
Approximate Hyperparameter Marginalisation for Gaussian Processes

Anonymous Author(s)

Affiliation

Address

email

Abstract

Kernel and covariance functions are often parameterised by hyperparameters, unknown *a-priori*. We propose a novel means of approximately integrating over these hyperparameters. This approach gives rise to a novel class of covariance functions, Marginal Generalised Exponentials, which are robust to hyperparameter uncertainty, avoiding the over-fitting that often accompanies type-II maximum-likelihood fitting. We demonstrate the efficacy of our covariances for Gaussian process regression on both synthetic and real data.

1 Introduction

Kernel and covariance functions are vital to much of machine learning and data mining [9]. In this paper, we consider Gaussian process [7] inference, which has at its heart just such a covariance function. The hyperparameters that define such functions are rarely known *a-priori*. Given training data, these hyperparameters are often set by type-II maximum likelihood, which can lead to overfitting. A more rigorous alternative is to integrate over, or marginalise, the hyperparameters, which is unfortunately usually intractable. While sampling approaches [4] have been developed to address this problem, they are prohibitively computationally expensive for large numbers of hyperparameters.

We introduce a Laplace approach to approximately marginalise covariance hyperparameters. This gives us a robust family of covariances that are tolerant to uncertainty in the hyperparameters. Our method is particularly appropriate when such uncertainty is likely to be significant, as when only small amounts of data are available. We are motivated by scenarios where obtaining data is expensive (computationally or otherwise), such as in global optimisation [5].

2 Gaussian Processes

Gaussian processes (GPs) constitute a powerful method for performing Bayesian inference about functions using a limited set of observations [7]. 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 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, or kernel, $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ 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:

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

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

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.

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. 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:

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

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

$$p(f_*|x_*, \mathbf{y}, \boldsymbol{\theta}) := \mathcal{N}(f_*; m_{\boldsymbol{\theta}}(f_*|x_*, \mathbf{y}), C_{\boldsymbol{\theta}}(f_*|x_*, \mathbf{y})) \quad (3)$$

where the posterior mean and covariance are:

$$m_{\boldsymbol{\theta}}(f_*|x_*, \mathbf{y}) := \mu_{\boldsymbol{\theta}}(x_*) + K_{\boldsymbol{\theta}}(x_*, \mathbf{x})V^{-1}(\mathbf{y} - \mu_{\boldsymbol{\theta}}(\mathbf{x}_*)) \quad (4)$$

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

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

2.1 Marginalising out the Hyperparameters

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

$$p(f_*|x_*, \mathbf{y}) = \frac{\int p(f_*|x_*, \mathbf{y}, \boldsymbol{\theta})p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}}{\int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}} \quad (6)$$

Unfortunately, such integrals are generally non-analytic, requiring numerical approximation. Randomized Monte Carlo techniques [1] form the most popular approaches to numerical integration, although Bayesian alternatives [8, 6] also exist. These techniques estimate the integral given the value of the integrand on a set of sample points, usually via a weighted mixture

$$p(f_*|x_*, \mathbf{y}) \simeq \sum_i \rho_i p(f_*|x_*, \mathbf{y}, \boldsymbol{\theta}_i) \quad (7)$$

for some weight vector $\boldsymbol{\rho}$. Unfortunately, the computational expense of evaluating the integrand at sufficient samples to estimate high-dimensional integrals is often prohibitive.

A less computationally demanding alternative is to select only a single sample. Type II maximum likelihood, or maximum marginal likelihood, approximates as

$$p(f_*|x_*, \mathbf{y}) \simeq p(f_*|x_*, \mathbf{y}, \boldsymbol{\theta}_{\text{ML}}) \quad (8)$$

where $\boldsymbol{\theta}_{\text{ML}}$ is the hyperparameter vector that maximises the marginal likelihood,

$$\boldsymbol{\theta}_{\text{ML}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \quad (9)$$

As per Figure 1a, (8) is equivalent to approximating the likelihood $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$ as the delta function $\delta(\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{ML}})$. This assumption is a poor representation of our true state of ignorance given only low numbers of data, for which the likelihood is typically broad and/or multi-modal. As a consequence, type II maximum likelihood can often lead to over-fitting. Nonetheless, the less onerous computational requirements of this approach have made it ubiquitous throughout machine learning.

We now present a computationally inexpensive approach to approximately marginalising hyperparameters.

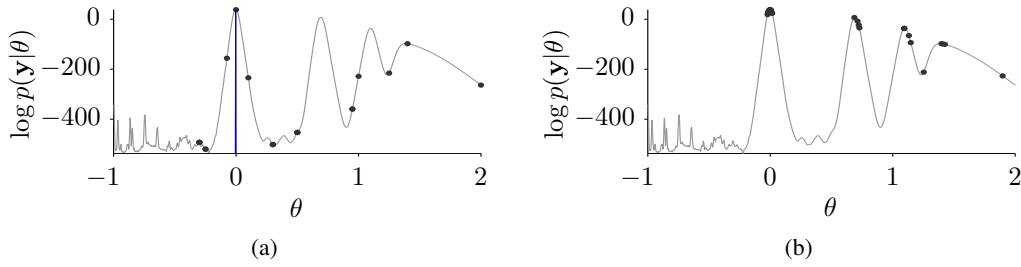


Figure 1: a) Samples (black dots) obtained by optimising the log-likelihood (grey) using a global optimiser, and in blue, the maximum likelihood approximation of the likelihood surface. b) Samples obtained by taking draws from the posterior using an MCMC method.

3 Approximate Hyperparameter Marginalisation

We now propose a general-purpose family of covariance functions that result from approximately marginalising the hyperparameters of an homogenous exponential covariance.

3.1 Generalised Exponential Covariance Functions

We consider the general form of covariance, which we term the *homogenous exponential*

$$K_{\theta}(x, x') = h^2 \exp \left(-\frac{1}{2} \frac{\Delta^2}{L^2} \right), \quad (10)$$

where h is a scaling parameter termed the *output scale*, and L is the *length scale*; the vector θ contains both hyperparameters. This form captures many of the most commonly used covariance functions. If $\Delta^2 = \sum_i (x_i - x'_i)^2$, (10) becomes the ubiquitous *squared exponential*. If $\Delta^2 = |x - x'|$, (10) is the exponential covariance. Many other covariances can be constructed by taking Δ^2 as other metrics, allowing for periodic covariances, warpings of the input space [10] and inference in more structured domains [2].

The output scale h can be tractably marginalised according to suitable priors [3]. Even if performing type-II maximum likelihood, the output scale is usually well-specified by even a moderate amount of data. We will not consider this hyperparameter further.

The length scale L affects how closely outputs co-vary as a function of Δ^2 . Our predictions will ultimately be very sensitive to the value of the length scale; a short length scale will favour rough functions, a long length scale will favour smoother functions. Unfortunately, the likelihood of the length scale is often highly multi-modal: Figure 1 illustrates the likelihood surface as a function of $\theta = L$ for a periodic covariance, where the multiple modes correspond to harmonics of the latent period and sampling frequency. That is, given limited data, the posterior for the input scale is likely to have large variance. As such, type-II maximum likelihood for this hyperparameter will lead to problematic over-fitting.

3.2 Marginal Generalised Exponential Covariance Function

Mike says: now we talk about our covariance, which I've called Marginal Generalised Exponential covariances

We propose to marginalise L (which is strictly positive) according to a Gaussian prior to its logarithm, $\beta := \log(L)$;

and define a Gaussian distribution over β so that $p(L)$ will have strictly positive support:

$$K_{\theta}(x, x') = h^2 \exp \left(-\frac{1}{2} \Delta^2 \exp(-2\beta) \right) \quad (11)$$

which we will write as K_β , and

$$p(\beta|\nu, \Lambda) = \mathcal{N}(\beta; \nu, \Lambda) = \frac{1}{\sqrt{2\pi\Lambda}} \exp\left(-\frac{(\beta - \nu)^2}{2\Lambda}\right) \quad (12)$$

where ν and Λ are hyper-hyperparameters.

As stated before, we marginalise out the uncertainty in β (or L), before fusing the prior with observations.

Jan says: Perhaps we should talk about the moment matching idea first.. As far as I can remember: Ie Taylor expand the marginalized process ... it then turns out the the kernel of the matched process equals the marginalized kernel you write down in (14) below..

The marginalised kernel is calculated by evaluating the following integral:

$$K_{\nu, \Lambda} = \int_{-\infty}^{+\infty} K_\beta p(\beta|\nu, \Lambda) d\beta \quad (13)$$

which is non-analytic given the current form of K_β . We approximate the solution to the integral by first approximating K_β using a second-order Taylor expansion around the mean of $p(\beta)$ given by ν . In order to ensure that the approximate kernel is positive semi-definite we force K_β to be in exponential form, hence we note the following identity and substitution:

$$K_\beta = \exp\left(-\ln \frac{1}{K_\beta}\right) = \exp(-A) \quad (14)$$

and we Taylor expand around A :

$$A \approx \ln \frac{1}{K_\beta} \Big|_\nu + (\beta - \nu) \frac{\partial A}{\partial \beta} \Big|_\nu + \frac{1}{2} (\beta - \nu)^2 \frac{\partial^2 A}{\partial^2 \beta} \Big|_\nu \quad (15)$$

Given (11), the derivatives of K_β can be calculated:

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

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

which can be used for calculating the partial derivatives of (15):

$$A = \ln \frac{1}{K_\beta} \quad (16)$$

$$\frac{\partial A}{\partial \beta} = -\frac{K'_\beta}{K_\beta} = -\Delta^2 \exp(-2\beta) \quad (17)$$

$$\frac{\partial^2 A}{\partial^2 \beta} = -\frac{K''_\beta}{K_\beta} + \frac{K'^2_\beta}{K^2_\beta} = 2\Delta^2 \exp(-2\beta) \quad (18)$$

Using (14) and substituting (16-18) into (15) we arrive at the following expression for K_β :

$$K_\beta \approx \exp\left(-\left(\ln \frac{1}{K_\nu} - (\beta - \nu)C + (\beta - \nu)^2 C\right)\right) = K_\nu \exp((\beta - \nu)C - (\beta - \nu)^2 C) \quad (19)$$

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

All that remains is to evaluate the integral in (13), which by using (12) and (19) is as follows:

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

$$= K_\nu \exp\left(\frac{\Lambda C^2}{2(1 + 2\Lambda C)}\right) \frac{1}{\sqrt{1 + 2\Lambda C}} \quad (20)$$

This result is a new kernel which is the product of the original kernel K_β evaluated at $\beta = \nu$, and an expression which is a function of the squared distance Δ^2 .

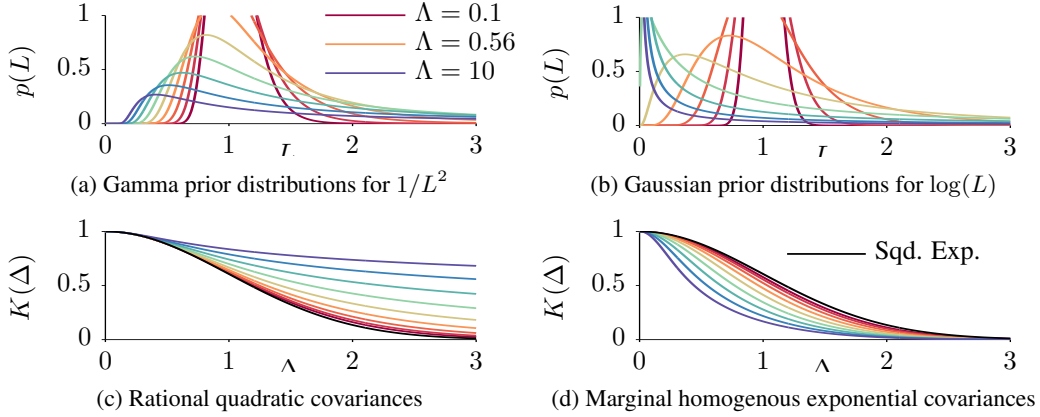


Figure 2: The different priors assumed by RQ and MHE covariances give rise to different behaviours as the variance in the input scale (Λ) increases. The priors of a assign more weight to very large input scales, whereas the priors of b assign more weight to small input scales. All figures have $\nu = 1$ and share the legends of a) and d).

3.3 Rational Quadratic Covariance Function

The Rational Quadratic covariance function represents one means of approximately marginalising the input scale of a squared exponential covariance. It assumes a gamma prior, which, problematically, puts excessive weight on very large input scales when we are very uncertain. See Figure 2.

the log-Gaussian prior dominates the inverse gamma prior

(21)

3.3.1 Proof of Positive Semi-Definiteness

In order to confirm that our new kernel $K_{\nu, \Lambda}$ is a legitimate covariance function, we must prove that it fulfils the necessary condition of kernels which is positive semi-definiteness. We first note that a sum of kernels is also a kernel, therefore:

$$K_{\nu, \Lambda} = \int K_{\beta} p(\beta) d\beta \quad (22)$$

should be a legitimate kernel if K_{β} is also a legitimate kernel as $p(\beta)$ just weights the contents of the integral. Using the expression for K_{β} in (19) and completing the square we arrive at the following expression:

$$K_{\beta} = K_{\nu} \exp \left(-\frac{1}{2} \left(\beta - \nu - \frac{1}{2} \right)^2 C \right) \exp \left(\frac{1}{4} C \right) \quad (23)$$

If we substitute in K_{ν} , where:

$$\begin{aligned} K_{\nu} &= h^2 \exp \left(-\frac{1}{2} \Delta^2 \exp(-2\nu) \right) \\ &= h^2 \exp \left(-\frac{C}{2} \right) \end{aligned} \quad (24)$$

we can write:

$$K_{\beta} = h^2 \exp \left(-\frac{C}{4} \right) \exp \left(-\frac{1}{2} \left(\beta - \nu - \frac{1}{2} \right)^2 C \right) \quad (25)$$

We note that the product of kernels is also a kernel. Both exponentials in (25) form legitimate kernels, which leaves us with the result that K_{β} is also a kernel and by definition is positive semi-definite as required.

4 Experiments

4.1 Posterior functions with differing length scales

When using standard kernels, functions drawn from the prior will have a fixed length scale; when using our kernel, the length scale of functions drawn from the prior can take a number a range of values as a result of the distribution over length scales. When our prior is fused with observations to form a posterior GP, we are able to describe functions that may have generated the observations over a range of length scales. When data is limited this is particularly crucial as a whole range of length scales may be permissible with regards to the data generated, even if one is favoured through type-II maximum likelihood methods. As data increases the range of permissible length scales reduces, eventually collapsing to the true length scale.

We demonstrate the above behaviour and its advantage over other methods in this section. We begin by drawing a function from a prior GP with known length scale. From this we can generate a training set of between X and X observations in the domain $\mathcal{X} \in [-3, 3]$. We train the hyperparameters or hyper-hyperparameters of our test methods using type-II maximum likelihood on the given training set. Now we go back to the original GP and draw a set of X functions from the posterior (so that the functions still go through the original training data) with length scales ranging from X to X . An example training set and resulting set of posterior functions can be seen in figure XXXX. A test set with 100 observations is generated from each function and a log likelihood is calculated for each method. This process is repeated with 100 different starting functions and training data.

We plot the mean and variance trials of the likelihood of the test data across all $X * 100$ against the log length scale of the posterior function.

4.2 Real world dataset

5 Conclusion

References

- [1] M.H. Chen, Q.M. Shao, and J.G. Ibrahim. *Monte Carlo methods in Bayesian computation*. Springer Verlag, 2000.
- [2] R. Garnett, MA Osborne, and SJ Roberts. Bayesian optimization for sensor set selection. In *Proceedings of the 9th ACM/IEEE International Conference on Information Processing in Sensor Networks*, pages 209–219. ACM, 2010.
- [3] M. Kennedy. Bayesian quadrature with non-normal approximating functions. *Statistics and Computing*, 8(4):365–375, 1998.
- [4] R.M. Neal. Monte carlo implementation of gaussian process models for bayesian regression and classification. *Arxiv preprint physics/9701026*, 1997.
- [5] M.A. Osborne, R. Garnett, and S.J. Roberts. Gaussian processes for global optimization. In *3rd International Conference on Learning and Intelligent Optimization (LION3)*, pages 1–15, 2009.
- [6] M.A. Osborne, R. Garnett, S.J. Roberts, C. Hart, S. Aigrain, N.P. Gibson, and S. Aigrain. Bayesian quadrature for ratios. *Fifteenth International Conference on Artificial Intelligence and Statistics*, 2012.
- [7] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press Cambridge, MA, USA, 2006.
- [8] C.E. Rasmussen and Z. Ghahramani. Bayesian monte carlo. *Advances in neural information processing systems*, 15:489–496, 2003.
- [9] J. Shawe-Taylor and N. Cristianini. *Kernel methods for pattern analysis*. Cambridge Univ Pr, 2004.
- [10] E. Snelson, C.E. Rasmussen, and Z. Ghahramani. Warped gaussian processes. *Advances in neural information processing systems*, 16:337–344, 2004.