

Preconditioning Kernel Matrices



Kurt Cutajar¹



John P. Cunningham²



Michael A. Osborne³



Maurizio Filippone¹

¹ EURECOM, Sophia Antipolis, France

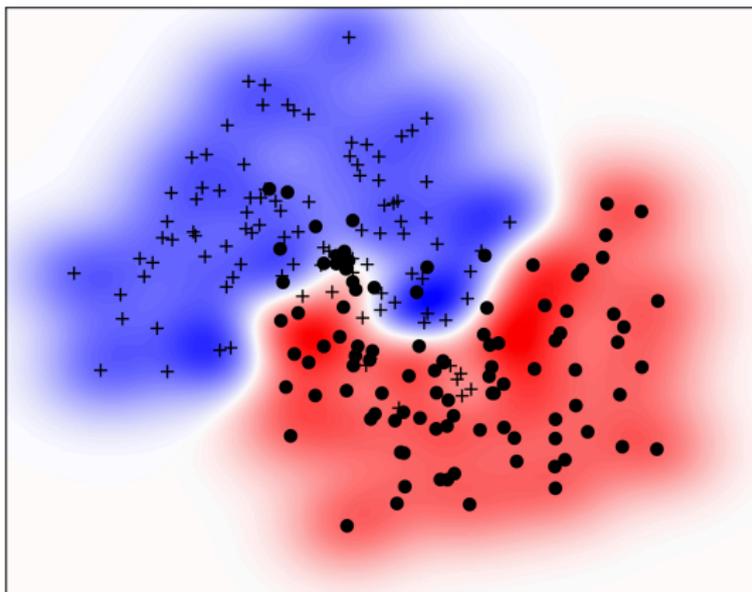
² Columbia University, New York, US

³ University of Oxford, Oxford, UK

ICML 2016, New York
June 22nd, 2016

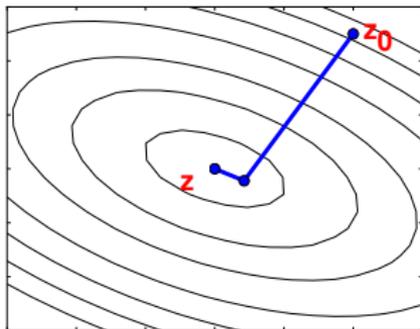
- Operate in a high-dimensional, implicit feature space
- Rely on the construction of an $n \times n$ Gram matrix K
- E.g. RBF : $k(\mathbf{x}_i, \mathbf{x}_j) = \sigma^2 \exp\left(-\frac{1}{2}d^2\right)$

where $d^2 = (\mathbf{x}_i - \mathbf{x}_j)^\top \Lambda (\mathbf{x}_i - \mathbf{x}_j)$

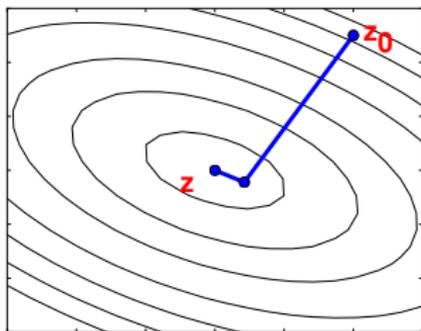


- Involve the solution of linear systems $Kz = v$
- **Cholesky Decomposition**
 - K must be stored in memory!
 - $\mathcal{O}(n^2)$ space and $\mathcal{O}(n^3)$ time - unfeasible for large n

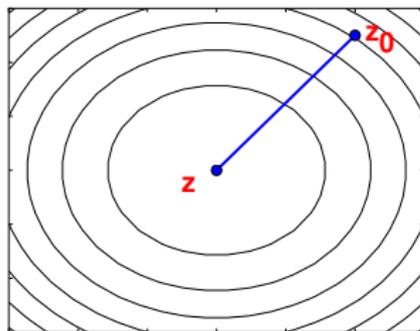
- Involve the solution of linear systems $Kz = v$
- **Cholesky Decomposition**
 - K must be stored in memory!
 - $\mathcal{O}(n^2)$ space and $\mathcal{O}(n^3)$ time - unfeasible for large n
- **Conjugate Gradient**
 - Numerical solution of linear systems
 - $\mathcal{O}(tn^2)$ for t CG iterations -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



- **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:

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

$$P^{-1} = \begin{bmatrix} \blacksquare & & & \\ & \blacksquare & & \\ & & \ddots & \\ & & & \blacksquare \end{bmatrix}^{-1} - \begin{bmatrix} \blacksquare & & & \\ & \blacksquare & & \\ & & \ddots & \\ & & & \blacksquare \end{bmatrix}^{-1} \begin{bmatrix} \blacksquare & & & \\ & \blacksquare & & \\ & & \blacksquare & \\ & & & \blacksquare \end{bmatrix} \left[\begin{bmatrix} \blacksquare & & & \\ & \blacksquare & & \\ & & \blacksquare & \\ & & & \blacksquare \end{bmatrix}^{-1} + \begin{bmatrix} \blacksquare & & & \\ & \blacksquare & & \\ & & \blacksquare & \\ & & & \blacksquare \end{bmatrix}^{-1} \begin{bmatrix} \blacksquare & & & \\ & \blacksquare & & \\ & & \blacksquare & \\ & & & \blacksquare \end{bmatrix} \right]^{-1} \begin{bmatrix} \blacksquare & & & \\ & \blacksquare & & \\ & & \blacksquare & \\ & & & \blacksquare \end{bmatrix}^{-1}$$

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:

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

$$P^{-1} = \begin{bmatrix} \times & & & \\ & \times & & \\ & & \times & \\ & & & \times \end{bmatrix}^{-1} - \begin{bmatrix} \times & & & \\ & \times & & \\ & & \times & \\ & & & \times \end{bmatrix}^{-1} \begin{bmatrix} \times & & & \\ & \times & & \\ & & \times & \\ & & & \times \end{bmatrix} \left[\begin{bmatrix} \times & & & \\ & \times & & \\ & & \times & \\ & & & \times \end{bmatrix}^{-1} + \begin{bmatrix} \times & & & \\ & \times & & \\ & & \times & \\ & & & \times \end{bmatrix}^{-1} \begin{bmatrix} \times & & & \\ & \times & & \\ & & \times & \\ & & & \times \end{bmatrix} \right]^{-1} \begin{bmatrix} \times & & & \\ & \times & & \\ & & \times & \\ & & & \times \end{bmatrix}^{-1}$$

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

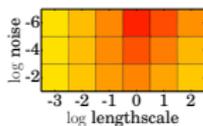
Block Jacobi $P = \text{bldiag}(K) + \lambda I$

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

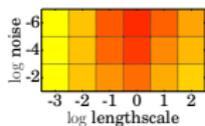
Regularization $P = K + \lambda I + \delta I$

Comparison of Preconditioners vs CG

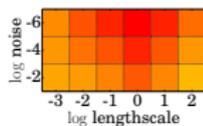
Concrete Dataset
($n = 1030, d = 8$)



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



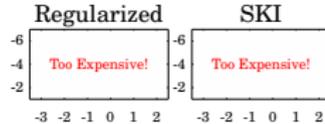
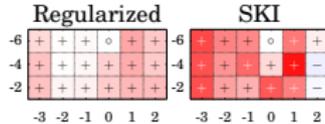
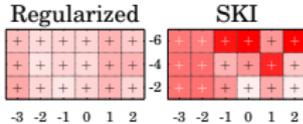
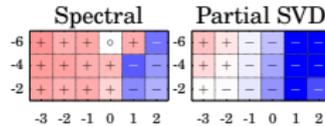
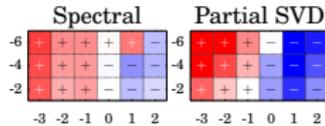
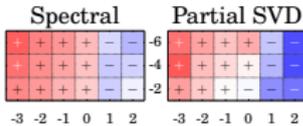
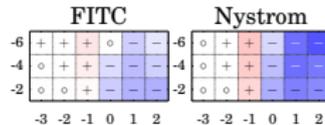
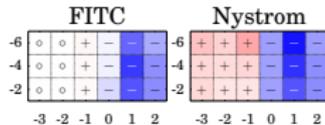
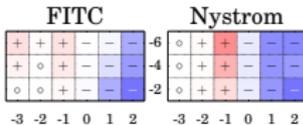
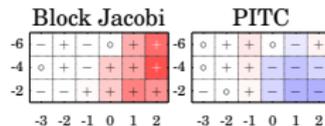
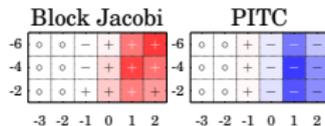
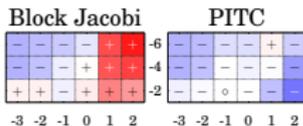
Protein Dataset
($n = 45730, d = 9$)



$\log_{10}(n_{it})$

Loss

Gain



- Marginal likelihood

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

- Derivatives wrt par

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

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

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

- Stochastic gradient

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

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

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

- Stochastic gradient

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

Linear systems only!

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

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

- Stochastic gradient

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

Linear systems only!

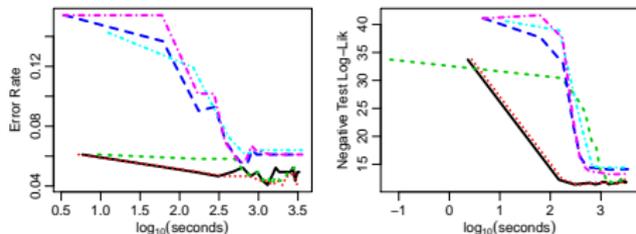
Also applicable to non-Gaussian likelihoods!

GP Kernel Parameter Optimization

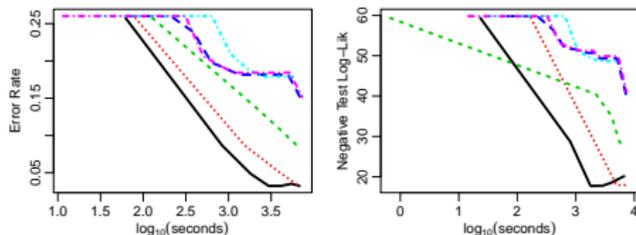
- Exact gradient-based optimization using **Cholesky** decomposition (CHOL)
- Stochastic gradient-based optimization (using **ADAGRAD**)
 - Linear systems solved with **CG** and **PCG**
- GP Approximations
 - Variational learning of inducing variables (**VAR**)
 - Fully Independent Training Conditional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

Classification

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

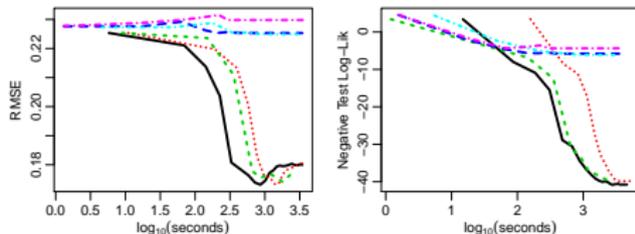


EEG ($n = 14979, d=14$)

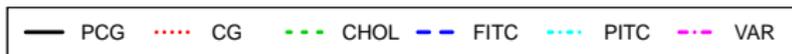
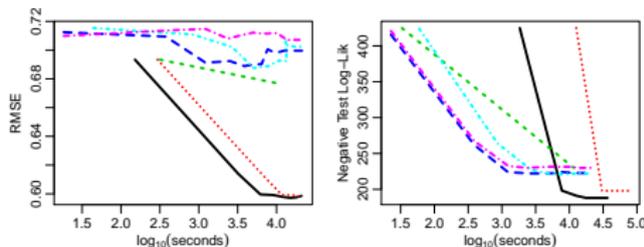


Regression

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



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



- Contributions
 - We provide an extensive discussion of preconditioning
 - We incorporate preconditioning in GP models with both Gaussian and non-Gaussian likelihoods
 - We carry out exact GP inference at a fraction of the cost

- Contributions
 - We provide an extensive discussion of preconditioning
 - We incorporate preconditioning in GP models with both Gaussian and non-Gaussian likelihoods
 - We carry out exact GP inference at a fraction of the cost
- Ongoing work
 - Extending this work to other kernel functions and models
 - Implementation on a distributed framework
 - Exploiting PCG in the solution of $f(K) \mathbf{z} = \mathbf{v}$

- Contributions
 - We provide an extensive discussion of preconditioning
 - We incorporate preconditioning in GP models with both Gaussian and non-Gaussian likelihoods
 - We carry out exact GP inference at a fraction of the cost
- Ongoing work
 - Extending this work to other kernel functions and models
 - Implementation on a distributed framework
 - Exploiting PCG in the solution of $f(K) \mathbf{z} = \mathbf{v}$

Thank you!