

Bayesian optimization in machine learning

José Jiménez Luna

March 9, 2017

Contents

1	Organization of this work	5
1.1	Introduction	5
1.2	Organization of the thesis	6
2	Gaussian Process regression	7
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	11
2.3.1	A toy example of Gaussian Process regression	13
2.3.2	Picking a winner	14
2.4	On covariance functions	16
2.4.1	Visualizing different covariance functions	18
2.5	Hyperparameter optimization	19
2.5.1	Type II Maximum Likelihood	20
2.5.2	Full Bayesian perspective	22
2.5.3	Cross validation	22
2.6	Further theoretical aspects	22
2.6.1	Gaussian processes as linear smoothers	22
2.6.2	Explicit basis functions	23
3	Bayesian optimization	25
3.1	Preliminaries	25
3.2	The bayesian optimization framework	25
3.3	On acquisition functions	26
3.3.1	Improvement-based policies	26
3.3.2	Optimistic policies	27
3.3.3	Information-based policies	27
3.3.4	Acquisition function portfolios	28
3.3.5	Visualizing the behaviour of an acquisition function	28
3.3.6	Why does Bayesian Optimization work?	29
3.4	Role of GP hyperparameters in optimization	31
3.5	Optimizing the acquisition function	32
3.6	Computational costs	33
3.6.1	Approximations to the analytical GP. Alternative surrogates.	33
3.6.2	Parallelization	34
3.7	Step-by-step examples	35
3.7.1	Optimizing (yet-again) the sine function	35
3.7.2	Optimizing the Rossenbrock function	35

4	Experiments	39
4.1	Benchmarking rules	39
4.1.1	Other strategies for hyper-parameter optimization	39
4.1.2	Evaluation metrics	40
4.1.3	Bayesian optimization setup	41
4.1.4	Machine-learning models used	41
4.2	The binding affinity dataset	41
5	pyGPGO: A simple Python Package for Bayesian Optimization	43
5.0.1	Installation	43
5.0.2	Usage	43
5.0.3	Examples	43
5.0.4	Features	43
5.0.5	Future work	43

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 master's thesis resembles a programming book or a manual. All the examples seen through this work were coded from scratch.
- 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 medical, biological or physical phenomena), applying such methodology and comparing to other common strategies applied for the same problem.
- Finally, to write a complete software package for users to apply Bayesian Optimization in their research. This comes in the form of a Python (>3.5) package named pyGPGO. The code can be accessed through <https://github.com/hawk31/pyGPGO>. The entire software package is MIT licensed. All the examples and code snippets throughout this manual are based on this software. While certainly there are a couple implementations of Global Optimization software in Python, the software described here is modular, very minimalistic and requires minimal dependencies, while still maintaining most of the other's functionality

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:

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

Most optimization procedures (local based ones such as gradient ascent, for example) assume that the function f is closed-form, that is, it has an analytical expression, 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 find the global minimum of non-necessarily convex, expensive functions that are expensive to evaluate. 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 (also known as statistical learning), we are usually interested in minimizing a loss function \mathcal{L} . These losses can take many forms. For instance, 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 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 Organization of the thesis

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.

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 point of view and from a functional point of 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 were used for other benchmarking purposes in other studies.

Chapter 5 has no theoretical content nor testing content. It will cover technical explanations of pyGPGO, the software developed alongside this manual. Software usage examples are also provided.

Chapter 2

Gaussian Process regression

In this chapter we will focus on regression problems. 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 our predictor matrix.

One can interpret a Gaussian Process in several ways. The most widely known is the function space view, which is the one we will cover second here and the one we will assume for the rest of the thesis. In this view, we consider a Gaussian Process to be a stochastic process, hence, 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, through the use of kernel functions.

2.1.1 Standard Bayesian linear regression

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

$$y = X^T \mathbf{w} + \epsilon \quad (2.1)$$

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:

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

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:

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

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:

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

where $A = \sigma^{-2}XX^T + \Sigma_p^{-1}$. Notice that a simple 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 due to its linearity. 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 through 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:

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

where we assume the same prior over \mathbf{w} as in Equation 2.3. 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:

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

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 Σ_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*:

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

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, the model $f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}$ 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 several covariance functions, 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:

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

It is straightforward to draw samples from a Gaussian Process. In particular, since we work with a finite number of points, choose an arbitrary number of them X_* and compute the squared exponential kernel. Then the procedure is simplified to just sampling from the following multivariate Gaussian:

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

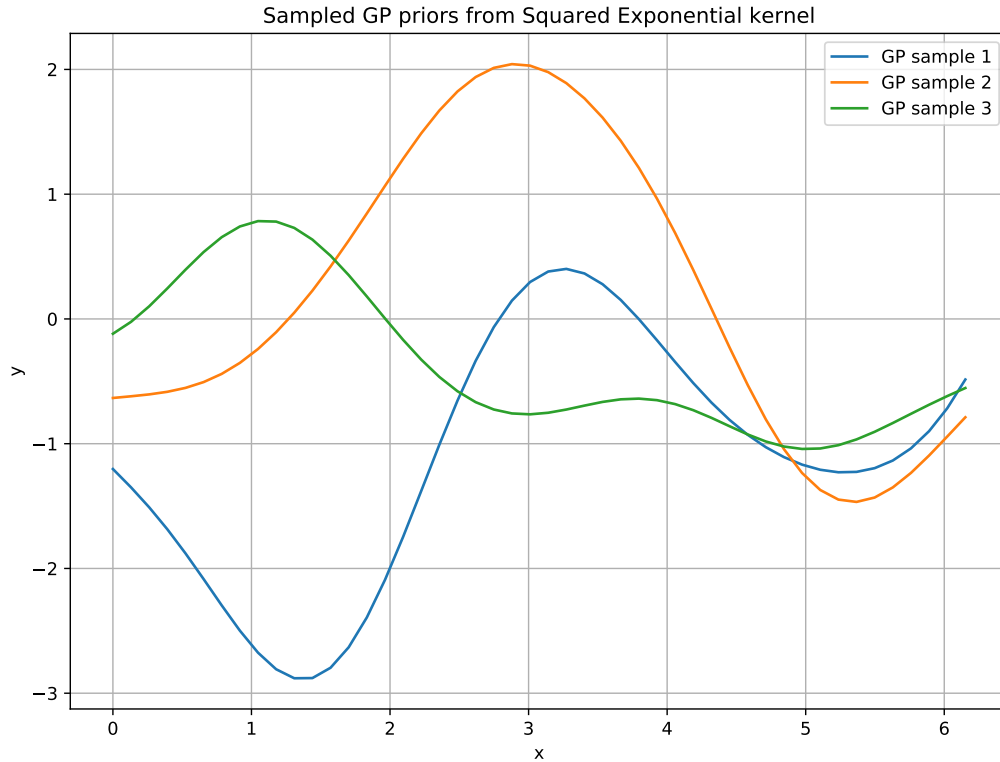
To illustrate this point, we will write a very simple Python script to draw samples from this function. For the moment, consider evenly spaced samples with a step of $\frac{\pi}{16}$. Some of the code here already uses the pyGPGO implementation of the covariance function, to simplify computation. This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

```

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



```

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.

2.3 Prediction using a Gaussian Process prior

In this particular section, arguably the most important one, we will learn how to incorporate the knowledge of training data $D = \{(\mathbf{x}_i, y_i) \mid i = 1, \dots, n\}$ 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 $\mathbf{x}|\mathbf{y} \sim (\boldsymbol{\mu}_x + CB^{-1}(\mathbf{y} - \boldsymbol{\mu}_y), A - CB^{-1}C^T)$

According to the prior distribution 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:

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

This covers all the basics for a Gaussian Process regression model. Notice that now we have a complete predictive distribution over test values \mathbf{f}_* , and this provides us with plenty of choices. 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 operations 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:

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

where now:

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

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

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 σ_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:    $\bar{f}_* \leftarrow \mathbf{k}_*^T \boldsymbol{\alpha}$ 
5:    $\mathbf{v} \leftarrow \text{linsolve}(L, \mathbf{k}_*)$ 
6:    $\mathbb{V}[f_*] \leftarrow k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^T \mathbf{v}$ 
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 σ_n^2 . It is only when you call the `fit` method on the function when you actually fit the model on some data X, \mathbf{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()

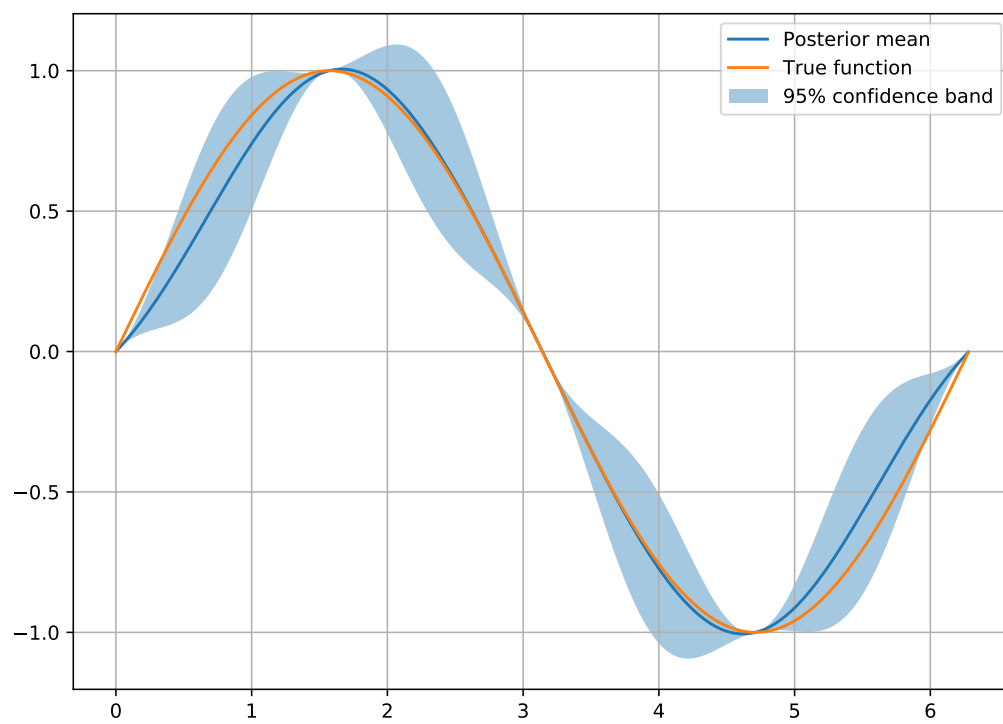
```

2.3.2 Picking a winner

In the previous section we have shown how to compute predictive posterior distribution for function outputs y_* given a new input \mathbf{x}_* . These are given by a Gaussian distribution with a certain mean and variance. In plenty of production settings, however, it is more common to provide a single value, or estimate y_{guess} that is *optimal* in some sense. To define a sense of optimality, define a loss function $\mathcal{L}(y_{\text{true}}, y_{\text{guess}})$. This, defines a penalty incurred by taking the decision to use y_{guess} when the true value is y_{true} . For example, this could be the mean square or mean absolute error function. In the Bayesian setting, there is no clear mention of a loss function in any stage. In the frequentist setting, however, a model is usually trained by minimizing this loss. Furthermore, there is a clear separation between loss and likelihood in the Bayesian setting, the latter used for training, with prior information. The loss function however only captures the consequences of making a single specific choice given a true state.

Again, we would like to somehow pick a *winner* y_{guess} that minimizes our loss. Without knowing y_{true} , our best choice is to minimize the expected loss, averaging with respect to our model:

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



$$\mathcal{R}_{\mathcal{L}}(y_{\text{guess}}|\mathbf{x}_*) = \int \mathcal{L}(y_*, y_{\text{guess}}) p(y_*|\mathbf{x}_*, \mathcal{D}) dy \quad (2.31)$$

Our optimal value is the one that minimizes this expected loss:

$$y_{\text{optimal}}|\mathbf{x}_* = \arg \min_{y_{\text{guess}}} \mathcal{R}_{\mathcal{L}}(y_{\text{guess}}|\mathbf{x}_*) \quad (2.32)$$

It can be proven that the value y_{guess} that minimizes Equation 2.32 for the absolute loss function is the median of $p(y_*|\mathbf{x}_*, \mathcal{D})$. For the squared loss function, it is the mean of the same distribution. Since in our case we are dealing with the Gaussian distribution, median and mean coincide, and the most reasonable *winner* will therefore be the specific value of the posterior mean.

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 all functions of \mathbf{x} and \mathbf{x}' can be defined as covariance function. Covariance functions (though not all) tend to satisfy different properties:

- *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 thesis. However, we will start providing examples of the most common covariance functions and their implementation in pyGPGO.

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:

$$k_{SE}(r) = \exp\left(-\frac{r^2}{2l^2}\right) \quad (2.33)$$

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:

$$k_{\text{Matern}}(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{l} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}r}{l} \right) \quad (2.34)$$

with $\nu, l > 0$ and K_ν 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         bessel = 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 * bessel
11        res[np.isnan(res)] = 1
12        return(res)

```

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

$$k(r) = \exp \left(- \left(\frac{r}{l} \right)^\gamma \right) \quad (2.35)$$

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:

$$k_{RQ}(r) = \left(1 + \frac{r^2}{2\alpha l^2} \right)^{-\alpha} \quad (2.36)$$

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 = 1
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:

$$k_{\text{arcSin}}(\mathbf{x}, \mathbf{x}') = \frac{2}{\pi} \sin^{-1} \left(\frac{2\mathbf{x}\Sigma\mathbf{x}'}{\sqrt{(1 + 2\mathbf{x}^T\Sigma\mathbf{x})(1 + 2\mathbf{x}'^T\Sigma\mathbf{x})}} \right) \quad (2.37)$$

```

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 Visualizing different covariance functions

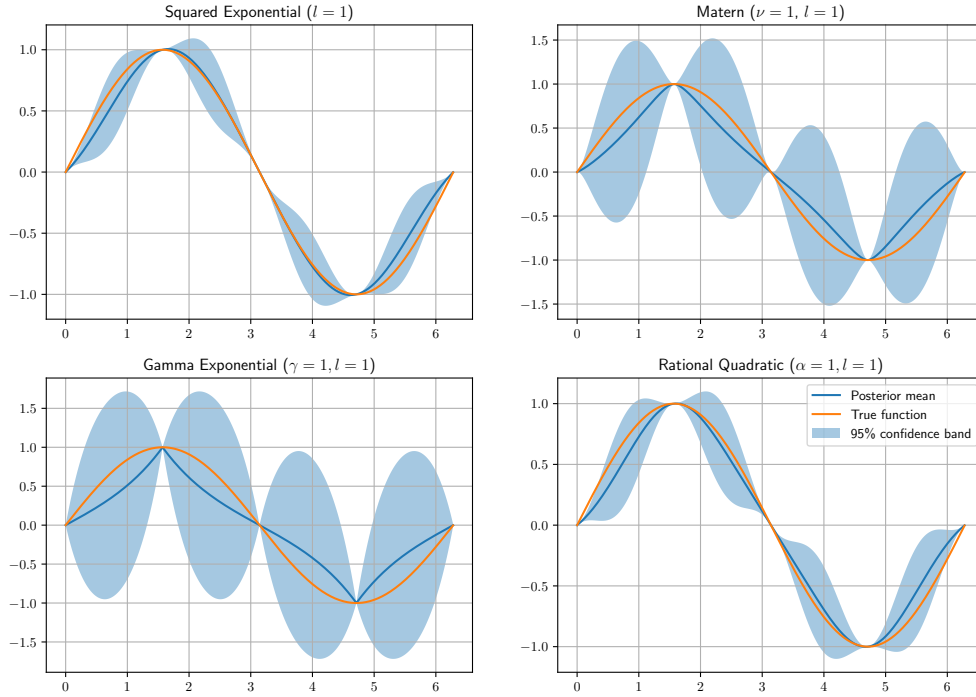
We have seen plenty of covariance function specifications in the last section. Remember that these control the degree of *similarity* between input points. As an exercise, it would be good to recreate the same sine function example that we saw before, using four different stationary different covariance 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)

```

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



```

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()

```

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 empirical Bayes 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:

$$p(\mathbf{w}|\mathbf{y}, X, \boldsymbol{\theta}, \mathcal{H}_i) = \frac{p(\mathbf{y}|X, \mathbf{w}, \mathcal{H}_i)p(\mathbf{w}|\boldsymbol{\theta}, \mathcal{H}_i)}{p(\mathbf{y}|X, \boldsymbol{\theta}, \mathcal{H}_i)} \quad (2.38)$$

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 \mathbf{w} space of the numerator in Equation 2.38.

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

$$p(\boldsymbol{\theta}|\mathbf{y}, X, \mathcal{H}_i) = \frac{p(\mathbf{y}|X, \boldsymbol{\theta}, \mathcal{H}_i)p(\boldsymbol{\theta}|\mathcal{H}_i)}{p(\mathbf{y}|X, \mathcal{H}_i)} \quad (2.39)$$

where now $p(\boldsymbol{\theta}|\mathcal{H}_i)$ is our prior over hyperparameters. We are interested however in optimizing the denominator in Equation 2.39 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:

$$\log p(\mathbf{y}|X) = -\frac{1}{2}\mathbf{y}^T(K + \sigma_n^2 I)^{-1}\mathbf{y} - \frac{1}{2}\log |K + \sigma_n^2 I| - \frac{n}{2}\log 2\pi \quad (2.40)$$

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.

$$\frac{\partial}{\partial \theta_j} \log p(\mathbf{y}|X, \boldsymbol{\theta}) = \frac{1}{2}\mathbf{y}^T K^{-1} \frac{\partial K}{\partial \theta_j} K^{-1} \mathbf{y} - \frac{1}{2} \text{tr} \left(K^{-1} \frac{\partial K}{\partial \theta_j} \right) \quad (2.41)$$

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 Full Bayesian perspective

The approach described in the previous subsection optimizes a marginal log-likelihood and therefore is not truly Bayesian. The full Bayesian treatment of hyper-parameters θ , we would marginalize over these and compute *integrated* acquisition functions of the type:

$$\hat{\alpha}(\mathbf{x}; \{\mathbf{x}_n, y_n\}) = \int \alpha(\mathbf{x}; \{\mathbf{x}_n, y_n\}, \theta) p(\theta | \{\mathbf{x}_n, y_n\}) d\theta \quad (2.42)$$

It can be shown [ref] that for both probability of improvement and expected improvement acquisition functions, this integral is the correct generalization to account for uncertainty in hyper-parameters. The calculations to compute such integrations however, are not analytically tractable, and one must resort to approximations, relying on MCMC or slice sampling, as described in [ref].

2.5.3 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 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 θ , 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 $\boldsymbol{\nu}$ the set of hyperparameter vectors $\boldsymbol{\nu}$ to test. Train your Gaussian Process regressor on D_T with some subset of hyperparameters from $\boldsymbol{\nu}$ 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.

2.6 Further theoretical aspects

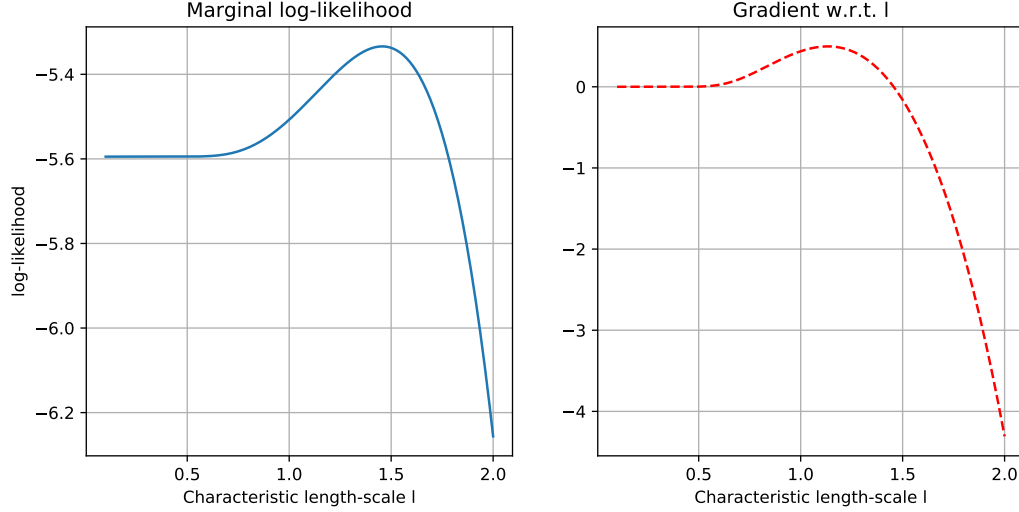
In this section, we will briefly mention other theoretical aspects of Gaussian Processes. This includes for example, how Gaussian Process can be seen as *linear smoothers*, by means of a spectral analysis or how to incorporate explicit basis functions into the model. Readers can safely skip this section if not particularly interested, as it will not be mentioned any other time throughout the text.

2.6.1 Gaussian processes as linear smoothers

As many machine learning algorithms, the main objective of a Gaussian Process regressor is to reconstruct the underlying signal f by removing noise ϵ . It does this by computing a weighted average of the values \mathbf{y} . In particular, as seen in 2.3, it can be written as:

$$\bar{f}(\mathbf{x}_*) = \mathbf{k}(\mathbf{x}_*)^T (K + \sigma_n^2 I)^{-1} \mathbf{y} \quad (2.43)$$

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 $l = 1.4$.



Therefore, one can see a Gaussian Process regressor as a linear smoother [ref]. We can study this smoothing in terms of spectral analysis [ref]. Again, for training points, predicted training points $\bar{\mathbf{f}}$ are:

$$\bar{\mathbf{f}} = \mathbf{K}(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} \quad (2.44)$$

Write \mathbf{K} using its eigenvalue decomposition $\mathbf{K} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^T$, with λ_i and \mathbf{u}_i its i -th eigenvalue and eigenvector respectively. Since \mathbf{K} is a covariance matrix, it is symmetric positive semidefinite, and therefore has positive eigenvalues. If we note $\gamma_i = \mathbf{u}_i^T \mathbf{y}$, then:

$$\bar{\mathbf{f}} = \sum_{i=1}^n \frac{\gamma_i \lambda_i}{\lambda_i + \sigma_n^2} \mathbf{u}_i \quad (2.45)$$

For the covariance functions we have studied in section 2.4, the eigenvalues are larger for slowly varying eigenvectors, so the more frequent items in \mathbf{y} get smoothed-out. The effective number of degrees of freedom in a Gaussian Process model can be defined as the number of used eigenvectors:

$$\text{df}(\mathbf{K}) = \text{tr}(\mathbf{K}(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}) = \sum_{i=1}^n \frac{\lambda_i}{\lambda_i + \sigma_n^2} \quad (2.46)$$

To make the explanation clearer, let us define $\mathbf{h}(\mathbf{x}_*) = (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{x}_*)$. So for a new given point, prediction is defined as $\bar{\mathbf{f}}(\mathbf{x}_*)^T \mathbf{y}$, that is, a linear combination of \mathbf{y} , with weights $\mathbf{h}(\mathbf{x}_*)$. A Gaussian Process regressor is a linear smoother, since the weight function \mathbf{h} does not depend directly on \mathbf{y} . While a regular linear model defines a linear combination of the inputs, a linear smoother defines a linear combination of the targets. This weight function depends directly on the specific location of the n training points, by means of the matrix inversion of $\mathbf{K} + \sigma_n^2 \mathbf{I}$, therefore observations close in input space are smoothed out.

2.6.2 Explicit basis functions

Notice that during the entire chapter, we have considered a Gaussian Process prior with mean $m(\mathbf{x}) = 0$ for simplicity reasons. One may want, however, to define a different mean value for the prior. On the other hand, imposing $m(\mathbf{x}) = 0$ is not a strong assumption, since the posterior is not constrained to be zero as well. With an explicit mean function $m(\mathbf{x}) \neq 0$, the prior becomes:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}^*)) \quad (2.47)$$

and the mean of the posterior predictive distribution then becomes, very naturally:

$$\mathbf{f}_* = \mathbf{m}(X_*) + k(X_*, X)K^{-1}(\mathbf{y} - \mathbf{m}(X)) \quad (2.48)$$

The variance of the posterior predictive distribution remains the same as in Equation 2.3. In practice, however, it may not be clear how to specify a prior mean function for the process. In some cases it may be useful to define a few parametric basis function, whose parameters β we have to estimate from training data. Formally:

$$g(\mathbf{x}) = f(\mathbf{x}) + \mathbf{h}(\mathbf{x})^T \beta \quad (2.49)$$

where $f(\mathbf{x})$ is a regular zero-mean Gaussian Process prior, $\mathbf{h}(\mathbf{x})$ are our chosen basis functions, and β are our parameters. For example, if we are interested in polynomial regression, then $\mathbf{h}(x) = (1, x, x^2, \dots)$. One could consider optimizing β the same way as with our kernel hyperparameters, but if we assume a Gaussian prior $\beta \sim \mathcal{N}(\mathbf{b}, B)$, we can solve analytically to obtain another Gaussian Process:

$$g(\mathbf{x}) \sim \mathcal{GP}(\mathbf{h}(\mathbf{x}^T \mathbf{b}, k(\mathbf{x}, \mathbf{x}^*) + \mathbf{h}(\mathbf{x})^T B \mathbf{h}(\mathbf{x}^*)) \quad (2.50)$$

Notice that now we have an extra term in the covariance function. This is caused by the uncertainty in the parameters of the mean. Now predictions are made by substituting these parameters into Equation 2.49. An explicit version for the mean and covariance is given by:

$$\bar{\mathbf{g}}(X_*) = H_*^T \hat{\beta} + K_*^T K^{-1}(\mathbf{y} - H^T \hat{\beta}) = \bar{\mathbf{f}}(X_*) + R^T \hat{\beta} \quad (2.51)$$

$$\text{Cov}(\mathbf{g}_*) = K_{**} + R^T (B^{-1} + H K^{-1} H^T)^{-1} R \quad (2.52)$$

where H and H_* are matrices containing the evaluation of the chosen basis functions over training and testing points respectively, $\hat{\beta} = (B^{-1} + H K^{-1} H^T)^{-1} (H K^{-1} \mathbf{y} + B^{-1} \mathbf{b})$ and $R = H_* - H K^{-1} K_*$. The posterior process parameters can be interpreted as such: $\hat{\beta}$ is a mean of the model linear parameters, a compromise between the prior and the likelihood provided by the data. The mean of the process is simply $\hat{\beta}$ plus our typical Gaussian Process prediction of the residuals. The covariance matrix is just the addition of our regular expression and a non-negative term.

Consider the limit of B^{-1} as it approaches O (O being a zero-filled matrix), that is when the prior is vague. We then get a predictive distribution independent of \mathbf{b} :

$$\bar{\mathbf{g}}(X_*) = \bar{\mathbf{f}}(X_*) + R^T \hat{\beta} \quad (2.53)$$

$$\text{Cov}(\mathbf{g}_*) = K_{**} + R^T (H K^{-1} H^T)^{-1} R \quad (2.54)$$

where now $\hat{\beta} = (H K^{-1} H^T)^{-1} H K^{-1} \mathbf{y}$.

We explore now the behaviour of the marginal log-likelihood under this model where we assume a Gaussian prior $\beta \sim \mathcal{N}(\mathbf{b}, B)$. Formally:

$$\log p(\mathbf{y}|X, \mathbf{b}, B) = -\frac{1}{2} \log |K + H^T B H| - \frac{n}{2} \log 2\pi \quad (2.55)$$

In the same way as before, exploring the limit $B^{-1} \rightarrow O$, the prior becomes irrelevant, so we can safely assume that $\mathbf{b} = 0$, yielding:

$$\log p(\mathbf{y}|X, \mathbf{b} = \mathbf{0}, B) = -\frac{1}{2} \mathbf{y}^T K^{-1} \mathbf{y} + \frac{1}{2} \mathbf{y}^T C \mathbf{y} \quad (2.56)$$

$$- \frac{1}{2} (\log |K| + \log |B| + \log |A| + n \log 2\pi) \quad (2.57)$$

where $A = B^{-1} + H K^{-1} H^T$ and $C = K^{-1} H^T A^{-1} H K^{-1}$.

Chapter 3

Bayesian optimization

3.1 Preliminaries

In this chapter we will deal with the main topic of this master's thesis, Bayesian Optimization. Here, we approach global optimization from the viewpoint of Bayesian theory, as a sequential problem. For the moment, imagine that we have a very expensive function to evaluate $f : \mathbb{R}^n \rightarrow \mathbb{R}$. This function, for the purposes of this work, will be the negative of a loss function in a machine learning problem, or any other fitness function that we wish to maximize. Formally, we wish to maximize over a compact set \mathcal{A} .

$$\max_{\mathbf{x} \in \mathcal{A}} f(\mathbf{x}) \quad (3.1)$$

For technical reasons, we also assume that the function is *Lipschitz-continuous*, that is, there exists some constant C such that $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}$:

$$\|f(\mathbf{x}_1) - f(\mathbf{x}_2)\| \leq C \|\mathbf{x}_1 - \mathbf{x}_2\| \quad (3.2)$$

We are also interested in global optimization instead of local, since loss functions do not have to be convex over hyperparameter space. That is, we can not assume that we can find a point \mathbf{x}^* such that:

$$f(\mathbf{x}^*) \geq f(\mathbf{x}), \forall \mathbf{x} \text{ s.t. } \|\mathbf{x}^* - \mathbf{x}\| < \epsilon \quad (3.3)$$

The function we are typically interested may not have an analytical expression that we can analyse, take derivatives etc. Most we will assume here is that we can just query the function over any point to evaluate $\mathbf{x} \in \mathcal{A}$ and some bounds to optimize over. This is normally called a *black box* function. Moreover, the function response can be noisy. This is the case when optimizing a loss or fitness function in machine learning on a holdout test, for example, only having an estimation of its real value.

Bayesian optimization has risen over the last few years as a very attractive method to optimize expensive to evaluate black box functions. [refs] It has grabbed the attention of machine learning researchers over simpler model hyperparameter optimization strategies, such as grid search [ref], random search [ref] or simulated annealing [ref]. Bayesian optimization uses prior information and evidence to define a posterior distribution over the space of functions. The model we will use to model this posterior is Gaussian Process regression, for which we have studied its basics in the previous chapter.

3.2 The bayesian optimization framework

Assume that we have sampled our function f to optimize a small number of times k . Notice this can be treated as a regression problem where \mathbf{x}_k is the k -th point we have sampled and y_k its (possibly noisy) function evaluation. We can fit a Gaussian Process regression model over the set of sampled points and evaluations. Remember from Section 2.3 that this gives us a posterior distribution over all possible values in \mathcal{A} . Basically, we will use this information to optimize the function efficiently. Note by

$$\mathcal{D}_n = \{\mathbf{x}_i, y_i, i = 1, \dots, n\}. \quad (3.4)$$

the set of training values.

Bayesian optimization is a sequential model-based approach for optimization. The posterior distribution facilitated by the Gaussian Process allows us to define what we will call an *acquisition function* α that will guide the search for the most promising point from \mathcal{A} to evaluate each step. Once we have sampled said point, we re-fit our Gaussian Process to update our posterior with the new information gathered and proceed the same way until convergence. The mentioned acquisition functions are both heuristic and myopic, in the sense that they define some behaviour given the posterior and only take the information available at a single step of the optimization. Typically, these functions trade-off exploration and exploitation of the target function, and their optima is close to where the posterior variance of the Gaussian Process is large (exploration) or where its posterior mean is high (exploitation). We will choose the next sampled point to evaluate by maximizing these acquisition functions. Algorithm 2 provides pseudo-code to implement a basic bayesian optimization module.

Algorithm 2 Bayesian optimization framework.

- 1: Sample a small number of points $\mathbf{x} \in \mathcal{A}$. Evaluate $f(\mathbf{x})$ to get \mathcal{D}_n
 - 2: **for** $n = 1, 2, \dots$ **do**
 - 3: Fit a GP regression model on \mathcal{D}_n
 - 4: $\mathbf{x}_{n+1} \leftarrow \arg \max_{\mathbf{x}} \alpha(\mathbf{x}, \mathcal{D}_n)$
 - 5: Evaluate $f(\mathbf{x}_{n+1}) = y_{n+1}$
 - 6: Augment data $\mathcal{D}_{n+1} = \{\mathcal{D}_n, (\mathbf{x}_{n+1}, y_{n+1})\}$
 - 7: **end for**
-

3.3 On acquisition functions

Thus far we have described the statistical model behind the optimization framework. The next natural step to ask is how we can define acquisition functions depending on its behaviour or the function we wish to maximize. In pyGPGO, the most common acquisition functions are implemented under the `Acquisition` class in the `acquisition` module. We can classify most of them in three main groups: improvement-based, optimistic, and information-based policies. We will start by analysing each of them:

3.3.1 Improvement-based policies

These acquisition functions' behaviour is to favour points that are in some way likely to improve upon the best observed value so far τ . Since any finite sample of a Gaussian Process is a multivariate Gaussian distribution, the most straightforward idea is to use an estimation of the *probability of improvement* of point evaluation ν w.r.t. τ .

$$\alpha_{\text{PI}}(\mathbf{x}, \mathcal{D}_n) = P(\nu > \tau) = \Phi\left(\frac{\mu_n(\mathbf{x}) - \tau}{\sigma_n(\mathbf{x})}\right) \quad (3.5)$$

where Φ denotes the standard normal cumulative density function and $\mu_n(\mathbf{x})$ and $\sigma_n(\mathbf{x})$ are the posterior mean and standard deviation of the fitted Gaussian Process at step n . In a sense, what this acquisition function is doing is just accumulating the posterior probability mass above τ at \mathbf{x} . The associated utility function is just an indicator of improvement $I(\mathbf{x}, \nu, \theta) = \mathbb{I}(\nu > \tau)$. While this is a very natural acquisition function to use, it has been shown [ref] that it behaves greedily if the best τ is not known.

Another very popular acquisition function is called *expected improvement*. This incorporates the amount of improvement over τ by weighing the probability of improvement over the difference $\nu - \tau$. Formally:

$$I(\mathbf{x}, \nu, \theta) = (\nu - \tau)\mathbb{I}(\nu > \tau) \quad (3.6)$$

Taking the expectation yields the expected improvement acquisition function:

$$\alpha_{\text{EI}}(\mathbf{x}, \mathcal{D}_n) = \mathbb{E}[I(\mathbf{x}, \nu, \theta)] = (\mu_n(\mathbf{x}) - \tau)\Phi\left(\frac{\mu_n(\mathbf{x}) - \tau}{\sigma_n(\mathbf{x})}\right) + \sigma_n(\mathbf{x})\phi\left(\frac{\mu_n(\mathbf{x}) - \tau}{\sigma_n(\mathbf{x})}\right) \quad (3.7)$$

where ϕ is in this case the standard normal density function. This acquisition function is by far the most used, since it has been empirically studied [ref] and proven convergence rates for [ref]. We have assumed that τ is the best observed value so far during the optimization procedure, but theoretical convergence is only guaranteed when τ is the best value f can take in \mathcal{A} . During practical research, however, this does not seem to be a concern.[ref]

3.3.2 Optimistic policies

Optimistic acquisition functions have their origins in the multi-armed bandit setting [ref]. These policies behave optimistically in the face of uncertainty, as a way to tradeoff exploration and exploitation. The most popular of methods in this class is the *Gaussian process upper confidence bound* (GP-UCB) [ref], with provable regret bounds. It works by taking a quantile of the posterior process, and since it is Gaussian, we can derive the result analytically:

$$\alpha_{\text{UCB}}(\mathbf{x}, \mathcal{D}_n) = \mu_n(\mathbf{x}) + \beta_n \sigma_n(\mathbf{x}) \quad (3.8)$$

where β_n controls the quantile we may be interested in. Theoretically motivated by the multi-armed bandits, there are guidelines to select and schedule β_n dynamically. Notice that if we choose $\beta = \beta_n$, to be one value or another, we will be encouraging the algorithm to exploit frequently by choosing points with high posterior mean (small β) or to explore frantically by choosing points with high posterior variance (high β).

3.3.3 Information-based policies

These are a newer class of methods that consider the posterior distribution over an unknown minimizer \mathbf{x}^* . One of the most popular policies in this categories is again motivated by the multi-armed bandit problem, Thompson sampling. [ref]. This very old strategy consists in randomly sampling rewards from the posterior distribution and picking the highest one. It is a randomized acquisition function in the sense:

$$\mathbf{x}_{n+1} \sim p_*(\mathbf{x}|\mathcal{D}_n) \quad (3.9)$$

This method, however, is not as simple to implement as the previously discussed one. It is not entirely clear how to sample in the continuous space of the Gaussian Process. There have been studies that solve this issue by using techniques like spectral sampling [ref]. We could define formally this acquisition function:

$$\alpha_{\text{TS}}(\mathbf{x}, \mathcal{D}_n) = f^{(n)}(\mathbf{x}) \quad (3.10)$$

where $f^{(n)} \sim GP(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$ by spectral sampling. It has been shown, however, that this method tends to perform greedily on high-dimensional spaces[ref]. Another new approach is entropy-based [ref]. They aim to reduce the uncertainty in location \mathbf{x}^* by choosing points likely to reduce the entropy in $p(\mathbf{x}|\mathcal{D}_n)$. The acquisition function can be defined as:

$$\alpha_{\text{ES}}(\mathbf{x}|\mathcal{D}_n) = H(\mathbf{x}^*|\mathcal{D}_n) - \mathbb{E}_{y|\mathcal{D}_n, \mathbf{x}}[H(\mathbf{x}^*|\mathcal{D}_n \cup \{(\mathbf{x}, y)\})] \quad (3.11)$$

where H notes the differential entropy function of the posterior distribution. As with Thompson sampling, the function is not tractable in continuous spaces. Several studies have been done approximating this quantity, either by using simple Monte Carlo sampling [ref] or a space discretization of \mathcal{A} . A recent paper [ref] introduced *predictive entropy search* (PES), a method to remove the need for discretization by rewriting Equation 3.11 as:

$$\alpha_{\text{PES}}(\mathbf{x}, \mathcal{D}_n) = H(y|\mathcal{D}_n, \mathbf{x}) - \mathbb{E}_{\mathbf{x}^*|\mathcal{D}_n} [H(y|\mathcal{D}_n, \mathbf{x}, \mathbf{x}^*)] \quad (3.12)$$

The expectation is approximated in the original paper by Monte Carlo with Thompson samples, with simplifying assumptions. This is arguably the state of the art in acquisition functions, according to the results reported in [ref].

3.3.4 Acquisition function portfolios

In a *no free lunch* fashion, it can be shown that no acquisition function will outperform the others in every single problem. In fact, it has been proven [ref] that the acquisition function to provide optimal performance can change even in different points of the optimization procedure. It is natural, therefore, to consider an ensemble of acquisition functions and act upon it. In general, this implies optimizing all of these functions at each optimization step and then choosing among candidate points using a meta-criteria. This higher order criteria can be seen as a second level acquisition function.

Earlier approaches rely on modifications of the Hedge algorithm [ref], again inspired by the multi-armed bandit problem. It is basically based on measuring past performance of points proposed by the different acquisition functions to predict future performance (or gain), via another objective function. However, this strategy tends to undervalue exploration, which also provides valuable information on the target. Another more recent approach [ref] is called *Entropy Search Portfolio*, that considers candidates by weighing the gain of information towards the optimum. Formally it is defined as:

$$\alpha_{\text{ESP}}(\mathbf{x}, \mathcal{D}_n) = -\mathbb{E}_{y|\mathcal{D}_n, \mathbf{x}} [H[\mathbf{x}_*|\mathcal{D}_n \cup \{(\mathbf{x}, y)\}]] \quad (3.13)$$

and then we try to maximize over the candidates provided by the k based acquisition functions $\mathbf{x}_{1:K,n}$.

$$\mathbf{x}_n = \arg \max_{\mathbf{x}_{1:K,n}} \alpha_{\text{ESP}}(\mathbf{x}|\mathcal{D}_n) \quad (3.14)$$

In other words, this method chooses the candidate that is expected to reduce the most the entropy about the minimizer \mathbf{x}_* .

3.3.5 Visualizing the behaviour of an acquisition function

To demonstrate the behaviour of different acquisition functions on a step of Bayesian optimization, we will create a small script with our sine function example. This will help us understand visually the trade-off between exploration and exploitation in each case. The code provided below produces Figure 3.1.

```

1 import numpy as np
2 from GPRRegressor import GPRRegressor
3 from acquisition import Acquisition
4 from covfunc import squaredExponential
5 from GPGO import GPGO
6
7
8 def plotGPGO(gpgo, param, index, new = True):
9     param_value = list(param.values())[0][1]
10    x_test = np.linspace(param_value[0], param_value[1], 1000).reshape((1000, 1))
11    y_hat, y_var = gpgo.GP.predict(x_test, return_std = True)
12    std = np.sqrt(y_var)
13    l, u = y_hat - 1.96 * std, y_hat + 1.96 * std
14    if new:
15        plt.figure()
```

```

16         plt.subplot(5, 1, 1)
17         plt.fill_between(x_test.flatten(), l, u, alpha = 0.2)
18         plt.plot(x_test.flatten(), y_hat)
19     plt.subplot(5, 1, index)
20     a = np.array([-gpgo._acqWrapper(np.atleast_1d(x)) for x in x_test]).flatten()
21     plt.plot(x_test, a, color = colors[index - 2], label = acq_titles[index - 2])
22     gpgo._optimizeAcq(method = 'L-BFGS-B', n_start = 1000)
23     plt.axvline(x = gpgo.best)
24     plt.legend(loc = 0)
25
26
27
28 if __name__ == '__main__':
29     def f(x):
30         return(np.sin(x))
31
32     acq_1 = Acquisition(mode = 'ExpectedImprovement')
33     acq_2 = Acquisition(mode = 'ProbabilityImprovement')
34     acq_3 = Acquisition(mode = 'UCB', beta = 0.5)
35     acq_4 = Acquisition(mode = 'UCB', beta = 1.5)
36     acq_list = [acq_1, acq_2, acq_3, acq_4]
37     sexp = squaredExponential()
38     param = {'x': ('cont', [0, 2 * np.pi])}
39     new = True
40     colors = ['green', 'red', 'orange', 'black']
41     acq_titles = [r'Expected improvement', r'Probability of Improvement',
42                  r'GP-UCB $\beta = .5$', r'GP-UCB $\beta = 1.5$']
43
44     for index, acq in enumerate(acq_list):
45         np.random.seed(200)
46         gp = GPRegressor(sexp)
47         gpgo = GPGO(gp, acq, f, param)
48         gpgo._firstRun(n_eval = 3)
49         plotGPGO(gpgo, param, index = index + 2, new = new)
50         new = False
51
52     plt.show()

```

3.3.6 Why does Bayesian Optimization work?

In this small section, we will consider very briefly why the Bayesian Optimization procedure is efficiently making use of the information provided by the mean and variance posterior distribution of the GP. The explanation is mostly visual by means of Figure 3.2. The plot shows a Gaussian Process regression model, as well as its confidence interval, comprised from a lower confidence bound $L = \mu(\mathbf{x}) - q\sigma(\mathbf{x})$ and an upper confidence bound $U = \mu(\mathbf{x}) + q\sigma(\mathbf{x})$ for a quantile $q > 0$.

In practice, when we are interested in doing hyperparameter search in machine-learning, the algorithms are very blunt, in the sense that the most common strategies explore the entire space, either exhaustively or randomly. The acquisition functions defined in this chapter before use information from the Gaussian Process prior to explore efficiently the space. The specifics of each acquisition function have already been discussed, specially their exploration/exploitation balance, but they all share some common logic.

Figure 3.1: Acquisition function behaviour for Expected Improvement, Probability of Improvement, GP-UCB ($\beta = .5$) and GP-UCB($\beta = 1.5$) in the sine function example.

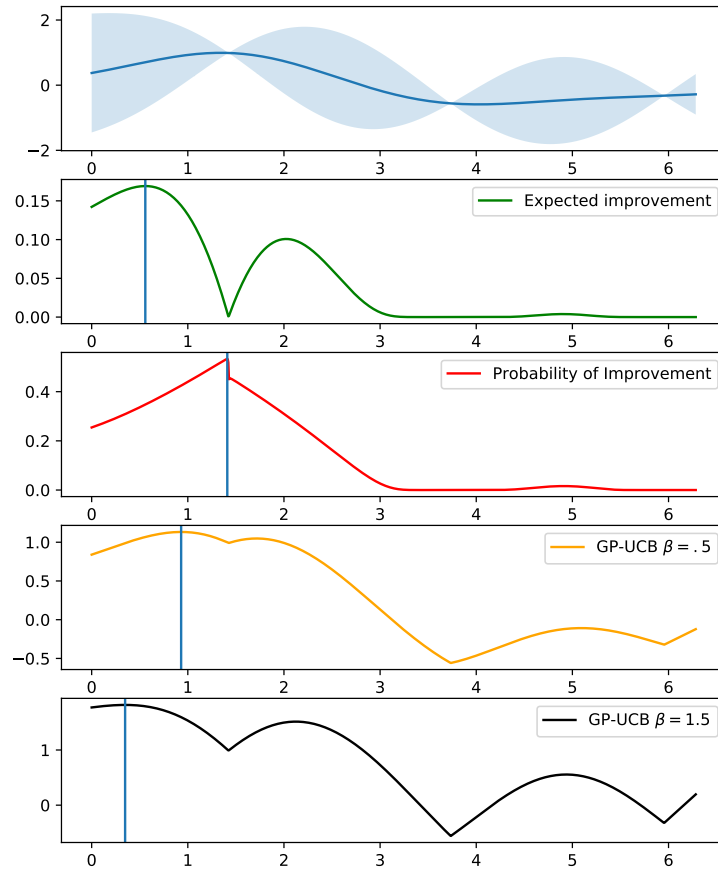
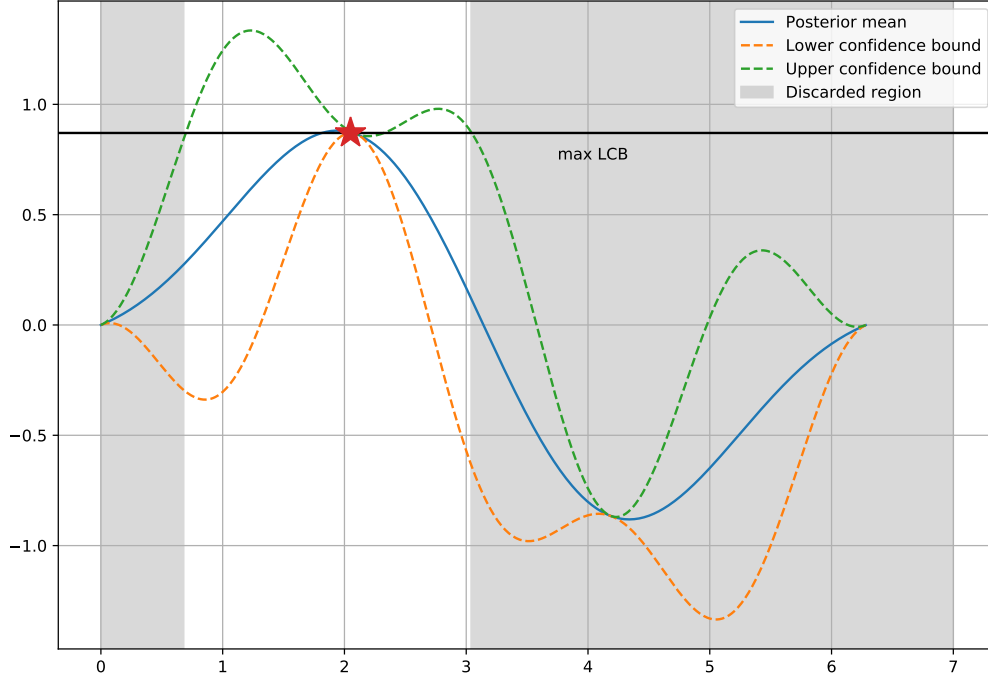


Figure 3.2: A visual explanation on why Bayesian optimization is efficient at exploring the space. It ignores all the input space where the UPB is lower than the point with maximum LCB.



Bayesian Optimization is efficient because it ignores all the space where the predicted upper confidence bound is lower than the maximum value of the lower confidence bound. It only uses the space that fulfills this criteria, according to the strategy selected by the chosen acquisition function. This remains a very reasonable assumption throughout the whole optimization procedure.

3.4 Role of GP hyperparameters in optimization

We already considered the role of hyperparameters in Gaussian Process regression in section 2.5. However, one may wonder how the estimation of this parameters, one way or another may affect the optimization procedure. So far, we have assumed that during the optimization procedure, parameters θ were given. Here we will consider two ways of handling hyperparameters during the optimization procedure, the one we presented, type II maximum likelihood estimation and approximate marginalization. For the moment, consider a generic acquisition function $\alpha : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$, where $\theta \in \Theta$ are our Gaussian Process hyperparameters. Naturally, one wishes to marginalize the uncertainty caused by θ with the following expression:

$$\alpha(\mathbf{x}) = \mathbb{E}_{\theta|\mathcal{D}_n} [\alpha(\mathbf{x}, \theta)] = \int_{\Theta} \alpha(\mathbf{x}|\theta) p(\theta|\mathcal{D}_n) d\theta. \quad (3.15)$$

The simplest way to do this is what we saw in 2.39, to optimize the marginal log-likelihood to obtain MAP estimates $\hat{\theta}_{\text{MAP}}$. Then, in each step of the optimization procedure, we simply maximize:

$$\hat{\alpha}(\mathbf{x}) = \alpha(\mathbf{x}, \hat{\theta}) \quad (3.16)$$

That is, we optimize the acquisition function defined by the *optimal* hyperparameters for the Gaussian Process determined in each step. Again, optimizing the marginal log-likelihood is a problem of its own, but it is common to use quasi-Newton methods such as L-BFGS-B methods.

A more Bayesian approach is to incorporate the uncertainty of θ into our model, since it may have an important role in guiding exploration. Point estimates are in a sense *winners* that may not capture the complexity of the response surface. The second approach we will be considering here, will be therefore to marginalize out hyperparameters using Markov Chain Monte Carlo (MCMC) sampling techniques. In practice, we will need to average M samples $\{\theta_n^{(i)}\}_{i=1}^M$ from the posterior distribution $p(\theta|\mathcal{D}_n)$.

$$\mathbb{E}_{\theta|\mathcal{D}_n} [\alpha(\mathbf{x}, \theta)] \approx \sum_{i=1}^M \alpha(\mathbf{x}, \theta_n^{(i)}) \quad (3.17)$$

Since it is not possible to have an analytical expression for the posterior distribution $p(\theta|\mathcal{D}_n)$, it is common to use MCMC techniques like Hamiltonian Monte Carlo [ref] to produce a sequence of samples whose stationary distribution is the posterior we are looking for. Once M valid samples are obtained, they are evaluated in the acquisition function and averaged. Quadrature methods can be used instead of MCMC techniques, yielding a weighted mixture:

$$\mathbb{E}_{\theta|\mathcal{D}_n} [\alpha(\mathbf{x}, \theta)] \approx \sum_{i=1}^M \omega_i \alpha(\mathbf{x}, \theta_n^{(i)}) \quad (3.18)$$

However, and to finish this section, in the problems that we will tackle in the experiments, it is usually a bad idea to estimate kernel hyperparameters. Estimating these hyperparameters with few function evaluations is a very challenging task, and can lead to disastrous results, as proven in [ref]. The marginal log-likelihood surface can easily fall into traps or be very flat, as seen by the (not cherry-picked) example in Figure 2.4. Even the more advanced MCMC or quadrature methods still suffer from this problem.

3.5 Optimizing the acquisition function

We have presented many acquisition functions in this chapter and provided a simple example to demonstrate basic functionality. However, so far, we have assumed that the acquisition function can be easily optimized. This is however, a problem of its own. The reader may be thinking that we have, in fact, changed one optimization problem (the one where we are interested in optimizing f) for another! (in which we now have to optimize α). This is technically true, but bear in mind that while f is very expensive to evaluate, α is very cheap, and it is reasonable to spend a bit more computational effort in evaluating α if it implies having to evaluate f less.

Maximizing α , however, is not an easy task. The acquisition function is often multi-modal and therefore non-convex, as it can be seen again in Figure 3.1. Theoretical convergence, furthermore, is only guaranteed when the optimal point \mathbf{x}^* in the acquisition function is found [ref]. At the end of the day, we encounter yet another global optimization problem that needs to be solved. From a practical point of view, there are many approaches the community has taken to solve this problem, from discretization [ref], adaptive-grids [ref]. If gradient information is available (rarely the case), a multi-start gradient ascent approach can be taken [ref]. Evolutionary approaches like CMA-ES can also be used [ref]. pyGPGO, in the GPGO module uses by default a multi-start quasi-Newton method (L-BFGS-B) to optimize the acquisition function, which in practice seems to work reasonably well. A CMA-ES optimizer interface is also available for the user.

Other methods have been proposed as alternatives to Bayesian Optimization in this aspect. [ref], sequentially building space-partitioning trees by splitting leaves with high function values or upper confidence bounds. This is called *Simultaneous Optimistic Optimization* (SOO) [ref]. Though these algorithms do not require any auxiliary information (smoothness) or optimization, it has been shown that they do not

perform quite as competitively as the Bayesian optimization framework, especially when prior knowledge is available.

3.6 Computational costs

Remember from chapter 2 that Gaussian Process regression has analytical expressions for the mean and variance of the posterior process. However, this exact inference is $O(n^3)$, where n is the number of training samples. This is caused by the inversion of K . pyGPGO uses a Cholesky decomposition that once computed, can reduce the cost of predicting to $O(n^2)$. If during the search procedure, however, we change K , for example, due to hyperparameter optimization, this $O(n^3)$ order is unavoidable for the traditional framework. Several approaches have been explored in the literature for approximating the output of the analytical solution. Mainly, there are two types of solutions, those that try to sparsify the process and those that use another type of surrogate model, such as Random Forests.

3.6.1 Approximations to the analytical GP. Alternative surrogates.

One of the first solutions is the *Sparse pseudo-input Gaussian Processes* (SPGP) [ref]. This is a straightforward approach to model large n using $m < n$ pseudo-inputs to reduce the rank of the covariance matrix to m . This method forces the interaction between the $\mathbf{x}_{1:n}$ data points and test points \mathbf{x}_* inducing m pseudo-inputs, achieving an approximate posterior in $O(nm^2 + m^3)$. Let \mathbf{f} and \mathbf{f}_* denote two sets of latent function (them being our training and testing points respectively). The assumption is that \mathbf{f} and \mathbf{f}_* are independent given a third set of variables \mathbf{u} , that is:

$$p(\mathbf{f}_*, \mathbf{f}) = \int p(\mathbf{f}_*, \mathbf{f}, \mathbf{u}) d\mathbf{u} \approx \int q(\mathbf{f}_*|\mathbf{u})q(\mathbf{f}|\mathbf{u})p(\mathbf{u}) d\mathbf{u} = q(\mathbf{f}, \mathbf{f}_*) \quad (3.19)$$

where \mathbf{u} is a vector representing the function values at the pseudo-inputs. Different pseudo-input Gaussian Process pseudo-input approximations specify their own form of $q(\mathbf{f}|\mathbf{u})$ and $q(\mathbf{f}_*|\mathbf{u})$ training and test conditionals. [ref]. How to choose the locations of these pseudo-inputs is another problem, is usually done by maximizing the marginal log-likelihood of the SPGP [ref]. Another approach uses variational inference [ref] to marginalize the pseudo-inputs to maximize the fidelity to the original Gaussian Process. It has been noted, however, that the computational savings of the pseudo-input methods impact heavily variance estimates. In the Bayesian optimization framework, this variance of the posterior predictive distribution is heavily used to guide exploration, so this behavior is undesirable.

Another approach is named *Sparse spectrum Gaussian Processes* (SSGP). Using the ideas of pseudo-input methods, these methods apply the same concept in kernel spectral space. [ref]. From basic spectral analysis, Bochner's theorem states that a stationary kernel $k(\mathbf{x}, \mathbf{x}_*)$ defines a positive finite Fourier spectrum $s(\boldsymbol{\omega})$:

$$k(\mathbf{x}) = \frac{1}{(2\pi)^d} \int e^{-i\boldsymbol{\omega}^T \mathbf{x}} s(\boldsymbol{\omega}) d\boldsymbol{\omega} \quad (3.20)$$

We can normalize this spectrum to make it a valid probability density function, such as $p(\boldsymbol{\omega}) = \frac{s(\boldsymbol{\omega})}{\nu}$. Now evaluating the kernel is the same as the expectation of the Fourier basis with respect to $p(\boldsymbol{\omega})$:

$$k(\mathbf{x}, \mathbf{x}_*) = \nu \mathbb{E}_{\boldsymbol{\omega}} \left[e^{-i\boldsymbol{\omega}^T (\mathbf{x} - \mathbf{x}_*)} \right] \quad (3.21)$$

Monte Carlo techniques can be used to approximate this expectation using m samples from the spectral density, so that:

$$k(\mathbf{x}, \mathbf{x}_*) \approx \frac{\nu}{m} \sum_{i=1}^m e^{-i\boldsymbol{\omega}_{(i)}^T \mathbf{x}} e^{-i\boldsymbol{\omega}_{(i)}^T \mathbf{x}_*} \quad (3.22)$$

where $\omega_{(i)} \sim s(\omega)/\nu$. The computational cost in this approach can be reduced to $O(nm^2 + m^3)$. It has been noted [ref] for this approximation method that whereas the uncertainty estimates are smoother than with the pseudo-input methods, observations away from the observed values exhibit irregular variance estimates. This again is undesirable behavior in the Bayesian Optimization framework.

As an alternative, several authors [ref] have suggested using different surrogate models in the Bayesian Optimization framework. Special attention has been drawn towards the Random Forest regression model, in the context of sequential model-based algorithm configuration. Random Forests are a very popular choice in the machine-learning community with very successful results in practice. They were introduced in 2001 by [ref], being an ensemble bagging tree-based method. Random Forests allow for training using sub-samples of data, giving it the ability to scale to large evaluation budgets, where exact analytical Gaussian Process regression would be infeasible in practice. The exploration strategy requires an uncertainty estimation for prediction at test points to apply the Bayesian optimization framework. The empirical variance in the predictions across trees was proposed as a substitute. This heuristic has been shown to work in practice [ref].

Random Forests are known to provide very good estimates when doing *interpolation* of seen data (or in its neighbourhood), they are very poor *extrapolators*. On points far from training data, the predictions across all trees in the model are very similar, providing poor predictions and more importantly, providing extremely confident (but erroneous) variance estimates. While Gaussian Processes are also terrible extrapolators, they produce reliable variance estimations far from training data, yielding better estimates for exploration and exploitation in the framework.

3.6.2 Parallelization

The Bayesian Optimization framework is inherently a sequential decision problem, where each candidate point to sample is selected after fitting a Gaussian Process to the currently available data. However, and given the parallel nature of current CPU architectures, significant speed-ups in clock time can be achieved if they are made good use of. That implies evaluating the acquisition function in parallel. While not entirely parallel in nature, some authors have proposed approaches based on the imputation of yet-to-run evaluations. Given $\mathcal{D}_n = \{(\mathbf{x}_n, y_n)\}$ known data and $\mathcal{D}_p = \{\mathbf{x}_p\}$ remaining to evaluate data, an idea is to impute the latter, $\hat{\mathcal{D}}_p = \{(\mathbf{x}_p, \hat{y}_p)\}$, and then use the typical Bayesian optimization framework on the augmented data $\mathcal{D}_n \cup \hat{\mathcal{D}}_p$.

Simple strategies have been proposed over time. The *constant liar* strategy proposes $\hat{y}_p = c, \forall p$, where $c \in \mathbb{R}$ is a predefined constant. Other strategies like the *kriging believer* uses the Gaussian Process posterior predictive mean instead: $\hat{y}_p = \mu_n(\mathbf{x}_p)$. More complex approaches have been proposed instead, for example [ref] proposed the use of s fantasies sampled from each unfinished experiment (out of a total of J) from the full GP posterior predictive distribution. Then they are averaged in a Monte Carlo fashion:

$$\alpha(x, \mathcal{D}_n, \mathcal{D}_p) = \int_{\mathbb{R}^J} \alpha(x, \mathcal{D}_n \cup \hat{\mathcal{D}}_p) P(y_{1:J}, \mathcal{D}_n) dy_{p,1:J} \quad (3.23)$$

$$\approx \frac{1}{S} \sum_{i=1}^S \alpha(x, \mathcal{D}_n \cup \hat{\mathcal{D}}_p^{(s)}) \quad (3.24)$$

and $\hat{\mathcal{D}}_p^{(s)} \sim P(\hat{y}_{1:J}, \mathcal{D}_n)$. [ref] has demonstrated empirically that this approach works reasonably well when α is chosen to be the Expected Improvement acquisition function. Other attempts [ref] have also been made using GP-UCB. Note that while these approaches are valid, they are not parallel per se. A true parallel approach to Bayesian optimization would propose simultaneously a set of candidates, or would use the information of posterior Gaussian Processes computed in another thread to make the sequential decision.

3.7 Step-by-step examples

We have explored the basics of Gaussian Process regression and Bayesian Optimization by now. To provide a more thorough understanding on the logic behind Algorithm 2, we will provide two different examples in this section. While not very complex, they serve illustrative purposes on how the optimization framework works. We hope this section will clear any practical consideration on how the procedure selects the next point to sample, graphically.

3.7.1 Optimizing (yet-again) the sine function

Again, it is no surprise for the reader that we choose the sine function as our first example to provide an intuition on how the Bayesian Optimization framework works step-by-step. We will be optimizing our function within the boundaries $x \in [0, 2\pi]$, where we know an optima is exactly at $x^* = \pi/2$. This example is particularly interesting to see how the next point is chosen balancing exploration and exploitation of the acquisition function. Now, as in most practical cases throughout this thesis, we will use the Expected Improvement acquisition function. The step-by-step optimization procedure for 6 complete epochs can be checked in Figure 3.3.

It can be appreciated that the chosen acquisition function leverages achieves a very reasonable compromise between mean and variance in the first stages of optimization, but exploits heavily in steps 3 and 4. It tries to explore again in step 5, and comes back to exploit at the end of the procedure.

3.7.2 Optimizing the Rössenbrock function

In this example, we will try to optimize the two-dimensional Rössenbrock function:

$$f(x, y|a, b) = (a - x)^2 + b(y - x^2)^2 \quad (3.25)$$

For given hyperparameters $a = 1$, $b = 100$ and in a bounding box defined by $x, y \in [-1, 1]$. A biplot of this function can be checked in Figure 3.4. This particular version of the Rössenbrock function has four different local optima, therefore it is multimodal. We would expect our GPGO algorithm to converge to one of these optima, while still exploring the space at earlier stages of optimization.

Given the low-dimensional space we are optimizing, we will plot the posterior mean and variance of the Gaussian Process, as well as the acquisition function. Again, the Expected Improvement acquisition function is used. The complete step-by-step optimization procedure for 6 epochs can be seen in Figure 3.5. For a more visually appealing explanation of the same procedure, there is a video in the linked GitHub repository of this thesis.

It can be seen that during the previous steps before fitting a Gaussian Process, that a random evaluation is reasonably close to one of the optima. This causes the acquisition function to exploit that area in the first step. In the second step, however, it starts exploring the area that is close to a second local optima, mostly because it is a large-variance area. It exploits in the third step, and keeps exploring the rest of the space during the following stages, therefore efficiently balancing exploitation and exploration as intended in the optimization procedure.

Figure 3.3: Six complete optimization epochs in the Bayesian Optimization framework for the target sine function in $x \in [0, 2\pi]$.

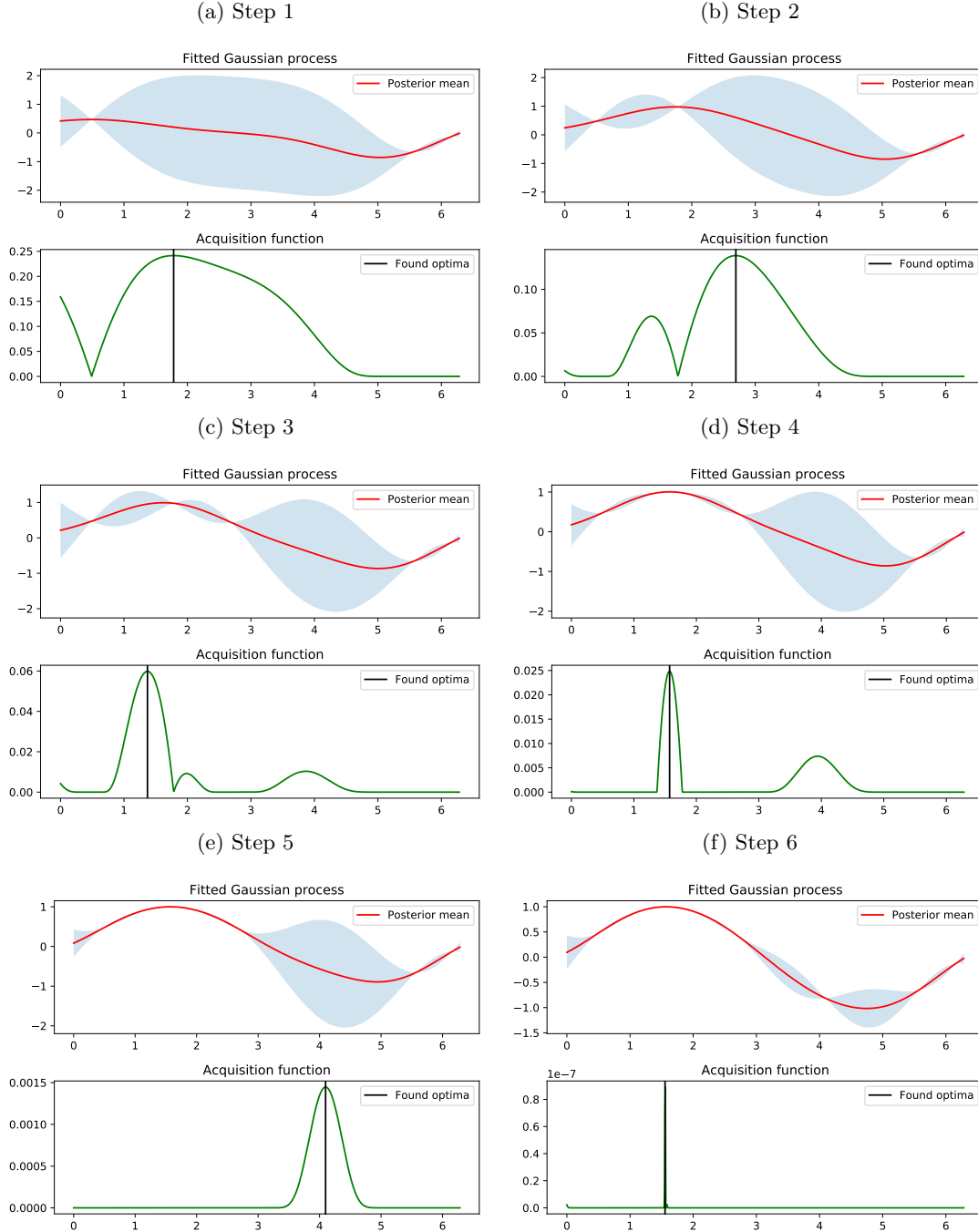


Figure 3.4: A 2D representation of the Rosenbrock function for $x, y \in [-1, 1]$.

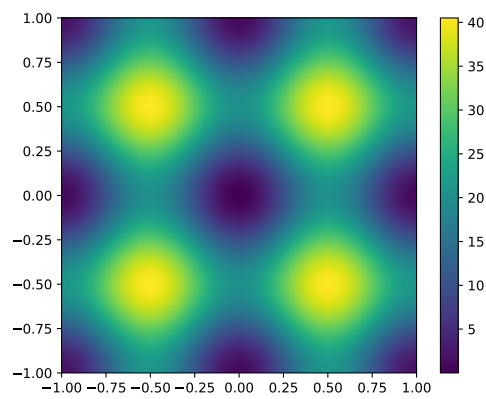
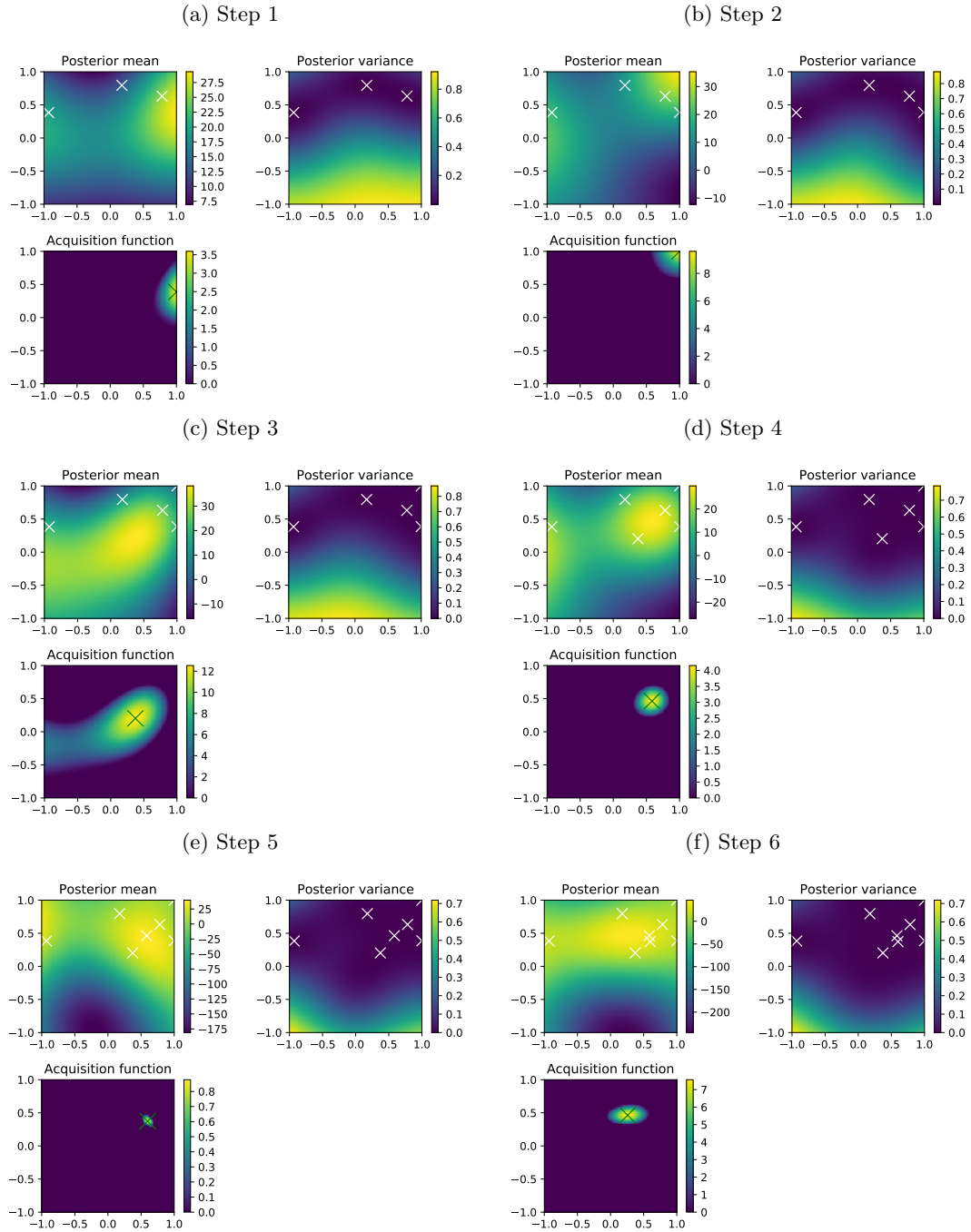


Figure 3.5: Six complete optimization epochs in the Bayesian Optimization framework for the target Rossenbrock 2D function.



Chapter 4

Experiments

In this chapter we will expose all the results using all implemented code for Gaussian Process regression and Bayesian Optimization in pyGPGO. First, we start by providing some benchmarking rules, how performance shall be evaluated across different models and datasets. In general, the models presented here hopefully span enough diversity: support vector machines, neural networks and bagging/boosting methods are among these candidates. Datasets come from very different research areas, spanning from biophysics to medicine.

4.1 Benchmarking rules

4.1.1 Other strategies for hyper-parameter optimization

Hyper-parameter optimization is usually the most tedious part in fitting a machine-learning algorithm. In practice, we are interested in models whose loss, evaluated on an independent part of data is as low as possible. Different hyper-parameter choices lead to different losses, therefore finding the optimal set is of importance. In reality, we will never know that the point found is the global optimum, but from a practical point of view, we are only interested in finding the model that works best in a production setting.

In the machine-learning community, hyper-parameter optimization is often overlooked, and in fact, some of the most famous models (e.g. Random Forests) have been shown to be somewhat insensitive to hyper-parameter optimization [ref]. Four common strategies for finding *better* hyper-parameters are presented here: grid search, random search, simulated annealing and gradient-based methods. Only random search and simulated annealing will be compared against our Bayesian Optimization algorithm, as implemented in pyGPGO. Details on how these work are laid down below:

- **Grid search.** Also known as parameter sweep. This is simply an exhaustive search in a user specified subset of parameter space. Search bounds have to be set manually. When multiple hyper-parameters have to be optimized over their sets, grid search considers the Cartesian product of all of them. This poses a problem when the number of hyper-parameters to optimize grows, also known as the curse of dimensionality. While suffering from this problem, grid search is still widely used for small to mid-sized problems, where function evaluations are not very expensive. Also, it is embarrassingly parallel: function evaluations can be distributed over in a simple way.
- **Random search.** Grid search is exhaustive since it considers all possible evaluations over a Cartesian product over parameter sets. Randomized search in hyper-parameter space is also a very popular method for the task at hand. In particular, this method has been shown to work considerably better in high-dimensional spaces than grid search [ref]. There is also evidence that often sometimes some hyper-parameters do not affect the loss significantly.
- **Simulated annealing.** While the other methods are search based, simulated annealing is a proper global optimization algorithm. This technique itself is a meta-heuristic to approximate the global minimum of a function f in a large search space. Its origins comes from physics, and it is in fact, an

adaptation of the Metropolis-Hastings algorithm [ref], a method to generate samples from a given thermodynamic system. Assume that our system is at a state s . Simulated annealing works by proposing a neighbouring candidate to evaluate s' , and decides whether to move towards it or not depending on a probabilistic criterion that is function of both s , s' , and an annealing temperature T . The initial annealing temperature T is cooled down according to some schema after each iteration of the algorithm. This also lowers the acceptance probability after each iteration.

In practice, how to choose to propose neighbouring candidates s' of s , the probabilistic criterion to accept the candidate and the temperature annealing scheme is left for the user to define. For all the experiments performed here the procedure was programmed as follows: the initial temperature T is set to 100. Neighbouring candidate points s' are randomly and uniformly sampled within hyper-parameter bounds, the acceptance probability is the standard Metropolis-Hastings acceptance criteria:

$$P(s, s', T) = \exp\left(-\frac{f(s') - f(s)}{T}\right) \quad (4.1)$$

For new states s' where $f(s') > f(s)$, $P(s, s', T)$ is automatically set to 1. This in practice means that if we find a state with a better value, we move to it, and if it is worse, we move to it with a certain probability given by Equation 4.1. We cool the temperature T in each iteration by a factor of 0.9. In particular, during the experiments, I tried different cooling schemas and found this one to yield the best results.

4.1.2 Evaluation metrics

Throughout all the experiments, we will be minimizing loss functions, in some form or the other. Every problem presented here is formulated either as a regression or classification task, therefore it is of importance to define early on what type of loss to use in each problem. For binary classification problems, we will use the logarithmic loss, also known as binary cross-entropy, defined by:

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{n} \sum_i (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)) \quad (4.2)$$

The log-loss is very popular among the deep learning folks [ref]. This particular metric does not only take into account whether a classifier makes the right decision given a threshold c , (like the 0/1 loss would), but also the confidence in predictions $\hat{\mathbf{y}}$. It is also very natural since it is just the negative log-likelihood of a Bernoulli random variable. The use of the logarithm both punishes erroneous extremely confident positive or negative predictions.

We generalize the previous metric for multi-class classification problems. The following expression is typically called categorical cross-entropy:

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{p}}) = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m y_{ij} \log(\hat{p}_{ij}) \quad (4.3)$$

where \hat{p}_{ij} is the predicted probability of a sample i belonging to class j , and m is the number of classes considered. For continuous regression problems, the loss that we use is the typical mean squared error, defined by:

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2 \quad (4.4)$$

We believe that the losses proposed here are very natural choices in both classification and regression problems.

While we have discussed the metrics to evaluate in each predictive problem, we still need to define how these losses will be evaluated in each step. To evaluate the performance of a given hyper-parameter optimization procedure, we have to balance both the performance in terms of loss and the number of evaluations necessary in order to get to a satisfactory solution. In practice, this is exactly how we will benchmark the different strategies, through a plot where the x -axis represent the number of function evaluations, and the y -axis the best log-loss found.

Evaluating the loss function itself can lead to multiple different values depending on the test values. One could choose an approach where this loss is evaluated on a single holdout test. This would lead to noisy estimates, however. We choose the more stable approach of performing a $k = 5$ cross-validation scheme to obtain a more reliable loss estimate. In practice, this means that we fit 5 models with the same architecture to different train/test splits and average the loss results in each.

4.1.3 Bayesian optimization setup

For all the tests performed with the Bayesian Optimization scheme, we choose the standard squared exponential kernel with a default characteristic length-scale $l = 1$. Different acquisition functions will be tested in each problem: Expected Improvement, GP-UCB ($\beta = .5$) and GP-UCB ($\beta = 1.5$).

4.1.4 Machine-learning models used

Since the shape of our objective function depends on both the dataset and the predictor, we try to span as many different types of machine-learning models as possible to provide the most extensive evaluation as possible. Except for rare occasions where we could not fit a model to a particular dataset for numerical conditions, all models are evaluated in all datasets with the same number of parameters and bounds to optimize over.

4.2 The binding affinity dataset

Chapter 5

pyGPGO: A simple Python Package for Bayesian Optimization

5.0.1 Installation

5.0.2 Usage

5.0.3 Examples

5.0.4 Features

5.0.5 Future work