

PDF hosted at the Radboud Repository of the Radboud University Nijmegen

The following full text is a preprint version which may differ from the publisher's version.

For additional information about this publication click this link.

<http://hdl.handle.net/2066/134607>

Please be advised that this information was generated on 2019-10-15 and may be subject to change.

Fast Laplace Approximation for Gaussian Processes with a Tensor Product Kernel

Perry Groot^a Markus Peters^b Tom Heskes^a Wolfgang Ketter^b

^a *Radboud University Nijmegen*

^b *Erasmus University Rotterdam*

Abstract

Gaussian processes provide a principled Bayesian framework, but direct implementations are restricted to small data sets due to the cubic time cost in the data size. In case the kernel function is expressible as a tensor product kernel and input data lies on a multidimensional grid it has been shown that the computational cost for Gaussian process regression can be reduced considerably. Tensor product kernels have mainly been used in regression with a Gaussian observation model since key steps in their algorithms do not easily translate to other tasks. In this paper we show how to obtain a scalable Gaussian process framework for gridded inputs and non-Gaussian observation models that factorize over cases. We empirically validate our approach on a binary classification problem and our results shows a major performance improvement in terms of run time.

1 Introduction

Gaussian processes (GPs) provide a rich, principled, and well-established Bayesian framework for tasks such as non-linear non-parametric regression and classification [8]. A direct implementation of GPs, however, limits their applicability to small data sets since the kernel matrix needs to be stored, costing $\mathcal{O}(N^2)$, and inverted, costing $\mathcal{O}(N^3)$, with N the number of data points. In recent years, several approaches have addressed this problem often selecting a subset of the training data (the active set) of size M reducing the computational complexity to $\mathcal{O}(M^2N)$ for $M \ll N$ [6, 10]. An alternative approach, is to exploit additional structure in the GP model to reduce the computational complexity.

The structural assumption that we exploit in this paper is the assumption of input data lying on a multidimensional grid.¹ Such data is often found in regression applications in time and space, for example, regular measurements over time and space from weather stations. This assumption allows the kernel function to be written as a tensor product kernel that allows for efficient computations using Kronecker products [5]. In recent years, the use of multidimensional grids has attracted increasing attention, and has been used in a variety of applications such as image reconstruction [3, 14], point process intensity estimation [2], network reconstruction [11], and density estimation [9].

The Kronecker product has, however, been limited to GP regression with spherical noise. The main reason for this is that a critical step is to compute the inverse $(\mathbf{K} + \sigma^2\mathbf{I})^{-1}$, which is not a Kronecker product even if \mathbf{K} is. In this particular case it is still possible to efficiently compute the inverse using an eigen decomposition $\mathbf{K} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ and the identity $(\mathbf{K} + \sigma^2\mathbf{I})^{-1}\mathbf{y} = \mathbf{Q}(\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1}\mathbf{Q}^T\mathbf{y}$. This identity, however, no longer holds with non-spherical noise and how to efficiently obtain an eigen decomposition in this case is in fact an open problem [4].

In this paper we propose a scalable GP framework for gridded inputs and non-Gaussian observation models that factorize over cases. Section 2 briefly introduces GPs, the Laplace approximation, and the Kronecker product. Thereafter we describe how to efficiently combine tensor product kernel functions in the Laplace approximation. We provide derivations and pseudocode for obtaining the mode and for hyperparameter learning in Sections 3 and 4. Section 5 reviews the memory and run time requirements of the algorithms. In Section 6 we empirically validate the new methods and Section 7 gives conclusions.

¹The grid does not necessarily have to be equispaced.

2 Background

We denote vectors \mathbf{x} and matrices \mathbf{K} with bold-face type and their components with regular type, i.e., x_i, K_{ij} . With \mathbf{x}^T we denote the transpose of the vector \mathbf{x} . Let $\mathbf{x} \in \mathbb{R}^D$ be an input, $\mathbf{y} \in \mathcal{Y}$ an output. We denote with \mathbf{X}, \mathbf{Y} the observed data. We denote with $|\mathbf{A}|$ the determinant of \mathbf{A} , $\mathbf{A} \otimes \mathbf{B}$ the Kronecker product, $\mathbf{A} \circ \mathbf{B}$ the Hadamard (element wise) product, and $\text{vec}(\mathbf{A})$ the vector obtained by stacking all the columns of matrix \mathbf{A} . Given a vector \mathbf{a} and matrix \mathbf{A} , $\text{diag}(\mathbf{a})$ gives a diagonal matrix with \mathbf{a} on its diagonal, and $\text{diag}(\mathbf{A})$ gives the vector of all diagonal elements of \mathbf{A} . Hyperparameters are denoted by $\boldsymbol{\theta} = \{\boldsymbol{\theta}^c, \boldsymbol{\theta}^l\}$, with $\boldsymbol{\theta}^c$ the kernel hyperparameters and $\boldsymbol{\theta}^l$ the likelihood hyperparameters.

2.1 Gaussian Processes

A Gaussian process (GP) is a collection of random variables $\{f(\mathbf{x}_i)\}_{i \in I}$ for some index set I , any finite number of which have a joint Gaussian distribution [8]. A GP $f \sim \mathcal{N}(\mathbf{m}, \mathbf{K})$ results in a finite multivariate Gaussian distribution where each element of the covariance matrix $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ is given by a kernel function $k(\cdot, \cdot)_{\boldsymbol{\theta}^c}$ with parameters $\boldsymbol{\theta}^c$. A GP effectively specifies a prior distribution over functions $f(\cdot)$ in a Bayesian framework in which the likelihood model $p(y_i|f(\mathbf{x}_i))$ is parameterized by $f(\cdot)$. Given a GP prior and likelihood, Bayes formula gives us the posterior distribution $p(\mathbf{f}|\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}) = p(\mathbf{Y}|\mathbf{f}, \boldsymbol{\theta})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta}^c)/p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})$. In the simplest case, observations are assumed to be generated by an unknown function possibly corrupted with Gaussian noise, i.e., $y_i = f(\mathbf{x}_i) + \epsilon_i$ with $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ which leads to a model in which inference is analytically tractable. In this case the predictive distribution for a test location \mathbf{x}_* is given by $f_*|\mathbf{X}, \mathbf{Y}, \mathbf{x}_* \sim \mathcal{N}(\mu_*, v_*)$ with $\mu_* = k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma^2\mathbf{I})^{-1}\mathbf{Y}$ and $v_* = k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma^2\mathbf{I})^{-1}k(\mathbf{X}, \mathbf{x}_*)$. For non-Gaussian likelihoods one needs to resort to sampling methods or approximations.

2.2 Laplace Approximation

With a non-Gaussian likelihood inference in Gaussian process models leads to analytically intractable integrals when making predictions. For example, the distribution of a latent variable given a test case is given by

$$p(f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*) = \int p(f_*|\mathbf{X}, \mathbf{x}_*, \mathbf{f})p(\mathbf{f}|\mathbf{X}, \mathbf{y})d\mathbf{f} \quad (1)$$

The Laplace approximation [13] gives an analytic approximation to these integrals by approximating the true posterior distribution $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ with a Gaussian $q(\mathbf{f})$ centered on the mode of the posterior. The variance is obtained through a second-order Taylor expansion around the posterior mode $\hat{\mathbf{f}}$, which results in

$$q(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, (\mathbf{K}^{-1} + \mathbf{W})^{-1}) \quad (2)$$

with $\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} p(\mathbf{f}|\mathbf{X}, \mathbf{y}, \boldsymbol{\theta})$ and $\mathbf{W} = -\nabla \nabla_{\mathbf{f}} \log p(\mathbf{y}|\mathbf{f})|_{\mathbf{f}=\hat{\mathbf{f}}}$ which is diagonal since the likelihood factorizes over cases.

2.3 Kronecker Products

In this paper we assume a tensor product kernel function, i.e., $k(\mathbf{x}_i, \mathbf{x}_j) = \prod_{d=1}^D k_d(\mathbf{x}_i^d, \mathbf{x}_j^d)$. Assuming a product kernel function together with input data on a multidimensional grid leads to a kernel matrix that decomposes into a Kronecker product of matrices of lower dimensions. The Kronecker product has a convenient algebra [5]:

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix} \quad \begin{aligned} (\mathbf{A} \otimes \mathbf{B})\text{vec}(\mathbf{X}) &= \text{vec}(\mathbf{B}\mathbf{X}\mathbf{A}^T) \\ \mathbf{A}\mathbf{B} \otimes \mathbf{C}\mathbf{D} &= (\mathbf{A} \otimes \mathbf{C})(\mathbf{B} \otimes \mathbf{D}) \\ (\mathbf{A} \otimes \mathbf{B})^{-1} &= \mathbf{A}^{-1} \otimes \mathbf{B}^{-1} \end{aligned} \quad (3)$$

Algorithm 1 Laplace mode finding

```
1: function LAPLACEMODE(covariance matrix  $\mathbf{K}$ , observations  $\mathbf{y}$ , likelihood function  $p(\mathbf{y}|\mathbf{f})$ )
2:    $\mathbf{f} \leftarrow \mathbf{0}$  ▷ Initialization
3:   repeat ▷ Newton iteration
4:      $\mathbf{W} \leftarrow -\nabla\nabla \log p(\mathbf{y}|\mathbf{f})$ 
5:      $\mathbf{b} \leftarrow \mathbf{W}\mathbf{f} + \nabla \log p(\mathbf{y}|\mathbf{f})$ 
6:     Iteratively solve:  $(\mathbf{I} + \mathbf{W}^{\frac{1}{2}}\mathbf{K}\mathbf{W}^{\frac{1}{2}})\mathbf{v} = \mathbf{W}^{\frac{1}{2}}\mathbf{K}\mathbf{b}$ 
7:      $\mathbf{a} \leftarrow \mathbf{b} - \mathbf{W}^{\frac{1}{2}}\mathbf{v}$ 
8:      $\mathbf{f} \leftarrow \mathbf{K}\mathbf{a}$ 
9:   until convergence ▷ Objective  $-\frac{1}{2}\mathbf{a}^T\mathbf{f} + \log p(\mathbf{y}|\mathbf{f})$ 
10:  return  $\hat{\mathbf{f}} \leftarrow \mathbf{f}$  ▷ Mode
11: end function
```

The Kronecker product can be used to construct memory and computationally efficient algorithms. For example, for N data points a full $N \times N$ matrix needs $\mathcal{O}(N^2)$ memory storage and a matrix vector product is of order $\mathcal{O}(N^2)$ whereas for a Kronecker matrix in which each dimension d has cardinality N_d , this is $\mathcal{O}(\sum_{d=1}^D N_d^2)$ and $\mathcal{O}(N(\sum_{d=1}^D N_d))$, respectively [14]. Furthermore, the Cholesky and singular value decompositions both reduce from $\mathcal{O}(N^3)$ to $\mathcal{O}(\sum_{d=1}^D N_d^3)$ since these operations need only be applied to each component.

3 Maximum a Posteriori Estimation

The Laplace approximation approximates the posterior by a Gaussian centered on the mode of the posterior. The mode $\hat{\mathbf{f}}$ of the posterior is the maximizer of the log posterior, which can be found by setting the first derivative of the log posterior to zero and solving for $\hat{\mathbf{f}}$. We can find the maximum iteratively by using Newton's algorithm for which we need to compute in each iteration the following update step:

$$\begin{aligned} \mathbf{f}^{new} &= (\mathbf{K}^{-1} + \mathbf{W})^{-1}\mathbf{b} \\ &= \mathbf{K}(\mathbf{I} - \mathbf{W}^{\frac{1}{2}}(\mathbf{I} + \mathbf{W}^{\frac{1}{2}}\mathbf{K}\mathbf{W}^{\frac{1}{2}})^{-1}\mathbf{W}^{\frac{1}{2}}\mathbf{K})\mathbf{b} \end{aligned} \quad (4)$$

with $\mathbf{b} = \mathbf{W}\mathbf{f} + \nabla \log p(\mathbf{y}|\mathbf{f})$. In the last line, we used the matrix inversion lemma to rewrite the Newton iteration in a numerically advantageous way. Conjugate Gradient (CG) can now be used iteratively to solve the linear system $(\mathbf{I} + \mathbf{W}^{\frac{1}{2}}\mathbf{K}\mathbf{W}^{\frac{1}{2}})\mathbf{v} = \mathbf{W}^{\frac{1}{2}}\mathbf{K}\mathbf{b}$, avoiding the need to directly compute any expensive matrix inverses. The mode is then found by setting $\mathbf{a} = \mathbf{b} - \mathbf{W}^{\frac{1}{2}}\mathbf{v}$ and $\mathbf{f} = \mathbf{K}\mathbf{a}$. Pseudocode is given in Algorithm 1. In line 6 one can obtain an exact answer for a quadratic system if the iterated solver is run for N steps, however, reasonable results are usually already obtained in a small fraction of iterations. Since Kronecker matrix vector products are efficient this greatly reduces the cost of mode finding compared to a standard Cholesky implementation.

4 Model Selection

After obtaining the mode $\hat{\mathbf{f}}$ for a given hyperparameter θ , the next problem is to learn the value of θ . This can be done by minimizing the negative marginal likelihood

$$-\log p(\mathbf{y}|\mathbf{X}, \theta) \approx \frac{1}{2}\hat{\mathbf{f}}^T \mathbf{K}^{-1}\hat{\mathbf{f}} - \log p(\mathbf{y}|\hat{\mathbf{f}}) + \frac{1}{2} \log |\mathbf{B}| \quad (5)$$

with $\mathbf{B} = \mathbf{I} + \mathbf{K}\mathbf{W}$. The first two terms are easily computed using the results of Newton's algorithm discussed in the previous paragraph, but the $\log |\mathbf{B}|$ and its gradients present difficulties. The \mathbf{B} matrix is not a Kronecker product and a direct implementation of the log determinant has order $\mathcal{O}(N^3)$.

To reduce the computational complexity, we propose to approximate the prior covariance matrix

Algorithm 2 Laplace negative approximate marginal likelihood

```

1: function LAPLACEMARGLIK(covariance matrix  $\mathbf{K}$ , observations  $\mathbf{y}$ , likelihood function  $p(\mathbf{y}|\mathbf{f})$ )
2:    $[\hat{\mathbf{f}}, \mathbf{a}] \leftarrow$  Obtain from Algorithm 1
3:   for  $d \leftarrow 1, \dots, D$  do
4:      $\mathbf{Q}_d, \mathbf{S}_d \leftarrow$  SCHUR( $\mathbf{K}_d$ ) ▷ Eigen decomposition  $\mathbf{K}_d = \mathbf{Q}_d \mathbf{S}_d \mathbf{Q}_d^T$ 
5:      $\mathbf{Q}_d, \mathbf{S}_d \leftarrow$  REDUCERANK( $\mathbf{Q}_d, \mathbf{S}_d$ ) ▷ Select largest eigenvalues
6:   end for
7:    $\mathbf{W} \leftarrow -\nabla \nabla \log p(\mathbf{y}|\mathbf{f})$ 
8:    $\mathbf{\Lambda}_1 \leftarrow \text{diag}(\text{diag}(\mathbf{K}) - (\mathbf{Q} \circ \mathbf{Q}) \cdot \text{diag}(\mathbf{S}))$  ▷ Eq. (6)
9:    $\mathbf{\Lambda}_2 \leftarrow \mathbf{I} + \mathbf{W}^{\frac{1}{2}} \mathbf{\Lambda}_1 \mathbf{W}^{\frac{1}{2}}$ 
10:   $\mathbf{\Lambda}_3 \leftarrow \mathbf{W}^{\frac{1}{2}} \mathbf{\Lambda}_2^{-1} \mathbf{W}^{\frac{1}{2}}$ 
11:   $\mathbf{L} \leftarrow \text{CHOLESKY}(\mathbf{S}^{-1} + \mathbf{Q}^T \mathbf{\Lambda}_3 \mathbf{Q})$  } ▷ Eq. (9)
12:   $\log \det B \leftarrow \sum_i \log \Lambda_{2ii} + \sum_i \log S_{ii} + 2 * \sum_i \log L_{ii}$ 
13:  return  $-\log q(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) \leftarrow \frac{1}{2} \mathbf{a}^T \hat{\mathbf{f}} - \log p(\mathbf{y}|\hat{\mathbf{f}}) + \frac{1}{2} \log \det B$  ▷ neg. approx. log. marg. lik.
14: end function

```

using a reduced-rank approximation:

$$\mathbf{K} \approx \mathbf{Q} \mathbf{S} \mathbf{Q}^T + \mathbf{\Lambda}_1, \text{ with } \mathbf{Q} = \bigotimes_{d=1}^D \mathbf{Q}_d, \mathbf{S} = \bigotimes_{d=1}^D \mathbf{S}_d, \mathbf{\Lambda}_1 = \text{diag}(\text{diag}(\mathbf{K}) - \text{diag}(\mathbf{Q} \mathbf{S} \mathbf{Q}^T)) \quad (6)$$

and where \mathbf{K} is an $N \times N$ matrix, \mathbf{Q} is an $N \times R$ matrix, and \mathbf{S} is an $R \times R$ diagonal matrix with $R \ll N$. The $\mathbf{\Lambda}_1$ matrix is a diagonal matrix that preserves the exact full-rank diagonal to obtain a better approximation without raising the computational costs. The sum of a diagonal matrix and a low rank matrix allows for efficient computations using the Woodbury formula and matrix determinant lemma

$$(\mathbf{A} + \mathbf{U} \mathbf{W} \mathbf{V})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{W}^{-1} + \mathbf{V} \mathbf{A}^{-1} \mathbf{U})^{-1} \mathbf{V} \mathbf{A}^{-1} \quad (7)$$

$$|\mathbf{A} + \mathbf{U} \mathbf{W} \mathbf{V}^T| = |\mathbf{W}^{-1} + \mathbf{V}^T \mathbf{A}^{-1} \mathbf{U}| |\mathbf{W}| |\mathbf{A}| \quad (8)$$

and has been used widely before in statistics and machine learning [1, 7]. The *key insight* here is that the Kronecker product allows the decomposition $\mathbf{Q} \mathbf{S} \mathbf{Q}^T$ to be *efficiently* obtained in $\mathcal{O}(\sum_{d=1}^D N_d^3)$.

4.1 Evaluating the Marginal Likelihood

With the reduced-rank approximation the log determinant in the negative marginal likelihood can be evaluated more efficiently using the relation:

$$\begin{aligned} |\mathbf{B}| &= |\mathbf{I} + \mathbf{W}^{\frac{1}{2}} \mathbf{K} \mathbf{W}^{\frac{1}{2}}| \approx |\mathbf{I} + \mathbf{W}^{\frac{1}{2}} \mathbf{\Lambda}_1 \mathbf{W}^{\frac{1}{2}} + \mathbf{W}^{\frac{1}{2}} \mathbf{Q} \mathbf{S} \mathbf{Q}^T \mathbf{W}^{\frac{1}{2}}| = |\mathbf{\Lambda}_2 + \mathbf{W}^{\frac{1}{2}} \mathbf{Q} \mathbf{S} \mathbf{Q}^T \mathbf{W}^{\frac{1}{2}}| \\ &= |\mathbf{\Lambda}_2| |\mathbf{S}| |\mathbf{S}^{-1} + \mathbf{Q}^T \mathbf{W}^{\frac{1}{2}} \mathbf{\Lambda}_2^{-1} \mathbf{W}^{\frac{1}{2}} \mathbf{Q}| = |\mathbf{\Lambda}_2| |\mathbf{S}| |\mathbf{S}^{-1} + \mathbf{Q}^T \mathbf{\Lambda}_3 \mathbf{Q}| \end{aligned} \quad (9)$$

with $\mathbf{\Lambda}_2 = \mathbf{I} + \mathbf{W}^{\frac{1}{2}} \mathbf{\Lambda}_1 \mathbf{W}^{\frac{1}{2}}$ and $\mathbf{\Lambda}_3 = \mathbf{W}^{\frac{1}{2}} \mathbf{\Lambda}_2^{-1} \mathbf{W}^{\frac{1}{2}}$. The final formula can be evaluated efficiently since it consists of two diagonal matrices and one low-rank matrix. To evaluate the last term we use the Cholesky decomposition $\mathbf{L} = \text{CHOLESKY}(\mathbf{S}^{-1} + \mathbf{Q}^T \mathbf{\Lambda}_3 \mathbf{Q})$, which can be reused in the gradients. Pseudocode is given in Algorithm 2.

4.2 Gradients

We can optimize the marginal likelihood with respect to the hyperparameters $\boldsymbol{\theta}$. For this we need the partial derivatives $\partial q(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})/\theta_j$. Since $\hat{\mathbf{f}}$ and \mathbf{W} are implicit functions of $\boldsymbol{\theta}$ we need to take care of both the explicit and implicit derivatives [8].

$$\frac{\partial \log q(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_j} = \left. \frac{\partial \log q(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_j} \right|_{\text{explicit}} + \sum_{i=1}^N \frac{\partial \log q(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})}{\partial \hat{f}_i} \frac{\partial \hat{f}_i}{\partial \theta_j} \quad (10)$$

4.2.1 Explicit Derivatives – Kernel Hyperparameters

The explicit derivatives are given by:

$$\left. \frac{\partial \log q(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_j^c} \right|_{\text{explicit}} = \frac{1}{2} \hat{\mathbf{f}}^T \mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial \theta_j^c} \mathbf{K}^{-1} \hat{\mathbf{f}} - \frac{1}{2} \text{tr} \left((\mathbf{W}^{-1} + \mathbf{K})^{-1} \frac{\partial \mathbf{K}}{\partial \theta_j^c} \right). \quad (11)$$

The first term is easily computed using Kronecker matrix-vector products, but we approximate the second term using the low-rank decomposition. Since

$$\begin{aligned} (\mathbf{W}^{-1} + \mathbf{K})^{-1} &= \mathbf{W}^{\frac{1}{2}} (\mathbf{I} + \mathbf{W}^{\frac{1}{2}} \mathbf{K} \mathbf{W}^{\frac{1}{2}})^{-1} \mathbf{W}^{\frac{1}{2}} \\ &\approx \mathbf{W}^{\frac{1}{2}} (\mathbf{I} + \mathbf{W}^{\frac{1}{2}} \boldsymbol{\Lambda}_1 \mathbf{W}^{\frac{1}{2}} + \mathbf{W}^{\frac{1}{2}} \mathbf{Q} \mathbf{S} \mathbf{Q}^T \mathbf{W}^{\frac{1}{2}})^{-1} \mathbf{W}^{\frac{1}{2}} \\ &= \mathbf{W}^{\frac{1}{2}} (\boldsymbol{\Lambda}_2 + \mathbf{W}^{\frac{1}{2}} \mathbf{Q} \mathbf{S} \mathbf{Q}^T \mathbf{W}^{\frac{1}{2}})^{-1} \mathbf{W}^{\frac{1}{2}} \\ &= \mathbf{W}^{\frac{1}{2}} (\boldsymbol{\Lambda}_2^{-1} - \boldsymbol{\Lambda}_2^{-1} \mathbf{W}^{\frac{1}{2}} \mathbf{Q} (\mathbf{S}^{-1} + \mathbf{Q}^T \boldsymbol{\Lambda}_3^{-1} \mathbf{Q})^{-1} \mathbf{Q}^T \mathbf{W}^{\frac{1}{2}} \boldsymbol{\Lambda}_2^{-1}) \mathbf{W}^{\frac{1}{2}} \\ &= \boldsymbol{\Lambda}_3 - \boldsymbol{\Lambda}_3 \mathbf{Q} (\mathbf{S}^{-1} + \mathbf{Q}^T \boldsymbol{\Lambda}_3 \mathbf{Q})^{-1} \mathbf{Q}^T \boldsymbol{\Lambda}_3 \end{aligned} \quad (12)$$

we can rewrite the trace term as follows

$$\begin{aligned} \text{tr} \left((\mathbf{W}^{-1} + \mathbf{K})^{-1} \frac{\partial \mathbf{K}}{\partial \theta_j^c} \right) &= \text{tr} \left(\boldsymbol{\Lambda}_3 \frac{\partial \mathbf{K}}{\partial \theta_j^c} \right) - \text{tr} \left(\boldsymbol{\Lambda}_3 \mathbf{Q} (\mathbf{S}^{-1} + \mathbf{Q}^T \boldsymbol{\Lambda}_3 \mathbf{Q})^{-1} \mathbf{Q}^T \boldsymbol{\Lambda}_3 \frac{\partial \mathbf{K}}{\partial \theta_j^c} \right) \\ &= \text{tr} \left(\boldsymbol{\Lambda}_3 \frac{\partial \mathbf{K}}{\partial \theta_j^c} \right) - \text{tr} \left(\mathbf{V}_1^T \frac{\partial \mathbf{K}}{\partial \theta_j^c} \mathbf{V}_1 \right) \end{aligned} \quad (13)$$

with $\mathbf{V}_1 = \boldsymbol{\Lambda}_3 \mathbf{Q} \text{CHOLESKY} [(\mathbf{S}^{-1} + \mathbf{Q}^T \boldsymbol{\Lambda}_3 \mathbf{Q})^{-1}]$. In this last line we also used the cyclic property of the trace operator $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA})$ to obtain a trace of a low-rank $R \times R$ matrix.

4.2.2 Implicit Derivatives – Kernel Hyperparameters

To evaluate the implicit derivatives we need the following expression:

$$\frac{\partial \log q(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})}{\partial \hat{f}_i} = -\frac{1}{2} [(\mathbf{K}^{-1} + \mathbf{W})^{-1}]_{ii} \frac{\partial^3}{\partial f_i^3} \log p(\mathbf{y}|\hat{\mathbf{f}}) \quad (14)$$

The troublesome part can be rewritten using the low-rank decomposition as follows

$$\begin{aligned} (\mathbf{K}^{-1} + \mathbf{W})^{-1} &= (\mathbf{K}^{-1} + \mathbf{W}^{\frac{1}{2}} \mathbf{I} \mathbf{W}^{\frac{1}{2}})^{-1} \\ &= \mathbf{K} - \mathbf{K} \mathbf{W}^{\frac{1}{2}} (\mathbf{I} + \mathbf{W}^{\frac{1}{2}} \mathbf{K} \mathbf{W}^{\frac{1}{2}})^{-1} \mathbf{W}^{\frac{1}{2}} \mathbf{K} \\ &= \mathbf{K} - \mathbf{K} (\mathbf{W}^{-1} + \mathbf{K})^{-1} \mathbf{K} \\ &\approx \mathbf{K} - \mathbf{K} \boldsymbol{\Lambda}_3 \mathbf{K} + \mathbf{K} \boldsymbol{\Lambda}_3 \mathbf{Q} (\mathbf{S}^{-1} + \mathbf{Q}^T \boldsymbol{\Lambda}_3 \mathbf{Q})^{-1} \mathbf{Q}^T \boldsymbol{\Lambda}_3 \mathbf{K} \\ &= \mathbf{K} - \mathbf{K} \boldsymbol{\Lambda}_3 \mathbf{K} + \mathbf{V}_2 \mathbf{V}_2^T \end{aligned} \quad (15)$$

with $\mathbf{V}_2 = \mathbf{K} \boldsymbol{\Lambda}_3 \mathbf{Q} \text{CHOLESKY} [(\mathbf{S}^{-1} + \mathbf{Q}^T \boldsymbol{\Lambda}_3 \mathbf{Q})^{-1}] = \mathbf{K} \mathbf{V}_1$. The diagonal $\text{diag}[(\mathbf{K}^{-1} + \mathbf{W})^{-1}]$ can be easily obtained using Eq. (15). The remaining term

$$\frac{\partial \hat{\mathbf{f}}}{\partial \theta_j^c} = (\mathbf{I} + \mathbf{K} \mathbf{W})^{-1} \frac{\partial \mathbf{K}}{\partial \theta_j^c} \nabla \log p(\mathbf{y}|\hat{\mathbf{f}}) \quad (16)$$

can be solved using conjugate gradient avoiding the need to directly compute the expensive matrix inverse.

4.2.3 Explicit Derivatives – Likelihood Hyperparameters

The explicit derivatives of the likelihood are given by:

$$\left. \frac{\partial \log q(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_j^l} \right|_{\text{explicit}} = \frac{\partial}{\partial \theta_j^l} \log p(\mathbf{y}|\mathbf{f}) + \frac{1}{2} [(\mathbf{K}^{-1} + \mathbf{W})^{-1}]_{ii} \frac{\partial}{\partial \theta_j^l} \nabla^2 \log p(\mathbf{y}|\mathbf{f}) \quad (17)$$

in which the second term can be approximated as outlined in Eq. (15).

Algorithm 3 Laplace gradients

```
1: function LAPLACEGRADIENT(covariance matrix  $\mathbf{K}$ , observations  $\mathbf{y}$ , likelihood function  $p(\mathbf{y}|\mathbf{f})$ )
2:    $[\hat{\mathbf{f}}, \mathbf{Q}, \mathbf{L}, \mathbf{\Lambda}_3] \leftarrow$  Obtain from Algorithm 1 and Algorithm 2
3:    $\mathbf{V}_1 \leftarrow \mathbf{\Lambda}_3 \mathbf{Q} / \mathbf{L}$   $\triangleright \mathbf{V}_1 = \mathbf{\Lambda}_3 \mathbf{Q} \text{CHOLESKY} [(\mathbf{S}^{-1} + \mathbf{Q}^T \mathbf{\Lambda}_3 \mathbf{Q})^{-1}]$ 
4:    $\mathbf{V}_2 \leftarrow \mathbf{K} \mathbf{V}_1$ 
5:    $\mathbf{C} \leftarrow \text{diag}(\mathbf{K}) - (\mathbf{K} \circ \mathbf{K}) \cdot \text{diag}(\mathbf{\Lambda}_3) + \text{SUMROWS}(\mathbf{V}_2 \circ \mathbf{V}_2)$   $\triangleright$  Eq. (15)
6:    $s_2 \leftarrow \frac{1}{2} \mathbf{C} \cdot \nabla^3 \log p(\mathbf{y}|\mathbf{f})$   $\triangleright$  Eq. (14)
7:   for  $j \leftarrow 1, \dots, \dim(\boldsymbol{\theta}^c)$  do  $\triangleright$  loop over kernel hyperparameters
8:      $s_1 \leftarrow \frac{1}{2} \mathbf{a}^T \frac{\partial \mathbf{K}}{\partial \theta_j^c} \mathbf{a} - \frac{1}{2} \text{diag}(\mathbf{\Lambda}_3)^T \text{diag}(\frac{\partial \mathbf{K}}{\partial \theta_j^c}) + \frac{1}{2} \text{tr}(\mathbf{V}_1^T \frac{\partial \mathbf{K}}{\partial \theta_j^c} \mathbf{V}_1)$   $\triangleright$  Eqs. (11–13)
9:     Iteratively solve:  $(\mathbf{I} + \mathbf{K} \mathbf{W}) \mathbf{s}_3 = \mathbf{K} \nabla \log p(\mathbf{y}|\mathbf{f})$   $\triangleright$  Eq. (16) using conjugate gradient
10:     $\nabla_j - \log p(\mathbf{y}|\mathbf{f}) \leftarrow -s_1 - \mathbf{s}_2^T \mathbf{s}_3$   $\triangleright$  Gradient  $-\frac{\partial}{\partial \theta_j^c} \log p(\mathbf{y}|\mathbf{f})$ 
11:  end for
12:  for  $j \leftarrow 1, \dots, \dim(\boldsymbol{\theta}^l)$  do  $\triangleright$  loop over likelihood hyperparameters
13:     $s_1 \leftarrow \frac{\partial}{\partial \theta_j^l} \log p(\mathbf{y}|\mathbf{f}) + \frac{1}{2} \mathbf{C} \frac{\partial}{\partial \theta_j^l} \nabla^2 \log p(\mathbf{y}|\mathbf{f})$   $\triangleright$  Eq. (17)
14:    Iteratively solve:  $(\mathbf{I} + \mathbf{K} \mathbf{W}) \mathbf{s}_3 = \mathbf{K} \frac{\partial}{\partial \theta_j^l} \nabla \log p(\mathbf{y}|\mathbf{f})$   $\triangleright$  Eq. (18) using conjugate gradient
15:     $-\nabla_j \log p(\mathbf{y}|\mathbf{f}) \leftarrow -s_1 - \mathbf{s}_2^T \mathbf{s}_3$   $\triangleright$  Gradient  $-\frac{\partial}{\partial \theta_j^l} \log p(\mathbf{y}|\mathbf{f})$ 
16:  end for
17:  return  $-\nabla \log p(\mathbf{y}|\mathbf{f})$   $\triangleright$  Gradient  $-\frac{\partial}{\partial \boldsymbol{\theta}} \log p(\mathbf{y}|\mathbf{f})$ 
18: end function
```

4.2.4 Implicit Derivatives – Likelihood Hyperparameters

The implicit derivatives of the likelihood (cf. Eq. (10)) are computed by combining Eq. (14) with

$$\frac{\partial \hat{\mathbf{f}}}{\partial \theta_j^l} = (\mathbf{I} + \mathbf{K} \mathbf{W})^{-1} \mathbf{K} \frac{\partial}{\partial \theta_j^l} \nabla \log p(\mathbf{y}|\mathbf{f}) \quad (18)$$

of which the latter is solved using conjugate gradient.

5 Memory and Run Time Complexity of the Algorithms

Algorithms 1, 2, and 3 are both *memory* and *run time* efficient. Care is taken to never evaluate a full $N \times N$ matrix, which may not fit into memory. The matrices $\mathbf{W}, \mathbf{b}, \mathbf{\Lambda}_1, \mathbf{\Lambda}_2, \mathbf{\Lambda}_3$ are diagonal each requiring $\mathcal{O}(N)$ storage, \mathbf{K} is a Kronecker matrix requiring $\mathcal{O}(\sum_{d=1}^D N_d^2)$ storage, $\mathbf{S} = \otimes \mathbf{S}_d$ requires at most $\mathcal{O}(R)$ storage, \mathbf{L} requires $\mathcal{O}(R^2)$ storage, and $\mathbf{Q} = \otimes \mathbf{Q}_d, \mathbf{V}_1, \mathbf{V}_2$ require (at most) $\mathcal{O}(RN)$ storage. Because all matrix and vector operations operate on small, i.e., $R \times N$, diagonal, or Kronecker matrices, instead of full $N \times N$ matrices, an increase in efficiency is obtained when N becomes large. Note that in Algorithm 2, line 8 and Algorithm 3, line 5 we made use of the identity

$$\text{diag}(\mathbf{A} \mathbf{B} \mathbf{C}^T) = (\mathbf{A} \circ \mathbf{C}) \cdot \text{diag}(\mathbf{B}) \quad (19)$$

which holds when \mathbf{B} is a diagonal matrix to speed up computations. For a Kronecker matrix \mathbf{K} , the Hadamard product $\mathbf{K} \circ \mathbf{K}$ can be evaluated efficiently and is again a Kronecker matrix.

6 Experiments

In order to compare the run time complexity of a standard Gaussian process with a Gaussian process having a Kronecker product kernel we implemented the algorithms in the publicly available Gaussian process toolbox GPstuff [12] and ran the algorithms on a synthetic dataset. As input we use gridded data on the D -dimensional hypercube $\mathcal{X} = \prod_{d=1}^D [0, 1]$. We consider the binary classification task of a datapoint belonging to the D -dimensional unit ball in \mathcal{X} .

For the likelihood model we use a probit model $p(y_i|f_i) = \Phi(y_i f_i)$ with $y_i \in \{-1, 1\}$. We generated inputs with $D = 2$ using varying grid sizes and ran both the standard Laplace approximation and the

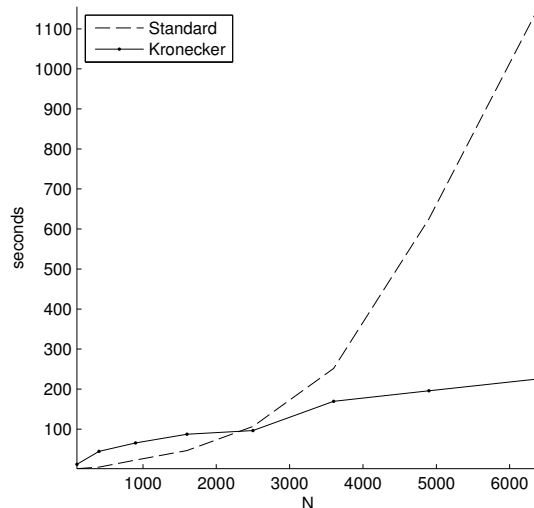


Figure 1: Run time of the standard Laplace algorithm versus the Kronecker Laplace algorithm on a 2D classification problem with a varying number of input data points.

Kronecker Laplace approximation and measured their run-time averaged over five runs. The results are shown in Figure 1. For the low-rank approximation we only selected the eigenvalues higher than $1E-6$, but if necessary only select from the highest eigen values such that the rank is below 10% of the number of input data points. The tolerance in the conjugate gradient iterations was set to $1E-10$. Clearly, the increase in computational complexity with increasing data size is much lower for the Kronecker Laplace approximation than for the standard Laplace approximation.

7 Conclusions

In this paper we have presented a scalable Gaussian process framework for gridded inputs and non-Gaussian observation models that factorize over cases. We used the Laplace approximation to perform approximate inference and showed it can efficiently be combined with tensor product kernels using fast matrix-vector operations and an efficiently obtained low-rank approximation of the kernel matrix. We provided derivations and pseudocode of the algorithms and empirically validated our approach on a binary classification problem showing a major performance improvement in terms of run time.

Directions for further work include among others handling missing observations and a more detailed analysis of the computational complexity. The standard Laplace approximation is quite straightforward, but the Kronecker Laplace approximation provides more options that each influence the computational complexity and quality of the algorithms. These include, among others, the rank of the low-rank approximation, the tolerances used in the Conjugate Gradients Iterated solvers, and the number of hyper-parameters.

References

- [1] N. Cressie and G. Johannesson. Fixed rank kriging for very large spatial data sets. *Journal of the Royal Statistical Society*, 70:209–226, 2008.
- [2] J. P. Cunningham, K. V. Shenoy, and M. Sahani. Fast Gaussian process methods for point process intensity estimation. In *Proceedings of the 25th international conference on Machine learning*, pages 192–199, 2008.
- [3] E. Gilboa, Y. Saatchi, and J. P. Cunningham. Scaling multidimensional Gaussian processes using projected additive approximations. In *International Conference on Machine Learning*, 2013.

- [4] A. Knutson and T. Tao. Honeycombs and sums of Hermitian matrices. *Notices Amer. Math. Soc.*, pages 175–186, 2000.
- [5] C. F. van Loan. The ubiquitous Kronecker product. *Journal of Computational and Applied Mathematics*, 123(1–2):85–100, 2000.
- [6] J. Quiñero-Candela and C. E. Rasmussen. A unifying view of sparse approximate Gaussian process regression. *Journal of Machine Learning Research*, 6:1939–1959, December 2005.
- [7] J. Quiñero-Candela, C.E. Rasmussen, and C.K.I. Williams. Approximation methods for Gaussian process regression. *Advances in Neural Information Processing Systems*, pages 203–223, 2007.
- [8] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [9] J. Riihimäki and A. Vehtari. Laplace approximation for logistic Gaussian process density estimation. *Bayesian Analysis*, in press, 2014.
- [10] E. Snelson and Z. Ghahramani. Sparse Gaussian processes using pseudo-inputs. *Advances in Neural Information Processing Systems*, 2006.
- [11] O. Stegle, C. Lippert, J. M. Mooij, N. D. Lawrence, and K. Borgwardt. Efficient inference in matrix-variate Gaussian models with iid observation noise. In *Advances in Neural Information Processing Systems 24*, pages 630–638, 2011.
- [12] J. Vanhatalo, J. Riihimäki, J. Hartikainen, P. Jylänki, V. Tolvanen, and A. Vehtari. GPstuff: Bayesian modeling with Gaussian processes. *Journal of Machine Learning Research*, 14(Apr):1175–1179, 2013.
- [13] C. K. I. Williams and D. Barber. Bayesian classification with Gaussian processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(12):1342–1351, 1998.
- [14] A. G. Wilson, E. Gilboa, A. Nehorai, and J. P. Cunningham. GPatt: Fast multidimensional pattern extrapolation with Gaussian processes. *arXiv*, 2013.