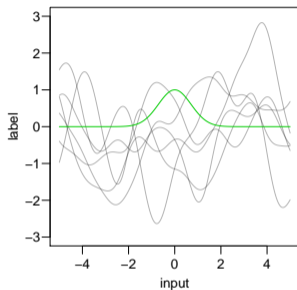
Gaussian Processes - Posterior Distribution over Functions

$$K_y = \begin{bmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare \end{bmatrix} + \begin{bmatrix} \blacksquare & & & \\ & \blacksquare & & \\ & & \blacksquare & \\ & & & \blacksquare \end{bmatrix}$$

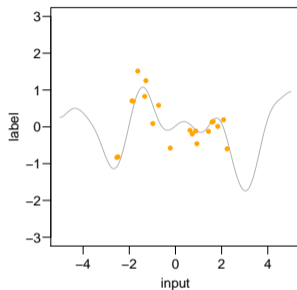


Gaussian Processes

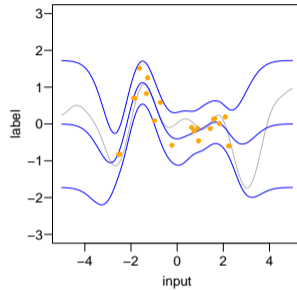
GP prior

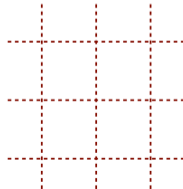


GP regression example

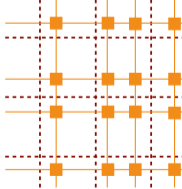


Inference result



$$K_{\infty} =$$


A diagram illustrating the kernel matrix K_{∞} . It consists of a grid of dashed lines, representing a kernel matrix where all elements are non-zero, indicating a fully connected prior.

$$K_{\infty} =$$


A diagram illustrating the kernel matrix K_{∞} . It consists of a grid of dashed lines, with orange squares at the intersections, representing a kernel matrix where all elements are non-zero, indicating a fully connected prior.

$$K_y =$$


A diagram illustrating the kernel matrix K_y . It consists of a 4x4 grid of orange squares, representing a kernel matrix where all elements are non-zero, indicating a fully connected prior. To the right of the grid is a plus sign followed by a diagonal matrix, representing the noise term.

Bayesian Learning vs Deep Learning

- **Deep Learning**

- + Scalable to very large datasets
- + Increased model flexibility/capacity
- Frequentist approaches make only point estimates
- Less robust to overfitting

- **Bayesian Learning**

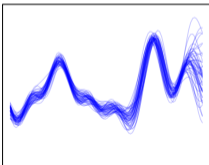
- + Incorporates uncertainty in predictions
- + Works well with smaller datasets
- Lack of conjugacy necessitates approximation
- Expensive computational and storage requirements

Bayesian Learning vs Deep Learning - Deep Gaussian Processes

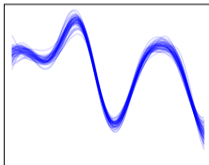
- Deep probabilistic models
- Composition of functions

$$f(\mathbf{x}) = \left(h^{(N_h-1)} \left(\theta^{(N_h-1)} \right) \circ \dots \circ h^{(0)} \left(\theta^{(0)} \right) \right) (\mathbf{x})$$

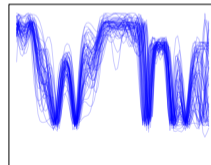
$h^{(0)}(\mathbf{x})$



$h^{(1)}(\mathbf{x})$



$h^{(1)}(h^{(0)}(\mathbf{x}))$



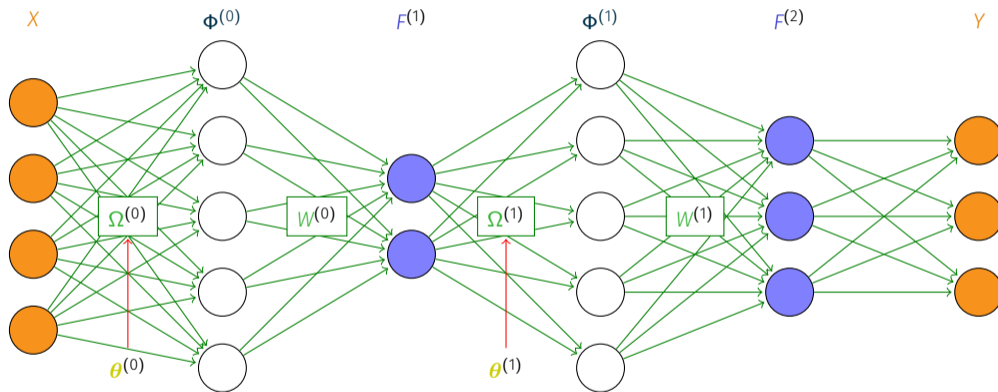
Bayesian Learning vs Deep Learning - Deep Gaussian Processes

- Inference requires calculating the marginal likelihood:

$$p(Y|X, \theta) = \int p(Y|F^{(N_h)}, \theta^{(N_h)}) \times \\ p(F^{(N_h)}|F^{(N_h-1)}, \theta^{(N_h-1)}) \times \dots \times \\ p(F^{(1)}|X, \theta^{(0)}) dF^{(N_h)} \dots dF^{(1)}$$

- Very challenging!

Bayesian Learning vs Deep Learning - Deep Gaussian Processes



Cutajar et al., *Random Feature Expansions for Deep Gaussian Processes*, ICML 2017
Yarin Gal, *Bayesian Deep Learning*, PhD Thesis

Scalable Gaussian Processes

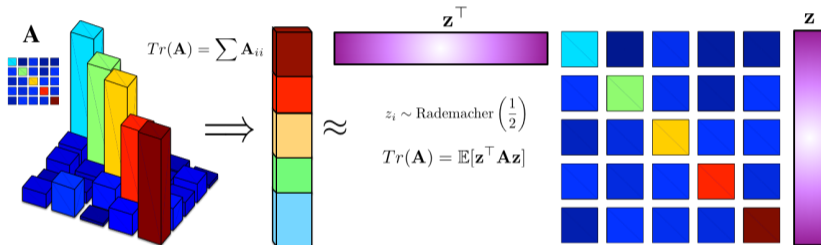
- Marginal likelihood

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

- Derivatives wrt \mathbf{par}

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

Gaussian Processes - Stochastic Trace Estimation



Taken from Shakir Mohamed's Machine Learning Blog

- Stochastic estimate of the trace - assuming $\mathbb{E}[\mathbf{r}\mathbf{r}^T] = \mathbf{I}$, then

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

- Stochastic gradient

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

- Stochastic estimate of the trace - assuming $\mathbb{E}[\mathbf{r}\mathbf{r}^T] = \mathbf{I}$, then

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

- Stochastic gradient

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

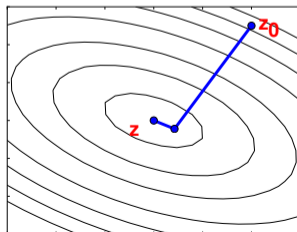
Linear systems only!

Solving Linear Systems

- Involve the solution of linear systems $Kz = v$
- **Cholesky Decomposition**
 - K must be stored in memory!
 - $\mathcal{O}(n^2)$ space and $\mathcal{O}(n^3)$ time - unfeasible for large n

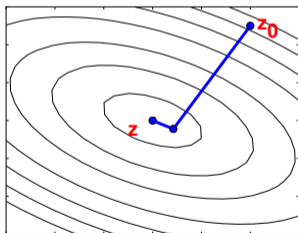
Solving Linear Systems

- Involve the solution of linear systems $Kz = v$
- **Cholesky Decomposition**
 - K must be stored in memory!
 - $\mathcal{O}(n^2)$ space and $\mathcal{O}(n^3)$ time - unfeasible for large n
- **Conjugate Gradient**
 - Numerical solution of linear systems
 - $\mathcal{O}(tn^2)$ for t CG iterations - in theory $t = n$ (possibly worse!)

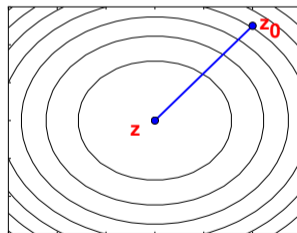


Solving Linear Systems

- **Preconditioned Conjugate Gradient** (henceforth **PCG**)
- Transforms linear system to be better conditioned, improving convergence
- Yields a new linear system of the form $P^{-1}Kz = P^{-1}v$
- $\mathcal{O}(tn^2)$ for t PCG iterations - in practice $t \ll n$



CG



PCG

Preconditioning Approaches

- Suppose we want to precondition $K_y = K + \lambda I$
- Our choice of preconditioner, P , should:
 - Approximate K_y as closely as possible
 - Be easy to invert

Preconditioning Approaches

- Suppose we want to precondition $K_y = K + \lambda I$
- Our choice of preconditioner, P , should:
 - Approximate K_y as closely as possible
 - Be easy to invert
- For low-rank preconditioners we employ the **Woodbury** inversion lemma:

$$K_y = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix} + \begin{bmatrix} \times & & & & \\ & \times & & & \\ & & \times & & \\ & & & \times & \\ & & & & \times \end{bmatrix} \quad P = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix} + \begin{bmatrix} \times & & & & \\ & \times & & & \\ & & \times & & \\ & & & \times & \\ & & & & \times \end{bmatrix}$$

$$P^{-1} = \begin{bmatrix} \times & & & & \\ & \times & & & \\ & & \times & & \\ & & & \times & \\ & & & & \times \end{bmatrix}^{-1} - \begin{bmatrix} \times & & & & \\ & \times & & & \\ & & \times & & \\ & & & \times & \\ & & & & \times \end{bmatrix}^{-1} \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix}^{-1} \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix}^{-1}$$

Preconditioning Approaches

- Suppose we want to precondition $K_y = K + \lambda I$
- Our choice of preconditioner, P , should:
 - Approximate K_y as closely as possible
 - Be easy to invert
- For low-rank preconditioners we employ the **Woodbury** inversion lemma:

$$K_y = \begin{bmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \end{bmatrix} + \begin{bmatrix} \blacksquare & & & & \\ & \blacksquare & & & \\ & & \blacksquare & & \\ & & & \blacksquare & \\ & & & & \blacksquare \end{bmatrix} \quad P = \begin{bmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \end{bmatrix} + \begin{bmatrix} \blacksquare & & & & \\ & \blacksquare & & & \\ & & \blacksquare & & \\ & & & \blacksquare & \\ & & & & \blacksquare \end{bmatrix}$$

$$P^{-1} = \begin{bmatrix} \blacksquare & & & & \\ & \blacksquare & & & \\ & & \blacksquare & & \\ & & & \blacksquare & \\ & & & & \blacksquare \end{bmatrix}^{-1} - \begin{bmatrix} \blacksquare & & & & \\ & \blacksquare & & & \\ & & \blacksquare & & \\ & & & \blacksquare & \\ & & & & \blacksquare \end{bmatrix}^{-1} \begin{bmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \end{bmatrix}^{-1} \begin{bmatrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare & \blacksquare \end{bmatrix}^{-1}$$

- For other preconditioners we solve **inner linear systems** once again using CG!

Preconditioning Approaches

Nyström $P = K_{XU}K_{UU}^{-1}K_{UX} + \lambda I$ where $U \subset X$

FITC $P = K_{XU}K_{UU}^{-1}K_{UX} + \text{diag}(K - K_{XU}K_{UU}^{-1}K_{UX}) + \lambda I$

PITC $P = K_{XU}K_{UU}^{-1}K_{UX} + \text{bldiag}(K - K_{XU}K_{UU}^{-1}K_{UX}) + \lambda I$

Spectral $P_{ij} = \frac{\sigma^2}{m} \sum_{r=1}^m \cos[2\pi \mathbf{s}_r^\top (\mathbf{x}_i - \mathbf{x}_j)] + \lambda I_{ij}$

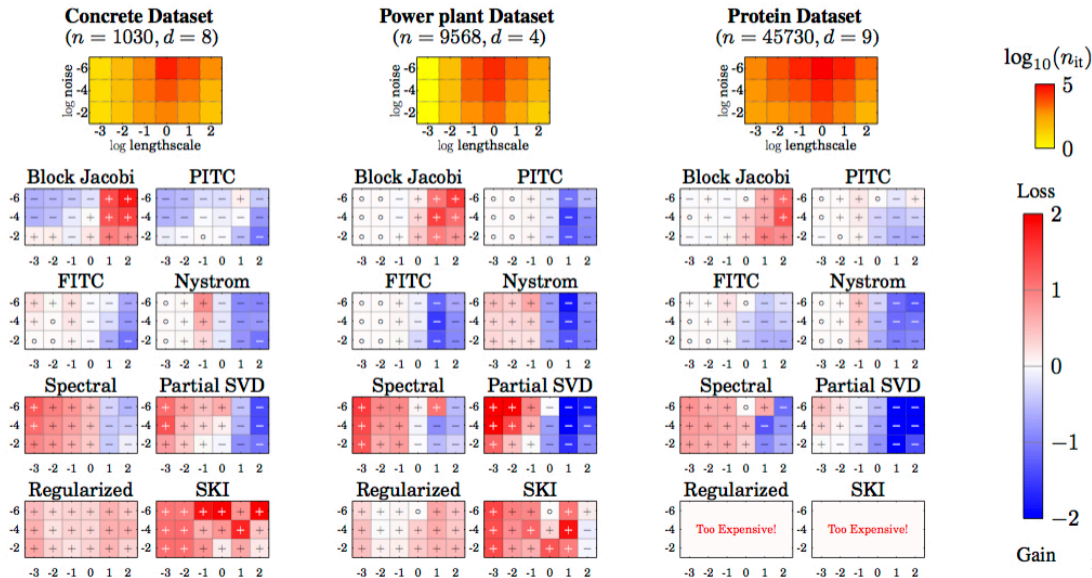
Partial SVD $K = A\Lambda A^\top \Rightarrow P = A_{[:,1:m]}\Lambda_{[1:m,1:m]}A_{[1:m,:]}^\top + \lambda I$

Block Jacobi $P = \text{bldiag}(K) + \lambda I$

SKI $P = WK_{UU}W^\top + \lambda I$ where K_{UU} is Kronecker

Regularization $P = K + \lambda I + \delta I$

Comparison of Preconditioners vs CG



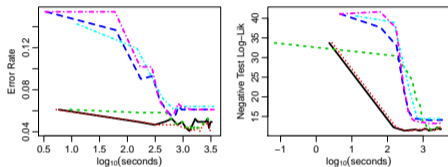
Experimental Setup - GP Kernel Parameter Optimization

- Exact gradient-based optimization using **Cholesky** decomposition (CHOL)
- Stochastic gradient-based optimization
 - Linear systems solved with **CG** and **PCG**
- GP Approximations
 - Variational learning of inducing variables (**VAR**)
 - Fully Independent Training Conditional (**FITC**)
 - Partially Independent Training Conditional (**PITC**)

Results - ARD Kernel

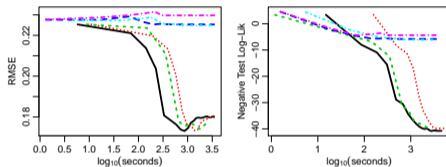
Classification

Spam ($n = 4061$, $d=57$)

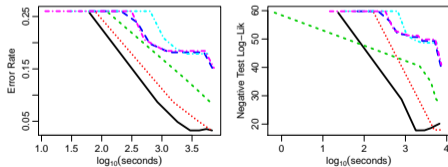


Regression

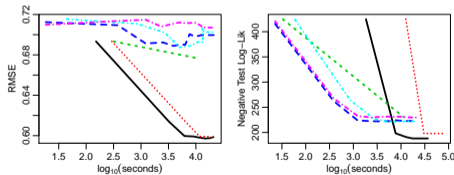
Power plant ($n = 9568$, $d=4$)



EEG ($n = 14979$, $d=14$)



Protein ($n = 45730$, $d=9$)



— PCG CG - - - CHOL - - - FITC PITC - - - VAR

- **Faster Kernel Ridge Regression Using Sketching and Preconditioning**
Avron et al. (2017)
- **FALKON: An Optimal Large Scale Kernel Method**
Rosasco et al. (2017)
- **Large Linear Multi-output Gaussian Process Learning for Time Series**
Feinberg et al. (2017)
- **Scaling up the Automatic Statistician: Scalable Structure Discovery using Gaussian Processes**
Kim et al. (2017)





... but what's left to do now?

Probabilistic Numerics

Probabilistic Numerics - Mission Statement

We deliver a call to arms for *probabilistic numerical methods*: algorithms for numerical tasks, including linear algebra, integration, optimization and solving differential equations, that return uncertainties in their calculations.

Such uncertainties, arising from the loss of precision induced by numerical calculation with limited time or hardware, are important for much contemporary science and industry.

Hennig et al., 2015

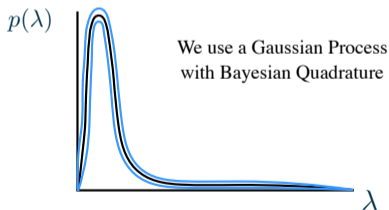
- Quadrature
- Linear Algebra
- Optimization
- Ordinary Differential Equations
- Partial Differential Equations
- Monte Carlo Sampling

- Quadrature
- Linear Algebra
- Optimization
- Ordinary Differential Equations
- Partial Differential Equations
- Monte Carlo Sampling

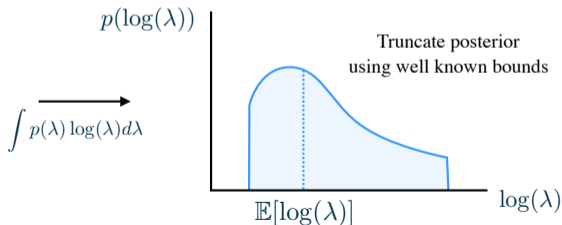
Bayesian Inference of Log Determinants

- Standard computation of $\log |A|$ is $\mathcal{O}(n^3)$
- Existing approximations do not give uncertainty estimates

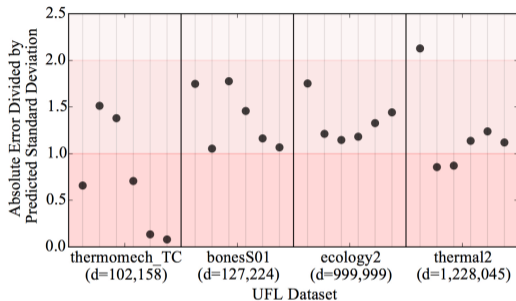
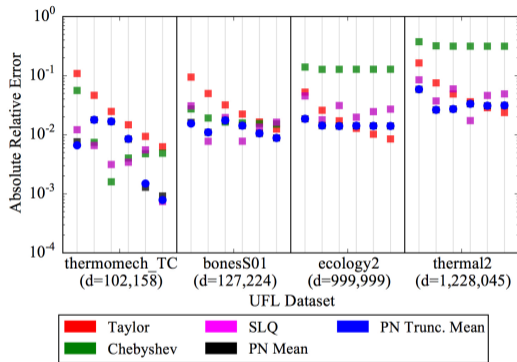
Infer Eigenspectrum



Distribution of Geometric Mean



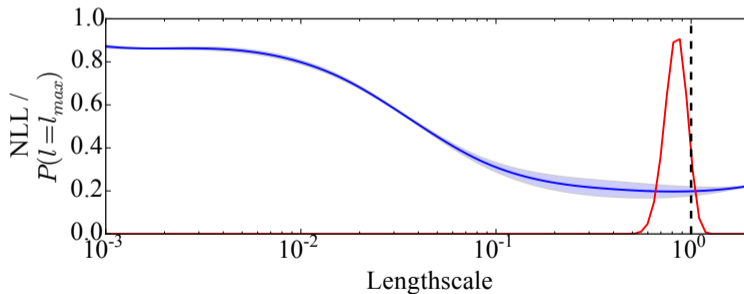
Bayesian Inference of Log Determinants



- Methods compared on a variety of UFL Sparse Datasets. For each dataset, the matrix was approximately raised up to the power of 5, 10, 15, 20, 25 and 30 (left to right) using stochastic trace estimation.

Bayesian Inference of Log Determinants - Application to DPPs

$$\log p(Z) = \log \frac{|L_J|}{|L + I|} \quad \forall J \subset X$$



Can we include this computational uncertainty within the full pipeline of GP inference?

Can we include this computational uncertainty within the full pipeline of GP inference?

- Better calibrated uncertainties
- Performance tunable to computational budget
- Applications to Bayesian optimisation

- **Preconditioning Kernel Matrices (ICML 2016)**
Kurt Cutajar, John Cunningham, Michael Osborne, Maurizio Filippone
- **Bayesian Inference of Log Determinants (UAI 2017)**
Jack Fitzsimons, Kurt Cutajar, Michael Osborne, Stephen Roberts, Maurizio Filippone
- **Random Feature Expansions for Deep Gaussian Processes (ICML 2017)**
Kurt Cutajar, Edwin V. Bonilla, Pietro Michiardi, Maurizio Filippone
- **AutoGP: Exploring the capabilities and limitations of Gaussian processes (UAI 2017)**
Karl Krauth, Edwin V. Bonilla, Kurt Cutajar, Maurizio Filippone
- **Entropic Trace Estimates for Log Determinants (ECML/PKDD 2017)**
Jack Fitzsimons, Diego Granziol, Kurt Cutajar, Maurizio Filippone, Michael Osborne, Stephen Roberts

Thank you!

<https://github.com/mauriziofilippone>

kurt.cutajar@eurecom.fr