
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 multi-variate 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 : 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:

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

$$K_\beta = h^2 \exp\left(-\frac{1}{2} \Delta^2 \exp(-2\beta)\right) \quad (16)$$

$$K'_\beta = K_\beta (\Delta^2 \exp(-2\beta)) \quad (17)$$

$$K''_\beta = K_\beta (\Delta^4 \exp(-4\beta) - \Delta^2 \exp(-2\beta)) \quad (18)$$

$$A = \ln \frac{1}{K_\beta} \quad (19)$$

$$\frac{\partial A}{\partial \beta} = -\frac{K'_\beta}{K_\beta} = \Delta^2 \exp(-2\beta) \quad (20)$$

$$\frac{\partial^2 A}{\partial^2 \beta} = -\frac{K''_\beta}{K_\beta} + \frac{K'^2_\beta}{K^2_\beta} = 2 \Delta^2 \exp(-2\beta) \quad (21)$$

$$C = \Delta^2 \exp(-2\beta_m)$$

$$K_\beta \approx \exp\left(-\left(\ln \frac{1}{K_{\beta_m}} + (\beta - \beta_m)C + (\beta - \beta_m)^2 C\right)\right) \quad (22)$$

$$= K_{\beta_m} \exp\left((\beta - \beta_m)C + (\beta - \beta_m)^2 C\right) \quad (23)$$

$$K_{\beta_m, \beta_v} = \int_{-\infty}^{+\infty} K_\beta p(\beta | \beta_m, \beta_v) d\beta \quad (24)$$

$$= K_{\beta_m} \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi\beta_v}} \exp\left(-\frac{(\beta - \beta_m)^2}{2\beta_v}\right) ((\beta - \beta_m)C + (\beta - \beta_m)^2 C) \quad (25)$$

$$= K_{\beta_m} \exp\left(\frac{\beta_v C^2}{2(1 + 2\beta_v C)}\right) \frac{1}{\sqrt{1 + 2\beta_v C}} \quad (26)$$

3.1 Proof of Positive Semi-Definiteness

++*+*+ I'll tidy this up!

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

$$K_{\beta_m, \beta_v} = \int K_\beta p(\beta) d\beta \quad (27)$$

should be a legitimate kernel if K_β is also a legitimate kernel as $p(\beta)$ just weights the contents of the integral.

$$K_\beta = K_{\beta_m} \exp\left((\beta - \beta_m)C + (\beta - \beta_m)^2 C\right) \quad (28)$$

completing the square:

$$K_\beta = K_{\beta_m} \exp\left(-\frac{1}{2}(\beta - \beta_m - 1)^2 C\right) \exp\left(\frac{1}{2}C\right) \quad (29)$$

$$(30)$$

Product of kernels is also a kernel. Therefore if all three parts of the above equation are kernel, then K_β is also a covariance function. First two parts are kernels, the last part isn't. However:

$$K_{\beta_m} = h^2 \exp \left(-\frac{1}{2} \Delta^2 \exp(-2\beta_m) \right) \quad (31)$$

$$= h^2 \exp \left(-\frac{C}{2} \right) \quad (32)$$

$$K_{\beta_m} \exp \left(\frac{1}{2} C \right) = h^2 \exp \left(-\frac{C}{2} \right) \exp \left(\frac{1}{2} C \right) \quad (33)$$

$$= h^2 \exp(0) \quad (34)$$

so K_β is a kernel.

++*+*+ I need to check this, as the result has changed since I did my substitution...

4 Experiments

5 Related Work

6 Conclusion

Acknowledgments

Do we have any? Aladdin / Orchid?

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.