

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs
Laplace
Approximation
Kronecker product

MAP
estimation

Model
Selection

Experiment

Fast Laplace Approximation for Gaussian Processes with a Tensor Product Kernel

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

Radboud University Nijmegen^a

Erasmus University Rotterdam^b

Introduction

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

Gaussian process models

- + rich, principled Bayesian framework
- scalability $\mathcal{O}(N^3)$

Approaches addressing scalability

- approximations / subset selection $\mathcal{O}(M^2N)$
- exploit additional structure
 - Tensor product kernels
 - + Efficient use of Kronecker products on grid data
 - Limited to standard regression, since if $\mathbf{K} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ then $(\mathbf{K} + \sigma^2\mathbf{I})^{-1}\mathbf{y} = \mathbf{Q}(\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1}\mathbf{Q}^T\mathbf{y}$.

Gaussian Processes

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

A **Gaussian process** (GP) is collection of random variables with the property that the joint distribution of any finite subset is a Gaussian.

A GP specifies a probability distribution over functions $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$ and is fully specified by its mean function $m(\mathbf{x})$ and covariance (or kernel) function $k(\mathbf{x}, \mathbf{x}')$.

Typically $m(\mathbf{x}) = \mathbf{0}$, which gives

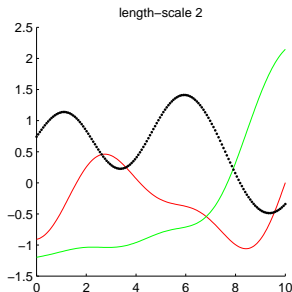
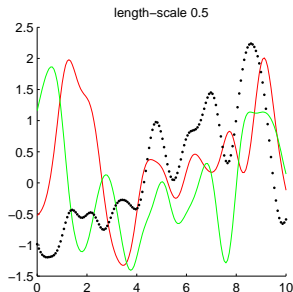
$$\{f(\mathbf{x}_1), \dots, f(\mathbf{x}_l)\} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}) \text{ with } K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

Gaussian Processes - Covariance function

Squared exponential (or Gaussian) covariance function:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\theta^2} \sum_{n=1}^N (x_n - x'_n)^2\right)$$

where θ is a length-scale parameter denoting how quickly the functions are to vary.



Gaussian Processes - 1D demo

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

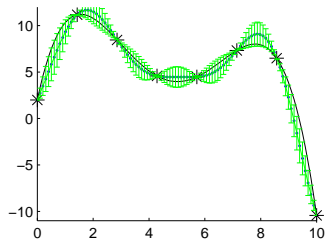
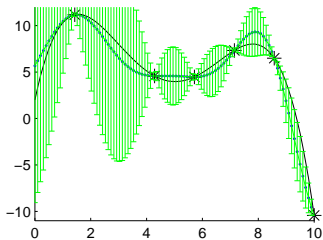
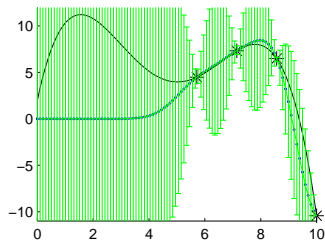
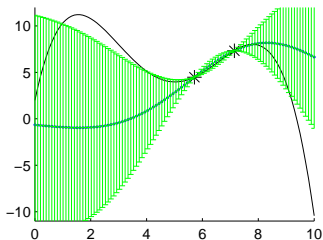
MAP

estimation

Model

Selection

Experiment



Laplace Approximation

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace
Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

For Gaussian process models with non-gaussian likelihood models we need approximations.

Laplace: approximate true posterior $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ with a Gaussian $q(\mathbf{f})$ centered on the mode of the posterior:

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

with $\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} p(\mathbf{f}|\mathbf{X}, \mathbf{y}, \theta)$ and $\mathbf{W} = -\nabla \nabla_{\mathbf{f}} \log p(\mathbf{y}|\mathbf{f})|_{\mathbf{f}=\hat{\mathbf{f}}}$.

Kronecker product

Assumptions

- Tensor product kernel function:
 $k(\mathbf{x}_i, \mathbf{x}_j) = \prod_{d=1}^D k_d(\mathbf{x}_i^d, \mathbf{x}_j^d)$.
- data on multi-dimensional grid

This results in a kernel matrix that decomposes into a Kronecker product of matrices of lower dimensions.

$$\mathbf{K} = \mathbf{K}_1 \otimes \cdots \otimes \mathbf{K}_D$$

where

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix}$$

Kronecker product

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

The Kronecker product has a convenient algebra

$$(\mathbf{A} \otimes \mathbf{B})\text{vec}(\mathbf{X}) = \text{vec}(\mathbf{B}\mathbf{X}\mathbf{A}^T)$$

$$\mathbf{A}\mathbf{B} \otimes \mathbf{C}\mathbf{D} = (\mathbf{A} \otimes \mathbf{C})(\mathbf{B} \otimes \mathbf{D})$$

$$(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$$

Operation	Standard	Kronecker
Storage	$\mathcal{O}(N^2)$	$\mathcal{O}(\sum_{d=1}^D N_d^2)$
Matrix vector product	$\mathcal{O}(N^2)$	$\mathcal{O}(N(\sum_{d=1}^D N_d))$
Cholesky / SVD	$\mathcal{O}(N^3)$	$\mathcal{O}(\sum_{d=1}^D N_d^3)$

MAP estimation

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP
estimation

Model
Selection

Experiment

Let $\mathbf{b} = \mathbf{W}\mathbf{f} + \nabla \log p(\mathbf{y}|\mathbf{f})$. We need

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

Repeat until convergence:

- iteratively solve $(\mathbf{I} + \mathbf{W}^{\frac{1}{2}} \mathbf{K} \mathbf{W}^{\frac{1}{2}}) \mathbf{v} = \mathbf{W}^{\frac{1}{2}} \mathbf{K} \mathbf{b}$
- $\mathbf{a} = \mathbf{b} - \mathbf{W}^{\frac{1}{2}} \mathbf{v}$
- $\mathbf{f} = \mathbf{K} \mathbf{a}$

Marginal Likelihood

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

learn the value of θ , which can be done by minimizing the negative marginal likelihood

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

with $\mathbf{B} = \mathbf{I} + \mathbf{K}\mathbf{W}$.

Reduced-rank approximation

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

We use a reduced-rank approximation:

$$\mathbf{K} \approx \mathbf{Q}\mathbf{S}\mathbf{Q}^T + \Lambda_1$$

with

$$\mathbf{Q} = \bigotimes_{d=1}^D \mathbf{Q}_d, \quad N \times R \text{ matrix}$$

$$\mathbf{S} = \bigotimes_{d=1}^D \mathbf{S}_d, \quad R \times R \text{ matrix}$$

$$\Lambda_1 = \text{diag}(\text{diag}(\mathbf{K}) - \text{diag}(\mathbf{Q}\mathbf{S}\mathbf{Q}^T))$$

Evaluating Marginal Likelihood

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

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

Gradients

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

Need the gradients wrt θ of

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

which is given by

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

Explicit Derivatives – Kernel Hyperparameters

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

Given by

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

Explicit Derivatives – Kernel Hyperparameters

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

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

Explicit Derivatives – Kernel Hyperparameters

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

$$\begin{aligned} & \text{tr} \left(\left(\mathbf{W}^{-1} + \mathbf{K} \right)^{-1} \frac{\partial \mathbf{K}}{\partial \theta_j^c} \right) \\ &= \text{tr} \left(\Lambda_3 \frac{\partial \mathbf{K}}{\partial \theta_j^c} \right) - \text{tr} \left(\Lambda_3 \mathbf{Q} \left(\mathbf{S}^{-1} + \mathbf{Q}^T \Lambda_3 \mathbf{Q} \right)^{-1} \mathbf{Q}^T \Lambda_3 \frac{\partial \mathbf{K}}{\partial \theta_j^c} \right) \\ &= \text{tr} \left(\Lambda_3 \frac{\partial \mathbf{K}}{\partial \theta_j^c} \right) - \text{tr} \left(\mathbf{V}_1^T \frac{\partial \mathbf{K}}{\partial \theta_j^c} \mathbf{V}_1 \right) \end{aligned}$$

$$\text{with } \mathbf{V}_1 = \Lambda_3 \mathbf{Q} \text{chol} \left[\left(\mathbf{S}^{-1} + \mathbf{Q}^T \Lambda_3 \mathbf{Q} \right)^{-1} \right].$$

Remaining Gradients

Perry Groot,
Markus Peters
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

Similar arguments for computing the remaining gradients

Linear Algebra Shortcuts

- $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA})$
- $\text{tr}(\mathbf{AB}^T) = \text{SUMROWS}(\mathbf{A} \circ \mathbf{B})$
- $\text{diag}(\mathbf{ABC}^T) = (\mathbf{A} \circ \mathbf{C}) \cdot \text{diag}(\mathbf{B})$ when \mathbf{B} is diagonal

Experiment

Perry Groot,
Markus Peters,
Tom Heskes,
Wolfgang
Ketter

Introduction

Background

GPs

Laplace

Approximation

Kronecker product

MAP

estimation

Model

Selection

Experiment

Artificial classification data generated on $\mathcal{X} = [0, 1]^2$ with various grid sizes.

