A Framework for Optimization under Limited Information

Tansu Alpcan Technical University Berlin Deutsche Telekom Laboratories *alpcan@sec.t-labs.tu-berlin.de*

Abstract

In many real world problems, optimization decisions have to be made with limited information. The decision maker may have no a priori or posteriori data about the often nonconvex objective function except from on a limited number of points that are obtained over time through costly observations. This paper presents an optimization framework that takes into account the information collection (observation), estimation (regression), and optimization (maximization) aspects in a holistic and structured manner. Explicitly quantifying the information acquired at each optimization step using the entropy measure from information theory, the (nonconvex) objective function to be optimized (maximized) is modeled and estimated by adopting a Bayesian approach and using Gaussian processes as a state-of-the-art regression method. The resulting iterative scheme allows the decision maker to solve the problem by expressing preferences for each aspect quantitatively and concurrently.

1 Introduction

In many real world problems, optimization decisions have to be made with limited information. Whether it is a static optimization or dynamic control problem, obtaining detailed and accurate information about the problem or system can often be a costly and time consuming process. In some cases, acquiring extensive information on system characteristics may be simply infeasible. In others, the observed system may be so nonstationary that by the time the information is obtained, it is already outdated due to system's fast-changing nature. Therefore, the only option left to the decision-maker is to develop a strategy for collecting information efficiently and choose a model to estimate the "missing portions" of the problem in order to solve it satisfactorily and according to a given objective.

To make the discussion more concrete, consider the problem of maximizing a (Lipschitz) continuous *nonconvex* objective function, which is unknown except from its value at only a small number of data points. The decision maker may have no a priori information about the function and start with zero data points. Furthermore, only a limited number of –possibly noisy– observations may be available before making a decision on the maximum value and its location. The function itself, however, remains unknown even after the decision is made. *What is the best strategy to address this problem*?

The decision making framework presented in this paper captures the posed problem by taking into account the information collection (observation), estimation (regression), and (multi-objective) optimization aspects in a holistic and structured manner. Hence, the framework enables the decision maker to solve the problem by expressing preferences for each aspect quantitatively and concurrently. It explicitly incorporates many concepts that have been implicitly considered by heuristic schemes, and builds upon many results from seemingly disjoint but relevant fields such as information theory, machine learning, and optimization and control theories. Specifically, it combines concepts from these fields by

- explicitly quantifying the information acquired using the entropy measure from information theory,
- modeling and estimating the (nonconvex) function or (nonlinear) system adopting a Bayesian approach and using Gaussian processes as a state-of-the-art regression method,
- using an iterative scheme for observation, learning, and optimization,
- capturing all of these aspects under the umbrella of a multi-objective "meta" optimization formulation.

Despite methods and approaches from machine (statistical) learning are heavily utilized in this framework, the problem at hand is very different from many classical machine learning ones, even in its learning aspect. In most classical application domains of machine learning such as data mining, computer vision, or image and voice recognition, the difficulty is often in handling significant amount of data in contrast to lack of it. Many methods such as Expectation-Maximization (EM) inherently make this assumption, except from "active learning" schemes [3]. Information theory plays plays an important role in evaluating scarce (and expensive) data and developing strategies for obtaining it. Interestingly, data scarcity converts at the same time the disadvantages of some methods into advantages, e.g. the scalability problem of Gaussian processes.

It is worth noting that the class of problems described here are much more frequently encountered in practice than it may first seem. For example, the class of black-box methods known as "kriging" [10] have been applied to such problems in geology and mining as well as to hydrology since mid-1960s. In addition, the solution framework proposed is applicable to a wide variety of fields due to its fundamental nature. One example is decentralized resource allocation decisions in networked and complex systems, e.g. wired and wireless networks, where parameters change quickly and global information on network characteristics are not available at the local decision-making nodes. Another example is security-related decisions where opponents spend a conscious effort to hide their actions. A related area is security and information technology risk management in large-scale organizations, where acquiring information on individual subsystems and processes can be very costly. Yet another example application is in biological systems where individual organisms or subsystems operate autonomously (even if they are part of a larger system) under limited local information.

2 Problem Definition and Approach

A concrete definition of the motivating problem mentioned in the introduction section is helpful for describing the multiple aspects of the limited information decision making framework. Without loss of any generality, let

$$\mathcal{X} \subseteq \Psi \subset \mathbb{R}^d$$

be a nonempty, convex, and compact (closed and bounded) subset of the original problem domain Ψ of d dimensions. The original domain Ψ does not have to be convex, compact, or even fully known. However, adopting a "divide and conquer" approach, the subset \mathcal{X} provides a reasonable starting point. Define next the objective function to be maximized

$$f: \mathcal{X} \to \mathbb{R},$$

which is unknown except from on a finite number of points (possibly imperfectly) observed. As a simplifying assumption, let f be Lipschitz continuous on \mathcal{X} . One of the main distinguishing characteristics of this problem is the limitations on set of observations

$$\Omega_n := \{x_1, \dots, x_n : x_i \in \mathcal{X} \, \forall i, n \ge 1\},\$$

due to cost of obtaining information or non-stationarity of the underlying system. Assume for now that the cost of observing the value of the objective function f(x) is the same for any $x \in \mathcal{X}$. Then, a basic search problem is defined as follows:

Problem 1 (Basic Search Problem) Consider a Lipschitz-continuous objective function $f : \mathcal{X} \to \mathbb{R}$ on the d-dimensional nonempty, convex, and compact set $\mathcal{X} \subset \mathbb{R}^d$. The function is unknown except from on a finite number of observed data points. What is the best search strategy

$$\Omega_N := \{x_1, \dots, x_N : x_i \in \mathcal{X} \; \forall i, \; N \ge 1\}$$

that solves

$$\max_{\Omega_N} f(x)$$

for a given N?

The number of observations, N, in Problem 1 may be imposed by the nature of the specific application domain. In many problems, where there is no time constraint, adopting an iterative (one-by-one) approach, and hence choosing N = 1 is clearly beneficial as it allows for usage of incoming new information at each step. Alternatively, the assumption on the equal observation cost can be relaxed and be formulated as a constraint

$$\sum_{x \in \Omega_n} c_o(x) \le C,$$

where $c_o(x) : \mathcal{X} \to \mathbb{R}$ is the observation cost function, and the scalar C is the total "exploration budget". It is also possible to define this cost iteratively based on the (distance from) previous observation, e.g. $c_o(x_n, x_{n-1})$. In such cases, a location-based iterative search scheme can be considered.

The simplest (both conceptually and computationally) strategy to solve Problem 1 is random search on the domain \mathcal{X} . As such no attempt is made to "learn" the properties of the function f. Unless, f is "algorithmically random" [14], which is rarely the case, this strategy wastes the information collected on f. A slightly more complicated and very popular set of strategies combine random search with simple modeling of the function through gradient methods. In this case, the collected information is used to model frudimentarily using derived gradients to "define slopes" in a heuristic manner. Then, these slopes of f are explored step-by-step in the upwards direction to find a local maximum, after which the search algorithm randomly jumps to another location. It is also possible to randomize the gradient climbing scheme for additional flexibility [24].

The framework presented in this paper takes one further step and explicitly models the (entire) objective function f (on the set \mathcal{X}) using the information collected instead of heuristically describing only the slopes. The function \hat{f} , which models, approximates, and estimates f, belongs to a certain class functions such that $\hat{f} \in \mathcal{F}$. The selection and properties of this class is based on "a priori" information available and can be interpreted as the "world view" of the decision maker. These properties can often be expressed using meta-parameters which are then updated based on the observations through a separate optimization process. Likewise, a slower time-scale process can be used for model selection if processing capabilities permit a multi-model approach.

This model-based search process, which lies at the center of the framework, is fundamentally a manifestation of the Bayesian approach [18]. It first imposes explicit and a priori modeling assumptions by choosing \hat{f} from a certain class of functions, \mathcal{F} , and then infers (learns, updates) \hat{f} in a structured manner as more information becomes available through observations.

From a computational point of view, the decision making framework with limited information lies at one end of the computation vs. observation spectrum, while random search is at the opposite end. The framework tries to utilize each piece of information to the maximum possible extent almost regardless of the computational cost. The underlying assumption here is: **observation is very costly whereas computation is rather cheap**. This assumption is not only valid for a wide variety of problems from different fields ranging from networking and security to economics and risk management, but also inspired from biological systems. In many biological organisms, from single cells to human beings, operating close to this end of the computation-observation spectrum is more advantageous than doing random search.

When doing random search on the domain \mathcal{X} , at each stage i.e. given the previous observations, each remaining candidate data point provides equivalent amount of information. However, this is not the case when doing model-based search. Depending on the model adopted and previous information collected, different unexplored points provide different amount of information. This information can be exactly quantified using the definition of entropy and information from the field of (Shannon) information theory. Accordingly, the scalar quantity $\mathcal{I}(\hat{f}, \Omega_n)$ denotes the aggregate information obtained from the set of observations Ω_n within the model represented by \hat{f} . A related issue is the reliability and possibly noisy nature of observations, which will be discussed in further detail in the next section.

An extension of Problem 1 that captures the aspects discussed above is defined next.

Problem 2 (Model-based Search Problem) Let $f : \mathcal{X} \to \mathbb{R}$ be an objective function on the d-dimensional nonempty, convex, and compact set $\mathcal{X} \subset \mathbb{R}^d$, which is unknown except from on a finite number of observed data points. Further let $\hat{f}(x)$ be an estimate of the objective function obtained using an a priori model and observed data. What is the best search strategy $\Omega_N := \{x_1, \ldots, x_N : x_i \in \mathcal{X} \; \forall i, N \geq 1\}$ that solves the multi-objective problem with the following components?

- Objective 1: $\max_{\Omega_N} f(x)$ given $\hat{f}(x)$
- Objective 2: $\arg \min_{\Omega_N} R\left(f(x), \hat{f}(x)\right), \ \hat{f} \in \mathcal{F}$
- Objective 3: $\max_{\Omega_N} \mathcal{I}(\hat{f}, \Omega_n)$

Here, $R(\cdot, \cdot)$ is a risk or expected loss function quantifying the mismatch between actual and estimated functions on the observation data [23]. The scalar quantity \mathcal{I} is the aggregate information obtained from the set of observations Ω_N within the model represented by \hat{f} . The cardinality of Ω_N , N, can be either given, e.g. N = 1, or defined as an additional constraint $\sum_{x \in \Omega_n} c_o(x) \leq C$, where $c_o(x) : \mathcal{X} \to \mathbb{R}$ is the observation cost function, and the scalar C is the total "exploration budget".



Figure 1: The three fundamental aspects of decision making with limited information.

It is important to observe here that the three objectives defined in Problem 2 are (almost) independent from and orthogonal to each other despite being closely related. *Objective 1* purely aims to maximize the unknown objective function f using the best estimate (model) \hat{f} . *Objective 2* focuses on minimizing the error between the estimate \hat{f} and the real unknown function f based on the observations made. *Objective 3* tries to maximize the amount of information provided by each (costly) observation or experiment. It is worth noting that *Objective 3* is independently formulated from *Objective 2*, in other words, exploration is done independently from estimation. In contrast, ensuring a balance between *Objective 1* and 2 is necessary to ensure that solution is robust. These objectives and the fundamental aspects of decision making with limited information are visually depicted in Figure 1.

Table 1: Fundamental Trade-offs

Exploration		Exploitation
Observation	versus	Computation
Robustness		Optimization

There are multiple trade-offs that are inherent to this problem as listed in Table 1. The first one, exploration versus exploitation, puts exploration or obtaining more observations against exploitation, i.e. trying to achieve the given objective. Observation versus computation captures the tradeoff between building sophisticated models using the available information to the fullest extend and making more observations. Robustness versus optimization puts risk avoidance against optimization with respect to the original objective as in exploitation.

3 Methodology

This section presents the methods that are utilized within the framework which addresses the problem defined in the previous one. First, the regression model and Gaussian Processes (GP) are presented. Subsequently, modeling and measurement of information is discussed based on (Shannon) information theory.

3.1 Regression and Gaussian Processes (GP)

Problem 2 presented in the previous section involves inferring or learning the function f using the set of observed data points. This is known as the *regression* problem in machine learning and is a supervised learning method since the observed data constitutes at the same time the learning data set. This learning process involves selection of a "model", where the learned function \hat{f} is, for example, expressed in terms of a set of parameters and specific basis functions, and at the same time minimization of an error measure between the functions f and \hat{f} on the learning data set. Gaussian processes (GP) provide a nonparametric alternative to this but follow in spirit the same idea.

The main goal of regression involves a trade-off. On the one hand, it tries to minimize the *observed* error between f and \hat{f} . On the other, it tries to infer the "real" shape of f and make good estimations using \hat{f} even at unobserved points. If the former is overly emphasized, then one ends up with "over fitting", which means \hat{f} follows f closely at observed points but has weak predictive value at unobserved ones. This delicate balance is usually achieved by balancing the prior "beliefs" on the nature of the function, captured by the model (basis functions), and fitting the model to the observed data.

This paper focuses on Gaussian Process [23] as the chosen regression method within the framework developed without loss of any generality. There are multiple reasons behind this preference. Firstly, GP provides an elegant mathematical method for easily combining many aspects of the framework. Secondly, being a nonparametric method GP eliminates any discussion on model degree. Thirdly, it is easy to implement and understand as it is based on well-known Gaussian probability concepts. Fourthly, noise in observations is immediately taken into account if it is modeled as Gaussian. Finally, one of the main drawbacks of GP namely being computational heavy, does not really apply to the problem at hand since the amount of data available is already very limited.

It is not possible to present here a comprehensive treatment of GP. Therefore, a very rudimentary overview is provided next within the context of the decision making problem. Consider a set of M data points

$$\mathcal{D} = \{x_1, \ldots, x_M\},\$$

where each $x_i \in \mathcal{X}$ is a *d*-dimensional vector, and the corresponding vector of scalar values is $f(x_i)$, $i = 1, \ldots, M$. Assume that the observations are distorted by a zero-mean Gaussian noise, *n* with variance $\sigma \sim \mathcal{N}(0, \sigma)$. Then, the resulting observations is a vector of Gaussian $y = f(x) + n \sim \mathcal{N}(f(x), \sigma)$.

A GP is formally defined as a collection of random variables, any finite number of which have a joint Gaussian distribution. It is completely specified by its mean function m(x) and covariance function $C(x, \tilde{x})$, where

$$m(x) = E[\hat{f}(x)]$$
 and $C(x, \tilde{x}) = E[(\hat{f}(x) - m(x))(\hat{f}(\tilde{x}) - m(\tilde{x}))], \forall x, \tilde{x} \in \mathcal{D}.$

Let us for simplicity choose m(x) = 0. Then, the GP is characterized entirely by its covariance function $C(x, \tilde{x})$. Since the noise in observation vector y is also Gaussian, the covariance function can be defined as the sum of a *kernel function* $Q(x, \tilde{x})$ and the diagonal noise variance

$$C(x,\tilde{x}) = Q(x,\tilde{x}) + \sigma I, \ \forall x, \tilde{x} \in \mathcal{D},$$
(1)

where I is the identity matrix. While it is possible to choose here any (positive definite) kernel $Q(\cdot, \cdot)$, one classical choice is

$$Q(x,\tilde{x}) = \exp\left[-\frac{1}{2} \left\|x - \tilde{x}\right\|^2\right].$$
(2)

Note that GP makes use of the well-known *kernel trick* here by representing an infinite dimensional continuous function using a (finite) set of continuous basis functions and associated vector of real parameters in accordance with the *representer theorem* [26].

The (noisy)¹ training set (\mathcal{D}, y) is used to define the corresponding GP, $\mathcal{GP}(0, C(\mathcal{D}))$, through the $M \times M$ covariance function $C(\mathcal{D}) = Q + \sigma I$, where the conditional Gaussian distribution of any point outside the training set, $\bar{y} \in \mathcal{X}, \bar{y} \notin \mathcal{D}$, given the training data (\mathcal{D}, t) can be computed as follows. Define the vector

$$k(\bar{x}) = [Q(x_1, \bar{x}), \dots Q(x_M, \bar{x})]$$
(3)

and scalar

$$\kappa = Q(\bar{x}, \bar{x}) + \sigma. \tag{4}$$

Then, the conditional distribution $p(\bar{y}|y)$ that characterizes the $\mathcal{GP}(0, C)$ is a Gaussian $\mathcal{N}(\hat{f}, v)$ with mean \hat{f} and variance v,

$$\hat{f}(\bar{x}) = k^T C^{-1} y \text{ and } v(\bar{x}) = \kappa - k^T C^{-1} k.$$
 (5)

This is a key result that defines GP regression as the mean function $\hat{f}(x)$ of the Gaussian distribution and provides a prediction of the objective

¹The special case of perfect observation without noise is handled the same way as long as the kernel function $Q(\cdot, \cdot)$ is positive definite

function f(x). At the same time, it belongs to the well-defined class $\hat{f} \in \mathcal{F}$, which is the set of all possible sample functions of the GP

$$\mathcal{F} := \{ \hat{f}(x) : \mathcal{X} \to \mathbb{R} \text{ such that } \hat{f} \in \mathcal{GP}(0, C(\mathcal{D})), \ \forall \mathcal{D}, C \},\$$

where $C(\mathcal{D})$ is defined in (1) and \mathcal{GP} through (3), (4), and (5), above. Furthermore, the variance function v(x) can be used to measure the uncertainty level of the predictions provided by \hat{f} , which will be discussed in the next subsection.

3.2 Quantifying Information in Observations

In the framework presented, each observation provides a data point to the regression problem (estimating f by constructing \hat{f}) as discussed in the previous subsection. Many works in the learning literature consider the "training" data used in regression available (all at once or sequentially) and do not discuss the possibility of the decision maker influencing or even optimizing the data collection process. The *active learning* problem defined in Section 2 requires, however, exactly addressing the question of "how to quantify information obtained and optimize the observation process?". Following the approach discussed in [17, 18], the framework here provides a precise answer to this question.

Making any decision on the next (set of) observations in a principled manner necessitates first *measuring the information obtained from each observation within the adopted model.* It is important to note that the information measure here is dependent on the chosen model. For example, the same observation provides a different amount of information to a random search model than a GP one.

Shannon information theory readily provides the necessary mathematical framework for measuring the information content of a variable. Let p be a probability distribution over the set of possible values of a discrete random variable A. The **entropy** of the random variable is given by $H(A) = \sum_{i} p_i \log_2(1/p_i)$, which quantifies the amount of uncertainty. Then, the information obtained from an observation on the variable, i.e. reduction in uncertainty, can be quantified simply by taking the difference of its initial and final entropy,

$$\mathcal{I} = H_0 - H_1.$$

It is important here to avoid the common conceptual pitfall of equating entropy to information itself as it is sometimes done in communication theory literature.² Within this framework, (Shannon) information is defined as a measure of the decrease of uncertainty after (each) observation (within a given model). This can be best explained with the following simple example.

3.2.1 Example: Bisection

Choose a number between 1 and 64 randomly with uniform probability (prior). What is the best searching strategy for finding this number? Let the random variable A represent this number. In the beginning the entropy of A is

$$H_0(A) = \sum_{i=1}^{64} \frac{1}{64} \log_2\left(\frac{1}{64}\right) = 6$$
 (bits).

The information maximization problem is defined as

$$\max \mathcal{I} = \max H_0 - H_1 = \min H_1,$$

since H_0 , the entropy before the action (obtaining information) is constant. The entropy H_1 is the one after information is obtained, and hence is directly affected by the specific action chosen. Now, define the action as setting a threshold 1 < t < 64 to check whether the chosen number is less or higher than this threshold t. To simplify the analysis, consider a continuous version of the problem by defining p as the probability of the chosen number being less than the threshold. Thus, in this uniform prior case, the problem simplifies to

$$\min_{p} H_1 = \min_{p} p \log(p) + (1-p) \log(1-p),$$

which has the derivative

$$\frac{dH_1}{dp} = \log(p) - \log(1-p).$$

Clearly, the threshold $p^* = 0.5$ is the global minimum, which roughly corresponds to t = 32 (ignoring quantization and boundary effects). Thus, bisection from the middle is the optimal search strategy for the uniform prior. In this example, the number can be found in the worst-case in 6 steps, each

²Since this issue is not of great importance for the class of problems considered in communication theory, it is often ignored. However, the difference is of conceptual importance in this problem. See http://www.ccrnp.ncifcrf.gov/~toms/information.is.not.uncertainty.html for a detailed discussion.

providing one bit of information. Nonuniform probabilities (priors) can be handled in a similar way.

If this search process (bisection) is repeatedly applied without any feedback, then it results in the optimal quantization of the search space both in the uniform case above and for the nonuniform probabilities. If feedback is available, i.e. one learns after each bisection whether the number is larger or less than the boundary, then this is as shown the best search strategy.

4 Model

The model adopted in the framework for decision making with limited information builds on the methods presented in the previous section and addresses the problem introduced in Section 2. The model consists of three main parts: observation, update of GP for regression, and optimization to determine next action. These three steps, shown in Figure 2 are taken iteratively to achieve the objectives in Problem 2. As a result of its iterative nature, this approach can be considered in a sense similar to the well-known Expectation-Maximization algorithm [3].



Figure 2: The main parts of the underlying model of the decision making framework.

Observations, given that they are a scarce resource in the class of problems considered, play an important role in the model. Uncertainties in the observed quantities can be modeled as additive noise. Likewise, properties (variance or bias) of additive noise can be used to model the reliability of (and bias in) the data points observed. GPs provide a straightforward mathematical structure for incorporating these aspects to the model under some simplifying assumptions.

The set of observations collected provide the (supervised) training data for GP regression in order to estimate the characteristics of the function or system at hand. This process relies on the GP methods described in Subsection 3.1. Thus, at each iteration an up-to-date description of the function or system is obtained based on the latest observations. Specifically, \hat{f} provides an estimate of the original function f^{3} . Assuming an additive Gaussian noise model, the noise variance σ can be used to model uncertainties, e.g. older and noisy data resulting in higher σ values.

The final and most important part of the model provides a basis for determining the next action after an optimization process that takes into account all three objectives in Problem 2. The information aspect of these objectives is already discussed in Subsection 3.2. An important issue here is the fact that there are infinitely many candidate points in this optimization process, but in practice only a finite collection of them can be evaluated.

4.1 Sampling Solution Candidates

When making a decision on the next action through multi-objective optimization, there are (infinitely) many candidate points. A pragmatic solution to the problem of finding solution candidates is to (adaptively) sample the problem domain \mathcal{X} to obtain the set

$$\Theta := \{x_1, \dots, x_T : x_i \in \mathcal{X}, \, x_i \notin \mathcal{D}, \, \forall i\}$$

that does not overlap with known points. In low (one or two) dimensions, this can be easily achieved through grid sampling methods. In higher dimensions, (Quasi) Monte Carlo schemes can be utilized. For large problem domains, the current domain of interest \mathcal{X} can be defined around the last or most promising observation in such a way that such a sampling is computationally feasible. Likewise, multi-resolution schemes can also be deployed to increase computational efficiency.

Although such a solution may seem restrictive at first glance, it is in spirit not very different from other schemes such as simulated annealing, which are widely used to address nonconvex optimization problems. However, a major difference between this and other schemes is the fact that the candidate sampling and evaluation are done here "a priori" due to experimentation being costly while other methods rely on abundance of information.

 $^{^{3}}$ See [23, Chap 7.2] for a discussion on asymptotic analysis of GP regression. It should not be noted, however, that asymptotic properties are of little relevance to the problem at hand.

A natural question that arises is: whether and under what conditions does such a sampling method give satisfactory results. The following result from [30, 31] provides an answer to this question in terms of number of samples required.

Theorem 1 Define a multivariate function f(x) on the convex, compact set \mathcal{X} , which admits the maximum $x^* = \arg \max_{x \in \mathcal{X}} f(x)$. Based on a set of N random samples $\Theta = \{x_1, \ldots, x_N : x_i \in \mathcal{X} \forall i\}$ from the entire set \mathcal{X} , let $\hat{x} := \arg \max_{x \in \Theta} f(x)$ be an estimate of the maximum x^* .

Given an $\varepsilon > 0$ and $\delta > 0$, the minimum number of random samples N which guarantees that

$$Pr\left(Pr[f(x^*) > f(\hat{x})] \le \varepsilon\right) \ge 1 - \delta,$$

i.e. the probability that 'the probability of the real maximum surpassing the estimated one being less than ε ' is larger than $1 - \delta$, is

$$N \ge \frac{\ln 1/\delta}{1/(1-\varepsilon)}$$

Furthermore, this bound is tight if the function f is continuous on \mathcal{X} .

It is interesting and important to note that this bound is independent of the sampling distribution used (as long as it covers the whole set \mathcal{X} with nonzero probability), the function f itself, as well as the properties and dimension of the set \mathcal{X} .

4.2 Quantifying Information in GP

The information measurement and GP approaches in Section 3 can be directly combined. Let the zero-mean multivariate Gaussian (normal) probability distribution be denoted as

$$p(x) = \frac{1}{\sqrt{2\pi |C_p(x)|}} \exp\left(-\frac{1}{2}[x-m]^T |C_p(x)|^{-1}[x-m]\right), \ x \in \mathcal{X},$$
(6)

where $|\cdot|$ is the determinant, m is the mean (vector) as defined in (5), and $C_p(x)$ is the covariance matrix as a function of the newly observed point $x \in \mathcal{X}$ given by

$$C_p(x) = \begin{bmatrix} C(\mathcal{D}) & k(x)^T \\ k(x) & \kappa \end{bmatrix}.$$
 (7)

Here, the vector k(x) is defined in (3) and κ in (4), respectively. The matrix $C(\mathcal{D})$ is the covariance matrix based on the training data \mathcal{D} as defined in (1).

The entropy of the multivariate Gaussian distribution (6) is [1]

$$H(x) = \frac{d}{2} + \frac{d}{2}\ln(2\pi) + \frac{1}{2}\ln|C_p(x)|,$$

where d is the dimension. Note that, this is the entropy of the GP estimate at the point x based on the available data \mathcal{D} . The aggregate entropy of the function on the region \mathcal{X} is given by

$$H^{agg} := \int_{x \in \mathcal{X}} \frac{1}{2} \ln |C_p(x)| dx.$$
(8)

The problem of choosing a new data point \hat{x} such that the information obtained from it within the GP regression model is maximized can be formulated in a way similar to the one in the bisection example:

$$\hat{x} = \arg\max_{\tilde{x}} \mathcal{I} = \arg\max_{\tilde{x}} \int_{x \in \mathcal{X}} \left[H_0 - H_1 \right] \, dx = \arg\min_{\tilde{x}} \int_{x \in \mathcal{X}} \frac{1}{2} \ln |C_q(x, \tilde{x})| dx,$$
(9)

where the integral is computed over all $x \in \mathcal{X}$, and the covariance matrix $C_q(x, \tilde{x})$ is defined as

$$C_q(x,\tilde{x}) = \begin{bmatrix} C(\mathcal{D}) & k^T(\tilde{x}) & k^T(x) \\ k(\tilde{x}) & \tilde{\kappa} & Q(x,\tilde{x}) \\ k(x) & Q(x,\tilde{x}) & \kappa \end{bmatrix},$$
(10)

and $\tilde{\kappa} = Q(\tilde{x}, \tilde{x}) + \sigma$. Here, $C(\mathcal{D})$ is a $M \times M$ matrix and C_q is a $(M+2) \times (M+2)$ one, whereas κ and $Q(x, \tilde{x})$ are scalars and k is a $M \times 1$ vector. This result is summarized in the following proposition.

Proposition 1 As a maximum information data collection strategy for a Gaussian Process with a covariance matrix $C(\mathcal{D})$, the next observation \hat{x} should be chosen in such a way that

$$\hat{x} = \arg \max_{\tilde{x}} \mathcal{I} = \arg \min_{\tilde{x}} \int_{x \in \mathcal{X}} \ln |C_q(x, \tilde{x})| dx,$$

where $C_q(x, \tilde{x})$ is defined in (10).

An Approximate Solution to Information Maximization

Given a set of (candidate) points Θ sampled from \mathcal{X} , the result in Proposition 1 can be revisited. The problem in (9) is then approximated [31] by

$$\max_{\tilde{x}} \mathcal{I} \approx \min_{\tilde{x}} \sum_{x \in \Theta} \ln |C_q(x, \tilde{x})|$$

$$\Rightarrow \hat{x} = \arg\min_{\tilde{x} \in \Theta} \prod_{x \in \Theta} |C_q(x, \tilde{x})|,$$
(11)

using monotonicity property of the natural logarithm and the fact that the determinant of a covariance matrix is non-negative. Thus, the following counterpart of Proposition 1 is obtained:

Proposition 2 As an approximately maximum information data collection strategy for a Gaussian Process with a covariance matrix $C(\mathcal{D})$ and given a collection of candidate points Θ , the next observation $\hat{x} \in \Theta$ should be chosen in such a way that

$$\hat{x} = \arg\min_{\tilde{x}\in\Theta} \prod_{x\in\Theta} |C_q(x,\tilde{x})| \approx \arg\max_{\tilde{x}\in\Theta} \mathcal{I},$$

where $C_q(x, \tilde{x})$ is given in (10).

Although it is an approximation, finding a solution to the optimization problem in Proposition 2 can still be computationally costly. Therefore, a greedy algorithm is proposed as a computationally simpler alternative. Let, $x^* \in \Theta$ be defined as

$$x^* := \arg\max_{x \in \Theta} |C_p(x)| = |C(\mathcal{D})| |\kappa(x) - k(x)C^{-1}(\mathcal{D})k^T(x)|,$$

where the matrix C_p is given by (7) [21]. The first term above, $|C(\mathcal{D})|$ is fixed and the second one,

$$|\kappa(x) - k(x)C^{-1}(\mathcal{D})k^T(x)|,$$

is the same as the GP variance v(x) in (5). Hence, the sample x^* is one of those with the maximum variance in the set Θ , given current data \mathcal{D} .

It follows from (10) and basic matrix theory that if $\tilde{x} = x$ for a given x then $|C_q(x, \tilde{x})|$ is minimized. As a simplification, ignore the dependencies

between $C_q(x, \tilde{x})$ matrices for different $x \in \Theta$. Then, choosing the maximum variance \hat{x} as

$$\hat{x} = \arg\max_{\tilde{x}\in\Theta} v(\tilde{x}) \approx \arg\min_{\tilde{x}\in\Theta} \prod_{x\in\Theta} |C_q(x,\tilde{x})|,$$

leads to a large (possibly largest) reduction in $\prod_{x \in \Theta} |C_q(x, \hat{x})|$, and hence provides a rough approximate solution to (11) and to the result in Proposition 1. This result is consistent with widely-known heuristics such as "maximum entropy" or "minimum variance" methods [28] and a variant has been discussed in [17].

Proposition 3 Given a Gaussian Process with a covariance matrix $C(\mathcal{D})$ and a collection of candidate points Θ , an approximate solution to the maximum information data collection problem defined in Proposition 1 is to choose the sample point(s) \tilde{x} in such a way that it has (they have) the maximum variance within the set Θ .

5 Optimization with Limited Information

Let $f : \mathcal{X} \to \mathbb{R}$ be the unknown Lipschitz-continuous function of interest on the *d*-dimensional nonempty, convex, and compact set $\mathcal{X} \subset \mathbb{R}^d$. The amount of information about this function available to the decision maker is limited to a finite number of possibly noisy observations. Since the observations are costly, the goal of the decision maker is to find the maximum of f, estimate f as accurately as possible using available observations, and select the most informative data points, at the same time. This naturally calls for an iterative and myopic optimization procedure since each new observation provides a new data point that concurrently affects the maximization, function estimation (regression), and information quantity.

The first and basic objective is the maximization of the function f(x) on $x \in \mathcal{X}$. As a simplification, observations are assumed to be sequential, one at a time. Since f is basically unknown, this problem has to be formulated as

$$\max_{\tilde{x}\in\mathcal{X}}F_1(\tilde{x})=f(\tilde{x}),$$

where \hat{f} is the best estimate obtained through GP regression (5) using the current data set \mathcal{D} . Data uncertainty (observation errors) is modeled through additive Gaussian noise with variance σ as a first approximation. The second objective is to minimize the difference (estimation error) between \hat{f} and f. Define $e(x) = \hat{f}(x) - f(x), \forall x \in \mathcal{X}$. Given the set of noisy observations

$$\mathcal{O} = \{ f(x_i) + n(x_i) : x \in \mathcal{D}, \forall i \},\$$

where $n \sim \mathcal{N}(0, \sigma)$ denotes zero mean Gaussian noise, it is possible to use another GP regression (5) to estimate this error function, $\hat{e}(\mathcal{D}, x)$, on the entire set \mathcal{X} . Thus, the second objective is to ensure that the next observation \tilde{x} solves

$$\min_{\tilde{x}\in\mathcal{X}}F_2(\tilde{x}) = \int_{\tau\in\mathcal{X}} |\hat{e}(\tilde{x},\mathcal{D},\tau)| \, d\tau$$

Note that, F_2 here corresponds to a risk or loss estimate function.

The third objective is to maximize the amount of information obtained with each observation \tilde{x} , or

$$\max_{\tilde{x}\in\mathcal{X}}F_3(\tilde{x}) = \mathcal{I}(\tilde{x},\hat{f}) = \int_{x\in\mathcal{X}} \ln|C_q(x,\tilde{x})| dx,$$

given the best estimate of the original function, \hat{f} . This objective has already been discussed in Section 3.2 in detail.

The values of the three objectives, F_1 , F_2 , F_3 , cannot be evaluated numerically on the entire set \mathcal{X} . Therefore, a sampling method is used as described in Section 4 to obtain a set of solution candidates Θ , which replaces \mathcal{X} in the maximization and minimization problems above. Next, specific problem formulations are presented based on such a sampling of solution candidates. The overall structure of the framework is visualized in Figure 3.

5.1 Solution Approaches

The most common approach to multi-objective optimization is the **weighted** sum method [19, 9]. The three objectives discussed above can be combined to obtain a single objective using the respective weights $[w_1, w_2, w_3]$, $\sum_{i=1}^{3} w_i = 1, \ 0 \le w_i \le 1 \ \forall i$. Assuming a single data point is chosen from and observed among the candidates Θ at each step, i.e. $\tilde{x} = \Omega_1$, a specific weighted sum formulation to address Problem 2 is obtained.

Proposition 4 The solution, $\tilde{x} \in \Theta$, to the optimization problem

$$\max_{\tilde{x}\in\Theta} F(\tilde{x}) = \sum_{i=1}^{3} F_i(\tilde{x}) = w_1 \hat{f}(\tilde{x}) - w_2 \frac{1}{N} \sum_{\tau\in\Theta} |\hat{e}(\tilde{x}, \mathcal{D}, \tau)| + w_3 \mathcal{I}(\tilde{x}, \hat{f}), \quad (12)$$



Figure 3: The decision making framework for static optimization with limited information.

constitutes the best search strategy for this weighted sum formulation of Problem 2.

As discussed in Subsection 3.2 and stated in Proposition 2, the information objective, F_3 , in (12) can be approximated by substituting it with GP variance v(x) in (5) to decrease computational load. Thus, an approximation to the solution in Proposition 4 is:

Proposition 5 The solution, $\tilde{x} \in \Theta$, to the optimization problem

$$\max_{\tilde{x}\in\Theta} F(x) = \sum_{i=1}^{3} F_i(\tilde{x}) = w_1 \hat{f}(x) - w_2 \frac{1}{N} \sum_{\tau\in\Theta} |\hat{e}(\tilde{x}, \mathcal{D}, \tau)| + w_3 v(\tilde{x}), \quad (13)$$

where $v(\tilde{x})$ is defined in (5), approximates the search strategy in Proposition 4.

The weighting scheme described is only meaningful if the three objectives are of the same order of magnitude. Therefore, the original objective functions, F_i , i = 1, 2, 3, have to be transformed or "normalized". There are many different approaches to perform such a transformation [19, 9]. The most common one, which coincidentally is known as normalization, aims to map each objective function to a predefined interval, e.g. [0, 1]. To do this, estimate first an upper F_i^U and lower F_i^L bound on each individual objective $F_i(x)$. Then, the i^{th} normalized objective is

$$F_i^N(x) = rac{F_i(x) - F_i^L}{F_i^U - F_i^L}.$$

The main issue in normalization is to determine the appropriate upper and lower bounds, which is a very problem-dependent one. In the case of Proposition 5, the estimated functions \hat{f} and \hat{e} on the set Θ as well as the existing observations \mathcal{D} , can be utilized to obtain these values. The specific bounds for the respective objectives $F_1^U = \max_{x \in \Theta} \hat{f}(x), F_1^L = \min_{x \in \Theta} \hat{f}(x), F_2^U = \max_{x \in \Theta} |\hat{e}(x, \mathcal{D})|, F_2^L = 0, F_3^U = \max_{x \in \Theta} \kappa(x)$, and $F_3^U = 0$ provide a suitable starting estimate and can be combined with a prior domain knowledge if necessary. Thus, a normalized version of the formulation in Proposition 5 is obtained.

Proposition 6 The solution, $\tilde{x} \in \Theta$, to the optimization problem

$$\max_{\tilde{x}\in\Theta} F(x) = \sum_{i=1}^{3} F_{i}^{N}(\tilde{x}) = \frac{w_{1}}{\Delta_{1}} \left(\hat{f}(x) - F_{1}^{L} \right) - \frac{w_{2}}{\Delta_{2}} \frac{1}{N} \sum_{\tau\in\Theta} |\hat{e}(\tilde{x}, \mathcal{D}, \tau)| + \frac{w_{3}}{\Delta_{3}} v(\tilde{x}),$$
(14)

where $\Delta_i = F_i^U - F_i^L$ i = 1, 2, 3, provides an approximation to the best search strategy for solving the normalized weighted-sum formulation of Problem 2.

The **bounded objective function** method provides a suitable alternative to the weighted sum formulation above in addressing the multi-objective problem defined. The bounded objective function method minimizes the single most important objective, in this case $F_1(x)$, while the other two objective functions $F_2(x)$ and $F_3(x)$ are converted to form additional constraints. Such constraints are in a sense similar to QoS ones that naturally exist in many real life problems [20, 2, 29]. As an advantage, in the bounded objective formulation there is no need for normalization.

The bounded objective counterpart of the result in Proposition 5 is as follows.

Proposition 7 The solution, $\tilde{x} \in \Theta$, to the constrained optimization problem

$$\max_{\tilde{x}\in\Theta} f(x)$$
(15)
such that $0 \le F_2(\tilde{x}) = \frac{1}{N} \sum_{\tau\in\Theta} |\hat{e}(\tilde{x}, \mathcal{D}, \tau)| \le b_1,$
and $0 \le F_3(\tilde{x}) = v(\tilde{x}) \le b_2,$

where b_1 and b_2 are given (predetermined) scalar bounds on F_2 and F_3 , respectively, provides an approximate best search strategy for a bounded-objective formulation of Problem 2.

The advantage of the bounded objective function method is that it provides a bound on the information collection and estimation objectives while maximizing the estimated function. This leads in practice to an initial emphasis on information collection and correct estimation of the objective function. In that sense, the method is more "classical", i.e. follows the common method of learn first and maximize later. Furthermore, it does not require normalization, i.e. it is easier to deploy. The method has, however, a significant disadvantage which makes its usage prohibitive. In large-scale or high-dimensional problems, the space to explore to satisfy any bound on information is simply immense. Therefore, one does not have the luxury of identifying the function first to maximize it later as it would take too many samples to do this. In such cases, it makes more sense to deploy the weighted sum method, possibly along with a cooling scheme to modify the weights as part of a cooling scheme to balance depth-first vs. breadth-first search.

Until now, it has been (implicitly) assumed that the static optimization problem at hand is stationary. However, in a variety of problems this is not the case and the function f(x,t) changes with time. The decision making framework allows for modeling such systems in the following way. Let

$$\mathcal{O}(t) = \{ f(x_i, t_i) + n(x_i, t_i) : x_i \in \mathcal{D}, t_i \le t, \forall i \},\$$

be the set of noisy or unreliable past observations until time t, where $n(x,t) \sim \mathcal{N}(0,\sigma(t))$ is the zero mean Gaussian "noise" term at time t. Now, the deterioration in the past information due to change in f(x,t) can be captured by increasing the variance of the noise term, $\sigma(t)$, with time. For example, a simple linear dynamic can be defined as

$$\frac{d\sigma(t)}{dt} = \eta$$

where $\eta > 0$ captures the level of stationarity, e.g. a large η indicates a rapidly changing system and function f(x, t).

5.2 Algorithm

An algorithmic summary of the solution approaches discussed above for a specific set of choices is provided by Algorithm 1, which describes both weighted-sum and bounded objective variants.

Algorithm 1 Optimization with Limited Information

1:	Input: Function domain, \mathcal{X} , GP meta-parameters, objective weights
	$[w_1, w_2, w_3]$ or bounds b_1, b_2 , initial data set (\mathcal{D}, y) .
2:	Use GP with a Gaussian kernel and specific expected error variances for
	function \hat{f} and error function \hat{e} estimation.
3:	while Search budget available, $1 \le n \le N_{max}$. do
4:	Sample domain \mathcal{X} to obtain $\Theta(n)$. In some cases, $\Theta(n) = \Theta \ \forall n$.
5:	Estimate \hat{f} and \hat{e} based on observed data (\mathcal{D}, y) on $\Theta(n)$ using GPs.
6:	Compute variance, $v(x)$, of $\hat{f}(5)$ on $\Theta(n)$ as an estimate of $\mathcal{I}(\hat{f})$.
7:	if Weighted-sum method then
8:	Next action maximizes a normalized and weighted sum of objectives
	$\sum_{i=1}^{3} F_i^N$ as stated in Proposition 6.
9:	else if Bounded objective method then
10:	Next action is solution to the constrained problem in Proposition
	7.
11:	end if
12:	Update the observed data (\mathcal{D}, y) .
13:	end while

5.3 Numerical Analysis

The Algorithm 1 is illustrated next with multiple numerical examples. It is worth reminding that the main issue here is to solve the optimization problems with minimum data using active learning. In all examples, a uniform grid is used to sample the solution space rather than resorting to a more sophisticated method since the examples are chosen to be only one or two dimensional for visualization purposes.

Example 1

The first numerical example aims to visualize the presented framework and algorithm. Hence, the chosen function is only one dimensional, f(x) = sin(5x)/x on the interval $\mathcal{X} = [0.1, 3.9]$. The interval is linearly sampled to obtain a grid with a distance of 0.01 between points, i.e. $\Theta = \{x_i \in \mathcal{X} \forall i : x_1 = 0.1, x_2 = 0.11, \ldots, x_N = 3.9\}$. A Gaussian kernel with variance 0.1 is chosen for estimating both \hat{f} and \hat{e} . The weights are equal to one, w = [1, 1, 1], in the weighted-sum method. The bounds are $b_1 = 0.5$ for the error bound and $b_2 = 0.2$ for the bound on maximum variance estimate in the bounded objective method. The initial data consists of a single point, x = 0.1.

Figure 4 shows the results based on the normalized weighted-sum method in Proposition 6 after 5 iterations (6 samples in total, together with the initial data point). The variance here is v(x) of the estimated function \hat{f} using data points \mathcal{D} . Clearly, the estimated peak is not the one of the real function f.

Next, Figure 5 shows that after 11 iterations (12 data points in \mathcal{D}), the function and the location of its peak is estimated correctly. The sequence of points selected during the iteration process are:

 $\mathcal{D} = \{0.47, 3.22, 1.17, 1.66, 2.43, 2.06, 3.9, 2.83, 3.6, 0.82, 1.42\}.$



Figure 4: Optimization result using the weighted-sum method with 6 data points.

The amount of information obtained during the iterative optimization is of particular interest. Figure 6 depicts the mean variance v and entropy \mathcal{I} of the estimated function \hat{f} on Θ at each iteration step. In this specific example, the two quantities are very well correlated. Note, however, that this correlation is a function of the relative weights between information collection and other objectives.

Finally, Figure 7 depicts the results of the bounded objective method with the given bounds. The number of iterations is 11 as before, which



Figure 5: Optimization result using the weighted-sum method with 12 data points.



Figure 6: Mean variance v and entropy \mathcal{I} on Θ at each iteration step.



Figure 7: Optimization result using the bounded objective method with 12 data points.

gives an opportunity of direct comparison with the weighted-sum method. The sequence of points selected during the iteration process are:

 $\mathcal{D} = \{0.47, 3.22, 1.17, 1.66, 2.43, 2.06, 3.9, 2.83, 3.6, 0.82, 1.42\}.$

Example 2

The objective function in the second numerical example is the Goldstein&Price function [8], which is shown in Figure 8 in its inverted form to ensure consistency with the maximization formulation in this paper. The problem domain consists of the two dimensional rectangular region $\mathcal{X} = [-2, 2] \times [-2, 2]$, which is linearly sampled to obtain a uniform grid with a 0.05 interval between sample points. A Gaussian kernel with variance 0.5 and 0.1 is chosen for estimating \hat{f} and \hat{e} , respectively. The weighted-sum method is utilized in Algorithm 1 with the weights w = [4, 2, 3]. The search budget is chosen as 50 before stopping the algorithm (for the search space of approx. 6400 samples in the grid). The real global minimum (peak) of the (inverted) Goldstein&Price function is at (0, -1) and the location found by the algorithm using the 50 data points is (-0.15, -1.05). Figure 9 depicts the estimated function, the data points as well as the optimum found. Although the real optimum value is -3 (in the inverted version) while the obtained one is -9.75, the result is still very satisfactory considering that the simple sampling scheme used and the Goldstein&Price function takes values in a range of 1 million, i.e. the error is less than 0.001 percent of the range. Finally, Figure 10 depicts the mean variance v and entropy \mathcal{I} of the estimated function \hat{f} on Θ at each iteration step.



Figure 8: The inverted Goldstein&Price function [8].

Example 3

The third example uses the same setup as the second one but this time with the (inverted) Brain function [6] shown in Figure 11. The rectangular problem domain $\mathcal{X} = [-5, 10] \times [0, 15]$ is sampled uniformly to obtain a grid of points with a 0.2 interval. The real global minimums (peaks) of the (inverted) Branin function are at (9.4, 2.47), ($-\pi$, 12.28), and (π , 2.28) whereas the locations found by the algorithm are (9, 2.6), (-3.2, 12), and (3, 2.2). The values at these locations found vary between -4.3 and -0.5compared to the real global value of -0.4 (of the inverted function). Thus, the algorithm again performs satisfactorily. Figure 9 shows the computed location of one optimum, the data points, as well as the estimated function based on the data points.



Figure 9: Optimization of the inverted Goldstein&Price function [8] using the weighted-sum method with 50 data points.



Figure 10: Mean variance v and entropy \mathcal{I} on Θ at each iteration step.



Figure 11: The inverted Branin function [6].



Figure 12: Optimization of the inverted Branin function [6] using the weighted-sum method with 50 data points.

Example 4

The fourth example is based on the six-hump camel function [7] (see Figure 13) on the domain $\mathcal{X} = [-2, 2] \times [-2, 2]$, which is sampled uniformly with a 0.05 interval. All of the parameters are chosen to be the same as before. Figure 14 shows the computed location of two optimums, the 50 data points, as well as the estimated function based on the data points. The optimum locations found are (0, 0.65) and (0.05, -0.6) with respective values of 0.98 and 1.06, whereas the real locations are (-0.09, 0.71) and (0.09, -0.71) with the value 1.03.



Figure 13: The inverted six-hump camel function.

6 Literature Review

Decision making with limited information is related to search theory. The idea of using information (theory) in this context is hardly new as evidenced by the article "A New Look at the Relation Between Information Theory and Search Theory" from 1979 [22]. The subject is further studied in [11]. The topic of optimal search is more recently revisited by [35], which contains substantial historical notes and studies problems where the search target distribution in itself is unobservable.



Figure 14: Optimization of the inverted six-hump camel function [7] using the weighted-sum method with 50 data points.

The book [18] provides important and valuable insights into the relationship between information theory, inference, and learning. Measuring information content of experiments using Shannon information is explicitly mentioned and a slightly informal version of the bisection example in Subsection 3.2 is discussed. However, focusing mainly on more traditional coding, communication, and machine learning topics, the book does not discuss the type of decision making problems presented in this paper.

Learning plays an important role in the presented framework, especially *regression*, which is a classical machine (or statistical) learning method. A very good introduction to the subject can be found in [3]. A complementary and detailed discussion on kernel methods is in [26]. Another relevant topic is Bayesian inference [33, 18], which is in the foundation of the presented framework. In machine learning literature, Gaussian processes (GPs) are getting increasingly popular due to their various favorable characteristics. The book [23] presents a comprehensive treatment of GPs. Additional relevant works on the subject include [18, 26, 16], which also discuss GP regression.

Convex optimization [4] is a well-understood topic that is often easy to handle even if available information is limited. Optimizing nonconvex functions, however, is still a research subject [12]. It is interesting to note that the method known as *kriging* in global optimization is almost the same as GP regression in machine learning. The field *stochastic programming* focuses on optimization under uncertainty but assumes a certain amount of prior knowledge on the problem at hand and models the uncertainty probabilistically [25]. The popular heuristic method *simulated annealing* [24] is essentially based on iterative random search. Another popular heuristic scheme particle swarm optimization [13] is also based on random search but parallel in nature as a distinguishing characteristic rather than iterative.

Gaussian processes have been recently applied to the area of optimization and regression [5] as well as system identification [32]. While the latter mentions active learning, neither work discusses explicit information quantification or builds a connection with Shannon information theory. The recent articles [15, 34], which utilize GP regression for optimization in a setting similar to the one in this paper and for state-space inference and learning, respectively, do not consider information-theoretic aspects of the problem, either. Likewise, the article [10] on stochastic black box optimization, which considers a problem similar to the one here, does not take into account explicit measurement of information.

The area of active learning or experiment design focuses on data scarcity in machine learning and makes use of Shannon information theory among other criteria [28]. The paper [17] discusses objective functions which measure the expected informativeness of candidate measurements within a Bayesian learning framework. The subsequent study [27] investigates active learning for GP regression using variance as a (heuristic) confidence measure for test point rejection.

7 Discussion

The foundation of the approach adopted in this paper is Bayesian inference, where the main idea is to choose an a priori model and update it with actual experimental data observed (see [18, Chap. 2] for a beautiful introductory discussion on the subject). As long as the a priori model is close to the reality (of the problem at hand), this inference methodology works very efficiently as indicated by the numerical examples in Section 5.3. In many cases this background information, which is sometimes referred to as "domain knowledge", is already available. However, in others one has to explore the model domain and learn model meta-parameters in a time scale naturally longer than the one of actual optimization [16].

The GP regression adopted in the presented framework is only one method for function estimation and other, e.g. parametric, methods can easily replace GP for the regression part. In any case, the regression methodology here is consistent with the principle of "Occam's razor", more specifically its interpretation using Kolmogorov complexity [14]. A priori, the optimization problems at hand are more probable to be simple rather than complex to describe in accordance with *universal distribution* [14]. Hence, given a data set it is reasonable to start describing it with the simplest explanation. GP regression already incorporates this line of thinking by relying on a kernel-based approach and making use of the representer theorem [23, Chap. 6.2]. As a visual example, we refer Figures 4 and 5 for a comparison of function estimates with different sets of available data.

This paper considers a class of problems where data is scarce and obtaining it is costly. Information theory plays an especially important role in devising optimal schemes for obtaining new data points (active learning). The entropy measure from Shannon information theory provides the necessary metric for this purpose, which quantifies the "exploration" aspect of the problem. Using a multi-objective optimization formulation, the presented framework allows explicit weighting of *exploration* vs. *exploitation* aspects. This trade-off is also very similar to one between the well-known depth-first vs. breadth-first search algorithms in search theory.

The amount of information obtained from each data point is different here only because a specific a priori general model is utilized to explain the observed data (GP regression). Because of this the amount of information obtained is specific to the model. Otherwise, without this Bayesian approach, each data point would give the same information (inversely proportional to the total number of candidate points).

The illustrative examples discussed are low-dimensional, which makes it possible to use grids for sampling. However, in higher dimensions (i.e. when the problem is much more "difficult") this "luxury" is not affordable and one has to necessarily resort to Monte Carlo methods. In such cases, the trade-off between exploration and exploitation is even more emphasized. Possible methods to address this issue include, "cooling" approaches similar to those used in simulated annealing, multi-resolution sampling based on region of interest or using topological properties of Gaussian mixtures to intelligently estimate candidate points based on the current state.

The optimization approach presented here can also be interpreted from a biological perspective. If an analogy between the decision-maker and a biological organism is established, then the a-priori Bayesian model (meta parameters of the GP) that is refined over a long time scale corresponds to evolution of a species in an environment (problem domain). Each individual organism belogning to the species obtains new information to achieve its objective while preserving resources as much as possible. The existing evolutionary basis (GP model) gives them an advantage to find a solution much faster compared to random search. From the perspective of the species, it also makes sense for some of its members to explore the model (meta parameter) domain and further refine it through adaptation. Those with better meta parameters achieve then their objectives even more efficiently and obtain an evolutionary edge in natural selection (assuming competition).

8 Conclusion

The decision making framework presented in this paper addresses the problem of decision making under limited information by taking into account the information collection (observation), estimation (regression), and (multiobjective) optimization aspects in a holistic and structured manner. The methodology is based on Gaussian processes and active learning. Various issues such as quantifying information content of new data points using information theory, the relationship between information and GP variance as well as related approximation and multi-objective optimization schemes are discussed. The framework is demonstrated with multiple numerical examples.

The presented framework should be considered mainly as an initial step. Future research directions are abundant and include further investigation of the exploration-exploitation trade-off, adaptive weighting parameters, and random sampling methods for problems in higher dimensional spaces. Additional research topics are the relationship of the framework with genetic/evolutionary methods, dynamic control problems, and multi-person decision making, i.e. game theory.

Acknowledgements

This work is supported by Deutsche Telekom Laboratories. The author wishes to thank Lacra Pavel, Slawomir Stanczak, Holger Boche, and Kivanc Mihcak for stimulating discussions on the subject.

References

 N. Ahmed and D. Gokhale, "Entropy expressions and their estimators for multivariate distributions," *IEEE Transactions on Information Theory*, vol. 35, no. 3, pp. 688–692, May 1989.

- [2] T. Alpcan, X. Fan, T. Başar, M. Arcak, and J. T. Wen, "Power control for multicell CDMA wireless networks: A team optimization approach," *Wireless Networks*, vol. 14, no. 5, pp. 647–657, October 2008. [Online]. Available: papers/Alpcan-Winet.pdf
- [3] C. M. Bishop, Pattern Recognition and Machine Learning (Information Science and Statistics). Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.
- [4] S. Boyd and L. Vandenberghe, *Convex Optimization*. New York, NY, USA: Cambridge University Press, 2004.
- [5] P. Boyle, "Gaussian for regression and optimiprocesses sation," Ph.D. dissertation, Victoria University of Welling-Wellington, New Zealand, 2007.[Online]. Available: ton, http://researcharchive.vuw.ac.nz/handle/10063/421
- [6] F. H. Branin, "Widely convergent method for finding multiple solutions of simultaneous nonlinear equations," *IBM Journal of Research and Development*, vol. 16, pp. 504–522, September 1972. [Online]. Available: http://dx.doi.org/10.1147/rd.165.0504
- [7] L. C. W. Dixon and G. P. Szego, "The optimization problem: An introduction," in *Towards Global Optimization II*, L. C. W. Dixon and G. P. Szego, Eds. New York, NY, USA: North-Holland, 1978.
- [8] A. A. Goldstein and J. F. Price, "On descent from local minima," Mathematics of Computation, vol. 25, no. 115, pp. 569–574, July 1971.
- [9] O. Grodzevich and O. Romanko, "Normalization and other topics in multi-objective optimization," Proceedings of the Fields-MITACS Industrial Problems Workshop, 2006. [Online]. Available: http://www.maths-in-industry.org/miis/233/
- [10] D. Huang, T. Allen, W. Notz, and N. Zeng, "Global optimization of stochastic black-box systems via sequential kriging meta-models," *Jour*nal of Global Optimization, vol. 34, pp. 441–466, 2006.
- [11] E. T. Jaynes, "Entropy and search-theory," in Maximum-Entropy and Bayesian Methods in Inverse Problems, C. R. Smith and J. W. T. Grandy, Eds. Springer, 1985, p. 443. [Online]. Available: http://bayes.wustl.edu/etj/articles/search.pdf

- [12] D. R. Jones, "A taxonomy of global optimization methods based on response surfaces," *Journal of Global Optimization*, vol. 21, pp. 345–383, December 2001. [Online]. Available: http://dx.doi.org/10.1023/A:1012771025575
- [13] J. Kennedy and R. Eberhart, "Particle swarm optimization," in Proc. of IEEE Intl. Conf. on Neural Networks, vol. 4, November 1995, pp. 1942–1948.
- [14] M. Li and P. Vitanyi, An Introduction to Kolmogorov Complexity and Its Applications, 2nd ed., ser. Texts in Computer Science. New York, NY, USA: Springer, 1997.
- [15] D. Lizotte, Τ. Bowling, and D. Schuurmans, Wang, М. optimization," "Gaussian process regression for in NIPS Workshop 2005 Value of Information in Inference, on Learning and Decision-Making, 2005.[Online]. Available: http://domino.research.ibm.com/comm/research_projects.nsf/pages/nips05workshop.index.html
- [16] D. J. C. MacKay, "Introduction to Gaussian processes," in *Neural Networks and Machine Learning*, ser. NATO ASI Series, C. M. Bishop, Ed. Kluwer Academic Press, 1998, pp. 133–166.
- [17] -"Information-based functions for objective acselection," Neural tive data Computation, vol. 4, 4. 590-604.1992.[Online]. Available: no. pp. http://www.mitpressjournals.org/doi/abs/10.1162/neco.1992.4.4.590
- [18] —, Information Theory, Inference, and Learning Algorithms. Cambridge University Press, 2003. [Online]. Available: http://www.inference.phy.cam.ac.uk/mackay/itila/
- [19] R. T. Marler and J. S. Arora, "Survey of multi-objective optimization methods for engineering," *Structural and Multidisciplinary Optimization*, vol. 26, no. 6, pp. 369–395, 2004. [Online]. Available: http://www.springerlink.com/openurl.asp?genre=article&id=doi:10.1007/s00158-003-0368-6
- [20] Y. Pan, L. Pavel, and T. Alpcan, "A system performance approach to OSNR optimization in optical networks," *IEEE Transactions on Communications*, vol. 58, no. 4, pp. 1193–1200, April 2010. [Online]. Available: papers/TComm_preprint_v7.pdf

- "The [21] K. В. Petersen and М. S. Pedersen. macookbook." 2008.trix October [Online]. Available: http://www2.imm.dtu.dk/pubdb/p.php?3274
- [22] J. G. Pierce, "A new look at the relation between information theory and search theory," Office of Naval Research, Arlington, VA, USA, Tech. Rep., June 1978. [Online]. Available: http://handle.dtic.mil/100.2/ADA063845
- [23] C. E. Rasmussen and C. K. I. Williams, Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning). The MIT Press, 2005.
- [24] R. Rutenbar, "Simulated annealing algorithms: an overview," IEEE Circuits and Devices Magazine, vol. 5, no. 1, pp. 19–26, January 1989.
- [25] N. V. Sahinidis, "Optimization under uncertainty: state-of-the-art and opportunities," Computers & Chemical Engineering, vol. 28, no. 6-7, pp. 971–983, June 2004, fOCAPO 2003 Special issue. [Online]. Available: http://www.sciencedirect.com/science/article/B6TFT-49YH97T-1/2/f15875aad97740410effc5264
- [26] B. Scholkopf and A. J. Smola, Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. Cambridge, MA, USA: MIT Press, 2001.
- [27] S. Seo, M. Wallat, T. Graepel, and K. Obermayer, "Gaussian process regression: active data selection and test point rejection," in *Proc. of IEEE-INNS-ENNS Intl. Joint Conf. on Neural Networks IJCNN 2000*, vol. 3, July 2000, pp. 241–246.
- [28] B. Settles, "Active learning literature survey," University of Wisconsin– Madison, Computer Sciences Technical Report 1648, 2009.
- [29] R. Srikant, The Mathematics of Internet Congestion Control, ser. Systems & Control: Foundations & Applications. Boston, MA: Birkhauser, 2004.
- [30] R. Tempo, E. W. Bai, and F. Dabbene, "Probabilistic robustness analysis: Explicit bounds for the minimum number of samples," Systems & Control Letters, vol. 30, no. 5, pp. 237–242, 1997. [Online]. Available: http://www.sciencedirect.com/science/article/B6V4X-3SP7DCD-4/2/3dc655107eff50f12b326ea10

- [31] R. Tempo, G. Calafiore, and F. Dabbene, Randomized Algorithms for Analysis and Control of Uncertain Systems. London, UK: Springer-Verlag, 2005.
- [32] K. R. Thompson, "Implementation of gaussian process models for nonlinear system identification," Ph.D. dissertation, University of Glasgow, Glasgow, Scotland, 2009.
- [33] M. E. Tipping, "Bayesian inference: An introduction to principles and practice in machine learning," in Advanced Lectures on Machine Learning, 2003, pp. 41–62. [Online]. Available: http://springerlink.metapress.com/openurl.asp?genre=article{&}issn=0302-9743{&}volume=317
- [34] R. Turner, M. P. Deisenroth, and C. E. Rasmussen, "State-space inference and learning with gaussian processes," in *Proc. of 13th Intl. Conf.* on Artificial Intelligence and Statistics (AISTATS), Chia Laguna Resort, Sardinia, Italy, May 2010.
- [35] Q. Zhu and J. Oommen, "On the optimal search problem: the case when the target distribution is unknown," in *Proc. of XVII Intl. Conf.* of Chilean Computer Science Society, 1997, pp. 268–277.