ACTIVELY LEARNING HYPERPARAMETERS FOR GPS

Roman Garnett Washington University in St. Louis

12.10.2016

Joint work with Michael Osborne (University of Oxford) Philipp Hennig (MPI Tübingen)

INTRODUCTION

Learning hyperparameters

Problem

- Gaussian processes (GPs) are powerful models able to express a wide range of structure in nonlinear functions.
- This power is sometimes a *curse*, as it can be very difficult to determine appropriate values of *hyperparameters*, especially with small datasets.

Small datasets

- Small datasets are inherent in situations when the function of interest is very *expensive*, as is typical in *Bayesian optimization*.
- Success on these problems hinges on accurate modeling of *uncertainty*, and undetermined hyperparameters can contribute a great deal (*often hidden!*).
- The traditional approach in these scenarios is to spend some portion of the budget on *model-agnostic* initialization (Latin hypercubes, etc.)
- We present a *model-driven* approach here.

Motivating problem: Learning embeddings

- High-dimensionality has stymied the progress of model-based approaches to many machine learning tasks.
- In particular, *Gaussian processes* approaches remain *intractable* for large numbers of input variables.
- An old idea for combating this problem is to exploit *low-dimensional* structure in the function, the most simple example of which is a *linear embedding*.

Learning embeddings for GPs

• We want to learn a function

$$f: \mathbb{R}^D \to \mathbb{R},$$

where D is very large.

• We assume that f has low *intrinsic dimension*, that is, that there is a function $g \colon \mathbb{R}^d \to \mathbb{R}$ such that

$$f(x) = g(Rx),$$

where $R \in \mathbb{R}^{d \times D}$ is a matrix defining a *linear embedding*.

Example

- Here f: ℝ² → ℝ (D = 2), but only depends on a one-dimensional projection of x (d = 1).
- All function values are realized along the *black line*.



 x_1

7

If we knew the embedding R, modeling f would be *straightforward.* Our model for f given the embedding R is a zero-mean Gaussian process:

$$p(f \mid R) = \mathcal{GP}(f; 0, K),$$

with

$$K(x, x'; R) = \kappa(Rx, Rx'),$$

where κ is a covariance on $\mathbb{R}^d \times \mathbb{R}^d$.

If κ is the familiar squared-exponential, then

$$K(x, x'; R, \gamma) = \gamma^2 \exp\left[-\frac{1}{2}(x - x')^{\top} R^{\top} R(x - x')\right].$$

This is a low-rank *Mahalanobis* covariance, also known as a *factor analysis* covariance.

Our approach

- Our goal is to learn R (in general, any $\theta)$ as quickly as possible!
- Unlike previous approaches, which focus on *random embeddings* (Wang, et al. 2013), we focus on *learning the embedding directly.*

What can happen with random choices



Djolonga, et al. NIPS 2013

LEARNING THE HYPERPARAMETERS

Learning the hyperparameters

We maintain a *probabilistic belief* on θ . We start with a prior

 $p(\theta),$

and given data $\ensuremath{\mathcal{D}}$ we find the (approximate) posterior

 $p(\theta \mid \mathcal{D}).$

The uncertainty in θ (in particular, its *entropy*) measures our progress!



The prior is *arbitrary*, but here we took diffuse independent prior distribution on each entry:

$$p(\theta_i) = \mathcal{N}(\theta_i; 0, \sigma_i^2).$$

Could also use something more sophisticated.

Now, given observations \mathcal{D} , we *approximate* the posterior distribution on θ :

$$p(\theta \mid \mathcal{D}) \approx \mathcal{N}(\theta; \hat{\theta}, \Sigma).$$

The method of approximation is also *arbitrary*, but we took a Laplace approximation.

SELECTING INFORMATIVE POINTS

Active learning

Selecting informative points

- We wish to sequentially *sample the most informative point* about *θ*.
- We suggest maximizing the *mutual information* between the observed function value and the hyperparameters, particularly in the form known as *Bayesian active learning by disagreement* (BALD).¹

$$x^* = \arg\max_{x} H[y \mid x, \mathcal{D}] - \mathbb{E}_{\theta} \big[H[y \mid x, \mathcal{D}, \theta] \big].$$

¹Houlsby, et al. BAYESOPT 2011

Breaking this down, we want to find points with high *marginal uncertainty* (à la uncertainty sampling)...

$$x^* = \arg\max_{x} H[y \mid x, \mathcal{D}] - \mathbb{E}_{\theta} [H[y \mid x, \mathcal{D}, \theta]].$$

... but would have *low* uncertainty if we *knew the hyperparameters* θ :

$$x^* = \arg\max_{x} H[y \mid x, \mathcal{D}] - \mathbb{E}_{\theta} \Big[H[y \mid x, \mathcal{D}, \theta] \Big].$$

BALD

- That is, we want to find points where the competing models (one for each value of θ) are all certain, but disagree highly with each other.
- These points are the most informative points about the hyperparameters! (We can discard hyperparameters that were confident about the *wrong answer*).

Computation of BALD

How can we compute or approximate the BALD objective for our model?

$$x^* = \operatorname*{arg\,max}_{x} H[y \mid x, \mathcal{D}] - \mathbb{E}_{\theta} [H[y \mid x, \mathcal{D}, \theta]].$$

The first term (marginal uncertainty in y) is especially *troubling*...

LEARNING THE FUNCTION

Approximate marginalization of GP hyperparameters

Given data \mathcal{D} , and an input x^* , we wish to capture our belief about the associated latent value f^* , accounting for uncertainty in θ :

$$p(f^* \mid x^*, \mathcal{D}) = \int p(f^* \mid x^*, \mathcal{D}, \theta) p(\theta \mid \mathcal{D}) d\theta.$$

We provide an *approximation* called the "marginal GP" (MGP).

The MGP

The result is this:

$$p(f^* \mid x^*, \mathcal{D}) \approx \mathcal{N}(f^*; m_{\mathcal{D}}^*, C_{\mathcal{D}}^*),$$

where

$$m_{\mathcal{D}}^* = \mu_{\mathcal{D},\hat{\theta}}^*.$$

The approximate mean is the MAP posterior mean, and...

The MGP

$$C_{\mathcal{D}}^{*} = \frac{4}{3} V_{\mathcal{D},\hat{\theta}}^{*} + \frac{\partial \mu^{*}}{\partial \theta}^{\top} \Sigma \frac{\partial \mu^{*}}{\partial \theta} + (3V_{\mathcal{D},\hat{\theta}}^{*})^{-1} \frac{\partial V^{*}}{\partial \theta}^{\top} \Sigma \frac{\partial V^{*}}{\partial \theta}.$$

The variance is *inflated* according to how the posterior mean and posterior variance change with the hyperparameters. The MGP gives us a *simple approximation* to the BALD objective; we maximize the following simple objective:

$$\frac{C_{\mathcal{D}}^*}{V_{\mathcal{D},\hat{\theta}}^*}.$$

So we sample the point with maximal *variance inflation*. This is the point where the plausible hyperparameters *maximally disagree* under our approximation!

BALD and the MGP





Example

Consider a simple *one-dimensional* example (here R is simply an inverse length scale).

- The blue envelope shows the uncertainty given by the *MAP embedding*.
- The red envelope shows the additional uncertainty due to *not knowing the embedding.*
- We sample where the ratio of these is *maximized*.





The inset shows our belief over $\log R$, it *tightens* as we continue to sample.



Example



Example





We sample at a variety of separations to *further refine our* belief about R.



Example

Notice that we are *relatively uncertain* about many function values! Nonetheless, we are effectively learning R.



2d example



2d example



Results

- We have tested this approach on numerous synthetic and real-world regression problems up to dimension D = 318, and our performance was significantly superior to:
 - random sampling,
 - Latin-hypercube designs, and
 - uncertainty sampling.

Test setup

For each method/dataset, we:

- Began with a *single observation* of the function at the center of the (box-bounded) domain,
- Allowed each method to select a sequence of $n=100 \ {\rm observations},$
- Given the resulting training data, found the *MAP* hyperparameters, and
- Used these hyperparameters to test on a *held-out* set of 1000 points, measuring RMSE and negative log likelihood.

Results: RMSE

Choosing 100 observations, predicting on 1000 more.

dataset	D/d	RAND	LH	UNC	BALD
synthetic	10/2	0.412	0.371	0.146	0.138
synthetic	10/3	0.553	0.687	0.557	0.523
synthetic	20/2	0.578	0.549	0.551	0.464
synthetic	20/3	0.714	0.740	0.700	0.617
Branin	10/2	18.2	17.8	3.63	2.29
Branin	20/2	18.3	14.8	13.4	15.0
communities & crime	96/2	0.720		0.782	0.661
temperature	106/2	0.423	_	0.427	0.328
CT slices	318/2	0.878		0.845	0.767

Reminder

The framework we have presented for actively learning linear embeddings is *completely general;* we can use it for actively learning hyperparameters in *any GP model!*

Question

Both these approaches suggest a *two-stage approach* for optimization. Is this necessary? Can we use BALD to learn the embedding while simultaneously optimizing the function?

Code

github.com/rmgarnett/ active_gp_hyperlearning



For more details

UAI 2014

Actively Learning Linear Embeddings for Gaussian Processes, UAI 2014.

Active Learning of Linear Embeddings for Gaussian Processes

Roman Garnett University of Bonn Rösserstrade 164 53117 Benn, Germany

Michael A. Osheme University of Oxford Parks Read Oxford OX1 3P3, UK In DEEFT1064, OX1 3P3, UK Philipp Honaig MPI for Intelligent Systems Spenantotale 72876 Tubingen, Germany Phoenil 07110, apr. do

Abstract

We propose as active learning method for discovering low-dimensional structure in highdimensional Grannian persons (00%) tracks. Stack problems are incomingly frequent and important, but have history personal server penetroid difficulties. We darba intraches a novel tochmique for approximately marginalizing or hyperparamenters, yidding marginal problems (no binfers an efficient means of performing of Progravient, quadrature, of Bayesian optimization in high-dimensioning space.

1 INTRODUCTION

We propose an ended in a network has a statistication, then a fraction of the biomediane all problem (or the ender the statistication of the statistication of the ender of the statistication of the statistication of the ender the statistication of the statistication of the ender of the statistication of the statistication of the ender statistication of the statistication

as the possible ever a to quartery the uncertainty in the embedding distortant, by Ta darm the function, we a family process hyperpersenters (heading these parameters) process hyperpersenters (heading these parameters) enhances of the second productions mode to hyperparater micropolication (decision 3). This sub-algorithm is more generally applicable to may Gensine process tasks and to the marginalization of hyperparameters often than proper limit, the active advision, we astend previous well (headby et al., 2011) to solver existantion that management of headby and a sub-second second sec

Estimators for E in wale was include 1.4500 (Thibitizm) 1990) and the Daming subcore (Tanion & Tao, 2007), both of which assume 4 \sim 1. These are parative methods estimaing the linear embedding from a shead dataset. This paper develops an algorithm that actively learns *R* for the densiti at *L* dataset. The goal is not first function estiuations to intelligently explore and identity 12. Notice that dishengis the analysis for a low to be linear, the function of the function of the function of the function.

This problem in ordinal to, but datated from, dimensionaling reduction (Lensme, 2012), for which active learning has recordly been proposed (Jonta et al., 2013). Dimendending reductions in the lacowar as variantizations or Walan sources reportation, and its order using, e.g., predicable components analysive (Level, Stations and Stations), net consider the problem of finding in low-dimensional arguments and the reduction of the second stations. Net consider the problem of finding is low-dimensional arguments in the reduction of the second station of the second station of the second statistics of the problem of finding to reduction of the reduction of the second statistical statistics and the second reduction of the second statistics of the second statistics of the reduction of the second statistics of the second statistics of the reduction of the second statistics of the problem of finding statistics of the second statistics o

Extension: NIPS 2015

Extension to *model selection*, one step closer to fully automated Bayesian optimization!

Bayesian Active Model Selection with an Application to Automated Audiometry, NIPS 2015.

Bayesian Active Model Selection with an Application to Automated Audiometry

Jacob R. Gardner	Gustavo Malkomes	Roman Garnett
CS, Cornell University	CSE, WUSTL	CSE, WUSTL
Ithaca, NY 14850	St. Louis, MO 63130	St. Louis, MO 63130
jrg365@cornell.edu	luizgustavo@wustl.e	du garnett&wustl.edu
Kilian Q. Weinberger	Dennis Barbour	John P. Cunningham
CS, Cornell University	BME, WUSTL	Statistics, Columbia University
Ithaca, NY 14850	St. Louis, MO 63130	New York, NY 10027
kqw4@cornell.edu	dbarbour@wustl.edu	jpc2181@columbia.edu

Abstract

We introduce as novel information-theoretic approach for active model selection demonstrates in fortune selection of the selection of the selection of the structure for Gaussian process (cr) models with arbitrary observation likelihoods (cr), and a selection of the selection of the selection of the selection (cr), and a selection of the selection of the selection of the selection (cr), and a selection of the selection of the selection of the selection (cr), and a selection of the selection of the selection of the selection (cr), and a selection of the selection

1 Introduction

Personalized modicine has long been a critical application area for machine learning [1–3], in which automated decision making and diagnosis are key components. Beyood improving outputs of life, machine learning in diagnosis extratings is particularly important because collecting additional motical data chem itoms supplicant functial buotes, incre cost, and partical discontient, it machine learning (e.g., a blood tstd) licens some cost, but will, with hope, better inform diagnosis, treatment, and prognosis. By careful analysis, wan upy colimize this trade of

However, many diagnostic settingis in medicine do not involve feature selection, but rather involve querying a sample gaues to discriminate different model description gueinst artitristica. A puricular, clarifying example that motivates this work is notive-sheadowed hearing for IRONTL, a prevalent disorder affecting 20 million working agae adults in the United States alone [4] and affecting over half of preventable with simple, low-cost solutions (e.g., carplags). The critical requirement for prevention is effective early diagnosis.

To be tested for NHL, patients must complete a time-consuming audiometric exam that presents a series of sones at various frequencies and intensities; at each tone the patient indicates whether bolds hears the tone [5–7]. From the responses, the clinicitican inferst the patient studieb threshold on a set of discrete frequencies (the audiogram); this process requires the delivery of up to handreds of tones. Audiologists scan the audiogram for a bearing deficit with a characteristic note whape—

Extension: NIPS 2016

Bayesian optimization for automated model selection

Gustavo Malkomes¹, Chip Schaff¹, Roman Garnett Department of Computer Science and Engineering Washington University in St. Louis St. Louis, MO 63130 [luizgustavo, checkaff, garnett)@wustl.edu

Abstract

Despite heaves of kernel-based nonparameteris methods, hereal electron site people scattering and engines and the start of the start of the start exploitation function of the start of the start of the start of the start in the start of the likelihood. Our proposed search method is based on Bayesian optimization to and start of the start of t

1 Introduction

Over the part decades, ecoremons human efforts has been devoide to machine karming proposeding dense models attacked, the hypothesis of the decade structure of the structure of the structure of the decade structure of the off "black struct" is kernel methods in particular, the structure of the decade of the structure of the str

Recent work has begun to tack the kernel-selection problem in a systematic way. Duresaud et al. [1] and Grosse et al. [2] described appearing regrammes for enumering a countably infinite space of arbitrarily complex kernels via exploiting the closure of kernels (10) Gross ad attact, Duresaud et al. (2) described and a structure of the series of a work of the series of a structure of the series of a work of the series of a structure of the series of a work of the series of a structure of the series of the ser

In this work, we develop a more sophisticated mechanism for searching through this space. The greedy search described above only considers a given dataset by querying a model's evidence. Our search performs a *metadarning* precedure, which, conditional on a dataset, establishes similarities among the models in terms of the space of explanations they can offer for the data. With this viewpoint, we construct a novel kernel between models (a 'kernel kernel'). We then approach to prevent the second s

[†]These authors contributed equally to this work

Another extension to *model selection* with *fixed datasets*, one step closer to fully automated Bayesian optimization!

Bayesian optimization for automated model selection, NIPS 2016.

THANK YOU!

Questions?