

# Efficient Approximation of Covariance Matrix Inverse for Solving the Gaussian Process Maximum Likelihood Problem

*Daniil Kononenko*

June 1 2012

**DATADVANCE**

AN EADS COMPANY

- 1 Gaussian processes regression
  - The problem statement
  - Maximum likelihood method
- 2 Existing approaches to acceleration of training
  - Quadratic function minimization
  - Quasi-Newton approach
  - Determinant approximation
- 3 Proposed extensions to approximation method
- 4 Computational experiments
- 5 Conclusion

# Table of contents

- 1 Gaussian processes regression
  - **The problem statement**
  - Maximum likelihood method
- 2 Existing approaches to acceleration of training
  - Quadratic function minimization
  - Quasi-Newton approach
  - Determinant approximation
- 3 Proposed extensions to approximation method
- 4 Computational experiments
- 5 Conclusion

# Problem statement

- Unknown function  $f(x)$ ,  $x \in \mathbb{R}^n$ ,  $f \in \mathbb{R}$ .
- Training set  $\{(x_i, y_i = y(x_i))\}_{i=1}^N$ .
- Observations  $y(x_i) = f(x_i) + \varepsilon(x_i)$ ,  $\varepsilon(x) \sim \mathcal{N}(0, \sigma^2)$ .
- Let  $f(x, w)$  — gaussian random field with parametric covariance function

$$K_0(x, x') = K_0(x, x' | \Theta).$$

- Covariance function of process  $y(x)$

$$K(x, x' | \Theta, \sigma^2) = K_0(x, x' | \Theta) + \sigma^2 \delta(x, x').$$

# Posterior distribution

- Conditional distribution at new point

$$Law(y(x)|y(x_1), y(x_2), \dots, y(x_n)) \sim \mathcal{N}(\hat{f}(x), \hat{\sigma}^2(x)).$$

- Mean

$$\hat{f}(x) = \mathbf{k}(x)\mathbf{K}^{-1}\mathbf{y}.$$

- Variance

$$\hat{\sigma}^2(x) = K(x, x) - \mathbf{k}(x)\mathbf{K}^{-1}\mathbf{k}(x)^T.$$

- Here  $\mathbf{k}(x) = \{K(x, x_i)\}_{i=1}^N$ ,  $\mathbf{K} = \{K(x_i, x_j)\}_{i,j=1}^N$ ,  
 $\mathbf{y} = (y_1, \dots, y_N)^T$ .
- Unknown parameters are tuned using maximum likelihood method.

# Table of contents

- 1 Gaussian processes regression
  - The problem statement
  - **Maximum likelihood method**
- 2 Existing approaches to acceleration of training
  - Quadratic function minimization
  - Quasi-Newton approach
  - Determinant approximation
- 3 Proposed extensions to approximation method
- 4 Computational experiments
- 5 Conclusion

# Likelihood

Likelihood:

$$l(\mathbf{a}) = \log p(\mathbf{y}|\mathbf{x}, \mathbf{a}) = -\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} - \frac{1}{2}\log(\det(\mathbf{K})) - \frac{N}{2}\log 2\pi.$$

Derivatives w.r.t. hyperparameters:

$$\frac{\partial l}{\partial a_i} = -\text{tr}\left(\mathbf{K}^{-1}\frac{\partial\mathbf{K}}{\partial a_i}\right) + \mathbf{y}^T\mathbf{K}^{-1}\frac{\partial\mathbf{K}}{\partial a_i}\mathbf{K}^{-1}\mathbf{y}.$$

Here  $\mathbf{a}$  — vector of hyperparameters.

It is necessary to efficiently calculate values

$$\mathbf{K}^{-1}\mathbf{y}, \text{tr}\left(\mathbf{K}^{-1}\frac{\partial\mathbf{K}}{\partial a_i}\right), \log(\det(\mathbf{K})).$$

# Cholesky decomposition

- Cholesky decomposition  $K = U^T U$ , where  $U$  is upper triangular matrix. It is time-consuming:  $O(N^3)$  operations.
- Typical number of iterations for likelihood optimization using gradient method with multistart is 100 – 200. If  $N \gtrsim 10^3$  then learning becomes rather long, up to half an hour.
- Stochastic global optimization methods require much more iterations.



# Table of contents

- 1 Gaussian processes regression
  - The problem statement
  - Maximum likelihood method
- 2 Existing approaches to acceleration of training
  - **Quadratic function minimization**
  - Quasi-Newton approach
  - Determinant approximation
- 3 Proposed extensions to approximation method
- 4 Computational experiments
- 5 Conclusion

# Product computation

Skilling (1993), Gibbs and MacKay (1997) suggest the following approach.

- Minimize  $g(u) = \frac{1}{2}u^T K u - u^T \mathbf{y}$ . Exact minimum  $u^* = K^{-1}\mathbf{y}$ .
- Quadratic function is minimized using some simple algorithm, for example conjugate gradients.
- Trace is approximated via Monte Carlo method  
$$\text{tr} \left( K^{-1} \frac{\partial K}{\partial a_i} \right) = E \left[ d^T K^{-1} \frac{\partial K}{\partial a_i} d \right], d \sim N(0, 1), K^{-1}d \text{ is calculated using method above.}$$

# Discussion of method

## Advantages:

- Accurate and fast method of approximating  $K^{-1}\mathbf{y}$ .
- Simple termination conditions.

## Drawbacks:

- Accurate trace estimator needs significant number of seeds, for each seed estimating  $K^{-1}d$  needs several  $N^2$  operations, so calculation time is significant.
- It is not clear what to do with the determinant term.

# Table of contents

- 1 Gaussian processes regression
  - The problem statement
  - Maximum likelihood method
- 2 Existing approaches to acceleration of training
  - Quadratic function minimization
  - **Quasi-Newton approach**
  - Determinant approximation
- 3 Proposed extensions to approximation method
- 4 Computational experiments
- 5 Conclusion

# Quasi-Newton approximation

Zhang and Leithead (2007) suggest an alternative algorithm.

- Quadratic function  $g(u) = \frac{1}{2}u^T K u - u^T \mathbf{y}$  is minimized using Quasi-Newton method.
- Thus, the sequence of approximations of inverse matrix  $K^{-1}$  is built using BFGS formula

$$\begin{aligned}
 u_{j+1} &= u_j - \eta_j K_j^{-1} \nabla g(u_j), \\
 K_{j+1}^{-1} &= K_j^{-1} + \Delta K_j^{-1}, \\
 \Delta K_j^{-1} &= \left( 1 + \frac{q_j^T K_j^{-1} q_j}{q_j^T p_j} \right) \frac{p_j p_j^T}{p_j^T q_j} - \frac{p_j q_j^T K_j^{-1} + K_j^{-1} q_j p_j^T}{q_j^T p_j},
 \end{aligned}$$

here  $p_j = u_{j+1} - u_j$ ,  $q_j = \nabla g(u_{j+1}) - \nabla g(u_j)$ .

# Initialization and termination

- Initial approximation  $K_0^{-1}$  is the value derived in the previous point, where likelihood was calculated (with previous values of hyperparameters).
- Each iteration is  $O(N^2)$  operations.
- Maximum number of iterations is some constant  $J$ , so, the whole method of calculating inverse covariance matrix is  $O(N^2)$ .
- Stopping criterion

$$\|\nabla g(u_j)\|_\infty \leq \frac{\varepsilon}{N}.$$

# Quality criterion and scaling

- Quality criterion

$$\frac{|\text{tr}(\tilde{\mathbf{K}}^{-1}\mathbf{K}) - N|}{N} < \varepsilon_1,$$

$\tilde{\mathbf{K}}^{-1}$  — approximation to the inverse matrix.

- Restart technique.
- Random scaling:

$$s = \mathbf{K}_j^{-1}(i, :) \times \mathbf{K}(:, i), i \text{ — random coordinate,}$$

$$\mathbf{K}_j^{-1} := \frac{\mathbf{K}_j^{-1}}{s}, j \text{ — number of iteration.}$$

# Table of contents

- 1 Gaussian processes regression
  - The problem statement
  - Maximum likelihood method
- 2 Existing approaches to acceleration of training
  - Quadratic function minimization
  - Quasi-Newton approach
  - **Determinant approximation**
- 3 Proposed extensions to approximation method
- 4 Computational experiments
- 5 Conclusion



# Determinant approximation

Determinant approximation:

$$\log(\det(\mathbf{K})) = \text{tr}(\log(\mathbf{K})),$$
$$\log(\mathbf{K}) = - \sum_{i=1}^{\infty} \left( \frac{\mathbf{K}^i}{i} \right) \text{ — matrix logarithm.}$$

Two-level scheme of approximation:

- Some finite number of terms in series:  $\log(\mathbf{K}) = - \sum_{i=1}^L \left( \frac{\mathbf{K}^i}{i} \right)$
- Each element  $\text{tr}(\mathbf{K}^i)$  is approximated via Monte Carlo method.
- Authors suggest several compensation schemes for errors at both levels.

# Discussion of method

## Advantages:

- Simultaneous approximation of  $K^{-1}\mathbf{y}$  and  $K^{-1}$ .
- Quality control for trace computations.

## Drawbacks:

- Complicated and rather long scheme of determinant approximation.
- No quality control for determinant term.

# Table of contents

- 1 Gaussian processes regression
  - The problem statement
  - Maximum likelihood method
- 2 Existing approaches to acceleration of training
  - Quadratic function minimization
  - Quasi-Newton approach
  - Determinant approximation
- 3 **Proposed extensions to approximation method**
- 4 Computational experiments
- 5 Conclusion

# Quasi-Newton determinant updates

- Matrix determinant lemma:

$$\det(A + uv^T) = \det(A)(1 + v^T A^{-1}u).$$

- DFP update formula instead of BFGS:

$$\Delta K_j^{-1} = \frac{p_j p_j^T}{q_j^T p_j} - \frac{K_j^{-1} q_j q_j^T (K_j^{-1})^T}{q_j^T K_j^{-1} q_j},$$

$$\det(K_{j+1}^{-1}) = \det(K_j^{-1}) \frac{q_j^T p_j}{q_j^T K_j^{-1} q_j}.$$

- At each iteration current approximation to determinant is exactly the determinant of current approximation to inverse matrix.

# Further extensions

- More efficient scaling technique:

$$m = \text{mean}(\text{diag}(\mathbf{K}_j^{-1} \times \mathbf{K})),$$
$$\mathbf{K}_j^{-1} := \frac{\mathbf{K}_j^{-1}}{m}.$$

- If norm of gradient increases on successive iterations, the restart is made.
- Initial approximation  $\mathbf{K}_0^{-1}$  is the value derived at one of the previous points, the closest one to the current point in the sense of special metric. Metric represents the  $l_2$ -norm difference between diagonals of covariance matrixes.

# Table of contents

- 1 Gaussian processes regression
  - The problem statement
  - Maximum likelihood method
- 2 Existing approaches to acceleration of training
  - Quadratic function minimization
  - Quasi-Newton approach
  - Determinant approximation
- 3 Proposed extensions to approximation method
- 4 **Computational experiments**
- 5 Conclusion

# Computational experiments

- Large set of 211 functions, typical for testing optimization tasks.
- Dimensionality of functions — from 2 to 8.
- For each function training sets of 320 and 1000 points were generated.
- Comparison of MSE and time of training.
- MSE is calculated on large test sets of 10000 points.
- Results — Dolan-More curves.

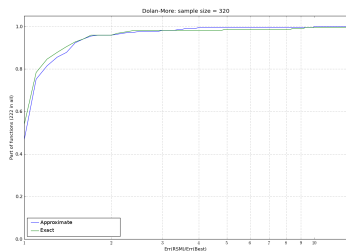
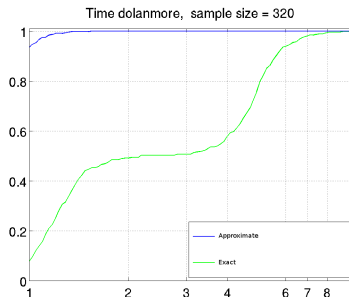
# Dolan-More curves

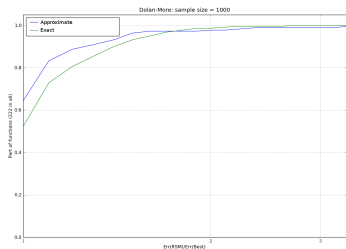
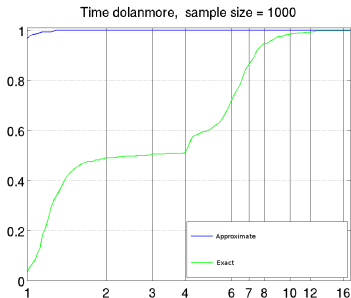
- $T$  tasks,  $A$  algorithms.
- $e_{ta}$  — MSE or time of training of  $a$ -th algorithm on  $t$ -th task.
- $\tilde{e}_t = \min_a e_{ta}$ .

$$p_a(h) = \frac{\#\{t : e_{ta} < h\tilde{e}_t\}}{T}$$

- If the curve is higher, then algorithm is better.
- $p_a(1)$  — part of tasks, on which algorithm  $a$  is the best.



Figure: MSE,  $N = 320$ Figure: Time of training,  $N = 320$

Figure: MSE,  $N = 1000$ Figure: Time of training,  
 $N = 1000$

# Comparison

- Gibbs and Mackay — only ideas without full algorithm.
- Zhang and Leithead — 1 restart at 6 – 7 iterations if  $N = 1000$ . But additional time for determinant approximation. Acceleration rate is less than 6.
- Proposed method: mean acceleration 4 – 6 times.

# Conclusion

## Advantages:

- We avoid complicated schemes for determinant approximation.
- Simultaneous approximation of inverse matrix and its determinant.
- Quality control for determinant term.
- Experiment results: significant acceleration without loss of quality.

Thank you for attention.