
Approximate Hyperparameter Marginalisation for Gaussian Processes

Anonymous Author(s)

Affiliation

Address

email

Abstract

Enter Abstract here

1 Introduction

2 Gaussian Processes

Gaussian processes (GPs) constitute a powerful method for performing Bayesian inference about functions using a limited set of observations [2]. A GP is defined as a distribution over the functions $f : \mathcal{X} \rightarrow \mathbb{R}$ such that the distribution over the possible function values on any finite set of \mathcal{X} is multivariate Gaussian. Given some arbitrary size n dataset, the observations $\mathbf{y} = \{y_1, \dots, y_n\}$ could be viewed as a single point sampled from a n -variate Gaussian distribution and can be partnered with a GP.

A GP is completely defined by its first and second moments: a mean function $\mu : \mathcal{X} \rightarrow \mathbb{R}$ which describes the overall trend of the function, and a symmetric positive semidefinite covariance function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ which describes how function values are correlated as a function of their locations in the domain. Given a function $f : \mathcal{X} \rightarrow \mathbb{R}$ which we would like to perform inference about, and a set of input points $\mathbf{x} \subseteq \mathcal{X}$, the Gaussian process prior distribution over the function values $\mathbf{f} = f(\mathbf{x})$ is given by:

$$p(\mathbf{f}|\mathbf{x}, \boldsymbol{\theta}, I) := \mathcal{N}(\mathbf{f}; \mu_{\boldsymbol{\theta}}(\mathbf{x}), K_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{x})) \quad (1)$$

$$:= \frac{1}{\sqrt{\det 2\pi K_{\mathbf{f}}}} \exp\left(-\frac{1}{2}(\mathbf{f} - \mu_{\mathbf{f}})^{\top} K_{\mathbf{f}}^{-1}(\mathbf{f} - \mu_{\mathbf{f}})\right) \quad (2)$$

where $\boldsymbol{\theta}$ is a vector containing any parameters required by μ and K , and which form the *hyperparameters* of the model, I . Due to the ubiquity of I we henceforth drop it from explicit representation for notational convenience. There exist a wide variety of mean and covariance functions which can be chosen in order to reflect any prior knowledge available about the function of interest. Note also that we need not place a GP directly on the function, for example a function known to be strictly positive might benefit from a GP over its logarithm.

Once we have observations of the function $(\mathbf{x}_s, \mathbf{y}_s)$ we can make predictions about the function value f_* at input x_* . As exact measurements of the function are often not available we assume a noise model, such that:

$$p(y|f, x, \sigma_n^2) := \mathcal{N}(y; f, \sigma_n^2) \quad (3)$$

which represents i.i.d Gaussian observation noise with variance σ_n^2 . As should be expected, the predictive distribution over f_* is Gaussian:

$$p(f_*|x_*, \mathbf{y}_s, \boldsymbol{\theta}) := \mathcal{N}(f_*; m_{\boldsymbol{\theta}}(f_*|x_*, \mathbf{y}_s), C_{\boldsymbol{\theta}}(f_*|x_*, \mathbf{y}_s)) \quad (4)$$

where the posterior mean and covariance are:

$$m_\theta(f_*|x_*, \mathbf{y}_s) := \mu_\theta(x_*) + K_\theta(x_*, \mathbf{x}_s)V^{-1}(\mathbf{y}_s - \mu_\theta(x_*)) \quad (5)$$

$$C_\theta(f_*|x_*, \mathbf{y}_s) := K_\theta(x_*, x_*) - K_\theta(x_*, \mathbf{x}_s)V^{-1}K_\theta(\mathbf{x}_s, x_*) \quad (6)$$

$$\text{where } V := K_\theta(\mathbf{x}_s, \mathbf{x}_s) + \sigma_n^2 \mathbf{I} \quad (7)$$

2.1 Dealing with Hyperparameters

The previous equations assume that the hyperparameters θ are known; in fact we can rarely be certain about θ *a priori*. This ignorance can be represented by a suitably uninformative prior distribution $p(\theta)$. Given such a *hyper-prior*, the hyperparameters should be marginalised to calculate the predictive distribution over f_* :

$$p(f_*|x_*, \mathbf{y}_s) = \frac{\int p(f_*|x_*, \mathbf{y}_s, \boldsymbol{\theta})p(\mathbf{y}_s|\mathbf{x}_s, \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}}{\int p(\mathbf{y}_s|\mathbf{x}_s, \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}} \quad (8)$$

Unfortunately such integrals are generally non-analytic, but can be well approximated by use of Bayesian Monte Carlo [3] techniques. This involves evaluating predictions for a range of hyperparameter samples $\{\theta_i : i \in S\}$, with a different mean $m_{\theta_i}(f_*|x_*, \mathbf{y}_s)$ and covariance $C_{\theta_i}(f_*|x_*, \mathbf{y}_s)$ for each, which are then combined in a weighted mixture:

$$p(f_*|x_*, \mathbf{y}_s) \simeq \sum_{i \in S} \rho_i N(f_*; m_{\theta_i}(f_*|x_*, \mathbf{y}_s), C_{\theta_i}(f_*|x_*, \mathbf{y}_s)) \quad (9)$$

with weights ρ as detailed in [1]. This approach gives a close approximation to full marginalisation but, in order to take a meaningful number of samples, suffers from high computational costs.

A far less demanding approach is to choose a single θ which maximises the marginal likelihood, calculated by marginalising over the function values \mathbf{f} :

$$p(\mathbf{y}_s|\mathbf{x}_s, \boldsymbol{\theta}) = \int p(\mathbf{y}_s|\mathbf{f}, \mathbf{x}_s, \boldsymbol{\theta})p(\mathbf{f}|\mathbf{x}, \boldsymbol{\theta})d\mathbf{f} \quad (10)$$

The value of θ which maximises (10) can be inserted into (4) in order to make predictions. This approximation, known as type II maximum likelihood, tends to the full marginalised prediction as the number of observations increases. Unfortunately the approximation can suffer from overfitting, especially when faced with many hyperparameters and low numbers of samples when the true uncertainty in θ will be high. The computational burden is however far less, hence the ubiquity of the method.

++*+*+*+*+*+ TEXT NEEDED!!! - Insert comment about what our paper does as a gap between
the two methods just discussed *+*+*+*+*+*+*+

3 Approximate Hyperparameter Marginalisation

3.1 Proof of Positive Semi-Definiteness

A sum of kernels is itself a kernel, which by definition fulfils the necessary condition of positive semi-definiteness. Therefore:

$$K = \int k(\beta)p(\beta)d\beta \quad (11)$$

4 Experiments

5 Related Work

6 Conclusion

Acknowledgments

Do we have any? Aladdin / Orchid?

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

References

- [1] M.A. Osborne, SJ Roberts, A. Rogers, SD Ramchurn, and N.R. Jennings. Towards real-time information processing of sensor network data using computationally efficient multi-output gaussian processes. In *Proceedings of the 7th international conference on Information processing in sensor networks*, pages 109–120. IEEE Computer Society, 2008.
- [2] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press Cambridge, MA, USA, 2006.
- [3] C.E. Rasmussen and Z. Ghahramani. Bayesian monte carlo. *Advances in neural information processing systems*, 15:489–496, 2003.