
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 [1]. 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. A vector of observations $\mathbf{y} = \{y_1, \dots, y_n\}$ could be viewed as a single point sampled from a n -variate Gaussian distribution.

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 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 $f : \mathcal{X} \rightarrow \mathbb{R}$ about which we wish to perform inference 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.

Once we have observations of the function $(\mathbf{x}_s, \mathbf{y})$ 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:

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

where the posterior mean and covariance are:

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

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

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

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

056
057
058
059

060
061
062

065

066
067
068

069
070
071

072

073

074

075

076
077

078
079
080
081
082

083
084

087

088
089

091

092
093
094

096
097

098

100
101
102

105

106

107

$$A = \ln \frac{1}{K_\beta} \quad (20)$$

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

$$\frac{\partial^2 A}{\partial^2 \beta} = -\frac{K''_\beta}{K_\beta} + \frac{K'^2_\beta}{K_\beta^2} = 2 \Delta^2 \exp(-2\beta) \quad (22)$$

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

$$K_\beta \approx \exp \left(- \left(\ln \frac{1}{K_\nu} + (\beta - \nu)C + (\beta - \nu)^2 C \right) \right) \quad (23)$$

$$= K_\nu \exp \left(((\beta - \nu)C + (\beta - \nu)^2 C) \right) \quad (24)$$

$$K_{\nu, \Lambda} = \int_{-\infty}^{+\infty} K_\beta p(\beta | \nu, \Lambda) d\beta \quad (25)$$

$$= K_\nu \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi\Lambda}} \exp \left(\frac{-(\beta - \nu)^2}{2\Lambda} \right) ((\beta - \nu)C + (\beta - \nu)^2 C) \quad (26)$$

$$= K_\nu \exp \left(\frac{\Lambda C^2}{2(1 + 2\Lambda C)} \right) \frac{1}{\sqrt{1 + 2\Lambda C}} \quad (27)$$

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:

$$K_{\nu, \Lambda} = \int K_\beta p(\beta) d\beta \quad (28)$$

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

$$K_\beta = K_\nu \exp \left(((\beta - \nu)C + (\beta - \nu)^2 C) \right) \quad (29)$$

completing the square:

$$K_\beta = K_\nu \exp \left(-\frac{1}{2}(\beta - \nu - 1)^2 C \right) \exp \left(\frac{1}{2}C \right) \quad (30)$$

$$(31)$$

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:

$$K_\nu = h^2 \exp \left(-\frac{1}{2} \Delta^2 \exp(-2\nu) \right) \quad (32)$$

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

$$K_\nu \exp \left(\frac{1}{2}C \right) = h^2 \exp \left(-\frac{C}{2} \right) \exp \left(\frac{1}{2}C \right) \quad (34)$$

$$= h^2 \exp(0) \quad (35)$$

so K_β is a kernel.

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

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

4 Experiments

5 Related Work

6 Conclusion

Acknowledgments

Do we have any? Aladdin / Orchid?

References

- [1] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press Cambridge, MA, USA, 2006.
- [2] C.E. Rasmussen and Z. Ghahramani. Bayesian monte carlo. *Advances in neural information processing systems*, 15:489–496, 2003.