Latent Variable Augmentation in Bayesian Inference

Applications for Gaussian Processes

vorgelegt von Dipl.-Ing. Théo Galy-Fajou geb. in Castres

von der Fakultät IV - Elektrotechnik und Informatik der Technischen Universität Berlin zur Erlangung des akademischen Grades

Doktor der Naturwissenschaften -Dr.-Ing.-

genehmigte Dissertation

Promotionsausschuss:

Vorsitzender: Prof. A

Gutachter: Prof. Manfred Opper

Gutachterin: Prof. C Gutachter: Prof. D

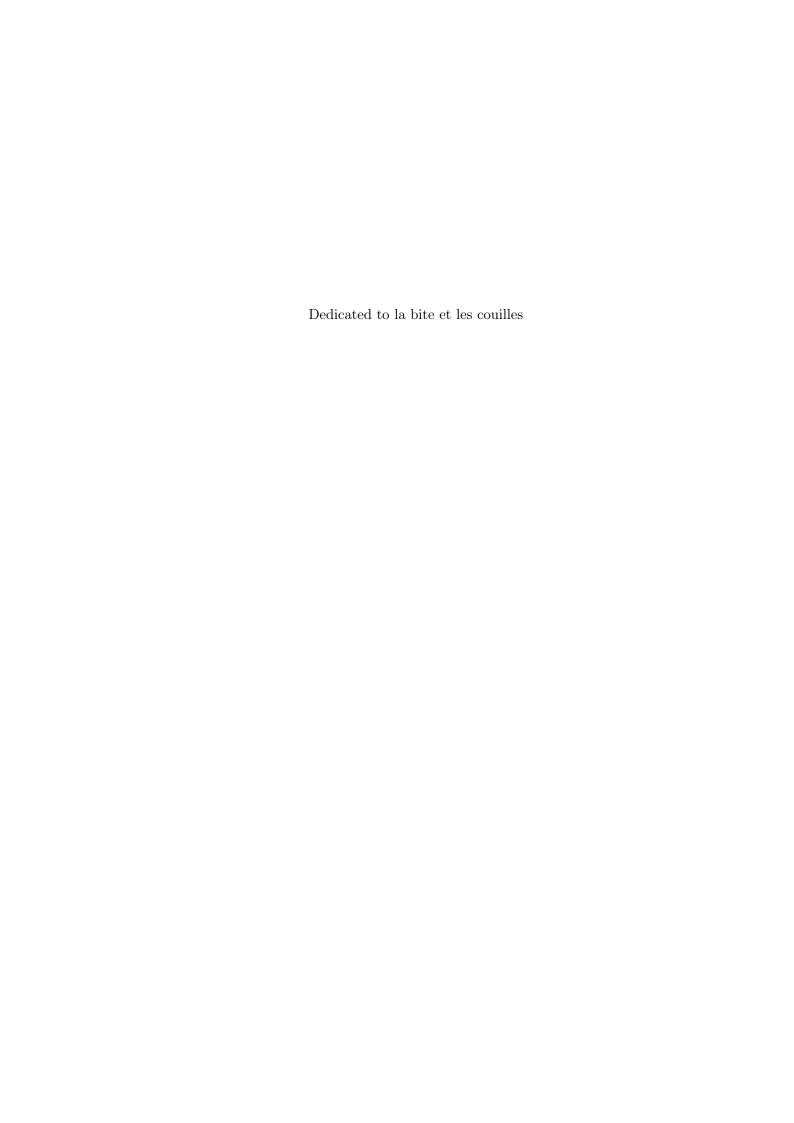
Tag der wissenschaftlichen Aussprache: XX. xxxx 2021

Zusammenfassung

Hier kommt der deutsche Abstrakt rein... ÜÖ sind ok.

Abstract

Put your abstract here...



Acknowledgements

I would like to acknowledge the thousands of individuals who have coded for open-source projects for free. It is due to their efforts that scientific work with powerful tools is possible.

Table of Contents

T	itle F	'age							1
\mathbf{Z}_{1}	usam	menfass	ung						iii
A	bstra	ıct							v
Li	ist of	Figures						3	xiii
Li	ist of	Tables							xv
A	bbre	viations						х	cvii
$\mathbf{S}_{\mathbf{y}}$	ymbo	ols						:	xix
1	Inti	coductio	n						1
	1.1	Bayesia	Machine Learning						1
	1.2	The unc	lerestimated importance of representation						1
	1.3	The use	of Gaussian Processes						2
	1.4	Thesis (Outline			•			2
2	Bac	kground	I						3
	2.1	Probabi	listic Bayesian Modeling						3
	2.2	Gaussia	n Processes						4
		2.2.1	Gaussian Process Regression						4
		2.2.2	Non-Conjugate Gaussian Processes						5
		2.2.3	Sparse Gaussian Processes						5
	2.3	Approx	mate Bayesian Inference						5
		2.3.1	Sampling						5
		2.3.2	Variational Inference						5
3	Effi tati		ussian Process Classification Using Polya-Gamma	Data	a A	.ue	gme	en-	. 9
,				- 4	_			4	
4			s Gaussian Process Classification Made Conjug la Data Augmentation	ate:	Ľ	iΠ	cie	nt	11
5	Aut	tomated	Augmented Conjugate Inference for Non-conjugate	gate	\mathbf{G}	au	ssi	an	
	\mathbf{Pro}	cess Mo	dels						13

TABLE OF CONTENTS

c	Westernal Country Destals Ele	1 -
b	Variational Gaussian Particle Flow	15
7	Discussion	17
\mathbf{A}	ppendix A Appendix A	19
\mathbf{R}	eferences	19

List of Figures

List of Tables

Abbreviations

This document is incomplete. The external file associated with the glossary 'abbreviations' (which should be called thesis.gls-abr) hasn't been created.

Check the contents of the file thesis.glo-abr. If it's empty, that means you haven't indexed any of your entries in this glossary (using commands like \gls or \glsadd) so this list can't be generated. If the file isn't empty, the document build process hasn't been completed.

Try one of the following:

• Add automake to your package option list when you load glossaries-extra.sty. For example:

\usepackage[automake]{glossaries-extra}

- Run the external (Lua) application: makeglossaries-lite.lua "thesis"
- Run the external (Perl) application: makeglossaries "thesis"

Then rerun LATEX on this document.

This message will be removed once the problem has been fixed.

Abbreviations

 \mathcal{GP} Gaussian Process

RKHS Reproducing Kernel Hilbert Space

 $\mathcal{GP}\mathbf{s}$ Gaussian Processes

 $\mathbf{MCMC}\,$ Markov Chain Monte Carlo

 ${f VI}$ Variational Inference

 \mathbf{VFE} Variational Free Energy

ELBO Evidence Lower BOund

 \mathbf{KL} Kullback-Leibler

 \mathbf{MF} Mean-Field

CAVI Coordinate Ascent Variational Inference

Symbols

This document is incomplete. The external file associated with the glossary 'symbolslist' (which should be called thesis.syi) hasn't been created.

Check the contents of the file thesis.syg. If it's empty, that means you haven't indexed any of your entries in this glossary (using commands like \gls or

\glsadd) so this list can't be generated. If the file isn't empty, the document build process hasn't been completed.

Try one of the following:

 Add automake to your package option list when you load glossaries-extra.sty. For example:

\usepackage[automake]{glossaries-extra}

- Run the external (Lua) application: makeglossaries-lite.lua "thesis"
- Run the external (Perl) application: makeglossaries "thesis"

Then rerun LATEX on this document.

This message will be removed once the problem has been fixed.

Introduction

Machine learning has become a wide field of research with a variety of sub-fields, each dedicated to solve various problems in different ways. One field in particular, usually called *probabilistic machine learning* aims at representing the statistical side of the different models.

- Motivate the idea of Bayesian machine learning
- Bring to the concept of relation between representation and inference
- Introduce each of the chapter properly

1.1 Bayesian Machine Learning

- Bayes is awesome
- Frequentist vs Bayesian

Bayesian Machine Learning is just another name for applied statistics on clean datasets.

In this thesis we are going to follow Bayesian principles. In a standard frequentist setting, one is interested in finding the best point estimate Using priors over latent parameters results in a posterior distribution instead of a point estimate. Posterior distributions have multiple advantages over point estimates: they are more robust to overfitting, they allow to compute prediction uncertainty and more. Of course they come at a higher computational cost: a distribution contains more information than a single point and finding analytical solutions is rare and often require manual derivations. However Bayesian methods gain an important edge when the datasets are small or when uncertainty about the predictions are uncertain. A typical example is in medicine, where data is scarce but the predictive outcome can have a dramatic effect (diagnosis, prognosis, etc...).

1.2 The underestimated importance of representation

• Different representation lead to very different results, efficiency etc

• Mention existing approaches

1.3 The use of Gaussian Processes

- All these things you can do with Gaussian processes
- why GPs vs other things

One of the strong

1.4 Thesis Outline

This thesis is constructed as follow:

- Chapter 2 will introduce in details all the common concepts to Bayesian inference and Gaussian Processes ($\mathcal{GP}s$). This background is generally introduced in each of the published articles, but this chapter allows to go more in-depth in the background theory. Bayesian inference will be properly introduced with a focus on variational inference and sampling.
- Chapter 3 introduces the paper [PUT PAPER NAME HERE], which was the first step in this work using augmentations to improve and scale up inference.
- Chapter 4 introduced the paper [PUT PAPER NAME HERE]. This paper brings new concepts of augmentation to a much more complex problem: multiclass classification.
- Chapter 5 introduces the paper [PUT PAPER NAME HERE]. This work was the first generalization of one type of augmentation and allowed to get a much better understanding of these concepts.
- Chapter 6 introduces a different way of representing variational inference with Gaussians using particle and defining optimal dynamics for those.

Background

A short introduction to the basic theory of Gaussian Processes and their scaling to larger is given in each paper. However, in this chapter we go a bit more into detail on how inference and prediction can be computed.

2.1 Probabilistic Bayesian Modeling

The Bayes' theorem is one of the simplest theorem in probabilities and its demonstration holds in one line, yet its implications are very important.

Let's give the very general modeling setting. Given a set of observed variables \boldsymbol{x} , a set of latent (unobserved) variables $\boldsymbol{\theta}$ with a prior distribution $p(\boldsymbol{\theta})$, and a likelihood function $p(\boldsymbol{x} \mid \boldsymbol{\theta})$, we can get the posterior distribution $p(\boldsymbol{\theta}|\boldsymbol{x})$:

$$p(\boldsymbol{\theta}|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\boldsymbol{x})} = \frac{p(\boldsymbol{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\boldsymbol{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}}$$
(2.1)

The interest of the posterior distribution is for making prediction on previously unseen data. Let's take the simple example of logistic regression: Given some input $x \in \mathbb{R}^D$ and binary label $y \in \{0,1\}$ we model the generative model as:

$$y \sim \text{Bernoulli}\left(\sigma(\boldsymbol{\theta}^{\top}\boldsymbol{x})\right),$$
 (2.2)

where $\boldsymbol{\theta} \in \mathbb{R}^D$ and σ is the logistic function $\sigma(x) = \frac{1}{1 + \exp(-x)}$. For simplicity we use an isotropic Normal prior on $\boldsymbol{\theta} : p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|0, I_D)$ and use the following likelihood function: $p(y_i|\boldsymbol{\theta}, \boldsymbol{x}_i) = \sigma\left(2(y_i - 1)\boldsymbol{\theta}^{\top}\boldsymbol{x}_i\right)$. Given the posterior $p(\boldsymbol{\theta}|\boldsymbol{y}, \boldsymbol{X})$ we can make predictions for new data using the following:

$$p(y^*|\mathbf{X}^*, \mathbf{y}, \mathbf{X}) = \int p(y^*, \boldsymbol{\theta} | \mathbf{X}^* \mathbf{y}, \mathbf{X}) d\boldsymbol{\theta} = \int p(y^*|\boldsymbol{\theta}, \mathbf{X}^*) p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) d\boldsymbol{\theta}.$$
(2.3)

The last term of the equation involves the posterior distribution $p(\theta|y, X)$. To solve this integral, we must either be able to know the posterior in closed form and solve the integral

numerically, or be able to sample from it and compute this integral with Monte-Carlo integration. Computing the posterior (2.1) in closed-form involves computing the integral $\int p(\boldsymbol{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}$, which is intractable for most non-trivial models. In Section 2.3, we mention methods which help solving this issue by introducing approximations or ways to sample directly from the posterior.

2.2 Gaussian Processes

Gaussian Process (\mathcal{GP}) are a class of stochastic processes used as non-parametric approximations to functions. A \mathcal{GP} is a stochastic process X_t , where the joint distribution on any collection of variables X_t follows a (multivariate) Gaussian distribution. This Gaussian nature is what make them attractive since operations on Gaussian variables tend to be easier and many calculus have closed-form solutions. The Gaussian distribution is to statistics what the harmonic oscillator is to physics. Although, $\mathcal{GP}s$ are defined to be a non-parametric model, one needs to define the covariance between each variable of the process. One of the most popular interpretation of \mathcal{GP} is as a prior on functions in the Reproducing Kernel Hilbert Space (RKHS). In practice the RKHS is infinite-dimensional, to be able to perform any computation one needs to project it into a finite-dimensional space. Considering a function f we wish to approximate with a \mathcal{GP} , we need some data \mathbf{X} to evaluate f on. We then consider the finite-dimensional vector \mathbf{f} where $f_i = f(X_i)$.

One resorts to kernel functions [NEED TO CITE THIS]. The kernel matrix K is defined by $K_{ij} = k(x_i, x_j)$. K is positive-definite, i.e. for $K \in \mathbb{R}^{D \times D}$, and $x \in \mathbb{R}^D$, $x^{\top}Kx > 0$.

2.2.1 Gaussian Process Regression

We now have a prior on the realisation of the function f on some data X, $p(f) = \mathcal{N}(f|\mu_0, K)$. We can now add information about some noisy observations y we got for X:

$$y_i = f(X_i) + \epsilon_i, \tag{2.4}$$

where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$. This leads to the likelihood $p(y_i|f_i) = \mathcal{N}(y_i|f_i, \sigma^2)$. Fortunately, multiplying Gaussian probability distributions together lead to another Gaussian distribution function. The posterior for \mathbf{f} is given by $p(\mathbf{f}|\mathbf{y}) = \mathcal{N}(\mathbf{f}|\mathbf{y}, \mathbf{K} + \sigma^2 I)$. The prediction of f^* on a new point \mathbf{x}^* can be done by computing:

$$p(f^*|\boldsymbol{x}^*, \boldsymbol{X}, \boldsymbol{y}) = \int p(f^*|\boldsymbol{f}, \boldsymbol{x}^*) p(\boldsymbol{f}|\boldsymbol{X}, \boldsymbol{y}) d\boldsymbol{f}.$$
 (2.5)

This integral is analytically solvable and results in $p(f^*|\boldsymbol{x}^*, \boldsymbol{X}, \boldsymbol{y}) = \mathcal{N}(f^*|m^*, s^*)$ where $m^* = K_{\boldsymbol{x}^*, \boldsymbol{X}} (K_{\boldsymbol{X}, \boldsymbol{X}} + \sigma^2 I)^{-1} y$ and $s^* = K_{\boldsymbol{x}^*, \boldsymbol{x}^*} - K_{\boldsymbol{x}^*, \boldsymbol{X}} (K_{\boldsymbol{X}, \boldsymbol{X}} + \sigma^2 I)^{-1} K_{\boldsymbol{X}, \boldsymbol{x}^*}$.

2.2.2 Non-Conjugate Gaussian Processes

A Gaussian prior is only conjugate¹ to a Gaussian likelihood Therefore $\mathcal{GP}s$ only give a Gaussian posterior with a Gaussian likelihood, for all other cases we talk about *Non-Conjugate Gaussian Processes*.

Since the posterior is not analytically tractable, one has to resort to some of the methods presented in Section 2.3.

2.2.3 Sparse Gaussian Processes

One of the largest issue of $\mathcal{GP}s$, regardless of if they are conjugate or not is the scalability with the number of observed samples.

2.3 Approximate Bayesian Inference

The posterior distribution in Eq.(2.1) cannot be computed in closed-form for non-trivial problems. To still be able to make predictions and render the model useful one can resort to different approximations. Out of a very large number of methods two of the most used are sampling and variational inference.

2.3.1 Sampling

When the posterior $p(\boldsymbol{\theta}|\boldsymbol{x})$ is not available in closed-form, it may be possible to draw samples from it. The set of methods is far too large to be even mentioned in this thesis, I will restrict the scope to methods tailored or adapted to $\mathcal{GP}s$. I will especially focus on Markov Chain Monte Carlo (MCMC) methods, where a chain of variable $\boldsymbol{\theta}^t$ is created with a Markovian assumption ($\boldsymbol{\theta}^t$ depends only of $\boldsymbol{\theta}^{t-1}$) and where the stationary distribution of $\boldsymbol{\theta}^t$ is the same as the target distribution (in our case the posterior $p(\boldsymbol{\theta}|\boldsymbol{x})$.

2.3.2 Variational Inference

Variational Inference (VI), sometimes called Variational Bayes, consists in approximating the posterior with another parametrized distribution. Given a family of distributions Q, parametrized by parameters φ one aims to solve the following optimization problem:

$$\varphi^* = \arg_{\omega} \min KL (q_{\varphi}(\boldsymbol{\theta}) || p(\boldsymbol{\theta} | \boldsymbol{x})), \qquad (2.6)$$

where the KL (Kullback-Leibler) divergence is defined (for continuous distributions as:

$$KL(q(x)||p(x)) = \int q(x) \log \frac{q(x)}{p(x)} dx$$
(2.7)

The objective of equation (2.6) is generally not tractable. Since computing $p(\boldsymbol{\theta}|\boldsymbol{x})$ involves the typically intractable normalization constant $p(\boldsymbol{x})$, one resort to a surrogate function, the Variational Free Energy (VFE) (or its negative counterpart the Evidence Lower BOund

¹A prior is said conjugate to a given likelihood when the resulting posterior is of the same family of the prior.

(ELBO)):

$$KL(q_{\varphi}(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\boldsymbol{x})) = \int q_{\varphi}(\boldsymbol{\theta}) (\log q_{\varphi}(\boldsymbol{\theta}) - \log p(\boldsymbol{\theta}|\boldsymbol{x})) d\boldsymbol{\theta}$$
(2.8)

$$= \int q_{\varphi}(\boldsymbol{\theta}) \left(\log q_{\varphi}(\boldsymbol{\theta}) - \log p(\boldsymbol{\theta}, \boldsymbol{x}) - \log p(\boldsymbol{x}) \right) d\boldsymbol{\theta}$$
 (2.9)

$$= \underbrace{-\log p(\boldsymbol{x})}_{\leq 0} + \int q_{\varphi}(\boldsymbol{\theta}) \left(\log q_{\varphi}(\boldsymbol{\theta}) - \log p(\boldsymbol{x}|\boldsymbol{\theta}) - \log p(\boldsymbol{\theta})\right) d\boldsymbol{\theta} \qquad (2.10)$$

$$\leq -\mathbb{E}_{q_{\varphi}}\left[\log p(\boldsymbol{x}|\boldsymbol{\theta})\right] + \mathrm{KL}\left(q_{\varphi}(\boldsymbol{\theta})||p(\boldsymbol{\theta})\right) = \mathcal{F}(\varphi) \tag{2.11}$$

By minimizing the VFE: $\mathcal{F}(\varphi)$ instead of the Kullback-Leibler (KL) divergence, we expect to find a solution close to the optimum of the problem stated in (2.6). A standard way is to perform gradient descent on the variational parameters φ

$$\boldsymbol{\varphi}^{t+1} = \boldsymbol{\varphi}^t - \epsilon \nabla_{\boldsymbol{\varphi}} \mathcal{F}(\boldsymbol{\varphi}^t). \tag{2.12}$$

Computing the gradient $\nabla_{\varphi} \mathcal{F}(\varphi)$ can be non-trivial but many methods were developed to tackle this problem.

[INTRODUCE DIFFERENT METHODS HERE]. One of the most important addition to the VI method is the Mean-Field (MF) approximation. MF is the assumption that the variational distribution $q(\theta)$ assumes every component of θ to be independent from each other. This way we can write

$$q_{\varphi}^{MF}(\boldsymbol{\theta}) = \prod_{i=1}^{D} q_{\varphi_i}(\theta_i)$$
 (2.13)

A more general method is also to consider blocks of variables instead.

Following the MF approach, it is sometimes possible to find the optimal parameters φ^* in closed-form. By solving:

$$\nabla_{\varphi_i} \mathcal{F}(\varphi)|_{\varphi_i = \varphi_i^*} = 0, \tag{2.14}$$

for each variable φ_i we can find a local optima, which with additional assumptions, can prove to be the local minima we are looking for. The advantage of this method is that one can also perform it independently for each latent variable θ_i . Concretely the updates are of the form:

$$q_{\varphi_i}^*(\theta_i) \propto \exp\left(\mathbb{E}_{q_{\varphi}(\boldsymbol{\theta}_{/i})} \left[\log p\left(\theta_i | \boldsymbol{\theta}_{/i}, \boldsymbol{x}\right)\right]\right)$$
 (2.15)

where $\theta_{/i}$ represent the collection of variables $\theta_{/i} = \{\theta_j | j \neq i\}$. When working with distribution coming from exponential families, it is straightforward to get the optimal variational parameters φ_i . By updating the parameters one after another we get a Coordinate Ascent Variational Inference (CAVI) scheme². Effectively, one update each variational parameter φ_i by its optimum

²The word ascent is used since the scheme was originally derived using the negative VFE or ELBO.

given the rest of the variational parameters $\pmb{arphi}_{/i}$ via closed-form functions:

$$\varphi_i^{t+1} = f_i \left(\varphi_{1:(i-1)}^{t+1}, \varphi_{(i+1):D}^t \right).$$
 (2.16)

The order of the updates do not matter as long as the variational parameters φ are initialized in their domain.

Efficient Gaussian Process Classification Using Polya-Gamma Data Augmentation

Multi-Class Gaussian Process
Classification Made Conjugate:
Efficient Inference via Data
Augmentation

Automated Augmented Conjugate Inference for Non-conjugate Gaussian Process Models

Variational Gaussian Particle Flow

Discussion