

# Bayesian optimization in machine learning

José Jiménez Luna

January 12, 2017



# Contents

<b>1</b>	<b>Organization of this work</b>	<b>5</b>
1.1	Introduction . . . . .	5
1.2	How do I read this? . . . . .	6
<b>2</b>	<b>Gaussian Process regression</b>	<b>7</b>
2.1	A weight space view for Gaussian Processes . . . . .	7
2.1.1	Standard Bayesian linear regression . . . . .	7
2.1.2	Kernel functions in feature space . . . . .	8
2.2	A function space view for Gaussian Processes . . . . .	8
2.3	Prediction using a Gaussian Process prior . . . . .	10
2.3.1	A toy example of Gaussian Process regression . . . . .	12
2.4	On covariance functions . . . . .	13
2.4.1	A zoo of covariance functions . . . . .	16
2.5	Hyperparameter optimization . . . . .	18
2.5.1	Type II Maximum Likelihood . . . . .	18
2.5.2	Cross validation . . . . .	20
<b>3</b>	<b>Bayesian optimization</b>	<b>23</b>
<b>4</b>	<b>Experiments</b>	<b>25</b>
<b>5</b>	<b>pyGPGO: A simple Python Package for Bayesian Optimization</b>	<b>27</b>



# Chapter 1

## Organization of this work

### 1.1 Introduction

This Master's Thesis provides an introduction to both Gaussian Processes and Bayesian Optimization. This work aims to be a multi-objective optimization task:

- The first objective of the thesis is to provide the reader with an introduction to Gaussian Process regression and Bayesian optimization. The work is written in such a way that it alternates very often between theoretical and practical (in terms of programming) background. This work is written like a programming book or a manual.
- To show the Bayesian Optimization framework works in several real-world machine learning tasks. I do so by selecting several datasets (most of them from biological/physical phenomena), applying such methodology and comparing to other common strategies applied for the same problem.
- Finally, and my personal favourite, to provide the user with usable and easy to use software to apply Bayesian Optimization in their research. This comes in the form of a Python (>3.5) package named pyGPGO. The code can be freely obtained in <https://github.com/hawk31/pyGPGO>. The software is MIT licensed. All the examples and code snippets throughout this manual are based on this software.

We begin by describing the title of this thesis. Bayesian Optimization focuses on the global optimization of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  over a compact set  $A$ . The problem can be formalized as:

$$\max_{\mathbf{x} \in A} f(\mathbf{x}) \quad (1.1)$$

Most optimization procedures (local based ones such as gradient ascent, for example) assume that the function  $f$  is closed-form, that is, can be written in a paper, that it is convex, with known first or second order derivatives or cheap to evaluate. Bayesian optimization focuses on all these problems proposing a very elegant solution. By the use of a surrogate model, a Gaussian Process, a Bayesian optimization procedure can help to find the global minimum of a non-necessarily convex, expensive functions. These methods shine also where there is no closed-form expression to evaluate and does not need any function derivatives.

Now is when the machine learning part of the title comes into play. In machine learning, we are usually interested in minimizing (or maximizing the opposite) a loss function  $\mathcal{L}$ . These losses can take many forms, for example, when doing regression, a typical loss might be the mean squared error between predictions and observed values on a holdout test set.

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2 \quad (1.2)$$

In binary classification, for example, a very popular choice is the logarithmic loss:

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{n} \sum_i (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)) \quad (1.3)$$

Notice that in any case, that these losses are typically defined in a subset of  $\mathbb{R}$ . We focus on the supervised setting of machine learning, and more specifically regression. Depending on the problem at hand, even evaluating these losses can be very expensive from a computational point of view. This may have to do with the machine learning algorithm used or the size of the dataset at hand. These machine learning algorithms typically have

*hyperparameters*, that is, parameters that have to be tuned in a sensible way to get the best performance possible out of these models. In the machine learning community it is common for practitioners to do an hyperparameter grid search or even randomize it. Since the training of a single model can already take substantial resources in terms of CPU cycles or memory, we would like to have a more efficient and cheap way to optimize these hyperparameters. Bayesian optimization will let us do that by proposing the next candidate point  $\mathbf{x}$  to evaluate according to several criteria.

## 1.2 How do I read this?

Before we dive directly into the topic at hand, it is mandatory to explain how this manual is intended to be read. Different chapters will cover different material, so if the reader is familiar with a topic in a chapter, for example, the chapter on Gaussian Processes, he can skip directly to the chapter on Bayesian Optimization or to the implementation if he wishes. There is no complicated or unnecessarily advanced material in any of the chapters, as this manual does not aim to be bleeding edge research material.

**Chapter 2** focuses on a swift but thorough introduction to regression problems using Gaussian Processes. These are the surrogate models we will use for Bayesian Optimization in Chapter 3. We will cover the theory behind them both from a weight space view and from a functional view. We will also explain different covariance functions and their role in these models, as well as methods for optimizing their hyperparameters. Finally we will provide the reader with usable code to fit a Gaussian Process in a regression problem. This chapter is heavily based on Carl E. Rasmussen excellent book *Gaussian Processes for Machine Learning*[ref], an more specially in chapters 2, 4 and 5. In fact, if the reader is interested in a more advanced, wider, slowly-paced introduction to Gaussian Processes, this is by far the best resource to go to.

**Chapter 3** is about the main topic in this work, Bayesian Optimization. Once we have laid down all the foundations of Gaussian Processes, we can start explaining the theory of Bayesian optimization using these as surrogate models. The algorithms, while simple, are very powerful. The role of several acquisition functions, that is to say, the functions that will propose the next point to evaluate will be discussed, as well as their advantages or disadvantages. The references on this chapter will be more diverse than on the previous chapter, as I will try to summarize several publications. Readers already familiar with Gaussian Processes should jump directly into this chapter.

**Chapter 4** covers experiments using the software provided alongside this manual. These are mostly mid-sized regression or classification problems where we will compare the performance of Bayesian Optimization of hyperparameters with several regressors/classifiers with other strategies, such as random search or simulated annealing. Most of these datasets are related to experimental sciences, and some of them are typically used for other benchmarking purposes in other studies.

**Chapter 5** has no theoretical content and is the shortest one. It will cover technical explanations of pyGPGO, the software developed alongside this manual. I will also try to provide practical examples on how to use the software.

## Chapter 2

# Gaussian Process regression

We will focus on regression problems in this chapter. Assume we have some labelled data  $D = \{(\mathbf{x}_i, y_i) \mid i = 1, \dots, n\}$ , where  $\mathbf{x}$  is a vector of covariates and  $y$  denotes a continuous objective variable. We wish to learn a predictive distribution over new values of  $\mathbf{x}$ , so that we can make predictions and inference over these. In practice, for simplicity we write that  $D = (X, \mathbf{y})$ , where  $X$  is a covariate matrix.

One can interpret a Gaussian Process in several ways, the most used one is the function space view, which is the one we will cover second here and the one we will assume for the rest for the manual. In this view, we consider a Gaussian Process to be a distribution over functions, instead of over values. Inference takes place directly in this space. For completeness, we will also provide a weight-space view first, that might be more appealing to readers familiar with Bayesian linear regression.

## 2.1 A weight space view for Gaussian Processes

In this section we will try to draw connections between Bayesian linear regression and our introduction to Gaussian Processes, by the use of kernel functions.

### 2.1.1 Standard Bayesian linear regression

A Bayesian linear regression model with Gaussian error can be formulated as:

$$y = f(\mathbf{x}) + \epsilon \quad (2.1)$$

where we typically assume  $\epsilon \sim N(0, \sigma_n^2)$ . This noise assumption directly implies a Gaussian likelihood, thus it can be easily proven that:

$$p(\mathbf{y} \mid X, \mathbf{w}) \sim N(X^T \mathbf{w}, \sigma_n^2 I) \quad (2.2)$$

Assume now a Gaussian prior on the weights  $\mathbf{w}$ :

$$\mathbf{w} \sim N(\mathbf{0}, \Sigma_p) \quad (2.3)$$

We are interested now on the posterior distribution of  $\mathbf{w}$ , given both  $X$  and  $\mathbf{y}$ , and assuming the model in Equation 2.1, that is:

$$p(\mathbf{w} \mid \mathbf{y}, X) = \frac{p(\mathbf{y} \mid X, \mathbf{w}) p(\mathbf{w})}{p(\mathbf{y} \mid X)} \quad (2.4)$$

One can solve this problem by means of sampling procedures like Markov Chain Monte Carlo, but in this particular case, there is a closed-form solution. It can be proven quite easily that:

$$p(\mathbf{w} \mid X, \mathbf{y}) \sim N\left(\frac{1}{\sigma_n^2} A^{-1} X \mathbf{y}, A^{-1}\right) \quad (2.5)$$

where  $A = \sigma^{-2} X X^T + \Sigma_p^{-1}$ . Notice that a very easy MAP (maximum a posteriori) estimate of the weights can be obtained by just computing the mean of this distribution. Now, to make predictions for a particular test case  $\mathbf{x}_*$ , we average over all possible parameter values, hence we get a whole predictive distribution. Again, it can be proven that:

$$f_*|x_*, X, \mathbf{y} \sim N\left(\frac{1}{\sigma_n^2} \mathbf{x}_*^T A^{-1} X \mathbf{y}, \mathbf{x}_*^T A^{-1} \mathbf{x}_*\right) \quad (2.6)$$

### 2.1.2 Kernel functions in feature space

We have presented a very simple Bayesian approach to linear regression in the previous section. While useful, it lacks expressiveness. A very simple idea is to project this data into a higher dimension, where it may be more easily separated by a linear model of this sort. This is called using the kernel trick. We can do this by the use of a covariance (or kernel) function  $\phi(\mathbf{x})$ . Note by  $\Phi(X)$  the aggregation of columns after computing this kernel function in the entire dataset at hand.

The model becomes now:

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w} \quad (2.7)$$

All the math presented in the previous section applies here, just placing  $\phi(\mathbf{x})$  instead of  $\mathbf{x}$ . The predictive distribution over  $y$  becomes now, for example:

$$f_*|x_*, X, \mathbf{y} \sim N\left(\frac{1}{\sigma_n^2} \phi(\mathbf{x}_*)^T A^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_*)^T A^{-1} \phi(\mathbf{x}_*)\right) \quad (2.8)$$

where for simplicity we have written  $\Phi = \Phi(X)$  and  $A = \sigma_n^{-2} \Phi \Phi^T + \Sigma_p^{-1}$ . The predictive distribution needs to invert  $N \times N$  matrix. Equation 2.8 can be rewritten as:

$$f_*|x_*, X, \mathbf{y} \sim N\left(\phi_*^T \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \mathbf{y}, \phi_*^T \Sigma_p \phi_* - \phi_*^T \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \Phi^T \Sigma_p \phi_*\right) \quad (2.9)$$

where we have again simplified notation by  $\phi_* = \phi(\mathbf{x}_*)$  and  $K = \Phi^T \Sigma_p \Phi$ . Now notice that the entries of  $K$  for both train and test set are of the form  $\phi(\mathbf{x}_*^T) \Sigma_p \phi(\mathbf{x}_*)$ . We have implicitly defined now a *covariance function* of the form  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x}_*^T) \Sigma_p \phi(\mathbf{x}_*)$ . This is in fact an inner product with respect to  $\Sigma_p$ . That is if we define  $\psi(\mathbf{x}) = \Sigma_p^{1/2}(\mathbf{x})$ , then a simple dot product representation of a covariance function is:

$$k(\mathbf{x}, \mathbf{x}') = \psi(\mathbf{x})^T \psi(\mathbf{x}') \quad (2.10)$$

where  $\Sigma_p^{1/2}$  can be defined by means of a singular value decomposition. Typically, we will replace the original feature vectors by these dot products, *lifting* to a higher space. This will make more sense in the following section.

## 2.2 A function space view for Gaussian Processes

We formally define a Gaussian Process as a collection of random variables, any finite number of which have a joint Gaussian distribution. This process is totally defined by two functions. The first one is its *mean function*:

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})] \quad (2.11)$$

The second one is its *covariance function*:

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))] \quad (2.12)$$

In practice, we say that  $f$  is a Gaussian Process with mean  $m(\mathbf{x})$  and covariance function  $k(\mathbf{x}, \mathbf{x}')$  and write:

$$f(\mathbf{x}) \sim \text{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (2.13)$$

In practice, for simplicity we will take  $m(\mathbf{x}) = 0$ , but this can be specified otherwise. As stated before, a Gaussian Process fulfils the marginalization property, that is to say that if the GP specifies  $(y_1, y_2) \sim N(\boldsymbol{\mu}, \Sigma)$  then this implies that  $y_1 \sim N(\mu_1, \Sigma_{11})$ . A Gaussian multivariate distribution is just a finite index set of a Gaussian process.

As seen in the previous section, a Gaussian Process can be viewed as a Bayesian regression model using a particular kernel, that is, for the model  $f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}$  and with the same prior as in equation 2.3 has mean and covariance functions:

$$\mathbb{E}[f(\mathbf{x})] = \phi(\mathbf{x})^T \mathbb{E}[\mathbf{w}] = 0 \quad (2.14)$$

$$\mathbb{E}[f(\mathbf{x})f(\mathbf{x}')] = \phi(\mathbf{x})^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi(\mathbf{x}') = \phi(\mathbf{x})^T \Sigma_p \phi(\mathbf{x}') \quad (2.15)$$



It is now a good time to start specifying our first covariance function, the *squared exponential* kernel, defined as:

$$k(x, x') = \exp\left(-\frac{1}{2}|x - x'|^2\right) \quad (2.16)$$

Where  $|\cdot|$  denotes the standard  $L_2$  norm. Most of the covariance functions that we will see here are a function of this norm, therefore it is much more comfortable to write  $r = |x - x'|$  and therefore the squared exponential kernel becomes:

$$k(x, x') = \exp\left(-\frac{1}{2}r^2\right) \quad (2.17)$$

The squared exponential covariance kernel is equivalent to a Bayesian linear model with an infinite number of basis functions. This also implies a distribution over functions. To see this, choose an arbitrary number of points  $X_*$  and compute the squared exponential kernel for those. Then just sample from the following multivariate Gaussian:

$$\mathbf{f}_* \sim N(\mathbf{0}, K(X_*, X_*)) \quad (2.18)$$

To illustrate this point, we will write a very simple Python script to draw (finite) samples from this function. For the moment, consider evenly spaced samples with a step of  $\frac{\pi}{16}$ . Some of the code here already uses the pyGPGO implementation of the covariance function, to simplify computation. This code produces Figure 2.1.

---

```

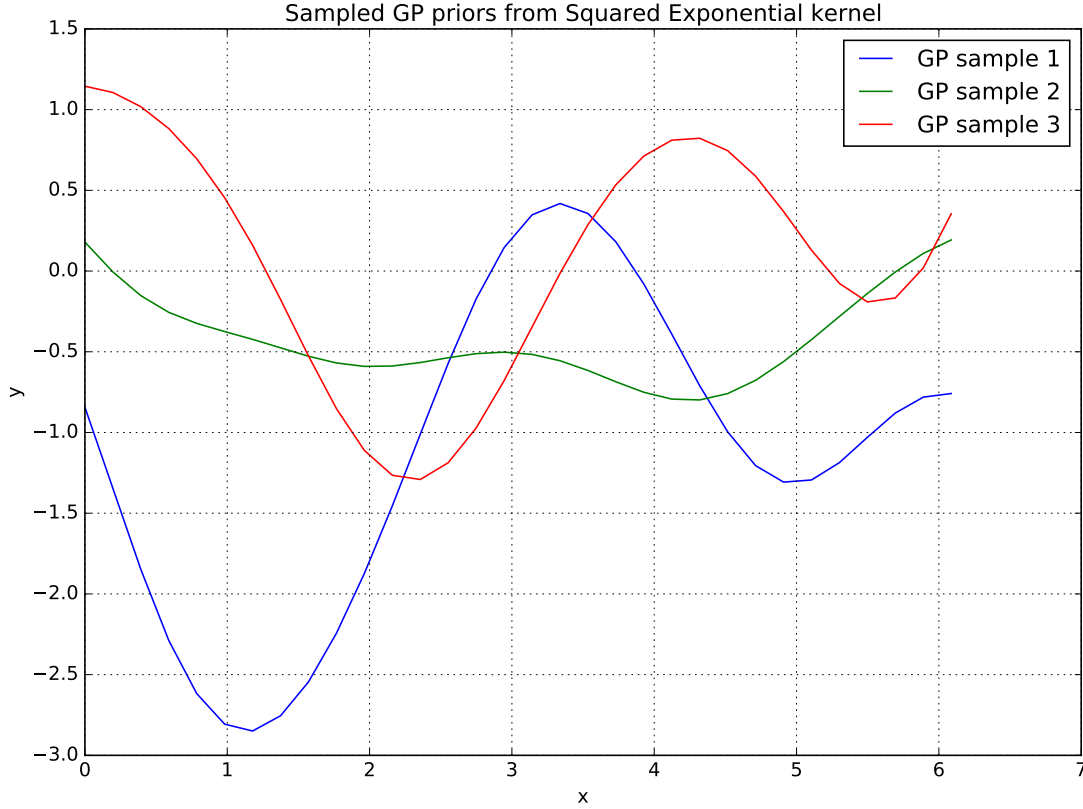
1 import numpy as np
2 from numpy.random import multivariate_normal
3 from covfunc import squaredExponential
4 import matplotlib.pyplot as plt
5
6 if __name__ == '__main__':
7     np.random.seed(93)
8     # Equally spaced values of Xstar
9     Xstar = np.arange(0, 2 * np.pi, step = np.pi/16)
10    Xstar = np.array([np.atleast_2d(x) for x in Xstar])[:, 0]
11    sexp = squaredExponential()
12    # By default assume mean 0
13    m = np.zeros(Xstar.shape[0])
14    # Compute squared-exponential matrix
15    K = sexp.K(Xstar, Xstar)
16
17    n_samples = 3
18    # Draw samples from multivariate normal
19    samples = multivariate_normal(m, K, size = n_samples)
20
21    # Plot values
22    x = Xstar.flatten()
23    plt.figure()
24    for i in range(n_samples):
25        plt.plot(x, samples[i], label = 'GP sample {}'.format(i + 1))
26    plt.xlabel('x')
27    plt.ylabel('y')
28    plt.title('Sampled GP priors from Squared Exponential kernel')
29    plt.grid()
30    plt.legend(loc = 0)
31    plt.show()

```

---

There's one important concept to explain before we move on. Notice in Figure 2.1 that the drawn functions seem to have a characteristic length-scale. This can be interpreted as the distance you have to move in input space before the function value changes significantly. By default, the squared exponential kernel uses a characteristic length-scale of 1 ( $l = 1$ ). To change this behaviour to another, it is sufficient to consider  $r/l$  instead of  $r$  in Equation 2.17. This can be thought as an hyperparameter to optimize, but we will return to this in another section.

Figure 2.1: Three sampled Gaussian Process priors using the Squared Exponential kernel.



### 2.3 Prediction using a Gaussian Process prior

This is probably the most important section in this chapter. We will learn how to incorporate the knowledge of training data  $D = \{(\mathbf{x}_i, y_i) | i = 1, \dots, n\}$  into our Gaussian Process to obtain a posterior predictive distribution. We will start considering the case that we have a noiseless function, that is to say, when  $\sigma_n^2 = 0$ . Let us define  $K(X, X_*)$ , the covariance function evaluated on train and test points,  $K(X, X)$  the covariance function evaluated at only the training points,  $K(X_*, X_*)$  equivalently defined for the test values. Notice the last two have to be square matrices by definition.

Let us also use the following theorem:

**Theorem 1** *Let  $\mathbf{x}$  and  $\mathbf{y}$  be jointly Gaussian:*

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim N \left( \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} A & C \\ C^T & B \end{bmatrix} \right) \quad (2.19)$$

*Then  $\mathbf{x}|\mathbf{y} \sim (\boldsymbol{\mu}_x + CB^{-1}(\mathbf{y} - \boldsymbol{\mu}_y), A - CB^{-1}C^T)$*

The reader might already be guessing what we are about to do now in terms of the training output points  $\mathbf{f}$  and the corresponding testing points  $\mathbf{f}_*$ . According to the prior chosen in 2.18, assume that  $\mathbf{f}$  and  $\mathbf{f}_*$  are jointly Gaussian:

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim N \left( \mathbf{0}, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right) \quad (2.20)$$

We are interested now in the distribution of  $\mathbf{f}_*|\mathbf{f}$ . Simply applying Theorem 1, we can obtain:

$$\mathbf{f}_*|\mathbf{f} \sim N \left( K(X_*, X)K(X, X)^{-1}\mathbf{f}, K(X_*, X_*) - K(X_*, X)K(X, X)^{-1}K(X, X_*) \right) \quad (2.21)$$

This is pretty much everything basic there is to know about Gaussian Process estimation for regression. Now we have a complete predictive distribution over test values  $\mathbf{f}_*$ , and we can do with it what we please. For

example, one could obtain an estimate of this function by drawing samples from a multivariate normal with the computed posterior parameters, or obtain a MAP estimate using the posterior mean.

Let us now consider the scenario where observations are not noise-free, that is, each time you query the function there is a i.i.d Gaussian error with mean 0 and variance  $\sigma_n^2 > 0$ . Assume now the following prior on the noisy observations:

$$\text{Cov}(\mathbf{y}) = K(X, X) + \sigma_n^2 I \quad (2.22)$$

Following the exact same math as before, but taking into account this new term, we got the following joint distribution:

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim N\left(\mathbf{0}, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right) \quad (2.23)$$

And conditioning again  $\mathbf{f}_*$  on  $\mathbf{f}$ , we obtain our final predictive distribution:

$$\mathbf{f}_* | \mathbf{f} \sim N(\bar{\mathbf{f}}_*, \text{Cov}(\mathbf{f}_*)) \quad (2.24)$$

where now:

$$\bar{\mathbf{f}}_* = K(X_*, X) (K(X, X) + \sigma_n^2 I)^{-1} \mathbf{y} \quad (2.25)$$

$$\text{Cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X) (K(X, X) + \sigma_n^2 I)^{-1} K(X, X_*) \quad (2.26)$$

It will probably be useful to note that a Gaussian Process model can be written easily in terms of a Bayesian hierarchical model, since:

$$\mathbf{y} | \mathbf{f} \sim N(\mathbf{f}, \sigma_n^2 I) \quad (2.27)$$

$$\mathbf{f} | X \sim N(\mathbf{0}, K(X, X)) \quad (2.28)$$

In fact, one can also assume other priors, even over  $\sigma_n^2$ . This representation may help us understand the introduction of the *marginal likelihood*. This marginal likelihood in a Gaussian Process setting is defined as:

$$p(\mathbf{y} | X) = \int p(\mathbf{y} | \mathbf{f}, X) p(\mathbf{f} | X) d\mathbf{f} \quad (2.29)$$

Using the results from Equations 2.27 and 2.28 we can derive the integral analytically to obtain:

$$\log p(\mathbf{y} | X) = -\frac{1}{2} \mathbf{y}^T (K + \sigma_n^2 I)^{-1} \mathbf{y} - \frac{1}{2} \log |K + \sigma_n^2 I| - \frac{n}{2} \log 2\pi \quad (2.30)$$

We have now all the necessary ingredients to lay down pseudo-code for your own implementation of a Gaussian Process regressor, as presented in Algorithm 1. It makes use of several tricks for computational stability, such as a Cholesky decomposition and several linear system of equations to avoid directly inverting matrices.

---

**Algorithm 1** Gaussian regressor pseudo-code.

---

```

1: function GPREGRESSOR( $X, \mathbf{y}, k, \sigma_n^2, \mathbf{x}_*$ )
2:    $L \leftarrow \text{chol}(K + \sigma_n^2 I)$ 
3:    $\boldsymbol{\alpha} \leftarrow \text{linsolve}(L^T, \text{linsolve}(L, \mathbf{y}))$ 
4:    $\bar{\mathbf{f}}_* \leftarrow \mathbf{k}_*^T \boldsymbol{\alpha}$ 
5:    $\mathbf{v} \leftarrow \text{linsolve}(L, \mathbf{k}_*)$ 
6:    $\mathbb{V}[\mathbf{f}_*] \leftarrow k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^T \mathbf{v}$ 
7:    $\log p(\mathbf{y} | X) \leftarrow -\frac{1}{2} \mathbf{y}^T \boldsymbol{\alpha} - \sum_i \log L_{ii} - \frac{n}{2} \log 2\pi$ 
8: end function
```

---

pyGPGO includes an implementation of a Gaussian Process regressor under the `GPRegressor` module. A minimal working skeleton of its code is shown below. We do now show the entirety of the module code since most of it is not relevant at this point.

---

```

1 import numpy as np
2 from scipy.linalg import cholesky, solve
3 from collections import OrderedDict
4 from scipy.optimize import minimize
5
6 class GPRegressor:
7     def __init__(self, covfunc, sigma = 0):
8         self.covfunc = covfunc
9         self.sigma = sigma
10    def fit(self, X, y):
11        self.X = X
12        self.y = y
13        self.nsamples = self.X.shape[0]
14        self.K = self.covfunc.K(self.X, self.X)
15        self.L = cholesky(self.K + self.sigma * np.eye(self.nsamples)).T
16        self.alpha = solve(self.L.T, solve(self.L, y))
17        self.logp = -.5 * np.dot(self.y, self.alpha) - np.sum(np.log(np.diag(self.L)))
18        - self.nsamples/2 * np.log(2 * np.pi)
19
20    def predict(self, Xstar, return_std = False):
21        Xstar = np.atleast_2d(Xstar)
22        kstar = self.covfunc.K(self.X, Xstar).T
23        fmean = np.dot(kstar, self.alpha)
24        v = solve(self.L, kstar.T)
25        fcov = self.covfunc.K(Xstar, Xstar) - np.dot(v.T, v)
26        if return_std:
27            fcov = np.diag(fcov)
28        return fmean, fcov
29
30    def update(self, xnew, ynew):
31        y = np.concatenate((self.y, ynew), axis = 0)
32        X = np.concatenate((self.X, xnew), axis = 0)
33        self.fit(X, y)

```

---

A few notes on the code, first notice that the `GPRegressor` class only takes as input the covariance function to use  $k$ , and the noise level  $\sigma_n^2$ . It is only when you call the `fit` method on the function when you actually fit the model on some data  $X, y$ . Quantities like the log-marginal likelihood are store within the object. The `predict` method not only can output the variance for a single prediction, but the covariance matrix for  $m$  new ones if asked to. For convenience, we also implement an `update` method to update the Gaussian Process parameters when new data is available to us.

### 2.3.1 A toy example of Gaussian Process regression

Now that we have both the algorithm and the tools at hand, it may be interesting how a Gaussian Process regressor behaves with a toy example. We try to approximate a simple sine function in the interval  $[0, 2\pi]$ , and plot both the posterior mean and a 95% confidence band using the posterior variance of the fitted process. The code shown below produces Figure 2.2.

---

```

1 import numpy as np
2 from GPRegressor import GPRegressor
3 from covfunc import squaredExponential
4 import matplotlib.pyplot as plt
5
6
7 if __name__ == '__main__':
8     # Build synthetic data (sine function)
9     x = np.arange(0, 2 * np.pi + 0.01, step = np.pi / 2)
10    y = np.sin(x)
11    X = np.array([np.atleast_2d(u) for u in x])[:, 0]
12
13    # Specify covariance function
14    sexp = squaredExponential()
15    # Instantiate GPRegressor class
16    gp = GPRegressor(sexp)
17    # Fit the model to the data
18    gp.fit(X, y)
19
20    # Predict on new data
21    xstar = np.arange(0, 2*np.pi, step = 0.01)
22    Xstar = np.array([np.atleast_2d(u) for u in xstar])[:, 0]
23    ymean, ystd = gp.predict(Xstar, return_std = True)
24
25    # Confidence interval bounds
26    lower, upper = ymean - 1.96 * ystd, ymean + 1.96 * ystd
27
28    # Plot values
29    plt.figure()
30    plt.plot(xstar, ymean, label = 'Posterior mean')
31    plt.plot(xstar, np.sin(xstar), label = 'True function')
32    plt.fill_between(xstar, lower, upper, alpha = 0.4, label = '95% confidence band')
33    plt.grid()
34    plt.legend(loc = 0)
35    plt.show()

```

---

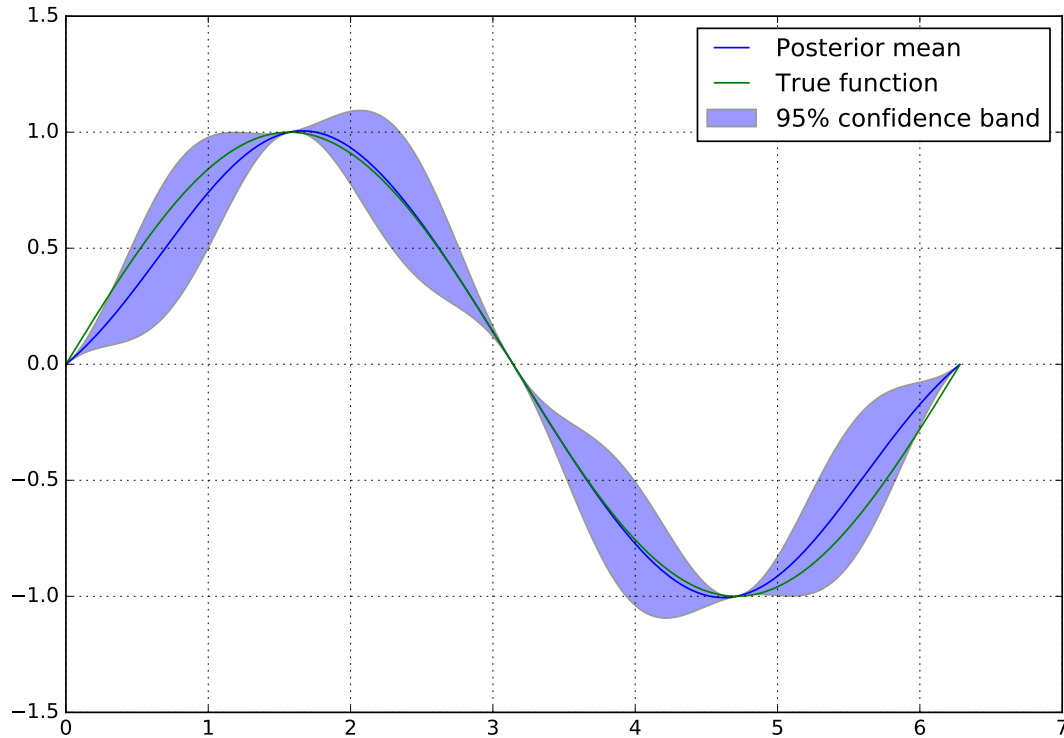
## 2.4 On covariance functions

A covariance function, like the squared exponential kernel that we have been using as an example throughout the chapter encodes our assumptions of similarity between inputs from  $\mathbf{x}$ . We assume *similar* items in input space to have similar values of the target value  $y$ . Not any function of  $\mathbf{x}$  and  $\mathbf{x}'$  can be considered a covariance function, in general should fulfil properties like the following:

- *Stationarity*. A covariance function is said to be stationary if it is a function of  $|\mathbf{x} - \mathbf{x}'|$ . That is to say that it is invariant to translations in the input space. Most of the covariance functions we will see fall into this category.
- *Isotropy*. A covariance function is said to be isotropic if it is *only* function of  $|\mathbf{x} - \mathbf{x}'|$ . Therefore, every isotropic covariance function is stationary.
- *Dot-product*. Some covariance functions are functionals of the dot-product  $|\mathbf{x}^T \mathbf{x}'|$ . These kernels, while invariant to rotations are not invariant to translations.

There is an excellent theoretical analysis of covariance functions in Chapter 4 of Carl E. Rasmussen's book *Gaussian Processes for Machine Learning*. We will not cover this here since it falls beyond the scope of this manual. However, we will start providing examples of the most common covariance functions and their implementation in pyGPGO.

Figure 2.2: A fitted Gaussian Process regressor to samples of the sine function.



The *Squared Exponential* covariance function is the one that we have been using so far. It is also arguably the most used in practice. It takes the general form:

$$k_{SE}(r) = \exp\left(-\frac{r^2}{2l^2}\right) \quad (2.31)$$

Where  $l$  is the parameter controlling its characteristic length-scale. It is useful to define these functions in terms of  $r$  since we can abstract this calculation to another function. A simple skeleton of its implementation in pyGPGO is the following:

---

```

1 import numpy as np
2 from scipy.special import gamma, kv
3 from scipy.spatial.distance import cdist
4
5 def l2norm_(X, Xstar):
6     return cdist(X, Xstar)
7
8 class squaredExponential:
9     def __init__(self, l = 1, sigma2f = 1, sigma2n = 0):
10         self.l = l
11         self.sigma2f = sigma2f
12         self.sigma2n = sigma2n
13
14     def K(self, X, Xstar):
15         r = l2norm_(X, Xstar)
16         return(np.exp(-.5 * r ** 2 / self.l ** 2))

```

---

The *Matern class* of covariance functions takes the form:

$$k_{\text{Matern}}(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \frac{\sqrt{2\nu}r}{l} \right)^\nu K_\nu \left( \frac{\sqrt{2\nu}r}{l} \right) \quad (2.32)$$

with  $\nu, l > 0$  and  $K_\nu$  is a modified Bessel function of the second kind. Notice that if we take limit  $\nu \rightarrow \infty$  we obtain our beloved squared exponential covariance function. pyGPGO implementation below.

---

```

1 class matern:
2     def __init__(self, v = 1, l = 1):
3         self.v, self.l = v, l
4
5     def K(self, X, Xstar):
6         r = l2norm_(x, xstar)
7         bessell = kv(self.v, np.sqrt(2 * self.v) * r / self.l)
8         f = 2 ** (1 - self.v) / gamma(self.v) *
9             (np.sqrt(2 * self.v) * r / self.l) ** self.v
10        res = f * bessell
11        res[np.isnan(res)] = 1
12        return(res)

```

---

The  $\gamma$ -exponential covariance function, of which the squared exponential is a special case. It takes the general form:

$$k(r) = \exp \left( - \left( \frac{r}{l} \right)^\gamma \right) \quad (2.33)$$

for  $0 < \gamma \leq 2$ . Its implementation in pyGPGO below:

---

```

1 class gammaExponential:
2     def __init__(self, gamma = 1, l = 1):
3         self.gamma = gamma
4         self.l = l
5
6     def K(self, X, Xstar):
7         r = l2norm_(X, Xstar)
8         return(np.exp(-(r / self.l) ** self.gamma))

```

---

The *Rational Quadratic* covariance function takes the form:

$$k_{RQ}(r) = \left( 1 + \frac{r^2}{2\alpha l^2} \right)^{-\alpha} \quad (2.34)$$

with  $\alpha, l > 0$ . And again its implementation:

---

```

1 class rationalQuadratic:
2     def __init__(self, alpha = 1, l = 1):
3         self.alpha = alpha
4         self.l = l
5
6     def K(self, X, Xstar):
7         r = l2norm_(X, Xstar)
8         return((1 + r**2/(2 * self.alpha * self.l **2))**(-self.alpha))

```

---

The *arcSin* kernel is an example of a dot product covariance function, therefore non-stationary:

$$k_{\text{arcSin}}(\mathbf{x}, \mathbf{x}') = \frac{2}{\pi} \sin^{-1} \left( \frac{2\mathbf{x}\Sigma\mathbf{x}'}{\sqrt{(1 + 2\mathbf{x}^T\Sigma\mathbf{x})(1 + 2\mathbf{x}'^T\Sigma\mathbf{x})}} \right) \quad (2.35)$$

And I promise this is the final covariance function implementation in this section.

---

```

1 class arcSin:
2     def __init__(self, n, sigma = None):
3         if sigma == None:
4             self.sigma = np.eye(n)
5         else:
6             self.sigma = sigma
7     def k(self, x, xstar):
8         num = 2 * np.dot(np.dot(x[np.newaxis, :], self.sigma), xstar)
9         a = 1 + 2 * np.dot(np.dot(x[np.newaxis, :], self.sigma), x)
10        b = 1 + 2 * np.dot(np.dot(xstar[np.newaxis, :], self.sigma), xstar)
11        res = num / np.sqrt(a * b)
12        return(res)

```

---

### 2.4.1 A zoo of covariance functions

We have seen plenty of covariance function specifications in the last section. Remember that these control the degree of *similarity* between input points. As an exercise, it would be good to recreate the same sine function example that we saw before, using four different stationary different covariance functions. The choice of parameters is the default one in pyGPGO. The script shown below produces Figure 2.3.

---

```

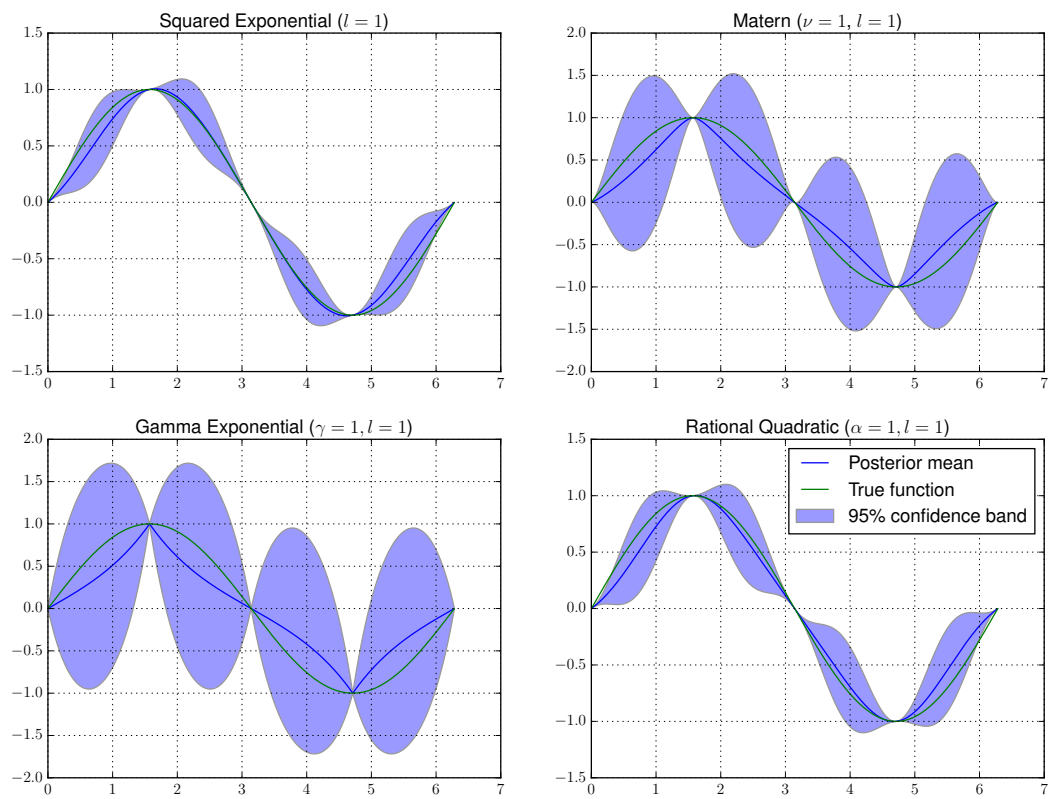
1 import numpy as np
2 from covfunc import *
3 from GPRegressor import GPRegressor
4 import matplotlib.pyplot as plt
5
6
7 if __name__ == '__main__':
8     # Build synthetic data (sine function)
9     x = np.arange(0, 2 * np.pi + 0.01, step = np.pi / 2)
10    y = np.sin(x)
11    X = np.array([np.atleast_2d(u) for u in x])[:, 0]
12
13    # Covariance functions to loop over
14    covfuncs = [squaredExponential(), matern(), gammaExponential(), rationalQuadratic()]
15    titles = [r'Squared Exponential ($l = 1$)', r'Matern ($\nu = 1$, $l = 1$)',
16             r'Gamma Exponential ($\gamma = 1$, $l = 1$)', r'Rational Quadratic ($\alpha = 1$, $l = 1$)']
17    plt.figure()
18    plt.rc('text', usetex=True)
19    for i, cov in enumerate(covfuncs):
20        gp = GPRegressor(cov)
21        gp.fit(X, y)
22        xstar = np.arange(0, 2 * np.pi, step = 0.01)
23        Xstar = np.array([np.atleast_2d(u) for u in xstar])[:, 0]
24        ymean, ystd = gp.predict(Xstar, return_std = True)
25
26        lower, upper = ymean - 1.96 * ystd, ymean + 1.96 * ystd
27        plt.subplot(2, 2, i + 1)
28        plt.plot(xstar, ymean, label = 'Posterior mean')
29        plt.plot(xstar, np.sin(xstar), label = 'True function')
30        plt.fill_between(xstar, lower, upper, alpha = 0.4,
31                        label = '95% confidence band')
32        plt.grid()
33        plt.title(titles[i])
34    plt.legend(loc = 0)
35    plt.show()

```

---



Figure 2.3: Behaviour of different stationary covariance functions with the default parameters in pyGPGO.



## 2.5 Hyperparameter optimization

As seen in the previous sections, different covariance functions have different *hyperparameters*. These control how the kernel measures similarity among different instances of  $\mathbf{x}$ . So far, we have chosen these hyperparameters according to those set default in pyGPGO, but one may want to optimize these according to the data they have at hand. Depending on the situation, this optimization may lead to better models, either in terms of accuracy or interpretability. There are several ways to select these hyperparameters, one more analytical, by optimizing the marginal log-likelihood, and simpler ones based on cross-validation.

### 2.5.1 Type II Maximum Likelihood

This is the Bayesian analytical approach to optimizing hyperparameters. One may quickly notice that Gaussian Processes are non-parametric models, in the sense that apart from the quantities set in the covariance functions, there is nothing else to optimize for. First we will provide a small background on Bayesian model selection. Assume that we have a model  $\mathcal{H}_i$  with *parameters*  $\mathbf{w}$ , *hyperparameters*  $\boldsymbol{\theta}$ , and we have some training data  $X, \mathbf{y}$ . The posterior over the parameters is easily given by:

$$p(\mathbf{w}|\mathbf{y}, X, \boldsymbol{\theta}, \mathcal{H}_i) = \frac{p(\mathbf{y}|X, \mathbf{w}, \mathcal{H}_i)p(\mathbf{w}|\boldsymbol{\theta}, \mathcal{H}_i)}{p(\mathbf{y}|X, \boldsymbol{\theta}, \mathcal{H}_i)} \quad (2.36)$$

where  $p(\mathbf{y}|X, \mathbf{w}, \mathcal{H}_i)$  is the likelihood,  $p(\mathbf{w}|\boldsymbol{\theta}, \mathcal{H}_i)$  our prior distribution over the parameters and  $p(\mathbf{y}|X, \boldsymbol{\theta}, \mathcal{H}_i)$  is called the *evidence* or marginal likelihood. Notice that this last quantity is nothing but the integral over parameter space of the numerator in Equation 2.36.

We can do the same at the next level of inference for the hyperparameters. The posterior of hyperparameters is defined as:

$$p(\boldsymbol{\theta}|\mathbf{y}, X, \mathcal{H}_i) = \frac{p(\mathbf{y}|X, \boldsymbol{\theta}, \mathcal{H}_i)p(\boldsymbol{\theta}|\mathcal{H}_i)}{p(\mathbf{y}|X, \mathcal{H}_i)} \quad (2.37)$$

where now  $p(\boldsymbol{\theta}|\mathcal{H}_i)$  is our prior over hyperparameters. We are interested however in optimizing the denominator in Equation 2.37 with respect to the hyperparameters. Typically, in Bayesian inference to perform the kind of integrals presented before, one has to resort to sampling procedures related to Markov Chain Monte Carlo, such as the Gibbs sampler. In the case of Gaussian processes, all computations are analytically tractable. In fact, the expression of the marginal likelihood was presented in Section 2.3. We reproduce the expression here for completeness:

$$\log p(\mathbf{y}|X) = -\frac{1}{2}\mathbf{y}^T(K + \sigma_n^2 I)^{-1}\mathbf{y} - \frac{1}{2}\log |K + \sigma_n^2 I| - \frac{n}{2}\log 2\pi \quad (2.38)$$

For typical local optimization methods to work fairly well, we may also need an specification of the derivative of the log-marginal likelihood w.r.t. the hyperparameters.

$$\frac{\partial}{\partial \theta_j} \log p(\mathbf{y}|X, \boldsymbol{\theta}) = \frac{1}{2}\mathbf{y}^T K^{-1} \frac{\partial K}{\partial \theta_j} K^{-1} \mathbf{y} - \frac{1}{2} \text{tr} \left( K^{-1} \frac{\partial K}{\partial \theta_j} \right) \quad (2.39)$$

where  $\frac{\partial K}{\partial \theta_j}$  denotes the derivative of the selected covariance function, evaluated at each pair of instances of the training set. For optimization, one may choose to make use of this expression or not, depending on both of the optimization algorithm (gradient ascent, L-BFGS-B...) or on the cost of evaluation of the derivative. In pyGPGO, most of the covariance functions are implemented with a method `gradK` to return the gradient, for example, the complete specification of the `squaredExponential` class is:

---

```

1 class squaredExponential:
2     def __init__(self, l = 1, sigma2f = 1, sigma2n = 0):
3         self.l = l
4         self.sigma2f = sigma2f
5         self.sigma2n = sigma2n
6
7     def K(self, X, Xstar):
8         r = l2norm_(X, Xstar)
9         return(np.exp(-.5 * r ** 2 / self.l ** 2))
10
11    def gradK(self, X, Xstar, param = 'l'):
```

---

```

12         if param == 'l':
13             r = l2norm_(X, Xstar)
14             num = r**2 * np.exp(-r**2 / (2 * self.l**2))
15             den = self.l ** 3
16             l_grad = num / den
17             return(l_grad)
18         else:
19             raise ValueError('Param not found')

```

---

By default, the `GPRegressor` class in `pyGPGO` does no hyperparameter optimization unless explicitly asked to. If the gradient of the covariance function is available, it will make use of it, if not, it will be estimated in a numerical fashion using the L-BFGS-B algorithm.

### Another toy example: Optimizing the characteristic length-scale

To illustrate the previous point, it may be a good idea to see the behaviour of the marginal log-likelihood and its gradient when we modify the characteristic length scale  $l$  in the squared exponential covariance function. The sine function will also serve as playground here. The code below produces Figure 2.4.

---

```

1 import numpy as np
2 from GPRegressor import GPRegressor
3 from covfunc import squaredExponential
4 import matplotlib.pyplot as plt
5
6 def gradient(gp, sexp):
7     alpha = gp.alpha
8     K = gp.K
9     gradK = sexp.gradK(gp.X, gp.X, 'l')
10    inner = np.dot(np.atleast_2d(alpha).T, np.atleast_2d(alpha)) - np.linalg.inv(K)
11    return(.5 * np.trace(np.dot(inner, gradK)))
12
13
14 if __name__ == '__main__':
15     x = np.arange(0, 2 * np.pi + 0.01, step = np.pi / 2)
16     X = np.array([np.atleast_2d(u) for u in x])[:, 0]
17     y = np.sin(x)
18
19     logp = []
20     grad = []
21     length_scales = np.linspace(0.1, 2, 1000)
22
23     for l in length_scales:
24         sexp = squaredExponential(l = l)
25         gp = GPRegressor(sexp)
26         gp.fit(X, y)
27         logp.append(gp.logp)
28         grad.append(gradient(gp, sexp))
29
30     plt.figure()
31     plt.subplot(2, 1, 1)
32     plt.plot(length_scales, logp)
33     plt.title('Marginal log-likelihood')
34     plt.xlabel('Characteristic length-scale l')
35     plt.ylabel('log-likelihood')
36     plt.grid()
37     plt.subplot(2, 1, 2)
38     plt.plot(length_scales, grad, '--', color = 'red')
39     plt.title('Gradient w.r.t. l')
40     plt.xlabel('Characteristic length-scale l')

```

```
41 plt.grid()
42 plt.show()
```

---

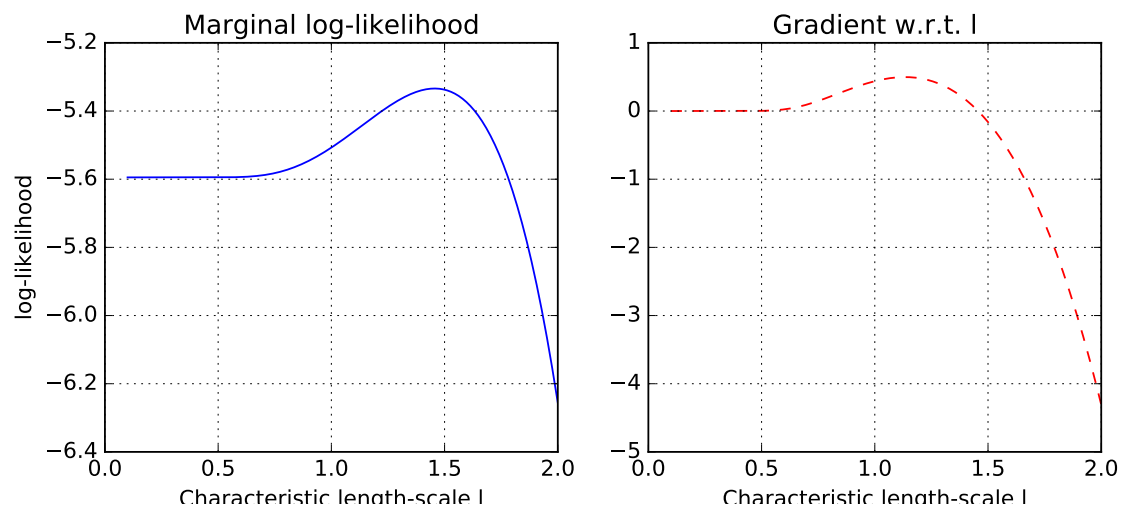
### 2.5.2 Cross validation

There is a very extensive analytical evaluation of cross validation using Gaussian Process in Rasmussen's book. Here, we will only lay down some very basic ideas related to model selection from a very basic machine learning perspective. The concepts presented here should not be new to the reader, but are more useful if one plans to use Gaussian Processes only as a regression model, not as a surrogate one towards another objective, as we will do in the following chapter.

To evaluate the performance of hyperparameters  $\theta$ , in  $D = (X, \mathbf{y})$  one could do the following:

- **Holdout test.** Consider  $D = \{D_T, D_V\}$  as a training and validation set from your data and  $\nu$  the set of hyperparameter vectors  $\nu$  to test. Train your Gaussian Process regressor on  $D_T$  with some subset of hyperparameters from  $\nu$  and test its performance according to some loss metric  $\mathcal{L}$  on  $D_V$ . Choose hyperparameters according to the lower loss obtained.
- **$k$ -fold cross validation.** Instead of considering a single test set  $D_V$ , partition  $D = D_1, \dots, D_k$ . Train your model iteratively on  $k - 1$  sets and test on the remaining one. Consider an average of losses for each hyperparameter to test.

Figure 2.4: Log-marginal likelihood and its gradient w.r.t to the characteristic length-scale. Notice there seems to be an optimal point at around  $l = 1.4$ .





## Chapter 3

# Bayesian optimization





## Chapter 4

# Experiments



## Chapter 5

# pyGPGO: A simple Python Package for Bayesian Optimization