

# Regret Bounds for Gaussian Process Bandits Without Observation Noise

by

Masrour Zoghi

B.Sc., The University of British Columbia, 2003

M.Sc., University of Toronto, 2004

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF  
THE REQUIREMENTS FOR THE DEGREE OF

MASTER OF SCIENCE

in

The Faculty of Graduate Studies

(Computer Science)

THE UNIVERSITY OF BRITISH COLUMBIA

(Vancouver)

August 2012

© Masrour Zoghi 2012

# Abstract

This thesis presents some statistical refinements of the bandits approach presented in [11] in the situation where there is no observation noise. We give an improved bound on the cumulative regret of the samples chosen by an algorithm that is related (though not identical) to the UCB algorithm of [11] in a complementary setting. Given a function  $f$  on a domain  $\mathcal{D} \subseteq \mathbb{R}^d$ , sampled from a Gaussian process with an anisotropic kernel that is four times differentiable at 0, and a lattice  $\mathcal{L} \subseteq \mathcal{D}$ , we show that if the points in  $\mathcal{L}$  are chosen for sampling using our branch-and-bound algorithm, the regret asymptotically decreases according to  $\mathcal{O}\left(e^{-\frac{\tau t}{(\ln t)^{d/4}}}\right)$  with high probability, where  $t$  is the number of observations carried out so far and  $\tau$  is a constant that depends on the objective function.

# Preface

This thesis grew out of a collaboration with Alex Smola and my supervisor, Nando de Freitas. An abridged version of the results presented here were presented in the Workshop on Bayesian Optimization at NIPS 2011. The results included in Section 3.1 and the Appendix were obtained by the other authors and the proofs are only included for the sake of completeness of the exposition. Everything else was proven and written by the author of this thesis.

# Table of Contents

- Abstract** . . . . . ii
- Preface** . . . . . iii
- Table of Contents** . . . . . iv
- List of Figures** . . . . . vi
- Acknowledgements** . . . . . vii
- Dedication** . . . . . viii
  
- 1 Introduction** . . . . . 1
  - 1.1 Global optimization . . . . . 2
  - 1.2 Costly evaluation . . . . . 4
  
- 2 Gaussian process bandits** . . . . . 7
  - 2.1 Gaussian processes . . . . . 7
  - 2.2 Surrogates for optimization . . . . . 8
  - 2.3 Our algorithm . . . . . 9
  
- 3 Analysis** . . . . . 11
  - 3.1 Approximation results . . . . . 11
  - 3.2 Finiteness of regret (statement) . . . . . 13
  - 3.3 Remarks on the main theorem . . . . . 16
    - 3.3.1 On the statement of Theorem 2 . . . . . 16
    - 3.3.2 Outline of the proof of Theorem 2 . . . . . 17

*Table of Contents*

---

3.3.3	Further remarks on the GP prior . . . . .	17
3.4	Finiteness of regret (proof) . . . . .	19
<b>4</b>	<b>Discussion</b> . . . . .	<b>27</b>
	<b>Bibliography</b> . . . . .	<b>28</b>
 <b>Appendices</b>		
	<b>Appendix: Auxilliary Lemmas</b> . . . . .	<b>30</b>

# List of Figures

1.1	Lipschitz hypothesis being used to discard pieces of the search space . . . . .	3
1.2	A sample application of the branch and bound algorithm . . . . .	5
2.1	Branch and Bound algorithm for a 2D function . . . . .	10
3.1	The elimination of other smaller peaks . . . . .	20
3.2	The shrinking of the relevant set . . . . .	24

# Acknowledgements

First and foremost, I would like to thank my supervisor for all the direction he provided throughout this thesis and for his patience when the proofs were not transpiring as effortlessly. I am also very grateful to Alex Smola for his collaboration and for providing me with an elegant proof of a key lemma when I was suffering through a much messier argument.

This work benefited greatly from many lengthy discussions with John Chia, who was working on the more experimental aspects of this project. I learned a lot from talking to both Matt Hoffman and Ziyu Wang, who were working on related topics, and who happened to sit at the two desks next to mine.

I would also like to thank our wonderful Group Assistant, Kath Imhira, for her incessant pleasantness, which is truly a scarce commodity.

Outside of the university, I would like to thank my many “sisters and brothers in arms” from the Mining Justice Alliance, who added some meaning to my life by allowing me to assist them in their efforts to shed some light on injustices emanating from Vancouver.

# Dedication

I would like to dedicate this thesis to my parents, without whose sacrifices over the last 33 years none of this would be possible, and to M.S. who contributed greatly to my intellectual cultivation during my short stay in Vancouver.



# Chapter 1

## Introduction

Let  $f : \mathcal{D} \rightarrow \mathbb{R}$  be a function on a compact subset  $\mathcal{D} \subseteq \mathbb{R}^d$ . We would like to address the global optimization problem

$$x_M = \operatorname{argmax}_{x \in \mathcal{D}} f(x).$$

Let us assume for the sake of simplicity that the objective function  $f$  has a unique global maximum (although it can have many local maxima)

The space  $\mathcal{D}$  might be the set of possible parameters that one could feed into a time-consuming algorithm or the locations where a sensor could be deployed, and the function  $f$  might be a measure of the performance of the algorithm (e.g. how long it takes to run) or whatever quantity it is that our given sensor measures (e.g. temperature). In this paper, our assumption is that once the function has been probed at point  $x \in \mathcal{D}$ , then the value  $f(x)$  can be observed with very high precision. This is the case when the deployed sensors are very accurate or if the algorithm whose parameters are being configured is deterministic.

The challenge in these applications is two-fold; the optimization problem is of a global nature without any convexity assumption on  $f$  and it is costly to evaluate  $f$  at any given point. We will discuss these two stumbling blocks below.

## 1.1 Global optimization

Global optimization is a difficult problem without any assumptions on the objective function  $f$ . The main complicating factor is the uncertainty over the extent of the variations of  $f$ , e.g. one could consider the characteristic function  $\chi_{\{x_M\}}$ , which is equal to 1 at  $x_M$  and 0 elsewhere, and none of the methods we mention here can optimize this function without exhaustively searching through every point in  $\mathcal{D}$ .

The way a large number of global optimization methods address this problem is by imposing some prior assumption on how fast the objective function  $f$  can vary. The most blatant manifestation of this remedy is the imposition of a Lipschitz assumption on the function  $f$ , which requires the change in the value of  $f(x)$ , as the point  $x$  moves around, to be smaller than a constant multiple of the distance traveled by  $x$  (cf. [6]). In fact, as pointed out in [3] (cf. Figure 3 therein), it is only important to have this kind of tight control over the function near its optimum: elsewhere in the space, we can have what they have dubbed a “weak Lipschitz” condition.

One way to relax these hard Lipschitz constraints is by putting a Gaussian Process (GP) prior on the function. Instead of restricting the function from oscillating too fast, a GP prior requires those fast oscillations to have low probability (cf. [5], Theorem 5). This is the setting of Bayesian optimization (aka GP bandits), with which we concern ourselves: a review of GPs and Bayesian optimization is provided in Section 2.

The main point of these derivative bounds (be it hard or soft, as outlined above) is to assist with the exploration-exploitation trade-off that global optimization algorithms have to grapple with. In the absence of any assumptions of convexity on the objective function, a global optimization algorithm is forced to explore enough until it reaches a point in the process when with some degree of certainty it can localize its search space and perform local optimization (aka exploitation). Derivative bounds such as the ones discussed here together with the boundedness of the search space, guaranteed by the compactness assumption on  $\mathcal{D}$ , provide us with such certainty by producing a useful upper bound that allows us to shrink the search space.

This is illustrated in Figure 1.1: suppose that we know that our function is Lipschitz with Lipschitz constant  $L$ , then given sample points as shown in the figure (the red dots), we can use the Lipschitz property to discard pieces of the search space (the shaded regions). This is done

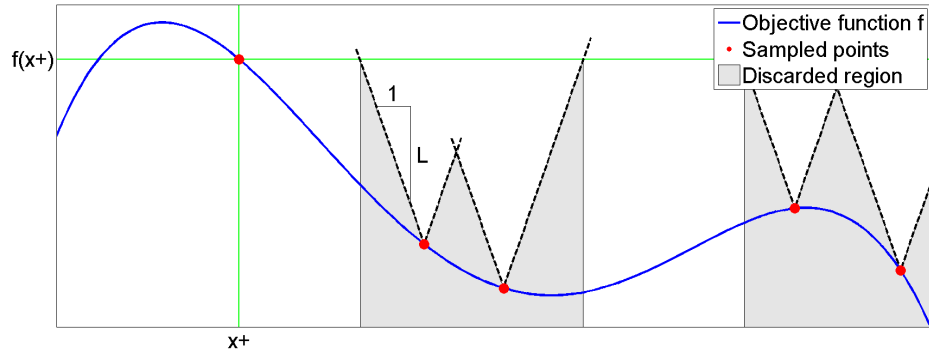


Figure 1.1: An example of the Lipschitz hypothesis being used to discard pieces of the search space.

by finding points in the search space where the function could not possibly be higher than the maximum value we've already observed without violating our Lipschitz hypothesis. Such points are found by placing cones at the sampled points with slope equal to  $L$  and seeing where those cones lie below the maximum observed value (the red dashed line in this case): since the function is guaranteed to be bounded from above by the cones, we know that there is nothing interesting to find in the shaded regions.

Note that this crude approach is often very wasteful because very often the slope of the function is much smaller than  $L$ , and as we will see below (cf. Figure 1.2), GPs do a better job of providing an upper bound that could be used to limit the search space, by essentially choosing Lipschitz constants that vary over space and time.

## 1.2 Costly evaluation

Now, let us also assume that the objective function  $f$  is costly to evaluate (e.g. time-wise or financially), in which case we would like to avoid probing  $f$  as much as possible and to get close to the optimum as quickly as possible. The solution to this problem is to approximate  $f$  with a *surrogate function* that provides a good upper bound for  $f$  and which is easier to calculate and optimize.

In the Lipschitz-based algorithms, this upper bound is provided by the Lipschitz constraint, whereas in the case of X-armed bandits, this bound is provided by the functions  $B_{h,i}(n)$  defined on page 1663 of [3].

Again, GPs provide a soft version of this by providing us with high probability bounds that could be modified depending on how certain we would like to be of our upper bounds. Moreover, the upper bounds provided by GPs have analytic expressions, which make evaluation and optimization relatively easy. We refer the reader to [2] for a general review of the literature on the various surrogate functions utilized in Bayesian optimization.

Let us also point out that in addition to being easier to evaluate, surrogate functions can also aid with global optimization by restricting the domain of interest. More precisely, if we know that  $f \leq g$  for some known function  $g$  and we have probed  $f$  at  $x_0$  to get the value  $y_0 = f(x_0)$ , then we can shrink our search space to the *relevant set*

$$\mathcal{R} = \{x \in \mathcal{D} | g(x) > y_0\},$$

which will be substantially smaller than  $\mathcal{D}$ , supposing that  $g$  is a relatively tight upper bound of  $f$  near the optimum: it can be atrociously inaccurate away from the optimum as long as it stays below  $\max f$ .

The surrogate function that we will make extensive use of here is called the Upper Confidence Bound (UCB), which is defined to be  $\mu + B\sigma$ , where  $\mu$  and  $\sigma$  are the posterior predictive mean and standard deviation of the GP and  $B$  is a constant to be chosen by the algorithm. The UCB surrogate function has the desirable property that it provides the tightest bound on  $f$  where the function has been probed the most: this is illustrated in Figure 1.2.

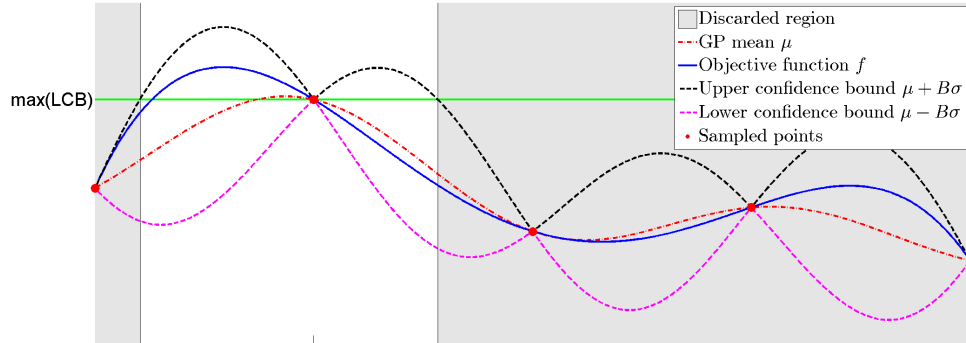


Figure 1.2: An example of our branch and bound maximization algorithm with UCB surrogate  $\mu + B\sigma$ , where  $\mu$  and  $\sigma$  are the mean and standard deviation of the GP respectively. The region consisting of the points  $x$  for which the upper confidence bound  $\mu(x) + B\sigma(x)$  is lower than the maximum value of the lower confidence bound  $\mu(x) - B\sigma(x)$  does not need to be sampled anymore. Note that the UCB surrogate function bounds  $f$  from above.

This surrogate function has been studied extensively in the literature and this paper relies heavily on the ideas put forth in the paper by Srinivas et al [11], in which the algorithm consists of repeated optimization of the UCB surrogate function after each sample.

One key difference between our setting and that of [11] is that, whereas we assume that the value of the function can be observed exactly, in [11] it is necessary for the noise to be non-trivial (and Gaussian) because the main quantity that is used in their estimates, namely information gain (cf. Equation (3) in [11]), becomes undefined when the variance of the observation noise ( $\sigma^2$  in their notation) is set to 0: cf. the expression for  $I(\mathbf{y}_A; \mathbf{f}_A)$  that was given in the paragraph following Equation (3).

Moreover, our algorithm can in some ways be thought of as a limiting case of the UCB algorithm in [11] when the parameter  $B$  goes to infinity because our main goal is to shrink  $\sigma$  as much as possible, even though we do use  $\mu + B\sigma$  to limit our search space. In subsequent work, we hope to prove that the algorithm presented here is inferior to the UCB algorithm, but in the mean time we thought it valuable to present this manifestly wasteful algorithm and show that the cumulative regret  $R_T$  one obtains from applying it belongs to  $O(1)$ , given a fixed function, in contrast to the previous bounds that grow as  $O(\sqrt{T})$ , where  $T$  is the number of samples collected. In fact, we show that the regret  $r_T$  decreases according to  $\mathcal{O}\left(e^{-\frac{\tau t}{(\ln t)^{d/4}}}\right)$ .

This asymptotic result is valuable to applications that are meant to run indefinitely because it unburdens the user from having to decide when to stop the bandit algorithm and just exploit the

optimum point found so far because the overall regret the algorithm can accrue is bounded from above.

The paper whose results are most similar to ours is [9], but there are some key differences in the methodology, analysis and obtained rates. For instance, we are interested in cumulative regret, whereas the results of [9] are proven for finite stop-time regret. In our case, the ideal application is the optimization of a function that is  $C^2$ -smooth and has an unknown non-singular Hessian at the maximum. We obtain a regret rate  $\mathcal{O}\left(e^{-\frac{\tau t}{(\ln t)^{d/4}}}\right)$ , whereas the DOO algorithm in [9] has regret rate  $\mathcal{O}(e^{-t})$  if the Hessian is known and the SOO algorithm has regret rate  $\mathcal{O}(e^{-\sqrt{t}})$  if the Hessian is unknown. In addition, the algorithms in [9] can handle functions that behave like  $-c\|x - x_M\|^\alpha$  near the maximum (cf. Example 2 therein). This problem was also studied by [14] and [4], but using the Expected Improvement surrogate instead of UCB. Our methodology and results are different, but complementary to theirs.

# Chapter 2

## Gaussian process bandits

### 2.1 Gaussian processes

As in [11], the objective function is distributed according to a Gaussian process prior:

$$f(x) \sim \text{GP}(m(\cdot), \kappa(\cdot, \cdot)). \quad (2.1)$$

For convenience, and without loss of generality, we assume that the prior mean vanishes, i.e.,  $m(\cdot) = 0$ . There are many possible choices for the covariance kernel. One obvious choice is the anisotropic kernel  $\kappa$  with a vector of known hyperparameters [10]:

$$\kappa(x_i, x_j) = \tilde{\kappa} \left( -(x_i - x_j)^\top \mathbf{D} (x_i - x_j) \right), \quad (2.2)$$

where  $\tilde{\kappa}$  is an isotropic kernel and  $\mathbf{D}$  is a diagonal matrix with positive hyperparameters along the diagonal and zeros elsewhere. Our results apply to squared exponential kernels and Matérn kernels with parameter  $\nu \geq 2$ . In this paper, we assume that the hyperparameters are fixed and known in advance.

We can sample the GP at  $t$  points by choosing points  $\mathbf{x}_{1:t} := \{x_1, \dots, x_t\}$  and sampling the values of the function at these points to produce the vector  $\mathbf{f}_{1:t} = [f(x_1) \cdots f(x_t)]^\top$ . The function values are distributed according to a multivariate Gaussian distribution  $\mathcal{N}(0, \mathbf{K})$ , with covariance entries  $\kappa(x_i, x_j)$ . Assume that we already have several observations from previous steps, and that we want to decide what action  $x_{t+1}$  should be considered next. Let us denote the value of the function at this arbitrary new point as  $f_{t+1}$ . Then, by the properties of GPs,  $\mathbf{f}_{1:t}$  and  $f_{t+1}$  are

jointly Gaussian:

$$\begin{bmatrix} \mathbf{f}_{1:t} \\ f_{t+1} \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \mathbf{K} & \mathbf{k}^\top \\ \mathbf{k} & \kappa(x_{t+1}, x_{t+1}) \end{bmatrix} \right),$$

where  $\mathbf{k} = [\kappa(x_{t+1}, x_1) \cdots \kappa(x_{t+1}, x_t)]^\top$ . Using the Schur complement, one arrives at an expression for the posterior predictive distribution:

$$P(f_{t+1} | \mathbf{x}_{1:t+1}, \mathbf{f}_{1:t}) = \mathcal{N}(\mu_t(x_{t+1}), \sigma_t^2(x_{t+1})),$$

where

$$\begin{aligned} \mu_t(x_{t+1}) &= \mathbf{k}^\top \mathbf{K}^{-1} \mathbf{f}_{1:t}, \\ \sigma_t^2(x_{t+1}) &= \kappa(x_{t+1}, x_{t+1}) - \mathbf{k}^\top \mathbf{K}^{-1} \mathbf{k} \end{aligned} \tag{2.3}$$

and  $\mathbf{f}_{1:t} = [f(x_1) \cdots f(x_t)]^\top$ .

## 2.2 Surrogates for optimization

When it is assumed that the objective function  $f$  is sampled from a GP, one can use a combination of the posterior predictive mean and variance given by Equations (2.3) to construct surrogate functions, which tell us where to sample next. Here we use the UCB combination, which is given by

$$\mu_t(x) + B_t \sigma_t(x),$$

where  $\{B_t\}_{t=1}^\infty$  is a sequence of numbers specified by the algorithm. This surrogate trades-off exploration and exploitation since it is optimized by choosing points where the mean is high (exploitation) and where the variance is large (exploration). Since the surrogate has an analytical expression that is easy to evaluate, it is much easier to optimize than the original objective function. Other popular surrogate functions constructed using the sufficient statistics of the GP include the Probability of Improvement, Expected Improvement and Thompson sampling. We refer the reader to [2, 7, 8] for details on these.



## 2.3 Our algorithm

The main idea of our algorithm (Algorithm 1) is to tighten the bound on  $f$  given by the UCB surrogate function by sampling the search space more and more densely and shrinking this space as more and more of the UCB surrogate function is “submerged” under the maximum of the Lower Confidence Bound (LCB). Figure 1.2 illustrates this intuition.

---

### Algorithm 1 Branch and Bound

---

Input: A compact subset  $\mathcal{D} \subseteq \mathbb{R}^d$ , a discrete lattice  $\mathcal{L} \subseteq \mathcal{D}$  and a function  $f : \mathcal{D} \rightarrow \mathbb{R}$ .

$\mathcal{R} \leftarrow \mathcal{D}$

$\delta \leftarrow 1$

**repeat**

**Sample Twice as Densely:**

- $\delta \leftarrow \frac{\delta}{2}$

- Sample  $f$  at enough points in  $\mathcal{L}$  so that every point in  $\mathcal{R}$  is contained in a simplex of size  $\delta$ .

**Shrink the Relevant Region:**

- Set

$$\tilde{\mathcal{R}} := \left\{ x \in \mathcal{R} \mid \mu_T(x) + \sqrt{\beta_T} \sigma_T(x) > \sup_{\mathcal{R}} \mu_T(x) - \sqrt{\beta_T} \sigma_T(x) \right\}.$$

$T$  is the number points sampled so far and  $\beta_T = 2 \ln \left( \frac{|\mathcal{L}|T^2}{\alpha} \right) = 4 \ln T + 2 \ln \frac{|\mathcal{L}|}{\alpha}$  with  $\alpha \in (0, 1)$ .

- Solve the following constrained optimization problem:

$$(x_1^*, x_2^*) = \underset{(x_1, x_2) \in \tilde{\mathcal{R}} \times \tilde{\mathcal{R}}}{\operatorname{argsup}} \|x_1 - x_2\|$$

- $\mathcal{R} \leftarrow B \left( \frac{x_1^* + x_2^*}{2}, \|x_1^* - x_2^*\| \right)$ , where  $B(p, r)$  is the ball of radius  $r$  centred around  $p$ .

**until**  $\mathcal{R} \cap \mathcal{L} = \emptyset$

---

More specifically, the algorithm consists of two iterative stages. During the first stage, the function is sampled along a lattice of points (the red crosses in Figure 2.1). In the second stage, the search space is shrunk to discard regions where the maximum is very unlikely to reside. Such regions are obtained by finding points where the UCB is lower than the LCB (the complement of the colored region in the same panel as before). The remaining set of relevant points is denoted by  $\tilde{\mathcal{R}}$ . In order to simplify the task of shrinking the search space, we simply find an enclosing ball, which is denoted by  $\mathcal{R}$  in Algorithm 1. Back to the first stage, we consider a lattice that is twice as dense as in the first stage of the previous iteration, but we only sample at points that lie within

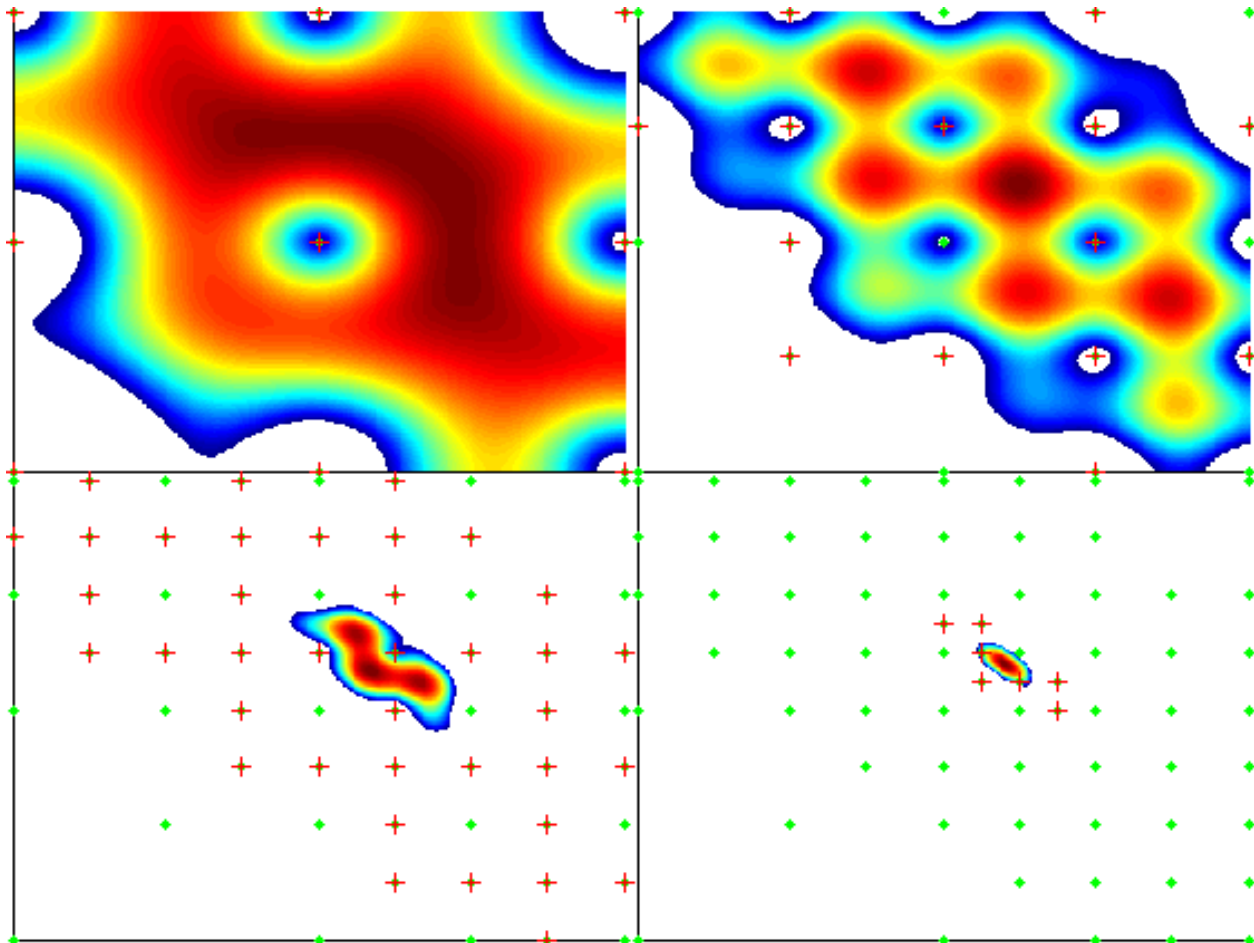


Figure 2.1: *Branch and Bound* algorithm for a 2D function. The colored region is the search space and the color-map, with red high and blue low, illustrates the value of the UCB. Four steps of the algorithm are shown; progressing from left to right and top to bottom. The green dots designate the points where the function was sampled in the previous steps, while the red crosses denote the freshly sampled points.

our new smaller search space.

In the second stage, the auxiliary step of approximating the relevant set  $\tilde{\mathcal{R}}$  with the ball  $\mathcal{R}$  introduces inefficiencies in the algorithm, since we only need to sample inside  $\tilde{\mathcal{R}}$ . This can be easily remedied in practice to obtain an efficient algorithm, for instance by using an ellipsoid instead of a ball. Our analysis will show that even without these improvements it is already possible to obtain very strong exponential convergence rates. Of course, practical improvements will result in better constants and ought to be considered seriously.

# Chapter 3

## Analysis

### 3.1 Approximation results

**Lemma 1 (Variance Bound)** *Let  $\kappa : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be a kernel that is four times differentiable along the diagonal  $\{(x, x) \mid x \in \mathbb{R}^d\}$ , with  $Q$  defined as in Lemma 7.2, and  $f \sim \text{GP}(0, \kappa(\cdot, \cdot))$  a sample from the corresponding Gaussian Process. If  $f$  is sampled at points  $x_{1:T} = \{x_1, \dots, x_T\}$  that form a  $\delta$ -cover of a subset  $\mathcal{D} \subseteq \mathbb{R}^d$ , then the resulting posterior predictive standard deviation  $\sigma_T$  satisfies*

$$\sup_{\mathcal{D}} \sigma_T \leq \frac{Q\delta^2}{4}$$

**Proof** Let  $\mathcal{H}$  be the RKHS corresponding to  $\kappa$  and  $h \in \mathcal{H}$  an arbitrary element, with  $g := (\mathbf{1} - P_{1:T})h$  the residual defined in Lemma 7.5. By Lemma 7.3, we know that  $\|\mathbf{1} - P_{1:T}\| \leq 1$  and so we have

$$\|g\|_{\mathcal{H}} \leq \|\mathbf{1} - P_{1:T}\| \|h\|_{\mathcal{H}} \leq \|h\|_{\mathcal{H}} \quad (3.1)$$

Moreover, by Lemma 7.2, we know that the second derivative of  $g$  is bounded by  $\|g\|_{\mathcal{H}} Q$ , and since by Lemma 7.5 we know that  $g$  vanishes at each  $x_i$ , we can use Lemma 9.2 and the inequality given by inequality (3.1) to conclude that

$$\begin{aligned} |h(x) - P_{1:T}h(x)| &:= |g(x)| \\ &\leq \frac{\|g\|_{\mathcal{H}} Q\delta^2}{4} \text{ by Lemma 9.2} \\ &\leq \frac{\|h\|_{\mathcal{H}} Q\delta^2}{4} \text{ by inequality (3.1)} \end{aligned}$$

and so for all  $x \in \mathcal{D}$  we have

$$|h(x) - P_{1:T}h(x)| \leq \frac{Q\delta^2}{4} \|h\|_{\mathcal{H}} \quad (3.2)$$

On the other hand, by Lemma 8, we know that for all  $x \in \mathcal{D}$  we have the following tight bound:

$$|h(x) - P_{1,T}h(x)| \leq \sigma_T(x) \|h\|_{\mathcal{H}}. \quad (3.3)$$

Now, given the fact that both inequalities (3.2) and (3.3) are bounding the same quantity and that the latter is a tight estimate, we necessarily have that

$$\sigma_T(x) \|h\|_{\mathcal{H}} \leq \frac{Q\delta^2}{4} \|h\|_{\mathcal{H}}.$$

Canceling  $\|h\|_{\mathcal{H}}$  gives the desired result. ■

## 3.2 Finiteness of regret (statement)

Having shown that the variance vanishes according to the square of the resolution of the lattice of sampled points, we now move on to show that this estimate implies an exponential asymptotic vanishing of the regret encountered by our Branch and Bound algorithm. This is laid out in our main theorem stated below and proven in the supplementary material.

The theorem considers a function  $f$ , which is a sample from a GP with a kernel that is four times differentiable along its diagonal. The global maximum of  $f$  can appear in the interior of the search space, with the function being twice differentiable at the maximum and with non-vanishing curvature. Alternatively, the maximum can appear on the boundary with the function having non-vanishing gradient at the maximum. Given a lattice that is fine enough, the theorem asserts that the regret asymptotically decreases in exponential fashion.

The main idea of the proof of this theorem is to use the bound on  $\sigma$  given by Proposition 1 to reduce the size of the search space. The key assumption about the function that the proof utilizes is the quadratic upper bound on the objective function  $f$  near its global maximum, which together with Proposition 1 allows us to shrink the relevant region  $\mathcal{R}$  in Algorithm 1 rapidly. The figures in the proof give a picture of this idea. The only complicating factor is the factor  $\sqrt{\beta_t}$  in the expression for the UCB that needs to be estimated. This is dealt with by modeling the growth in the number of points sampled in each iteration with a difference equation and finding an approximate solution of that equation.

Recall that  $\mathcal{D} \subseteq \mathbb{R}^d$  is assumed to be a non-empty compact subset and  $f$  a sample from the Gaussian Process  $\text{GP}(0, \kappa(\cdot, \cdot))$  on  $\mathcal{D}$ . Moreover, in what follows we will use the notation  $x_M := \underset{x \in \mathcal{D}}{\operatorname{argmax}} f(x)$ . Also, by convention, for any set  $\mathcal{S}$ , we will denote its interior by  $\mathcal{S}^\circ$ , its boundary by  $\partial\mathcal{S}$  and if  $S$  is a subset of  $\mathbb{R}^d$ , then  $\operatorname{conv}(S)$  will denote its convex hull. The following holds true:

**Theorem 2** *Suppose we are given:*

1.  $\alpha > 0$ , a compact subset  $\mathcal{D} \subseteq \mathbb{R}^d$ , and  $\kappa$  a stationary kernel on  $\mathbb{R}^d$  that is four times differentiable;
2.  $f \sim \text{GP}(0, \kappa)$  a continuous sample on  $\mathcal{D}$  that has a unique global maximum  $x_M$ , which satisfies

one of the following two conditions:

(†)  $x_M \in \mathcal{D}^\circ$  and  $f(x_M) - c_1\|x - x_M\|^2 < f(x) \leq f(x_M) - c_2\|x - x_M\|^2$  for all  $x$  satisfying  $x \in B(x_M, \rho_0)$  for some  $\rho_0 > 0$ ;

(‡)  $x_M \in \partial\mathcal{D}$  and both  $f$  and  $\partial\mathcal{D}$  are smooth at  $x_M$ , with  $\nabla f(x_M) \neq 0$ ;

3. any lattice  $\mathcal{L} \subseteq \mathcal{D}$  satisfying the following two conditions

$$\bullet \quad 2\mathcal{L} \cap \text{conv}(\mathcal{L}) \subseteq \mathcal{L} \tag{3.4}$$

$$\bullet \quad 2^{\lceil -\log_2 \frac{\rho_0}{\text{diam}(\mathcal{D})} \rceil + 1} \mathcal{L} \cap \mathcal{L} \neq \emptyset \tag{3.5}$$

if  $f$  satisfies (†)

Then, there exist positive numbers  $A$  and  $\tau$  and an integer  $T$  such that the points specified by the Branch and Bound algorithm,  $\{x_t\}$ , will satisfy the following asymptotic bound: For all  $t > T$ , with probability  $1 - \alpha$  we have

$$r(x_t) < Ae^{-\frac{\tau t}{(\ln t)^{d/4}}}.$$

We would like to make a few clarifying remarks about the theorem. First, note that for a random sample  $f \sim \text{GP}(0, \kappa)$  one of conditions (†) and (‡) will be satisfied almost surely if  $\kappa$  is a Matérn kernel with  $\nu > 2$  and the squared exponential kernel because the sample  $f$  is twice differentiable almost surely by [1, Theorem 1.4.2] and [12, §2.6]) and the vanishing of at least one of the eigenvalues of the Hessian is a co-dimension 1 condition in the space of all functions that are smooth at a given point, so it has zero chance of happening at the global maximum. Second, the two conditions (3.4) and (3.5) simply require that the lattice be “divisible by 2” and that it be fine enough so that the algorithm can sample inside the ball  $B(x_M, \rho_0)$  when the maximum of the function is located in the interior of the search space  $\mathcal{D}$ . Finally, it is important to point out that the rate decay  $\tau$  does not depend on the choice of the lattice  $\mathcal{L}$ , even though as stated, the statement of the theorem chooses  $\tau$  only after  $\mathcal{L}$  is specified. The theorem was written this way simply for the sake of readability.

Given the exponential rate of convergence we obtain in Theorem 2, we have the following finiteness conclusion for the cumulative regret accrued by our Branch and Bound algorithm:

**Corollary 3** *Given  $\kappa$ ,  $f \sim \text{GP}(0, \kappa)$  and  $\mathcal{L} \subseteq \mathcal{D}$  as in Theorem 2, the cumulative regret is bounded from above.*

**Remark 4** *It is worth pointing out the trivial observation that using a simple UCB algorithm with monotonically increasing and unbounded factor  $\sqrt{\beta_t}$ , without any shrinking of the search space as we do here, necessarily leads to unbounded cumulative regret since eventually  $\sqrt{\beta_t}$  becomes large enough so that at points  $x'$  far away from the maximum,  $\sqrt{\beta_t}\sigma_t(x')$  becomes larger than  $f(x_M) - f(x)$ . In fact, eventually the UCB algorithm will sample every point in the lattice  $\mathcal{L}$ .*

### 3.3 Remarks on the main theorem

This section includes a discussion of the assumptions placed on the objective function in Theorem 2 as well as an outline of the proof, the full details of which are included in the appendix.

#### 3.3.1 On the statement of Theorem 2

A few remarks on the assumptions and the conclusion of the main theorem are in order:

**A. Relationship between the local and global assumptions on  $f$ :** The theorem has two seemingly unrelated restrictions on the function  $f$ : the global GP prior and the local behaviour near the global maximum. However, in many circumstances of interest, the local condition follows almost surely from the global condition. Two such circumstances are if  $\kappa$  is a Matérn kernel with  $\nu > 2$  (including the squared exponential kernel) or if  $\kappa$  is six times differentiable. In either case, the sample  $f$  is twice differentiable almost surely, in the former case by [1, Theorem 1.4.2] and [12, §2.6]) and in the latter situation by [5, Theorem 5]. If the global maximum  $x_M$  lies in the interior of  $\mathcal{D}$ , the Hessian of  $f$  at  $x_M$  will almost surely be non-singular since the vanishing of at least one of the eigenvalues of the Hessian is a co-dimension 1 condition in the space of all functions that are smooth at a given point, hence justifying condition (†).

On the other hand, if  $x_M$  lies on the boundary of  $\mathcal{D}$ , then condition (‡) will be satisfied almost surely, since the additional event of the vanishing of  $\nabla f(x_M)$  is a codimension  $d$  phenomenon in the space of functions with global maximum at  $x_M$ .

**B. Uniqueness of the global maximum:** A randomly drawn continuous sample from a GP on a compact domain will almost surely have a unique global maximum: this is because the space of continuous functions on a compact domain that attain their global maximum at more than one point have codimension one in the space of all continuous functions on that domain.

**C. Assumptions on  $\mathcal{L}$ :** The two conditions (3.4) and (3.5) simply require that the lattice be “divisible by 2” and that it be fine enough so that the algorithm can sample inside the ball  $B(x_M, \rho_0)$  when the maximum of the function is located in the interior of the search space  $\mathcal{D}$ . One can simply choose  $\mathcal{L}$  to be the set of points in  $\mathcal{D}$  that have floating point coordinates: it’s just the points at which the algorithm is allowed to sample the function.

**D. On  $\tau$ ’s dependence:** Finally, it is important to point out that the decay rate  $\tau$  does not



depend on the choice of the lattice  $\mathcal{L}$ , even though as stated, the statement of the theorem chooses  $\tau$  only after  $\mathcal{L}$  is specified. The theorem was written this way simply for the sake of readability.

### 3.3.2 Outline of the proof of Theorem 2

The starting point for the proof is the observation that one can use the posterior predictive mean and standard deviation of the GP to obtain a high probability envelope around the objective function (cf. Lemma 10 in the appendix). Given the fact that the thickness of this envelope is determined by the height of the posterior predictive standard deviation,  $\sigma$ , we can use the bound given by Proposition 1 to show that asymptotically one can rapidly dispense with large portions of the search space, as illustrated in Figure 1.2.

One disconcerting component of Algorithm 1 is the step that requires sampling twice as densely in each iteration, since the number of samples can start to grow exponentially, hence killing any hope of obtaining exponentially decreasing regret. However, this is where the assumption on the local behaviour near the global maximum becomes relevant. Since Proposition 1 tells us that every time the function is sampled twice as densely,  $\sigma$  decreases by a factor of 4, and given our assumption that the function has quadratic behaviour near the global maximum, we can conclude that the radius of the search space is halved after each iteration and so the number of sampled points added in each iteration roughly remains constant. Of course, this assumes that the multiplicative factor  $\sqrt{\beta_t}$  remains constant in this process. However, the algorithm requires  $\sqrt{\beta_t}$  to grow logarithmically, and so to fill this gap we need to bound the growth of  $\sqrt{\beta_t}$ , which is tied to the number of samples needed in each iteration of the algorithm, which in turn is linked to the resolution of the lattice of sampled points  $\delta$  and the size of the relevant set  $\mathcal{R}$ , which in turn depends on the size of  $\sqrt{\beta_t}\sigma_t$ . This circular dependence gives rise to a difference equation, whose solutions we bound by solving the corresponding differential equation.

### 3.3.3 Further remarks on the GP prior

Let us step back for a moment and pose the question of whether it would be possible to carry out a similar line of reasoning under other circumstances. To answer this, one needs to identify the key ingredients of the proof, which are the following:

- A. A mechanism for calculating a high probability envelope around the objective function (cf. Lemma 10);
- B. An estimate showing that the thickness of the envelope diminishes rapidly as the function is sampled more and more densely (cf. Proposition 1), so that the search space can be shrunk under reasonable assumptions on the behaviour of the function near the peak.

The reason for our imposing a GP prior on  $f$  is that it gives us property A, while our smoothness assumption on the kernel guarantees property B. However, GPs are but one way one could obtain these properties and they do this essentially by coming up with local estimates of the Lipschitz constant based on the observed values of the objective function nearby. Perhaps one could explicitly incorporate similar local estimates on the Lipschitz constant into tree based approaches like HOO and SOO, cf. [3] and [9], in which case one would be able to dispense with the GP assumption and get similar performance. But, that is beyond the scope of this paper and will be left for future work.

### 3.4 Finiteness of regret (proof)

We will need the following two definitions in the proof of the theorem:

**Definition 5** *Given the above setup, the **regret** function is defined to be*

$$r(x) = \max_{\mathcal{D}} f - f(x).$$

**Definition 6 (Covering Number)** *Denote by  $\mathcal{B}$  a Banach space with norm  $\|\cdot\|$ . Furthermore denote by  $B \subseteq \mathcal{B}$  a set in this space. Then the covering number  $n_\epsilon(B, \mathcal{B})$  is defined as the minimum number of  $\epsilon$  balls with respect to the Banach space norm that are required to cover  $B$  entirely.*

**Proof** [Theorem 2] The proof consists of the following steps:

- **Global:** We first show that after a finite number of steps the algorithm zooms in on the neighbourhood  $B(x_M, \rho_0)$ . This is done by first showing that  $\epsilon$  can be chosen small enough to squeeze the set  $f^{-1}((f_M - \epsilon, f_M])$  into any arbitrarily small neighbourhood of  $x_M$  and that as the function is sampled more and more densely, the UBC-LCB envelope around  $f$  becomes arbitrarily tight, hence eventually fitting the relevant set inside a small neighbourhood of  $x_M$ . Please refer to Figure 3.1 for a graphical depiction of this process.

$G_I$ : Since  $\mathcal{D}$  is compact and  $f$  is continuous and has a unique maximum, for every  $\rho > 0$ , we can find an  $\epsilon = \epsilon(\rho) > 0$  such that

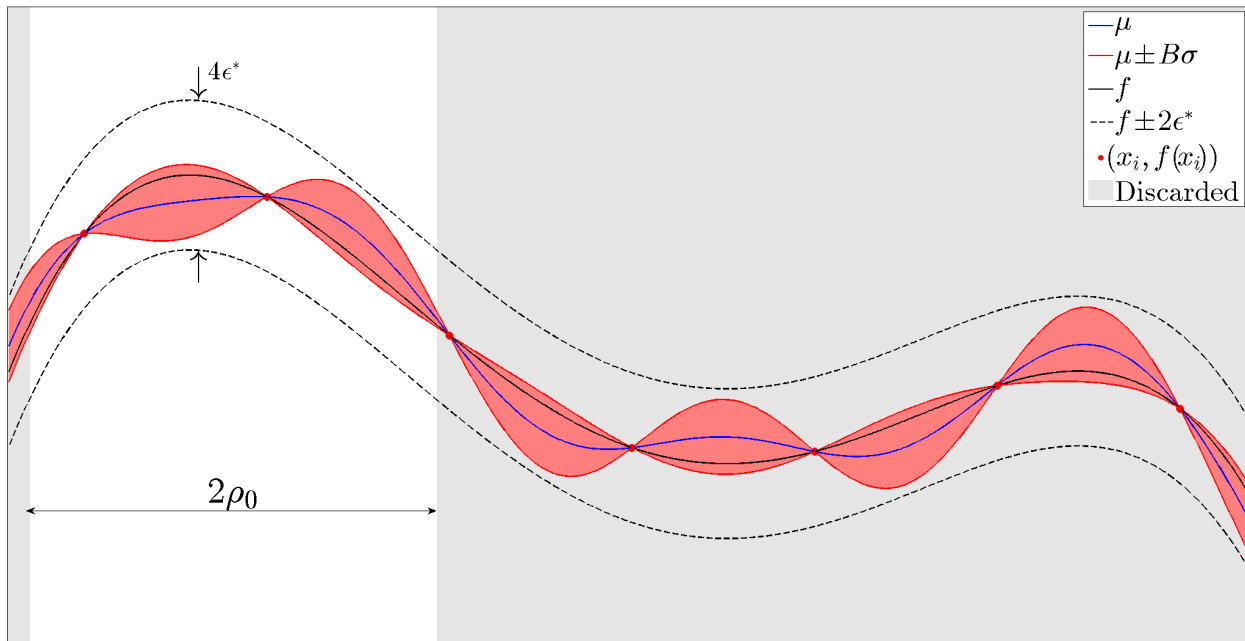
$$f^{-1}((f_M - \epsilon, f_M]) \subseteq B(x_M, \rho),$$

where  $f_M = \max f$ .

To see this, suppose on the contrary that there exists a radius  $\rho > 0$  such that for all  $\epsilon > 0$  we have

$$f^{-1}((f_M - \epsilon, f_M]) \not\subseteq B(x_M, \rho)$$

which means that there exists a point  $x \in \mathcal{D}$  such that  $f(x_M) - f(x) < \epsilon$  but  $\|x - x_M\| > \rho$ . Now, for each  $i \in \mathbb{N}$ , pick a point  $x^i \in f^{-1}((f_M - \frac{1}{i}, f_M]) \setminus B(x_M, \rho)$ : this gives us a sequence of points  $\{x^i\}$  in  $\mathcal{D}$ , which by the compactness of  $\mathcal{D}$  has a convergent


 Figure 3.1: *The elimination of other smaller peaks.*

subsequence  $\{x^{i_k}\}$ , whose limit we will denote by  $x^*$ . From the continuity of  $f$  and the fact that  $f(x_M) - f(x^i) < \frac{1}{i}$ , we can conclude that  $f(x_M) - f(x^*) = 0$ , which contradicts our assumption that  $f$  has a unique global maximum since we necessarily have  $x^* \notin B(x_M, \rho)$ .

**G<sub>II</sub>**: Define  $\epsilon^* := \frac{\epsilon(\rho_0)}{4}$ , with  $\rho_0$  as in Condition (†) of the statement of Theorem 2.

**G<sub>III</sub>**: For each  $T$ , define the “relevant set”  $\mathcal{R}_T \subseteq \mathcal{D}$  as follows:

$$\mathcal{R}_T = \left\{ x \in \mathcal{D} \mid \mu_T(x) + \sqrt{\beta_T} \sigma_T(x) > \sup_{\mathcal{R}} \mu_T(x) - \sqrt{\beta_T} \sigma_T(x) \right\}.$$

**G<sub>IV</sub>**: Choose  $\beta_T = b \ln(T)$ , with  $b$  chosen large enough to satisfy the conditions of Lemma 10.

Then, it is possible to sample  $f$  densely enough so that

$$\sqrt{\beta_T} \max_{x \in \mathcal{D}} \sigma_T(x) < \epsilon^*, \quad (3.6)$$

so that  $\mathcal{R}_T \subseteq B(x_M, \rho_0)$ . This is because as  $\mathcal{D}$  is sampled more and more densely we have  $\sigma = O(\delta^2)$ , where  $\delta$  is the distance between the points of the grid, and  $\beta = O(\ln \frac{1}{\delta^d}) = O(-\ln \delta)$  and so  $\sqrt{\beta} \sigma \rightarrow 0$  as  $\delta \rightarrow 0$ , and so there exists a  $\delta_0$  small enough so that

a lattice of resolution  $\delta_0$  would give us the bound given in inequality (3.6). The end point of this process is depicted in Figure 3.1, where the relevant set  $\mathcal{R}_T$  lies inside the non-shaded region: the reason for this inclusion and “thickness”  $4\epsilon^*$  is described below, in Step L<sub>1</sub> of the proof: cf. Equation (3.7).

- **Local:** Once the algorithm has localized attention to a neighbourhood of  $x_M$ , then we can show that the regret decreases exponentially; to do so, we will proceed by sampling the relevant set twice as densely and shrinking the relevant set, and repeating these two steps. The claim is that in each iteration, the maximum regret goes down exponentially and the number of the new points that are sampled in each **refining** iteration is asymptotically constant. To prove this, we will write down the equations governing the behaviour of the number of sampled points and  $\sigma$ . We will adopt the following notation to carry out this task:

- $\delta_\ell$  - the resolution of the lattice of sampled points at the end of the  $(\ell + 1)^{th}$  refining iteration inside  $\mathcal{R}_{\ell+1}$  (defined below).
- $\epsilon_\ell = \sup_{x \in \mathcal{R}_\ell} \sigma_{N_\ell}(x)$  at the end of the  $\ell^{th}$  iteration. Note that  $\epsilon_\ell \propto \delta_\ell^2$ . Also, note that  $\epsilon_0 \leq \epsilon^*$  by the choice of  $\delta_0$ .
- $N_\ell$  - number of points that have been sampled by the end of the  $\ell^{th}$  iteration.
- $\Delta N_\ell = N_{\ell+1} - N_\ell$ .
- $\mathcal{R}_\ell$  - the relevant set at the beginning of the  $\ell^{th}$  iteration. Note that  $\mathcal{R}_1 \subseteq B(x_M, \rho_0)$ .
- $\rho_\ell = \frac{\text{diam}(\mathcal{R}_\ell)}{2}$ . Note that  $\rho_1 < \rho_0$ .

L<sub>1</sub>: We have the following chain of inequalities:

$$\begin{aligned}
 N_1 &\leq N_0 + W n_{\delta_0} \left( \mathcal{R}_0, (\mathbb{R}^d, \|\cdot\|_2) \right) && \text{where } n_{\delta_0} \left( \mathcal{R}_0, (\mathbb{R}^d, \|\cdot\|_2) \right) \text{ is the } \delta_0\text{-covering number} \\
 &&& \text{as defined in Definition 6} \\
 &\leq N_0 + W \mathcal{N}(\rho_0, \delta_0) && \text{where } \mathcal{N}(\rho, \delta) := n_\delta \left( B(0, \rho), (\mathbb{R}^d, \|\cdot\|_2) \right) \\
 &\leq N_0 + W \mathcal{N} \left( \sqrt{\frac{4\epsilon_0 \sqrt{\beta_{N_0}}}{c_2}}, \delta_0 \right) \\
 &\leq N_0 + W \mathcal{N} \left( \sqrt{\frac{4\epsilon_0 \sqrt{b \ln N_0}}{c_2}}, \delta_0 \right) \\
 &= N_0 + W \mathcal{N} \left( c \sqrt{\epsilon_0} \sqrt[4]{\ln N_0}, \delta_0 \right) && \text{where } c := \sqrt{\frac{4\sqrt{b}}{c_2}}
 \end{aligned}$$

In the first line of the above chain of inequalities, the factor  $W$  multiplying the last term is any number greater than one that gives an upper bound on the algorithm's wastefulness when it comes to producing a  $\delta$ -covering of any set  $\mathcal{S} \subseteq \mathbb{R}^d$  as compared to the minimum number of points necessary, i.e.  $n_\delta(\mathcal{S}, (\mathbb{R}^d, \|\cdot\|_2))$ .

The expression  $\sqrt{\frac{4\epsilon_0 \sqrt{\beta_{N_0}}}{c_2}}$  comes about as follows: using the notations  $B = \sqrt{\beta_{N_0}}$  and  $\sigma = \sigma_{N_0}$  we know by Lemma 10 that  $f$  and  $\mu$  are intertwined with each other in the sense that both of the following chains of inequality hold:

$$\begin{aligned}
 \mu - B\sigma &\leq f \leq \mu + B\sigma \\
 f - B\sigma &\leq \mu \leq f + B\sigma,
 \end{aligned}$$

which, combined together, give us the following chain of inequalities

$$f - 2B\sigma \leq \mu - B\sigma \leq f \leq \mu + B\sigma \leq f + 2B\sigma. \quad (3.7)$$

Since, we also know that  $\sigma(x) \leq \epsilon_0$  for all  $x \in \mathcal{R}_0$ , we can conclude that

$$f - 2B\epsilon_0 \leq \mu - B\sigma \leq \mu + B\sigma \leq f + 2B\epsilon_0.$$

Moreover, if condition  $(\dagger)$  holds, we know that in  $\mathcal{R}_0$ , the function  $f$  satisfies  $-c_1\mathbf{r}^2 < f(x) - f(x_M) < -c_2\mathbf{r}^2$ , where  $\mathbf{r} = \mathbf{r}(x) := \|x - x_M\|$ , so we get that

$$f(x_M) - c_1\mathbf{r}^2 - 2B\epsilon_0 \leq \mu - B\sigma \leq \mu + B\sigma \leq f(x_M) - c_2\mathbf{r}^2 + 2B\epsilon_0.$$

Now, recall that  $\mathcal{R}_0$  is defined to consist of points  $x$  where  $\mu(x) + B\sigma(x) \geq \sup_{\mathcal{D}} \mu(x) - B\sigma(x)$ , but given the fact that we have the above outer envelope for  $\mu \pm B\sigma$ , we can conclude that

$$\begin{aligned} \mathcal{R}_0 &\subseteq \left\{ x \mid f(x_M) - c_2\mathbf{r}(x)^2 + 2B\epsilon_0 \geq \max f(x_M) - c_1\mathbf{r}(x)^2 - 2B\epsilon_0 \right\} \\ &= \left\{ x \mid f(x_M) - c_2\mathbf{r}(x)^2 + 2B\epsilon_0 \geq f(x_M) - 2B\epsilon_0 \right\} \\ &= \left\{ x \mid -c_2\mathbf{r}(x)^2 + 2B\epsilon_0 \geq -2B\epsilon_0 \right\} \\ &= \left\{ x \mid c_2\mathbf{r}(x)^2 \leq 4B\epsilon_0 \right\} \\ &= \left\{ x \mid \mathbf{r}(x) \leq \sqrt{\frac{4B\epsilon_0}{c_2}} \right\} \end{aligned}$$

Now, if, on the other hand,  $f$  satisfies condition  $(\ddagger)$ , then by the smoothness assumptions in  $(\ddagger)$ , we know that  $\nabla f(x_M)$  is perpendicular to  $\partial\mathcal{D}$  at  $x_M$  and so there exist positive numbers  $c_1$  and  $c_2$  such that in a neighbourhood of  $x_M$  we have

$$-c_1\mathbf{r} \leq f - f(x_M) \leq -c_2\mathbf{r}^2.$$

Note that in the argument above in the case of  $(\dagger)$ , the precise form of the lower bound on  $f$  was irrelevant, since all we are interested in is its maximum. So, the same argument goes through again.

This is depicted in Figure 3.2, where  $B := \sqrt{\beta_{N_0}} = \sqrt{b \ln N_0}$ .

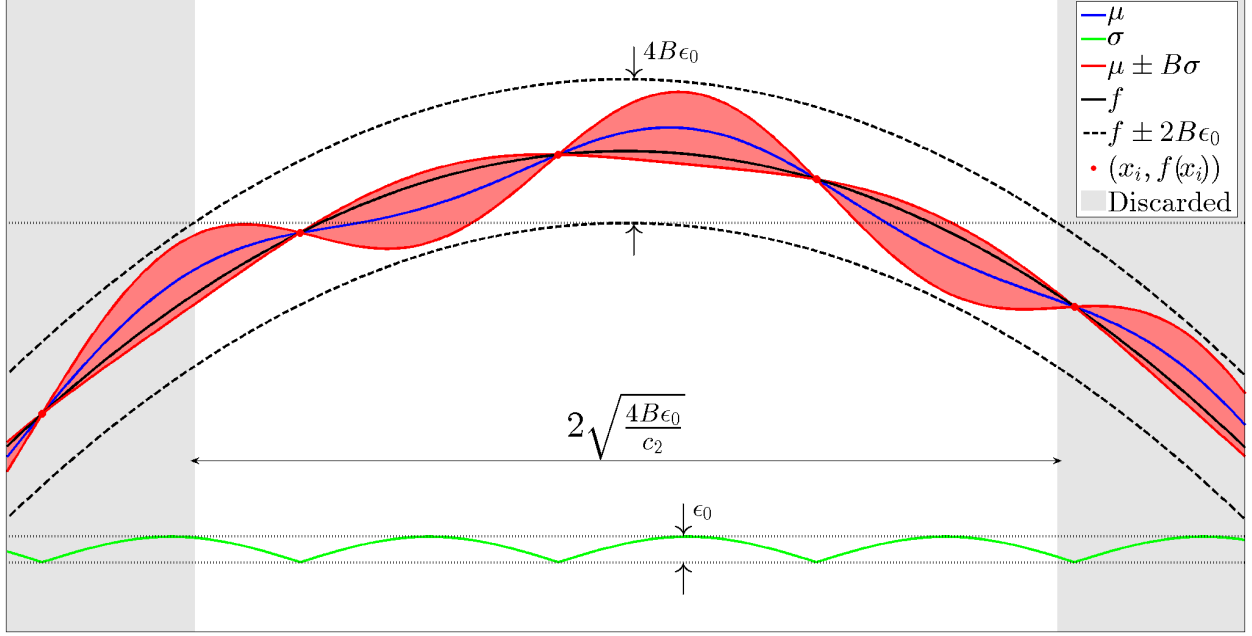


Figure 3.2: *The shrinking of the relevant set  $\mathcal{R}_\ell$ . Here,  $B = \sqrt{\beta_{N_0}}$*

$L_{\ell+1}$ : Now, let us suppose that we are the end of the  $\ell^{\text{th}}$  iteration. We have

$$\begin{aligned}
 N_{\ell+1} &\leq N_\ell + W\mathcal{N}(\rho_\ell, \delta_\ell) \\
 &= N_\ell + W\mathcal{N}\left(c\sqrt{\epsilon_\ell} \sqrt[4]{\ln N_\ell}, \delta_\ell\right) \\
 &\leq N_\ell + W\mathcal{N}\left(c\sqrt{\frac{\epsilon_0}{4^\ell}} \sqrt[4]{\ln N_\ell}, \frac{\delta_0}{2^\ell}\right) && \text{by Proposition 1} \\
 &= N_\ell + W\mathcal{N}\left(c\sqrt{\epsilon_0} \sqrt[4]{\ln N_\ell}, \delta_0\right) && \text{since } \mathcal{N}(2\rho, 2\delta) = \mathcal{N}(\rho, \delta) \text{ for any } \rho \text{ and } \delta \\
 &:= N_\ell + Wn_{\delta_0}\left(B(0, c\sqrt{\epsilon_0} \sqrt[4]{\ln N_\ell}), (\mathbb{R}^d, \|\cdot\|_2)\right) && \text{by the definition of } \mathcal{N} \\
 &\leq N_\ell + Wn_{\delta_0}\left(\left[-c\sqrt{\epsilon_0} \sqrt[4]{\ln N_\ell}, c\sqrt{\epsilon_0} \sqrt[4]{\ln N_\ell}\right]^d, (\mathbb{R}^d, \|\cdot\|_2)\right) \\
 &\leq N_\ell + W\left(\frac{2c\sqrt{\epsilon_0} \sqrt[4]{\ln N_\ell}}{\delta_0}\right)^d && \text{since a regular lattice of resolution } \delta_0 \text{ gives a } \delta_0\text{-covering} \\
 &\leq N_\ell + C(\ln N_\ell)^{\frac{d}{4}} && \text{where } C = W\left(\frac{2c\sqrt{\epsilon_0}}{\delta_0}\right)^d
 \end{aligned}$$

So, the number of samples needed by the branch and bound algorithm is governed by the difference inequation

$$\Delta N_\ell \leq C(\ln N_\ell)^{\frac{d}{4}}. \quad (3.8)$$



To study the solutions of this difference equation, we consider the corresponding differential equation:

$$\frac{dN}{d\ell} = C(\ln N)^{\frac{d}{4}}. \quad (3.9)$$

Since this equation is separable, we can write

$$\frac{dN}{(\ln N)^{\frac{d}{4}}} = C d\ell.$$

Now, letting  $\ell = L$  be a given number of iterations in the algorithm and  $N(L)$  the corresponding number of sampled points, we can integrate both sides of the above equation to get

$$\int_{N(0)}^{N(L)} \frac{dN}{(\ln N)^{\frac{d}{4}}} = \int_0^L C d\ell = CL.$$

Given the fact that the integral on the left can't be solved analytically, we will use the lower bound

$$\frac{N(L) - N(0)}{(\ln N(L))^{\frac{d}{4}}} \leq \int_{N(0)}^{N(L)} \frac{dN}{(\ln N)^{\frac{d}{4}}}$$

to get

$$\frac{N(L) - N(0)}{C(\ln N(L))^{\frac{d}{4}}} \leq L \quad (3.10)$$

Given a time  $t$ , we will denote by  $\ell_t$  the largest non-negative integer such that  $N_{\ell_t} < t$  or 0 if no such number exists. We illustrate this somewhat obtuse definition with the following example:



Now, by Lemma 11, for all  $t \gg N_0$  we have

$$\begin{aligned}
 r_t &\leq 2\sqrt{\beta_t} \max_{\mathcal{R}_{\ell_t}} \sigma_t \leq 2\sqrt{b \ln t} \epsilon_{\ell_t} \leq \frac{2\epsilon_0 \sqrt{b \ln t}}{4^{\ell_t}} \leq \frac{8\epsilon_0 \sqrt{b \ln t}}{4^{\ell_t+1}} \\
 &\leq 8\epsilon_0 \sqrt{b \ln t} \left(\frac{1}{4}\right)^{\frac{N_{\ell_t+1}-N_0}{C(\ln N_{\ell_t+1})^{d/4}}} \quad \text{by Equation 3.10} \\
 &\leq 8\epsilon_0 \sqrt{b \ln t} \left(\frac{1}{4}\right)^{\frac{DN_{\ell_t+1}}{(\ln N_{\ell_t+1})^{d/4}}} \quad \text{for some } D > 0 \text{ since } N_{\ell_t+1} > N_0 \\
 &\leq 8\epsilon_0 \sqrt{b \ln t} \left(\frac{1}{4}\right)^{\frac{Dt}{(\ln t)^{d/4}}} \quad \text{for } t \text{ satisfying } \ln t > \frac{d}{4} \text{ (see } \star \text{ below) since } t \leq N_{\ell_t+1} \\
 &\leq 8\epsilon_0 \sqrt{b} e^{-\frac{Et}{(\ln t)^{d/4}} + \frac{\ln \ln t}{2}} \\
 &\leq 8\epsilon_0 \sqrt{b} e^{-\frac{Et}{(\ln t)^{d/4}} + \frac{Et}{2(\ln t)^{d/4}}} \quad \text{for large enough } t \\
 &= A e^{-\frac{\tau t}{(\ln t)^{d/4}}} \quad \text{for } A = 8\epsilon_0 \sqrt{b} \text{ and } \tau = E/2.
 \end{aligned}$$

★ The reason for the specific criterion  $\ln t > \frac{d}{4}$  is that the function  $\frac{x}{(\ln x)^{d/4}}$  is increasing when this condition is satisfied, and so decreasing  $x$  from  $N_{\ell_t} + 1$  to  $t$  decreases its value, increasing the overall expression  $\left(\frac{1}{4}\right)^{\frac{x}{(\ln x)^{d/4}}}$ . To see that  $\frac{x}{(\ln x)^{d/4}}$  becomes increasing when  $\ln x > \frac{d}{4}$ , we simply need to calculate its derivative:

$$\begin{aligned}
 \frac{d}{dx} \frac{x}{(\ln x)^{d/4}} &= \frac{1}{(\ln x)^{d/4}} - \frac{d}{4} \frac{x}{x(\ln x)^{d/4+1}} \\
 &= \frac{\ln x - \frac{d}{4}}{(\ln x)^{d/4}}.
 \end{aligned}$$

Moreover, since  $N_{\ell_t+1} \geq t$ , if the derivative of  $\frac{x}{(\ln x)^{d/4}}$  is positive at  $t$ , it is also positive between  $t$  and  $N_{\ell_t+1}$  and so the function is indeed increasing in that interval. ■

# Chapter 4

## Discussion

In this thesis, we have put forth a modification of the UCB algorithm presented in [11], which has the additional step of discarding pieces of the search space where the global maximum is very unlikely to be. We show that this additional step leads to a considerable improvement of the regret accrued by the algorithm: in particular, the cumulative regret obtained by our Branch and Bound (BB) algorithm is bounded from above, whereas the cumulative regret bounds obtained in [11] are unbounded. This is because UCB’s regret is  $\mathcal{O}\left(\frac{1}{\sqrt{t}}\right)$  up to logarithmic factors, whereas for BB, the regret is  $\mathcal{O}\left(e^{-\frac{\tau t}{(\ln t)^{d/4}}}\right)$ .

The main realization that we had in the process of writing this paper was that what affects the rate of convergence the most is the possibility of dispensing with chunks of the search space. This observation can also be seen in the works involving hierarchical partitioning, e.g. [9], where regions of the space are deemed as less worthy of probing as time goes on.

The most obvious next step is to carry out similar analysis for the case when there is non-trivial observation noise. Another natural question is whether or not, in each iteration of the Branch and Bound algorithm, one could dispense with points that lie outside the current (or perhaps the previous) relevant set when calculating the covariance matrix  $\mathbf{K}$  in the posterior equations (2.3), and how much of an effect such a computational cost-cutting measure would have on the regret encountered by the algorithm.

Another obvious question is how does the performance of the Branch and Bound algorithm presented here compares with a modified version of UCB algorithm in which in addition to optimizing the UCB surrogate function, one introduces an occasional “cleansing” step of getting rid of regions in the search space where the maximum is very unlikely to be.

Finally, we believe that similar regret bounds should hold for contextual bandit problems, using similar methods, although the unpredictability of the context introduces new difficulties.

# Bibliography

- [1] Robert J. Adler and Jonathan E. Taylor. *Random Fields and Geometry*. Springer, 2007.
- [2] Eric Brochu, Vlad M Cora, and Nando de Freitas. A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. Technical Report TR-2009-023, arXiv:1012.2599v1, University of British Columbia, Department of Computer Science, 2009.
- [3] Sébastien Bubeck, Rémi Munos, Gilles Stoltz, and Csaba Szepesvari. X-armed bandits. *Journal of Machine Learning Research*, 12:1655–1695, 2011.
- [4] Adam D. Bull. Convergence rates of efficient global optimization algorithms. *Journal of Machine Learning Research*, 12:2879–2904, 2011.
- [5] Subhashis Ghosal and Anindya Roy. Posterior consistency of Gaussian process prior for non-parametric binary regression. *Ann. Stat.*, 34:2413–2429, 2006.
- [6] P. Hansen, B. Jaumard, and S. Lu. Global optimization of univariate Lipschitz functions: I. survey and properties. *Mathematical Programming*, 55:251–272, 1992.
- [7] Matthew Hoffman, Eric Brochu, and Nando de Freitas. Portfolio allocation for Bayesian optimization. In *Uncertainty in Artificial Intelligence*, pages 327–336, 2011.
- [8] Benedict May, Nathan Korda, Anthony Lee, and David Leslie. Optimistic Bayesian sampling in contextual-bandit problems. 2010.
- [9] Rémi Munos. Optimistic optimization of a deterministic function without the knowledge of its smoothness. In *Advances in Neural Information Processing Systems*, 2011.
- [10] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, 2006.

- [11] Niranján Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *International Conference on Machine Learning*, 2010.
- [12] Michael L. Stein. *Interpolation of Spatial Data: Some Theory for Kriging*. Springer, 1999.
- [13] Ingo Steinwart and Andreas Christmann. *Support Vector Machines*. Springer, 2008.
- [14] Emmanuel Vazquez and Julien Bect. Convergence properties of the expected improvement algorithm with fixed mean and covariance functions. *Journal of Statistical Planning and Inference*, 140:3088–3095, 2010.

# Apendix: Auxilliary Lemmas

We begin our analysis by showing that, given sufficient explored locations, the residual variance is small. More specifically, for any point  $x$  contained in the convex hull of a set of  $d$  points that are no further than  $\delta$  apart from  $x$ , we show that the residual is bounded by  $O(\|h\|_{\mathcal{H}} \delta^2)$ , where  $\|h\|_{\mathcal{H}}$  is the Hilbert Space norm of the associated function and that furthermore the residual variance is bounded by  $O(\delta^2)$ . This quantity is subsequently related to  $\ell_2$  and  $\ell_1$  covering numbers of the associated domain. We begin by relating residual variance, projection operators, and interpolation in Hilbert Spaces. The results are standard and we only present them here for the purpose of being self-contained.

**Lemma 7 (Hilbert Space Properties)** *Given a set of points  $x_{1:T} := \{x_1, \dots, x_T\} \in \mathcal{D}$  and a Reproducing Kernel Hilbert Space (RKHS)  $\mathcal{H}$  with kernel  $\kappa$  the following bounds hold:*

1. Any  $h \in \mathcal{H}$  is Lipschitz continuous with constant  $\|h\|_{\mathcal{H}} L$ , where  $\|\cdot\|_{\mathcal{H}}$  is the Hilbert space norm and  $L$  satisfies the following:

$$L^2 \leq \sup_{x \in \mathcal{D}} \partial_x \partial_{x'} \kappa(x, x')|_{x=x'} \text{ and for } \kappa(x, x') = \tilde{\kappa}(x - x') \text{ we have } L^2 \leq \partial_x^2 \tilde{\kappa}(x)|_{x=0}. \quad (1)$$

2. Any  $h \in \mathcal{H}$  has its second derivative bounded by  $\|h\|_{\mathcal{H}} Q$  where

$$Q^2 \leq \sup_{x \in \mathcal{D}} \partial_x^2 \partial_{x'}^2 \kappa(x, x')|_{x=x'} \text{ and for } \kappa(x, x') = \tilde{\kappa}(x - x') \text{ we have } Q^2 \leq \partial_x^4 \tilde{\kappa}(x)|_{x=0}. \quad (2)$$

3. The projection operator  $P_{1:T}$  on the subspace  $\text{span}\{\kappa(x_t, \cdot)\}_{t=1:T} \subseteq \mathcal{H}$  is given by

$$P_{1:T}h := \mathbf{k}^\top(\cdot) \mathbf{K}^{-1}(\mathbf{k}(\cdot), h) \quad (3)$$

where  $\mathbf{k}(\cdot) = \mathbf{k}_{1:T}(\cdot) := [\kappa(x_1, \cdot) \cdots \kappa(x_T, \cdot)]^\top$  and  $\mathbf{K} := [\kappa(x_i, x_j)]_{i,j=1:T}$ ; moreover, we have

$$\langle \mathbf{k}(\cdot), h \rangle := \begin{bmatrix} \langle \kappa(x_1, \cdot), h \rangle \\ \vdots \\ \langle \kappa(x_T, \cdot), h \rangle \end{bmatrix} = \begin{bmatrix} h(x_1) \\ \vdots \\ h(x_T) \end{bmatrix}.$$

Here  $P_{1:T}P_{1:T} = P_{1:T}$  and  $\|P_{1:T}\| \leq 1$  and  $\|\mathbf{1} - P_{1:T}\| \leq 1$ .

4. Given sets  $x_{1:T} \subseteq x_{1:T'}$  it follows that  $\|P_{1:T}h\|_{\mathcal{H}} \leq \|P_{1:T'}h\|_{\mathcal{H}} \leq \|h\|_{\mathcal{H}}$ .
5. Given tuples  $(x_i, h_i)$  with  $h_i = h(x_i)$ , the minimum norm interpolation  $\bar{h}$  with  $\bar{h}(x_i) = h(x_i)$  is given by  $\bar{h} = P_{1:T}h$ . Consequently its residual  $g := (\mathbf{1} - P_{1:T})h$  satisfies  $g(x_i) = 0$  for all  $x_i \in x_{1:T}$ .

**Proof** We prove the claims in sequence.

1. This follows from Corollary 4.36 in [13], with  $|\alpha| = 1$ .
2. Same as above, just with  $|\alpha| = 2$ .

3. For any operator  $V$  with full column rank the projection on the image of  $V$  is given by  $V(V^\top V)^{-1}V^\top$ . In our case,  $V$  is defined as

$$V : \mathbb{R}^T \longrightarrow \mathcal{H} \quad .$$

$$\mathbf{w} := \begin{bmatrix} w_1 \\ \vdots \\ w_T \end{bmatrix} \longmapsto \mathbf{k}^\top(\cdot)\mathbf{w}$$

To calculate the adjoint  $V^\top$ , let  $h \in \mathcal{H}$  and  $\mathbf{w} \in \mathbb{R}^T$  be arbitrary elements and consider the following chain of equalities:

$$\begin{aligned} \left\langle V^\top h, \mathbf{w} \right\rangle_{\mathbb{R}^T} &= \langle h, V\mathbf{w} \rangle_{\mathcal{H}} && \text{(cf. [13] Equation A.19)} \\ &= \left\langle h, \mathbf{k}^\top(\cdot)\mathbf{w} \right\rangle_{\mathcal{H}} = \left\langle h, \begin{bmatrix} k(x_1, \cdot) & \cdots & k(x_T, \cdot) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_T \end{bmatrix} \right\rangle_{\mathcal{H}} \\ &= \langle h, k(x_1, \cdot)w_1 + \cdots + k(x_T, \cdot)w_T \rangle_{\mathcal{H}} \\ &= \langle h, k(x_1, \cdot) \rangle w_1 + \cdots + \langle h, k(x_T, \cdot) \rangle w_T && \text{(by linearity of } \langle \cdot, \cdot \rangle_{\mathcal{H}} \text{)} \\ &= \begin{bmatrix} \langle h, k(x_1, \cdot) \rangle_{\mathcal{H}} & \cdots & \langle h, k(x_T, \cdot) \rangle_{\mathcal{H}} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_T \end{bmatrix} \\ &= \langle \mathbf{k}(\cdot), h \rangle_{\mathcal{H}}^\top \mathbf{w} && \text{(by the symmetry of } \langle \cdot, \cdot \rangle_{\mathcal{H}} \text{)} \\ &= \langle \langle \mathbf{k}(\cdot), h \rangle_{\mathcal{H}}, \mathbf{w} \rangle_{\mathbb{R}^T} && \text{(by the definition of } \langle \cdot, \cdot \rangle_{\mathbb{R}^T} \text{)} \end{aligned}$$

Since, this equality holds for all  $\mathbf{w}$ , we can conclude that for all  $h \in \mathcal{H}$

$$V^\top h = \langle \mathbf{k}(\cdot), h \rangle_{\mathcal{H}}.$$

Now, all we need to do is to calculate the expression  $V^\top V$  to complete the derivation of the



expression for  $P_{1:T}$ ; to this end let  $\mathbf{w} \in \mathbb{R}^T$  be arbitrary:

$$\begin{aligned}
 V^\top V \mathbf{w} &= \left\langle \mathbf{k}(\cdot), \mathbf{k}^\top(\cdot) \mathbf{w} \right\rangle_{\mathcal{H}} \\
 &= \left\langle \mathbf{k}(\cdot), \mathbf{k}^\top(\cdot) \right\rangle_{\mathcal{H}} \mathbf{w} = \left\langle \begin{bmatrix} \kappa(x_1, \cdot) \\ \vdots \\ \kappa(x_T, \cdot) \end{bmatrix}, \begin{bmatrix} \kappa(x_1, \cdot) & \cdots & \kappa(x_T, \cdot) \end{bmatrix} \right\rangle_{\mathcal{H}} \mathbf{w} \\
 &= \begin{bmatrix} \langle \kappa(x_1, \cdot), \kappa(x_1, \cdot) \rangle_{\mathcal{H}} & \cdots & \langle \kappa(x_1, \cdot), \kappa(x_T, \cdot) \rangle_{\mathcal{H}} \\ \vdots & \ddots & \vdots \\ \langle \kappa(x_T, \cdot), \kappa(x_1, \cdot) \rangle_{\mathcal{H}} & \cdots & \langle \kappa(x_T, \cdot), \kappa(x_T, \cdot) \rangle_{\mathcal{H}} \end{bmatrix} \mathbf{w} \\
 &= \begin{bmatrix} \kappa(x_1, x_1) & \cdots & \kappa(x_1, x_T) \\ \vdots & \ddots & \vdots \\ \kappa(x_T, x_1) & \cdots & \kappa(x_T, x_T) \end{bmatrix} \mathbf{w} \quad (\text{cf. [13] Definition 4.18 and Equation 4.14}) \\
 &= \mathbf{K} \mathbf{w}
 \end{aligned}$$

and so  $V^\top V = \mathbf{K}$ .

This provides us with  $P_{1:T}$ . The remaining claims follow from standard properties of projection operators.

4. Projection operators satisfy  $\|P_{1:T}\| \leq 1$ . This proves the second claim. The first claim can be seen from the fact that projecting on a subspace can only have a smaller norm than the superspace projection.
5. We first show that the projection is an interpolation. This follows from

$$\bar{h}(x_i) = P_{1:T} h(x_i) = \langle P_{1:T} h, \kappa(x_i, \cdot) \rangle = \langle h, P_{1:T} \kappa(x_i, \cdot) \rangle = \langle h, \kappa(x_i, \cdot) \rangle = h(x_i).$$

Correspondingly  $g(x_i) = h(x_i) - \bar{h}(x_i) = 0$  for all  $x_i \in x_{1:T}$ . By construction  $P_{1:T} h$  uses  $h$  only in evaluations  $h(x_i)$ , hence for any two functions  $h, h'$  with  $h(x_i) = h'(x_i)$  we have  $P_{1:T} h = P_{1:T} h'$ . Since  $\|P_{1:T}\| \leq 1$  it follows that  $\|P_{1:T} h\| \leq \|h\|_{\mathcal{H}}$ . Hence there is no interpolation with norm smaller than  $\|P_{1:T} h\|$ .

■

**Lemma 8 (Gaussian Process Variance)** *Under the assumptions of Lemma 7 it follows that*

$$|h(x) - P_{1:T}h(x)| \leq \|h\|_{\mathcal{H}} \sigma_T(x) \text{ where } \sigma_T^2(x) = \kappa(x, x) - \mathbf{k}_{1:T}^\top(x) \mathbf{K}^{-1} \mathbf{k}_{1:T}(x) \quad (4)$$

*and this bound is tight. Moreover,  $\sigma_T^2(x)$  is the residual variance of a Gaussian process with the same kernel.*

**Proof** To see the bound we again use the Cauchy-Schwartz inequality

$$\begin{aligned} |h(x) - P_{1:T}h(x)| &= |(\mathbf{1} - P_{1:T})h(x)| \\ &= |\langle (\mathbf{1} - P_{1:T})h, \kappa(x, \cdot) \rangle_{\mathcal{H}}| \quad (\text{by the defining property of } \langle \cdot, \cdot \rangle_{\mathcal{H}}, \text{ cf. [13] Def. 4.18}) \\ &= |\langle h, (\mathbf{1} - P_{1:T})\kappa(x, \cdot) \rangle_{\mathcal{H}}| \quad (\text{since } \mathbf{1} - P_{1:T} \text{ is an orthogonal projection and so self-adjoint}) \\ &\leq \|h\|_{\mathcal{H}} \|(\mathbf{1} - P_{1:T})\kappa(x, \cdot)\| \quad (\text{by Cauchy-Schwarz}) \end{aligned}$$

This inequality is clearly tight for  $h = (\mathbf{1} - P_{1:T})\kappa(x, \cdot)$  by the nature of dual norms. Next note that

$$\begin{aligned} \|(\mathbf{1} - P_{1:T})\kappa(x, \cdot)\|^2 &= \langle (\mathbf{1} - P_{1:T})\kappa(x, \cdot), (\mathbf{1} - P_{1:T})\kappa(x, \cdot) \rangle = \langle \kappa(x, \cdot), (\mathbf{1} - P_{1:T})\kappa(x, \cdot) \rangle \\ &= \kappa(x, x) - \langle \kappa(x, \cdot), P_{1:T}\kappa(x, \cdot) \rangle = \sigma_T^2(x). \end{aligned}$$

The second equality follows from the fact that  $\mathbf{1} - P_{1:T}$  is idempotent. The last equality follows from the definition of  $P_{1:T}$ . The fact that  $\sigma_T^2(x)$  is the residual variance of a Gaussian Process regression estimate is well known in the literature and follows, e.g. from the matrix inversion lemma. ■

**Lemma 9 (Approximation Guarantees)** *In the following we denote by  $x_{1:T} \subseteq \mathcal{D}$  a set of locations and assume that  $g(x_i) = 0$  for all  $x_i \in x_{1:T}$ .*

1. *Assume that  $g$  is Lipschitz continuous with bound  $L$ . Then  $g(x) \leq Ld(x, x_{1:T})$ , where  $d(x, x_{1:T})$  is the minimum distance  $\|x - x_i\|$  between  $x$  and any  $x_i \in x_{1:T}$ .*
2. *Assume that  $g$  has its second derivative bounded by  $Q'$ . Moreover, assume that  $x$  is contained inside the convex hull of  $x_{1:T}$  such that the smallest such convex hull has a maximum pairwise distance between vertices of  $d$ . Then we have  $g(x) \leq \frac{1}{4}Q'd^2$ .*

**Proof** The first claim is an immediate consequence of the Lipschitz property of  $g$ . To see the second claim we need to establish a number of issues: without loss of generality assume that the maximum within the convex hull containing  $x$  is attained at  $x$  (and that the maximum rather than the minimum denotes the maximum deviation from 0).

The maximum distance of  $x$  to one of its vertices is bounded by  $\delta/\sqrt{2}$ . This is established by considering the minimum enclosing ball and realizing that the maximum distance is achieved for the regular polyhedron.

To see the maximum deviation from 0 we exploit the fact that  $\partial_x g(x) = 0$  by the assumption of  $x$  being the maximum (we need not consider cases where  $x$  is on a facet of the polyhedral set since in this case we could easily reduce the dimensionality). In this case the largest deviation between  $g(x)$  and  $g(x_i)$  is obtained by making  $g$  a quadratic function  $g(x') = \frac{Q'}{2} \|x' - x\|^2$ . At distance  $\frac{\delta}{\sqrt{2}}$  the function value is bounded by  $\frac{\delta^2 Q'}{4}$ . Since the latter bounds the maximum deviation it does bound it for  $g$  in particular. This proves the claim. ■

For the sake of completeness of the exposition, we will also include the proofs of the following two lemmas, which first appeared in [11]:

**Lemma 10 (Lemma 5.1 of [11])** *Given any finite set  $\mathcal{L}$ , any sequence of points  $\{x_1, x_2, \dots\} \subseteq \mathcal{L}$  and  $f : \mathcal{L} \rightarrow \mathbb{R}$  a sample from  $\text{GP}(0, \kappa(\cdot, \cdot))$ , for all  $\alpha \in (0, 1)$ , we have*

$$P \left\{ \forall x \in \mathcal{L}, t \geq 1 : |f(x) - \mu_{t-1}(x)| \leq \sqrt{\beta_t} \sigma_{t-1}(x) \right\} \geq 1 - \alpha,$$

where  $\beta_t = 2 \ln \left( \frac{|\mathcal{L}| \pi_t}{\alpha} \right)$  and  $\{\pi_t\}$  is any positive sequence satisfying  $\sum \frac{1}{\pi_t} = 1$ . Here  $|\mathcal{L}|$  denotes the number of elements in  $\mathcal{L}$ .

**Proof** Looking at the complement of the set whose probability we're interested in, we have

$$\begin{aligned} & P \left\{ \forall x \in \mathcal{L}, t \geq 1 : |f(x) - \mu_{t-1}(x)| \leq \sqrt{\beta_t} \sigma_{t-1}(x) \right\} \\ &= 1 - P \left\{ \exists x \in \mathcal{L}, t \geq 1 \text{ s.t. } |f(x) - \mu_{t-1}(x)| > \sqrt{\beta_t} \sigma_{t-1}(x) \right\} \\ &= 1 - P \left( \bigcup_{t \geq 1} \bigcup_{x \in \mathcal{L}} |f(x) - \mu_{t-1}(x)| > \sqrt{\beta_t} \sigma_{t-1}(x) \right) \\ &\geq 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} P \left( |f(x) - \mu_{t-1}(x)| > \sqrt{\beta_t} \sigma_{t-1}(x) \right) \\ &= 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} 2P \left( f(x) > \mu_{t-1}(x) + \sqrt{\beta_t} \sigma_{t-1}(x) \right) \\ &= 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} 2P \left\{ r > \sqrt{\beta_t} \mid r \sim \mathcal{N}(0, 1) \right\} = 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} \frac{2}{\sqrt{2\pi}} \int_{\sqrt{\beta_t}}^{\infty} e^{-\frac{r^2}{2}} dr \\ &= 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} \frac{2}{\sqrt{2\pi}} \int_0^{\infty} e^{-\frac{(u+\sqrt{\beta_t})^2}{2}} du \quad (u = r - \sqrt{\beta_t}, du = dr, r = u + \sqrt{\beta_t}) \\ &= 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} \frac{2}{\sqrt{2\pi}} \int_0^{\infty} e^{-\frac{u^2 + 2u\sqrt{\beta_t} + \beta_t}{2}} du = 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} \frac{2e^{-\frac{\beta_t}{2}}}{\sqrt{2\pi}} \int_0^{\infty} e^{-u\sqrt{\beta_t}} e^{-\frac{u^2}{2}} du \\ &\geq 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} \frac{2e^{-\frac{\beta_t}{2}}}{\sqrt{2\pi}} \int_0^{\infty} e^{-\frac{u^2}{2}} du \quad (\because e^{-u\sqrt{\beta_t}} \leq 1 \quad \forall u \geq 0) \\ &= 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} e^{-\frac{\beta_t}{2}} \left[ \frac{2 \int_0^{\infty} e^{-\frac{u^2}{2}} du}{\sqrt{2\pi}} \right] = 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} e^{-\frac{\beta_t}{2}} \left[ \frac{\int_{-\infty}^{\infty} e^{-\frac{u^2}{2}} du}{\sqrt{2\pi}} \right] \\ &= 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} e^{-\frac{\beta_t}{2}} = 1 - \sum_{t \geq 1} \sum_{x \in \mathcal{L}} \frac{\alpha}{\pi_t |\mathcal{L}|} = 1 - \alpha \sum_{t \geq 1} \frac{1}{\pi_t} \sum_{x \in \mathcal{L}} \frac{1}{|\mathcal{L}|} = 1 - \alpha. \end{aligned}$$

■

**Lemma 11 (Lemma 5.2 in [11])** *Let  $\mathcal{L}$  a non-empty finite set and  $f : \mathcal{L} \rightarrow \mathbb{R}$  an arbitrary function. Also assume that there exist functions  $\mu, \sigma : \mathcal{L} \rightarrow \mathbb{R}$  and a constant  $\sqrt{\beta}$ , such that*

$$|f(x) - \mu(x)| \leq \sqrt{\beta}\sigma \quad \forall x \in \mathcal{L}. \quad (\dagger)$$

*Then, we have*

$$r(x) \leq 2\sqrt{\beta}\sigma(x) \leq 2\sqrt{\beta} \max_{\mathcal{L}} \sigma.$$

**Proof** We know from the definition of sup and the upper and the lower bounds given by  $(\dagger)$  that

$$\mu(x) - \sqrt{\beta}\sigma(x) \leq f(x) \leq \sup_{\mathcal{L}} f \leq \sup_{\mathcal{L}} (\mu + \sqrt{\beta}\sigma) = \mu(x) + \sqrt{\beta}\sigma(x).$$

So,  $f(x)$  and  $\sup_{\mathcal{L}} f$  lie in the interval  $[\mu(x) - \sqrt{\beta}\sigma(x), \mu(x) + \sqrt{\beta}\sigma(x)]$ , and so their distance from each other can be at most  $2\sqrt{\beta}\sigma(x)$ , and so

$$r(x) := \sup_{\mathcal{L}} f - f(x) \leq 2\sqrt{\beta}\sigma(x).$$

■