

BAYESIAN OPTIMIZATION

Roman Garnett

(To be published by) Cambridge University Press

CONTENTS

PREFACE	vii
NOTATION	xi
1 INTRODUCTION	1
1.1 Formalization of optimization	2
1.2 The Bayesian approach	5
2 GAUSSIAN PROCESSES	15
2.1 Definition and basic properties	16
2.2 Inference with exact and noisy observations	18
2.3 Overview of remainder of chapter	26
2.4 Joint Gaussian processes	26
2.5 Continuity	28
2.6 Differentiability	30
2.7 Existence and uniqueness of global maxima	33
2.8 Inference with non-Gaussian observations and constraints	35
3 MODELING WITH GAUSSIAN PROCESSES	45
3.1 The prior mean function	46
3.2 The prior covariance function	49
3.3 Notable covariance functions	51
3.4 Modifying and combining covariance functions	54
3.5 Modeling functions on high-dimensional domains	61
4 MODEL ASSESSMENT, SELECTION, AND AVERAGING	67
4.1 Models and model structures	68
4.2 Bayesian inference over parametric model spaces	70
4.3 Model selection via posterior maximization	73
4.4 Model averaging	74
4.5 Multiple model structures	78
4.6 Automating model structure search	81
5 DECISION THEORY FOR OPTIMIZATION	87
5.1 Introduction to Bayesian decision theory	89
5.2 Sequential decisions with a fixed budget	91
5.3 Cost and approximation of the optimal policy	99
5.4 Cost-aware optimization and termination as a decision	103
6 UTILITY FUNCTIONS FOR OPTIMIZATION	109
6.1 Expected utility of terminal recommendation	109
6.2 Cumulative reward	114
6.3 Information gain	115
6.4 Dependence on model of objective function	116
6.5 Comparison of utility functions	117

7 COMMON BAYESIAN OPTIMIZATION POLICIES	123
7.1 Example optimization scenario	124
7.2 Decision-theoretic policies	124
7.3 Expected improvement	127
7.4 Knowledge gradient	129
7.5 Probability of improvement	131
7.6 Mutual information and entropy search	135
7.7 Multi-armed bandits and optimization	141
7.8 Maximizing a statistical upper bound	145
7.9 Thompson sampling	148
7.10 Other ideas in policy construction	150
8 COMPUTING POLICIES WITH GAUSSIAN PROCESSES	157
8.1 Notation for objective function model	157
8.2 Expected improvement	158
8.3 Probability of improvement	167
8.4 Upper confidence bound	170
8.5 Approximate computation for one-step lookahead	171
8.6 Knowledge gradient	172
8.7 Thompson sampling	176
8.8 Mutual information with x^*	180
8.9 Mutual information with f^*	187
8.10 Averaging over a space of Gaussian processes	192
8.11 Alternatives to Gaussian processes	196
9 IMPLEMENTATION	201
9.1 Gaussian process inference	201
9.2 Optimizing acquisition functions	207
9.3 Starting and stopping optimization	210
10 THEORETICAL ANALYSIS	213
10.1 Regret	213
10.2 Useful function spaces for studying convergence	215
10.3 Relevant properties of covariance functions	220
10.4 Bayesian regret with observation noise	224
10.5 Worst-case regret with observation noise	231
10.6 The exact observation case	235
10.7 The effect of unknown hyperparameters	239
11 EXTENSIONS AND RELATED SETTINGS	243
11.1 Unknown observation costs	243
11.2 Constrained optimization and unknown constraints	247
11.3 Synchronous batch observations	250
11.4 Asynchronous observation with pending experiments	260
11.5 Multifidelity optimization	261
11.6 Multitask optimization	264
11.7 Multiobjective optimization	267
11.8 Gradient observations	274

11.9	Environmental variables	275
11.10	Incremental optimization of sequential procedures	276
11.11	Non-Gaussian observation models and active search	277
12	A BRIEF HISTORY OF BAYESIAN OPTIMIZATION	283
12.1	Historical precursors and optimal design	283
12.2	Sequential analysis and Bayesian experimental design	283
12.3	The rise of Bayesian optimization	285
12.4	Later rediscovery and development	286
12.5	Multi-armed bandits to infinite-armed bandits	288
12.6	What's next?	290
A	THE GAUSSIAN DISTRIBUTION	291
B	METHODS FOR APPROXIMATE BAYESIAN INFERENCE	297
C	GRADIENTS	303
D	ANNOTATED BIBLIOGRAPHY OF APPLICATIONS	309
	BIBLIOGRAPHY	327
	INDEX	345

PREFACE

My interest in Bayesian optimization began in 2007 at the start of my doctoral studies. I was frustrated that there seemed to be a Bayesian approach to every task I cared about, *except* optimization. Of course, as was often the case at that time (not to mention now!), I was mistaken in this belief, but one should never let ignorance impede inspiration.

Meanwhile, my labmate and soon-to-be frequent collaborator Mike Osborne had a fresh copy of RASMUSSEN and WILLIAMS's *Gaussian Processes for Machine Learning* and just would *not* stop talking about GPs at our lab meetings. Through sheer brute force of repetition, I slowly built a hand-wavy intuition for Gaussian processes – my mental model was the “sausage plot” – without even being sure about their precise definition. However, I was pretty sure that marginals were Gaussian (what else?), and one day it occurred to me that one could achieve Bayesian optimization by maximizing the probability of improvement. This was the algorithm I was looking for! In my excitement I shot off an email to Mike that kicked off years of fruitful collaboration:

Can I ask a dumb question about GPs? Let's say that I'm doing function approximation on an interval with a GP. So I've got this mean function $m(x)$ and a variance function $v(x)$. Is it true that if I pick a particular point x , then $p(f(x)) \sim \mathcal{N}(m(x), v(x))$? Please say yes.

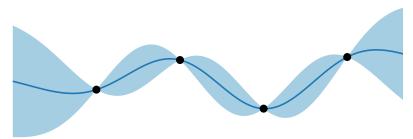
If this is true, then I think the idea of doing Bayesian optimization using GPs is, dare I say, trivial.

The hubris of youth!

Well, it turned out I was 45 years too late in proposing this algorithm,¹ and that it only seemed “trivial” because I had no appreciation for its theoretical foundation. However, truly great ideas are rediscovered many times, and my excitement did not fade. Once I developed a deeper understanding of Gaussian processes and Bayesian decision theory, I came to see them as a “Bayesian crank” I could turn to realize adaptive algorithms for *any* task. I have been repeatedly astonished to find that the resulting algorithms – seemingly by magic – *automatically* display intuitive emergent behavior as a result of their careful design. My goal with this book is to paint this grand picture. In effect, it is a gift to my former self: the book I wish I had in the early years of my career.

In the context of machine learning, Bayesian optimization is an ancient idea – KUSHNER's paper appeared only three years after the term “machine learning” was coined! Despite its advanced age, Bayesian optimization has been enjoying a period of revitalization and rapid progress over the past ten years. The primary driver of this renaissance has been advances in computation, which have enabled increasingly sophisticated tools for Bayesian modeling and inference.

Ironically, however, perhaps the most critical development was not Bayesian at all, but the rise of deep neural networks – another old idea



The first of many “sausage plots” to come.

¹ H. J. KUSHNER (1962). A Versatile Stochastic Model of a Function of Unknown and Time Varying Form. *Journal of Mathematical Analysis and Applications* 5(1):150–167.

² J. SNOEK et al. (2012). Practical Bayesian Optimization of Machine Learning Algorithms. *NeurIPS 2012*.

granted new life by modern computation. The extreme cost of training these models demands efficient routines for hyperparameter tuning, and in a timely and influential paper, SNOEK et al. demonstrated (dramatically!) that Bayesian optimization was up to the task.² Hyperparameter tuning proved to be a “killer app” for Bayesian optimization, and the surge of interest that followed has yielded a mountain of publications developing new algorithms and improving old ones, exploring countless variations on the basic setup, establishing theoretical guarantees on performance, and applying the framework to a huge range of domains.

Due to the nature of the computer science publication model, these recent developments are scattered across dozens of brief papers, and the pressure to establish novelty in a limited space can obscure the big picture in favor of minute details. This book aims to provide a self-contained and comprehensive introduction to Bayesian optimization, starting “from scratch” and carefully developing all the key ideas along the way. This bottom-up approach allows us to identify unifying themes in Bayesian optimization algorithms that may be lost when surveying the literature.

The intended audience is graduate students and researchers in machine learning, statistics, and related fields. However, it is also my sincere hope that practitioners from more distant fields wishing to harness the power of Bayesian optimization will also find some utility here.

For the bulk of the text, I assume the reader is comfortable with differential and integral calculus, probability, and linear algebra. On occasion the discussion will meander to more esoteric areas of mathematics, and these passages can be safely ignored and returned to later if desired. A good working knowledge of the Gaussian distribution is also essential, and I provide an abbreviated but sufficient introduction in appendix A.

The book is divided into three main parts. Chapters 2–4 cover theoretical and practical aspects of modeling with Gaussian processes. This class of models is the overwhelming favorite in the Bayesian optimization literature, and the material contained within is critical for several following chapters. It was daunting to write this material in light of the many excellent references already available, in particular the aforementioned *Gaussian Processes for Machine Learning*. However, I heavily biased the presentation in light of the needs of optimization, and even experts may find something new.

Chapters 5–7 develop the theory of sequential decision making and its application to optimization. Although this theory requires a model of the objective function and our observations of it, the presentation is agnostic to the choice of model and may be read independently from the preceding chapters on Gaussian processes.

These threads are unified in chapters 8–10, which discuss the particulars of Bayesian optimization with Gaussian process models. Chapters 8–9 cover details of computation and implementation, and chapter 10 discusses theoretical performance bounds on Bayesian optimization algorithms, where most results depend intimately on a Gaussian process model of the objective function or the associated reproducing kernel Hilbert space.

intended audience	
prerequisites	
chapters 2–4: modeling the objective function with Gaussian processes	
chapters 5–7: sequential decision making and policy building	
chapters 8–10: Bayesian optimization with Gaussian processes	

The nuances of some applications require modifications to the basic sequential optimization scheme that is the focus of the preceding chapters, and chapter 11 introduces several notable extensions to this basic setup. Each is systematically presented through the unifying lens of Bayesian decision theory to illustrate how one might proceed when facing a novel situation.

Finally, chapter 12 provides a brief and standalone history of Bayesian optimization. This was perhaps the most fun chapter for me to write, if only because it forced me to plod through old Soviet literature (in an actual library! what a novelty these days!). To my surprise I was able to antedate many Bayesian optimization policies beyond their commonly attested origin, including expected improvement, knowledge gradient, probability of improvement, and upper confidence bound. (A reader familiar with the literature may be surprised to learn the last of these was actually the first policy discussed by KUSHNER in his 1962 paper.) Despite my best efforts, there may still be stones left to be overturned before the complete history is revealed.

Dependencies between the main chapters are illustrated in the margin. There are two natural linearizations of the material. The first is the one I adopted and personally prefer, which covers modeling prior to decision making. However, one could also proceed in the other order, reading chapters 5–7 first, then looping back to chapter 2. After covering the material in these chapters (in either order), the remainder of the book can be perused at will. Logical partial paths through the book include:

- a minimal but self-contained introduction: chapters 1–2, 5–7
- a shorter introduction requiring leaps of faith: chapters 1 and 7
- a crash course on the underlying theory: chapters 1–2, 5–7, 10
- a head start on implementing a software package: chapters 1–9

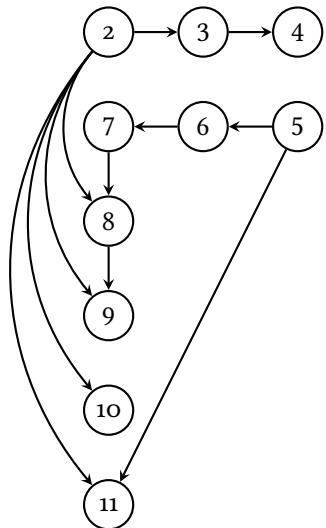
A reader already quite comfortable with Gaussian processes might wish to skip over chapters 2–4 entirely.

I struggled for some time over whether to include a chapter on applications. On the one hand, Bayesian optimization ultimately owes its popularity to its success in optimizing a growing and diverse set of difficult objectives. However, these applications often require extensive technical background to appreciate, and an adequate coverage would be tedious to write and tedious to read. As a compromise, I provide an annotated bibliography outlining the optimization challenges involved in notable domains of interest and pointing to studies where these challenges were successfully overcome with the aid of Bayesian optimization.

The sheer size of the Bayesian optimization literature – especially the output of the previous decade – makes it impossible to provide a complete survey of every recent development. This is especially true for the extensions discussed in chapter 11 and even more so for the bibliography on applications, where work has proliferated in myriad branching directions. Instead I settled for presenting what I considered to be the most important ideas and providing pointers to entry points

chapter 11: extensions

chapter 12: brief history of Bayesian optimization



A dependency graph for chapters 2–11. Chapter 1 is a universal dependency.

annotated bibliography of applications:
appendix D, p. 309

for the relevant literature. The reader should not read anything into any omissions; there is simply too much high-quality work to go around.

Additional information about the book, including a list of errata as they are discovered, may be found at the companion webpage:

bayesoptbook.com

I encourage the reader to report any errata or other issues to the companion GitHub repository for discussion and resolution:

github.com/bayesoptbook/bayesoptbook.github.io

Thank you!

Preparation of this manuscript was facilitated tremendously by numerous free and open source projects, and the creators, developers, and maintainers of these projects have my sincere gratitude. The manuscript was typeset in \LaTeX using the excellent and extremely flexible memoir class. The typeface is Linux Libertine. Figures were laid out in MATLAB and converted to TikZ/PGF/PGFPLOTS for further tweaking and typesetting via the `matlab2tikz` script. The colors used in figures were based on www.colorbrewer.org by Cynthia A. Brewer, and I endeavored to the best of my ability to ensure that the figures are colorblind friendly. The colormap used in heat maps is a slight modification of the Matplotlib viridis colormap where the “bright end” is pure white.

I would like to thank Eric Brochu, Nando de Freitas, Matt Hoffman, Frank Hutter, Mike Osborne, Bobak Shahriari, Jasper Snoek, Kevin Swersky, and Ziyu Wang, who jointly provided the activation energy for this undertaking. I would also like to thank Eytan Bakshy, Peter Frazier, Jake Gardner, Javier González, Frank Hutter, Mike Osborne, Matthias Poloczeck, Jonathan Scarlett, Bobak Shahriari, and Jasper Snoek for valuable discussions along the way, as well as David Tranah, Anna Scriven, and Abigail Walkington at Cambridge University Press for their support and patience. Special thanks are due to the students of two seminars run at Washington University covering the material in this book; their feedback was also instrumental in shaping the book’s content. This list is currently incomplete and there are a huge number of additional people to thank. If you are reading this draft, you are one of them!

Funding support was provided by the United States National Science Foundation under award number 1845434. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation.

This book took far more time than I initially anticipated, and I would especially like to thank my wife Marion and son Max (arg Max?) for their understanding and support during this long journey.

Roman Garnett
St. Louis, Missouri
September 2021

NOTATION

All vectors are column vectors and are denoted in lowercase bold: $\mathbf{x} \in \mathbb{R}^d$.
 Matrices are denoted in uppercase bold: \mathbf{A} . We adopt the “numerator layout” convention for matrix calculus: the derivative of a vector by a scalar is a (column) vector, whereas the derivative of a scalar by a vector is a row vector. This results in the chain rule proceeding from left-to-right – if the vector $\mathbf{x}(\theta)$ depends on a scalar parameter θ , then for a function $f(\mathbf{x})$, we have:

$$\frac{\partial f}{\partial \theta} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \theta}.$$

When an indicator function is required, we use the Iverson bracket notation. For a statement s , we have:

$$[s] = \begin{cases} 1 & \text{if } s \text{ is true;} \\ 0 & \text{otherwise.} \end{cases}$$

The statement may depend on a parameter: $[x \in A]$, $[x \geq 0]$, etc.

Logarithms are taken with respect to their natural base, e . Quantities in log units such as log likelihoods or entropy thus have units of *nats*, the base- e analogue of the more familiar base-2 bits.

vectors and matrices
 matrix calculus convention

chain rule

indicator functions

logarithms
 nats

SYMBOLS WITH IMPLICIT DEPENDENCE ON LOCATION

There is one notational innovation in this book compared with the Gaussian process and Bayesian optimization literature at large: we make heavy use of symbols for quantities that depend *implicitly* on a putative (arbitrary) input location x . Most importantly, to refer to the value of an objective function f at a given location x , we introduce the symbol $\phi = f(x)$. This avoids clash with the name of the function itself, f , while avoiding an extra layer of brackets. We use this scheme throughout the book, including variations such as:

$$\phi' = f(x'); \quad \phi = f(\mathbf{x}); \quad \gamma = g(x); \quad \text{etc.}$$

To refer to the outcome of a (possibly inexact) measurement at x , we use the symbol y ; the distribution of y presumably depends on ϕ .

We also allocate symbols to describe properties of the marginal predictive distributions for the objective function value ϕ and observed value y , all of which also have implicit dependence on x . These appear in the table below.

COMPREHENSIVE LIST OF SYMBOLS

A list of important symbols appears on the following pages, arranged roughly in alphabetical order.

NOTATION

symbol	description
\equiv	identical equality of functions; for example, for a constant c , $f \equiv c$ is a constant function
∇	gradient operator
\emptyset	termination option: the action of immediately terminating optimization
\prec	either Pareto dominance or the Löwner order: for symmetric \mathbf{A}, \mathbf{B} , $\mathbf{A} \prec \mathbf{B}$ if and only if $\mathbf{B} - \mathbf{A}$ is positive definite
$x \sim p(x)$	is sampled according to: x is a realization of a random variable with probability density $p(x)$
$\bigsqcup_i \mathcal{X}_i$	disjoint union of $\{\mathcal{X}_i\}$: $\bigsqcup_i \mathcal{X}_i = \bigcup_i \{(x, i) \mid x \in \mathcal{X}_i\}$
$ \mathbf{A} $	determinant of square matrix \mathbf{A}
$ x $	Euclidean norm of vector x ; $ x - y $ is thus the Euclidean distance between vectors x and y
$\ f\ _{\mathcal{H}_K}$	norm of function f in reproducing kernel Hilbert space \mathcal{H}_K
\mathbf{A}^{-1}	inverse of square matrix \mathbf{A}
\mathbf{x}^\top	transpose of vector x
$\mathbf{0}$	vector or matrix of zeros
\mathcal{A}	action space for a decision
$\alpha(x; \mathcal{D})$	acquisition function evaluating x given data \mathcal{D}
$\alpha_\tau(x; \mathcal{D})$	expected marginal gain in $u(\mathcal{D})$ after observing at x then making $\tau - 1$ additional optimal observations given the outcome
$\alpha_\tau^*(\mathcal{D})$	value of \mathcal{D} with horizon τ : expected marginal gain in $u(\mathcal{D})$ from τ additional optimal observations
α_{EI}	expected improvement
α_{f^*}	mutual information between y and f^*
α_{KG}	knowledge gradient
α_{PI}	probability of improvement
α_{x^*}	mutual information between y and x^*
α_{UCB}	upper confidence bound
α_{TS}	Thompson sampling “acquisition function:” a draw $f \sim p(f \mid \mathcal{D})$
β	confidence parameter in Gaussian process upper confidence bound policy
$\beta(\mathbf{x}; \mathcal{D})$	batch acquisition function evaluating \mathbf{x} given data \mathcal{D} ; may have modifiers analogous to α
\mathbf{C}	prior covariance matrix of observed values \mathbf{y} : $\mathbf{C} = \text{cov}[\mathbf{y}]$
$c(\mathcal{D})$	cost of acquiring data \mathcal{D}
$\text{chol } \mathbf{A}$	Cholesky decomposition of positive definite matrix \mathbf{A} : if $\Lambda = \text{chol } \mathbf{A}$, then $\mathbf{A} = \Lambda \Lambda^\top$
$\text{corr}[\omega, \psi]$	correlation of random variables ω and ψ ; with a single argument, $\text{corr}[\omega] = \text{corr}[\omega, \omega]$
$\text{cov}[\omega, \psi]$	covariance of random variables ω and ψ ; with a single argument, $\text{cov}[\omega] = \text{cov}[\omega, \omega]$
\mathcal{D}	set of observed data, $\mathcal{D} = (\mathbf{x}, \mathbf{y})$
$\mathcal{D}', \mathcal{D}_1$	set of observed data after observing at x : $\mathcal{D}' = \mathcal{D} \cup \{(x, y)\} = (\mathbf{x}', \mathbf{y}')$
\mathcal{D}_τ	set of observed data after τ observations
$D_{\text{KL}}[p \parallel q]$	Kullback–Leibler divergence between distributions with probability densities p and q
$\Delta(x, y)$	marginal gain in utility after acquiring observation (x, y) : $\Delta(x, y) = u(\mathcal{D}') - u(\mathcal{D})$
$\delta(x - a)$	Dirac delta distribution on x with point mass at a
$\text{diag } \mathbf{x}$	diagonal matrix with diagonal \mathbf{x}
\mathbb{E}, \mathbb{E}_x	expectation, expectation with respect to x
ε	measurement error associated with an observation at x : $\varepsilon = y - \phi$
f	objective function; $f: \mathcal{X} \rightarrow \mathbb{R}$
$f _{\mathcal{Y}}$	the restriction of f onto the subdomain $\mathcal{Y} \subset \mathcal{X}$
f^*	globally maximal value of the objective function: $f^* = \max f$
γ_τ	information capacity of an observation process given τ iterations

symbol	description
$\mathcal{GP}(f; \mu, K)$	Gaussian process on f with mean function μ and covariance function K
\mathcal{H}_K	reproducing kernel Hilbert space associated with kernel K
$\mathcal{H}_K[B]$	ball of radius B in \mathcal{H}_K : $\{f \mid \ f\ _{\mathcal{H}_K} \leq B\}$
$H[\omega]$	discrete or differential entropy of random variable ω
$H[\omega \mid \mathcal{D}]$	discrete or differential of random variable ω after conditioning on \mathcal{D}
$I(\omega; \psi)$	mutual information between random variables ω and ψ
$I(\omega; \psi \mid \mathcal{D})$	mutual information between random variables ω and ψ after conditioning on \mathcal{D}
\mathbf{I}	identity matrix
K	prior covariance function: $K = \text{cov}[f]$
$K_{\mathcal{D}}$	posterior covariance function given data \mathcal{D} : $K_{\mathcal{D}} = \text{cov}[f \mid \mathcal{D}]$
K_M	Matérn covariance function
K_{SE}	squared exponential covariance function
κ	cross covariance between f and observed values \mathbf{y} : $\kappa(x) = \text{cov}[\mathbf{y}, \phi \mid x]$
ℓ	either a length-scale parameter or the lookahead horizon
λ	output-scale parameter
\mathcal{M}	space of models indexed by the hyperparameter vector θ
\mathbf{m}	prior expected value of observed values \mathbf{y} , $\mathbf{m} = \mathbb{E}[\mathbf{y}]$
μ	either the prior mean function, $\mu = \mathbb{E}[f]$, or the predictive mean of ϕ : $\mu = \mathbb{E}[\phi \mid x, \mathcal{D}] = \mu_{\mathcal{D}}(x)$
$\mu_{\mathcal{D}}$	posterior mean function given data \mathcal{D} : $\mu_{\mathcal{D}} = \mathbb{E}[f \mid \mathcal{D}]$
$\mathcal{N}(\phi; \mu, \Sigma)$	multivariate normal distribution on ϕ with mean vector μ and covariance matrix Σ
\mathbf{N}	measurement error covariance corresponding to observed values \mathbf{y}
\mathcal{O}	“big O” notation: for nonnegative functions f, g of τ , $f = \mathcal{O}(g)$ if f/g is bounded as $\tau \rightarrow \infty$
\mathcal{O}^*	as above with logarithmic factors suppressed: $f = \mathcal{O}^*(g)$ if $f(\tau)(\log \tau)^k = \mathcal{O}(g)$ for some k
p	probability density
q	either an approximation to probability density p or a quantile function
$\Phi(z)$	standard normal cumulative density function: $\Phi(z) = \int_{-\infty}^z \phi(z) dz$
ϕ	value of the objective function at x : $\phi = f(x)$
$\phi(z)$	standard normal probability density function: $\phi(z) = (\sqrt{2\pi})^{-1} \exp(-\frac{1}{2}z^2)$
Pr	probability
\mathbb{R}	set of real numbers
R_{τ}	cumulative regret after τ iterations
$\bar{R}_{\tau}[B]$	worst-case cumulative regret after τ iterations on the RKHS ball $\mathcal{H}_K[B]$
r_{τ}	simple regret after τ iterations
$\bar{r}_{\tau}[B]$	worst-case simple regret after τ iterations on the RKHS ball $\mathcal{H}_K[B]$
\mathbf{P}	a correlation matrix
ρ	a scalar correlation
ρ_{τ}	instantaneous regret on iteration τ
s^2	predictive variance of y ; for additive Gaussian noise, $s^2 = \text{var}[y \mid x, \mathcal{D}] = \sigma^2 + \sigma_n^2$
Σ	a covariance matrix, usually the Gram matrix associated with \mathbf{x} : $\Sigma = K_{\mathcal{D}}(\mathbf{x}, \mathbf{x})$
σ^2	predictive variance of ϕ : $\sigma^2 = K_{\mathcal{D}}(x, x)$
σ_n^2	variance of measurement error at x : $\sigma_n^2 = \text{var}[\varepsilon \mid x]$
$\text{std}[\omega]$	standard deviation of random variable ω
$\mathcal{T}(\phi; \mu, \sigma^2, v)$	Student- t distribution on ϕ with v degrees of freedom, mean μ , and variance σ^2
$\mathcal{TN}(\phi; \mu, \sigma^2, I)$	truncated normal distribution, $\mathcal{N}(\phi; \mu, \sigma^2)$ truncated to interval I
τ	either decision horizon (in the context of decision making) or number of optimization iterations passed (in the context of asymptotic analysis)

NOTATION

symbol	description
θ	vector of hyperparameters indexing a model space \mathcal{M}
$\text{tr } \mathbf{A}$	trace of square matrix \mathbf{A}
$u(\mathcal{D})$	utility of data \mathcal{D}
$\text{var}[\omega]$	variance of random variable ω
x	putative input location of the objective function
\mathbf{x}	either a sequence of observed locations $\mathbf{x} = \{x_i\}$ or (when the distinction is important) a vector-valued input location
x^*	a location attaining the globally maximal value of f : $x^* \in \arg \max f; f(x^*) = f^*$
\mathcal{X}	domain of objective function
y	value resulting from an observation at x
\mathbf{y}	observed values resulting from observations at locations \mathbf{x}
z	z -score of measurement y at x : $z = (y - \mu)/s$

1

INTRODUCTION

Optimization is an innate human behavior. On an individual level, we all strive to better ourselves and our surroundings. On a collective level, societies struggle to allocate limited resources seeking to improve the welfare of their members, and optimization has been an engine of societal progress since the domestication of crops through selective breeding over 12 000 years ago – an effort that continues to this day.

Given its pervasiveness, it should perhaps not be surprising that optimization is also *difficult*. While searching for an optimal design, we must spend – sometimes quite significant – resources evaluating suboptimal alternatives along the way. This observation compels us to seek methods of optimization that, when necessary, can carefully allocate resources to identify optimal parameters as efficiently as possible. This is the goal of mathematical optimization.

Since the 1960s, the statistics and machine learning communities have refined a *Bayesian* approach to optimization that we will develop and explore in this book. Bayesian optimization routines rely on a statistical model of the objective function, whose beliefs guide the algorithm in making the most fruitful decisions. These models can be quite sophisticated, and maintaining them throughout optimization may entail significant cost of its own. However, the reward for doing so is unparalleled sample efficiency. For this reason, Bayesian optimization has found a niche in optimizing objectives that:

- are costly to compute, precluding exhaustive evaluation,
- lack a useful expression, causing them to function as “black boxes,”
- cannot be evaluated exactly, but only through some indirect or noisy mechanism, and/or
- offer no efficient mechanism for estimating their gradient.

Let us consider an example setting motivating the machine learning community’s recent interest in Bayesian optimization. Consider a data scientist crafting a complex machine learning model – say a deep neural network – from training data. To ensure success, the scientist must carefully tune the model’s hyperparameters, including the network architecture and details of the training procedure, which have massive influence on performance. Unfortunately, effective settings can only be identified via trial-and-error: by training several networks with different settings and evaluating their performance on a validation dataset.

The search for the best hyperparameters is of course an exercise in optimization. Mathematical optimization has been under continual development for centuries, and numerous off-the-shelf procedures are available. However, these procedures usually make assumptions about the objective function that may not always be valid. For example, we might assume that the objective is cheap to evaluate, that we can easily compute its gradient, or that it is convex, allowing us to reduce from global to local optimization.

INTRODUCTION

In hyperparameter tuning, all of these assumptions are invalid. Training a deep neural network can be extremely expensive in terms of both time and energy. When some hyperparameters are discrete – as many features of network architecture naturally are – the gradient does not even *exist*. Finally, the mapping from hyperparameters to performance may be highly complex and multimodal, so local refinement may not yield an acceptable result.

- 1 R. TURNER et al. (2021). Bayesian Optimization is Superior to Random Search for Machine Learning Hyperparameter Tuning: Analysis of the Black-Box Optimization Challenge 2020. *Proceedings of the Neurips 2020 Competition and Demonstration Track*.
- 2 A. GELMAN and A. VEHTARI (2021). What are the most important statistical ideas of the past 50 years? *Journal of the American Statistical Association*.

annotated bibliography of applications:
appendix D, p. 309

outline and reading guide: p. viii

The Bayesian approach to optimization allows us to relax all of these assumptions when necessary, and Bayesian optimization algorithms can deliver impressive performance even when optimizing complex “black box” objectives under severely limited observation budgets. Bayesian optimization has proven successful in settings spanning science, engineering, and beyond, including of course hyperparameter tuning.¹ In light of this broad success, GELMAN and VEHTARI identified adaptive decision analysis – and Bayesian optimization in particular – as one of the eight most important statistical ideas of the past 50 years.²

Covering all these applications and their nuances could easily fill a separate volume (although we do provide an overview of some important application domains in an annotated bibliography), so in this book we will settle for developing the mathematical foundation of Bayesian optimization underlying its success. In the remainder of this chapter we will lay important groundwork for this discussion. We will first establish the precise formulation of optimization we will consider and important conventions of our presentation, then outline and illustrate the key aspects of the Bayesian approach. The reader may find an outline of and reading guide for the chapters to come in the preface.

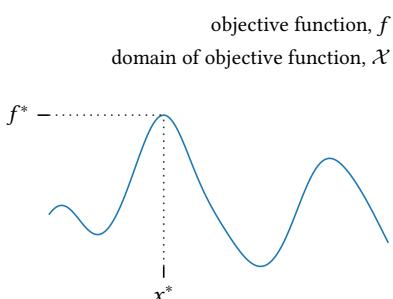
1.1 FORMALIZATION OF OPTIMIZATION

Throughout this book we will consider a simple but flexible formulation of sequential global optimization outlined below. There is nothing inherently Bayesian about this model, and countless solutions are possible.

We begin with a real-valued objective function defined on some domain \mathcal{X} ; $f: \mathcal{X} \rightarrow \mathbb{R}$. We make no assumptions regarding the nature of the domain. In particular, it need not be Euclidean but might instead, for example, comprise a space of complex structured objects. The goal of optimization is to systematically search the domain for a point x^* attaining the globally maximal value f^* :³

$$x^* \in \arg \max_{x \in \mathcal{X}} f(x); \quad f^* = \max_{x \in \mathcal{X}} f(x) = f(x^*). \quad (1.1)$$

Before we proceed, we note that our focus on maximization rather than minimization is entirely arbitrary; the author simply judges maximization to be the more optimistic choice. If desired, we can freely transform one problem to the other by negating the objective function. We caution the reader that some translation may be required when comparing expressions derived here to what may appear in parallel texts focusing on minimization.



- 3 A skeptical reader may object that, without further assumptions, a global maximum may not exist at all! We will sidestep this issue for now and pick it up again in § 2.7, p. 34.

```

input: initial dataset  $\mathcal{D}$                                 ▶ can be empty
repeat
     $x \leftarrow \text{POLICY}(\mathcal{D})$           ▶ select the next observation location
     $y \leftarrow \text{OBSERVE}(x)$                 ▶ observe at the chosen location
     $\mathcal{D} \leftarrow \mathcal{D} \cup \{(x, y)\}$       ▶ update dataset
until termination condition reached      ▶ e.g., budget exhausted
return  $\mathcal{D}$ 

```

In a significant departure from classical mathematical optimization, we do not require that the objective function have a known functional form or even be computable directly. Rather, we only require access to a mechanism revealing *some* information about the objective function at identified points on demand. By amassing sufficient information from this mechanism, we may hope to infer the solution to (1.1). Avoiding the need for an explicit expression for f allows us to consider so-called “black box” optimization, where a system is optimized through indirect measurements of its quality. This is one of the greatest strengths of Bayesian optimization.⁴

Optimization policy

Directly solving for the location of global optima is infeasible except in exceptional circumstances. The tools of traditional calculus are virtually powerless in this setting; for example, enumerating and classifying every stationary point in the domain would be tedious at best and perhaps even impossible. Mathematical optimization instead takes an indirect approach: we design a sequence of experiments to probe the objective function for information that, we hope, will reveal the solution to (1.1).

The iterative procedure in algorithm 1.1 formalizes this process. We begin with an initial (possibly empty) dataset \mathcal{D} that we grow incrementally through a sequence of observations of our design. In each iteration, an *optimization policy* inspects the available data and selects a point $x \in \mathcal{X}$ where we make our next observation.⁵ This action in turn reveals a corresponding value y provided by the system under study. We append the newly observed information to our dataset and finally decide whether to continue with another observation or terminate and return the current data. When we inevitably do choose to terminate, the returned data can be used by an external consumer as desired, for example to inform a subsequent decision.

We place no restrictions on how an optimization policy is implemented beyond mapping an arbitrary dataset to some point in the domain. A policy may be deterministic or stochastic, as demonstrated respectively by the prototypical examples of grid search and random search. In fact, these popular policies are *nonadaptive* and completely ignore the observed data. However, when observations only come at significant cost, we will naturally prefer policies that adapt their behavior in light of evolving information. The primary challenge in optimization

Algorithm 1.1: Sequential optimization.

⁴ Of course, we do not *require* but merely *allow* that the objective function act as a black box. Access to a closed-form expression does not preclude a Bayesian approach!

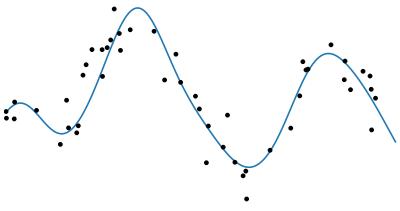
⁵ Here “policy” has the same meaning as in other decision-making contexts: it maps our state (indexed by our data, \mathcal{D}) to an action (the location of our next observation, x).

terminal recommendations: § 5.1, p. 90

INTRODUCTION

is designing policies that can *rapidly* optimize a broad class of objective functions, and intelligent policy design will be our focus for the majority of this book.

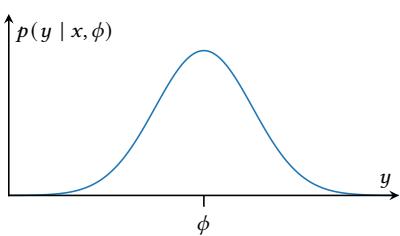
Observation model



Inexact observations of an objective function corrupted by additive noise.

measured value, y
observation location, x
objective function value, $\phi = f(x)$

conditional independence of observations given objective values



Additive Gaussian noise: the distribution of the value y observed at x is Gaussian, centered on the objective function value ϕ .

observation noise scale, σ_n
heteroskedastic noise: § 2.2, p. 25

For optimization to be feasible, the observations we obtain must provide information about the objective function that can guide our search and in aggregate determine the solution to (1.1). A near-universal assumption in mathematical optimization is that observations yield *exact* evaluations of the objective function at our chosen locations. However, this assumption is unduly restrictive: many settings feature inexact measurements due to noisy sensors, imperfect simulation, or statistical approximation. A typical example featuring additive observation noise is shown in the margin. Although the objective function is not observed directly, the noisy measurements nonetheless constrain the plausible options due to strong dependence on the objective.

We thus relax the assumption of exact observation and instead assume that observations are realized by a stochastic mechanism depending on the objective function. Namely, we assume that the value y resulting from an observation at some point x is distributed according to an observation model depending on the underlying objective function value $\phi = f(x)$:

$$p(y | x, \phi). \quad (1.2)$$

Through judicious design of the observation model, we may consider a wide range of observation mechanisms.

As with the optimization policy, we do not make any assumptions about the nature of the observation model, save one. Unless otherwise mentioned, we assume that a set of *multiple* measurements y are conditionally independent given the corresponding observation locations x and objective function values $\phi = f(x)$:

$$p(y | x, \phi) = \prod_i p(y_i | x_i, \phi_i). \quad (1.3)$$

This is not strictly necessary but is overwhelmingly common in practice and will simplify our presentation considerably.

One particular observation model will enjoy most of our attention in this book: *additive Gaussian noise*. Here we model the value y observed at x as

$$y = \phi + \varepsilon,$$

where ε represents measurement error. Errors are assumed to be Gaussian distributed with mean zero, implying a Gaussian observation model:

$$p(y | x, \phi, \sigma_n) = \mathcal{N}(y; \phi, \sigma_n^2). \quad (1.4)$$

Here the observation noise scale σ_n may optionally depend on x , allowing us to model both homoskedastic or heteroskedastic errors.

If we take the noise scale to be identically zero, we recover the special case of exact observation, where we simply have $y = \phi$ and the observation model collapses to a Dirac delta distribution:

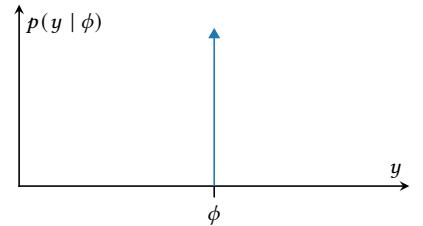
$$p(y | \phi) = \delta(y - \phi).$$

Although not universally applicable, many settings do feature exact observations such as optimizing the output of a deterministic computer simulation. We will sometimes consider the exact case separately as some results simplify considerably in the absence of measurement error.

We will focus on additive Gaussian noise as it is a reasonably faithful model for many systems and offers considerable mathematical convenience. This observation model will be most prevalent in our discussion on Gaussian processes in the next three chapters and on the explicit computation of Bayesian optimization policies with this model class in chapter 8. However, the general methodology we will build in the remainder of this book is not contingent on this choice, and we will occasionally address alternative observation mechanisms.

Termination

The final decision we make in each iteration of optimization is whether to terminate immediately or continue with another observation. As with the optimization policy, we do not assume any particular mechanism by which this decision is made. Termination may be deterministic – such as stopping after reaching a certain optimization goal or exhausting a preallocated observation budget – or stochastic, and may optionally depend on the observed data. In many cases, the time of termination may in fact not be under the control of the optimization routine at all but instead decided by an external agent. However, we will also consider scenarios where the optimization procedure can dynamically choose when to return based upon inspection of the available data.



Exact observations: every value measured equals the corresponding function value, corresponding to a Dirac delta observation model.

inference with non-Gaussian observations:
§ 2.8, p. 35

optimization with non-Gaussian observations: § 11.11, p. 277

optimal termination: § 5.4, p. 103
practical termination: § 9.3, p. 210

1.2 THE BAYESIAN APPROACH

Bayesian optimization does not refer to one particular algorithm but rather to a philosophical approach to optimization grounded in Bayesian inference from which an extensive family of algorithms have been derived. Although these algorithms display significant diversity in their details, they are bound by common themes in their design.

Optimization is fundamentally a sequence of decisions: in each iteration, we must choose where to make our next observation and then whether to terminate depending on the outcome. As the outcomes of these decisions are governed by the system under study and outside our control, the success of optimization rests entirely on effective decision making.

Increasing the difficulty of these decisions is that they must be made under *uncertainty*, as it is impossible to know the outcome of an observation before making it. The optimization policy must therefore design each

INTRODUCTION

observation with some measure of faith that the outcome will ultimately prove beneficial and justify the cost of obtaining it. The sequential nature of optimization further compounds the weight of this uncertainty, as the outcome of each observation not only has an immediate impact, but also forms the basis on which all future decisions are made. Developing an effective policy requires somehow addressing this uncertainty.

The Bayesian approach systematically relies on probability and Bayesian inference to reason about the uncertain quantities arising during optimization. This critically includes the objective function itself, which is treated as a random variable to be inferred in light of our prior expectations and any available data. In Bayesian optimization, this belief then takes an active role in decision making by guiding the optimization policy, which may evaluate the merit of a proposed observation location according to our belief about the value we might observe. We introduce the key ideas of this process with examples below, starting with a refresher on Bayesian inference.

Bayesian inference

⁶ The literature is vast. The following references are excellent, but no list can be complete:

D. J. C. MACKAY (2003). *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press.

A. O'HAGAN and J. FORSTER (2004). *Kendall's Advanced Theory of Statistics*. Vol. 2B: Bayesian Inference. Arnold.

J. O. BERGER (1985). *Statistical Decision Theory and Bayesian Analysis*. Springer-Verlag.

To frame the following discussion, we offer a quick overview of Bayesian inference as a reminder to the reader. This introduction is far from complete, but there are numerous excellent references available.⁶

Bayesian inference is a framework for inferring uncertain features of a system of interest from observations grounded in the laws of probability. To illustrate the basic ideas, we may begin by identifying some unknown feature of a given system that we wish to reason about. In the context of optimization, this might represent, for example, the value of the objective function at a given location, or the location x^* or value f^* of the global optimum (1.1). We will take the first of these as a running example: inferring about the value of an objective function at some arbitrary point x , $\phi = f(x)$. We will shortly extend this example to inference about the *entire* objective function.

In the Bayesian approach to inference, *all* unknown quantities are treated as random variables. This is a powerful convention as it allows us to represent beliefs about these quantities with probability distributions reflecting their plausible values. Inference then takes the form of an inductive process where these beliefs are iteratively refined in light of observed data by appealing to probabilistic identities.

As with any induction, we must start somewhere. Here we begin with a so-called *prior distribution* (or simply *prior*) $p(\phi | x)$, which encodes what we consider to be plausible values for ϕ before observing any data.⁷ The prior distribution allows us to inject our knowledge about and experience with the system of interest into the inferential process, saving us from having to begin “from scratch” or entertain patently absurd possibilities. The left panel of figure 1.1 illustrates a prior distribution for our example, indicating support over a range of values.

Once a prior has been established, the next stage of inference is to refine our initial beliefs in light of observed data. Suppose in our

prior distribution, $p(\phi | x)$

⁷ Here we assume the location of interest x is known, hence our conditioning the prior on its value.

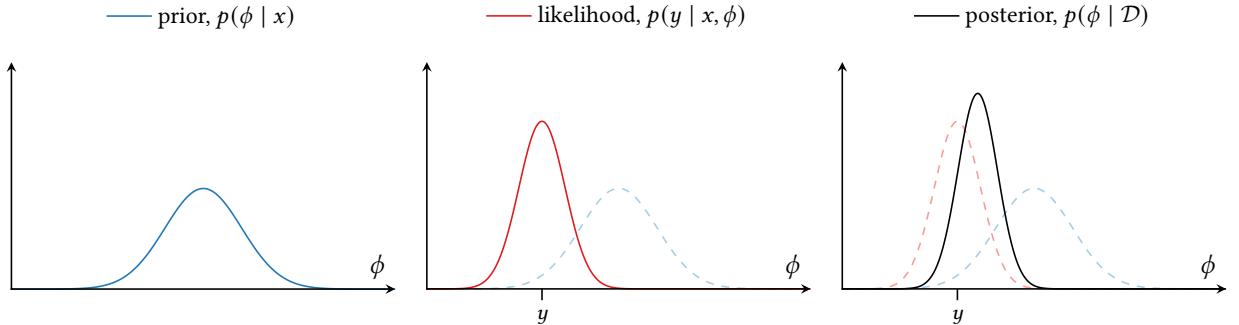


Figure 1.1: Bayesian inference for an unknown function value $\phi = f(x)$. Left: a prior distribution over ϕ ; middle: the likelihood of the marked observation y according to an additive Gaussian noise observation model (1.4) (prior shown for reference); right: the posterior distribution in light of the observation and the prior (prior and likelihood shown for reference).

example we make an observation of the objective function at x , revealing a measurement y . In our model of optimization, the distribution of this measurement is assumed to be determined by the value of interest ϕ through the observation model $p(y | x, \phi)$ (1.2). In the context of Bayesian inference, a distribution explaining the observed values (here y) in terms of the values of interest (here ϕ) is known as a *likelihood function* or simply a *likelihood*. The middle panel of figure 1.1 show the likelihood – as a function of ϕ – for a given measurement y , here assumed to be generated by additive Gaussian noise (1.4).

Finally, given the observed value y , we may derive the updated *posterior distribution* (or simply *posterior*) of ϕ by appealing to Bayes' theorem:

$$p(\phi | x, y) = \frac{p(\phi | x) p(y | x, \phi)}{p(y | x)}. \quad (1.5)$$

The posterior is proportional to the prior weighted by the likelihood of the observed value. The denominator is a constant with respect to ϕ that ensures normalization:

$$p(y | x) = \int p(y | x, \phi) p(\phi | x) d\phi. \quad (1.6)$$

The right panel of figure 1.1 shows the posterior resulting from the measurement in the middle panel. The posterior represents a compromise between our experience (encoded in the prior) and the information contained in the data (encoded in the likelihood).

Throughout this book we will use the catchall notation \mathcal{D} to represent all the information influencing a posterior belief; here the relevant information is $\mathcal{D} = (x, y)$, and the posterior distribution is then $p(\phi | \mathcal{D})$.

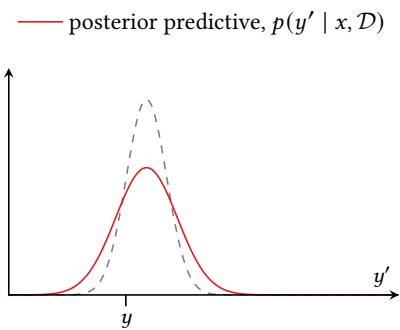
As mentioned previously, Bayesian inference is an inductive process whereby we can continue to refine our beliefs through additional observation. At this point, the induction is trivial: to incorporate a new

likelihood function (observation model),
 $p(y | x, \phi)$

posterior distribution, $p(\phi | x, y)$

data informing posterior belief, \mathcal{D}

INTRODUCTION



Posterior predictive distribution for a repeated measurement at x for our running example. The location of our first measurement y and the posterior distribution of ϕ are shown for reference. There is more uncertainty in y' than ϕ due to the effect of observation noise.

⁸ This expression takes the same form as (1.6), which is simply the (prior) predictive distribution evaluated at the actual observed value.

observation, what was our posterior serves as the prior in the context of the new information, and multiplying by the likelihood and renormalizing yields a new posterior. We may continue in this manner as desired.

The posterior distribution is not usually the end result of Bayesian inference but rather a springboard enabling follow-on tasks such as prediction or decision making, both of which are integral to Bayesian optimization. To address the former, suppose that after deriving the posterior (1.5), we wish to predict the result of an independent, *repeated* noisy observation at x, y' . Treating the outcome as a random variable, we may derive its distribution by integrating our posterior belief about ϕ against the observation model (1.2):⁸

$$p(y' | x, \mathcal{D}) = \int p(y' | x, \phi) p(\phi | x, \mathcal{D}) d\phi; \quad (1.7)$$

this is known as the *posterior predictive distribution* for y' . By integrating over all possible values of ϕ weighted by their plausibility, the posterior predictive distribution naturally accounts for uncertainty in the unknown objective function value; see the figure in the margin.

The Bayesian approach to decision making also relies on a posterior belief about unknown features affecting the outcomes of our decisions, as we will discuss shortly.

Bayesian inference of the objective function

At the heart of any Bayesian optimization routine is a probabilistic belief over the objective function. This takes the form of a *stochastic process*, a probability distribution over an infinite collection of random variables – here the objective function value at every point. The reasoning behind this inference is, in essence, the same as our single-point example above.

We begin by encoding any assumptions we may have about the objective function, such as smoothness or other features, in a *prior process* $p(f)$. Conveniently, we can specify a stochastic process via the distribution of the function values ϕ corresponding to an arbitrary *finite* set of locations \mathbf{x} :

$$p(\phi | \mathbf{x}). \quad (1.8)$$

The family of *Gaussian processes* – where these finite-dimensional distributions are multivariate Gaussian – is especially convenient and widely used in Bayesian optimization. We will explore this model class in depth in the following three chapters; here we provide a motivating illustration.

Figure 1.2 shows a Gaussian process prior on a one-dimensional objective function, constructed to reflect a minimal set of assumptions we will elaborate on later in the book:

- that the objective function is smooth (that is, infinitely differentiable),
- that correlations among function values have a characteristic scale, and
- that the function's expected behavior does not depend on location (that is, the prior process is *stationary*).

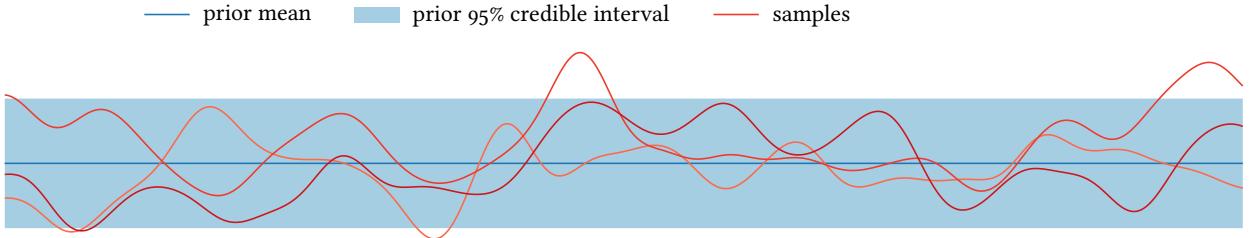


Figure 1.2: An example prior process for an objective defined on an interval. We illustrate the marginal belief at every location with its mean and a 95% credible interval and also show three example functions sampled from the prior process.

We summarize the marginal belief of the model, for each point in the domain showing the prior mean (dark blue) and 95% credible interval (light blue) for the corresponding function value. We also show three functions sampled from the prior process, each exhibiting the assumed behavior. We encourage the reader to become comfortable with this plotting convention, as we will use it throughout this book. In particular we eschew axis labels, as they are always the same: the horizontal axis represents the domain \mathcal{X} and the vertical axis the function value. Further, we do not mark units on axes to stress relative rather than absolute behavior, as scale is arbitrary in this illustration.

We can encode a vast array of information into the prior process and can model significantly more complex structure than in this simple example. We will explore the world of possibilities in chapter 3, including interaction at different scales, nonstationarity, low intrinsic dimensionality, and more.

With the prior process in hand, suppose we now make a set of observations at some locations \mathbf{x} , revealing corresponding values \mathbf{y} ; we aggregate this information into a dataset $\mathcal{D} = (\mathbf{x}, \mathbf{y})$. Bayesian inference accounts for these observations by forming the *posterior process* $p(f | \mathcal{D})$.

The derivation of the posterior process can be understood as a two-stage process. First we consider the impact of the data on the corresponding function values ϕ alone (1.5):

$$p(\phi | \mathcal{D}) \propto p(\phi | \mathbf{x}) p(\mathbf{y} | \mathbf{x}, \phi). \quad (1.9)$$

The quantities on the right-hand side are known: the first term is given by the prior process (1.8), and the second by the observation model (1.3), which serves the role of a likelihood. We now extend the posterior on ϕ to all of f :⁹

$$p(f | \mathcal{D}) = \int p(f | \mathbf{x}, \phi) p(\phi | \mathcal{D}) d\phi. \quad (1.10)$$

The posterior encapsulates our belief regarding the objective in light of the data, incorporating both the assumptions of the prior process and the information contained in the observations.

We illustrate an example posterior in figure 1.3, where we have conditioned our prior from figure 1.2 on three exact observations. As the

plotting conventions

nonstationarity, warping: § 3.4, p. 56
low intrinsic dimensionality: § 3.5, p. 61

observed data, $\mathcal{D} = (\mathbf{x}, \mathbf{y})$
objective function posterior, $p(f | \mathcal{D})$

⁹ The given expression sweeps some details under the rug. A careful derivation of the posterior process proceeds by finding the posterior of an arbitrary *finite*-dimensional vector $\phi_* = f(\mathbf{x}_*)$:

$$p(\phi_* | \mathbf{x}_*, \mathcal{D}) = \int p(\phi_* | \mathbf{x}_*, \mathbf{x}, \phi) p(\phi | \mathcal{D}) d\phi,$$

which specifies the process. The distributions on the right-hand side are known: the posterior on ϕ is in (1.9), and the posterior on ϕ_* given the *exact* function values ϕ can be found by computing their joint prior (1.8) and conditioning.

INTRODUCTION

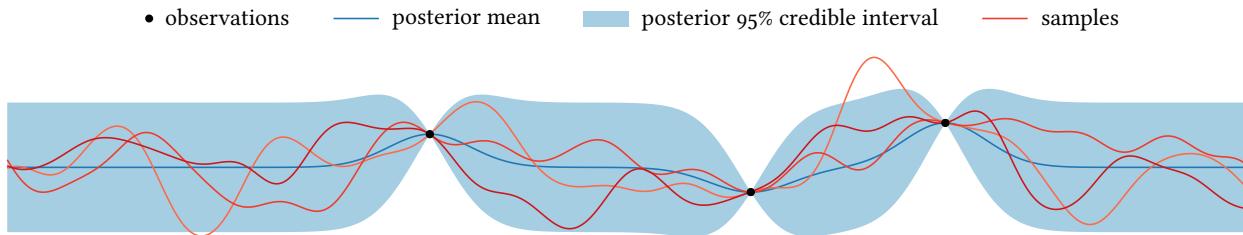


Figure 1.3: The posterior process for our example scenario in figure 2.1 conditioned on three exact observations.

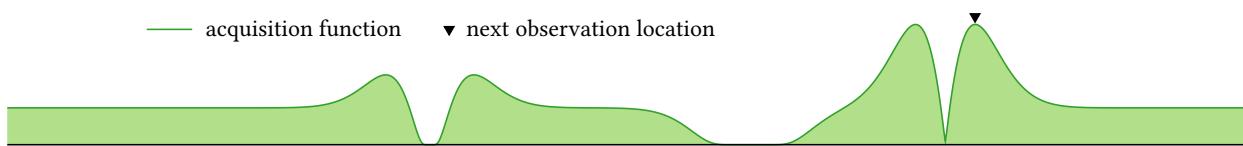


Figure 1.4: A prototypical acquisition function corresponding to our example posterior from figure 1.3.

observations are assumed to be exact, the objective function posterior collapses onto the observed values. The posterior mean interpolates through the data, and the posterior credible intervals reflect increased certainty regarding the function near the observed locations. Further, the posterior continues to reflect the structural assumptions encoded in the prior, demonstrated by comparing the behavior of the samples drawn from the posterior process to those drawn from the prior.

Uncertainty-aware optimization policies

Although Bayesian inference provides an elegant means of reasoning about an uncertain objective function, but the success of optimization is measured not by the fidelity of our beliefs but by the outcomes of our actions. These actions are determined by the optimization policy, which examines available data to design each successive observation location. Each of these decisions is fraught with uncertainty, as we must commit to each observation before knowing its result, which will form the context of all following decisions. Bayesian inference enables us to express this uncertainty, but effective decision making additionally requires us to establish preferences over outcomes and act to maximize those preferences.

To proceed we need to establish a framework for decision making under uncertainty, an expansive subject with a world of possibilities. A natural and common choice is *Bayesian decision theory*, the subject of chapters 5–6. We will discuss this and other approaches to policy construction at length in chapter 7 and derive popular optimization policies from first principles.

Ignoring the details in policy design, a thread running through all Bayesian optimization policies is a uniform handling of uncertainty in

chapter 5: decision theory for optimization,
p. 87

chapter 6: utility functions for optimization,
p. 109

chapter 7: common Bayesian optimization
policies, p. 123

the objective function and the outcomes of observations via Bayesian inference. Instrumental in connecting our beliefs about the objective function to decision making is the posterior predictive distribution (1.7), representing our belief about the outcomes of proposed observations. Bayesian optimization policies are designed with reference to this distribution, which guides the policy in discriminating between potential actions.

In practice, Bayesian optimization policies are defined indirectly by optimizing a so-called *acquisition function* assigning a score to potential observation locations commensurate with their perceived ability to benefit the optimization process. Acquisition functions tend to be cheap to evaluate with analytically tractable gradients, allowing the use of off-the-shelf optimizers to efficiently design each observation. Numerous acquisition functions have been proposed for Bayesian optimization, each derived from different considerations. However, all notable acquisition functions address the classic tension between *exploitation* – sampling where the objective function is expected to be high – and *exploration* – sampling where we are uncertain about the objective function to inform future decisions. These opposing concerns must be carefully balanced for effective global optimization.

An example acquisition function is shown in figure 1.4, corresponding to the posterior from figure 1.3. Consideration of the exploitation–exploration tradeoff is apparent: this example acquisition function attains relatively large values both near local maxima of the posterior mean and in regions with significant marginal uncertainty. Local maxima of the acquisition function represent optimal compromises between these concerns. Note that the acquisition function vanishes at the location of the current observations: as the objective function values at these locations are already known, observing there would be pointless. Maximizing the acquisition function determines the policy; here the policy chooses to search around the local optimum on the right-hand side.

Figure 1.5 demonstrates an entire session of Bayesian optimization, beginning from the belief and initial decision from figure 1.4 and progressing iteratively following algorithm 1.1. The red function shows the true (unknown) objective function, whose maximum is near the center of the domain. The running marks below each posterior show the locations of each measurement made, progressing in sequence from top to bottom, and we show the objective function posterior at four waypoints.

Dynamic consideration of the exploitation–exploration tradeoff is evident in the algorithm’s behavior. The first two observations map out the neighborhood of the initially best-seen point, exhibiting exploitation. Once sufficiently explored, the policy continues exploitation around the second best-seen point, discovering and refining the global optimum in iterations 7–8. Finally, the policy switches to exploration in iterations 13–19, systematically covering the domain to ensure nothing has been missed. At termination, there is clear bias in the collected data toward higher objective values, and all remaining uncertainty is in regions where the credible intervals indicate the optimum is unlikely to reside.

acquisition functions: § 5, p. 88

example and discussion

INTRODUCTION

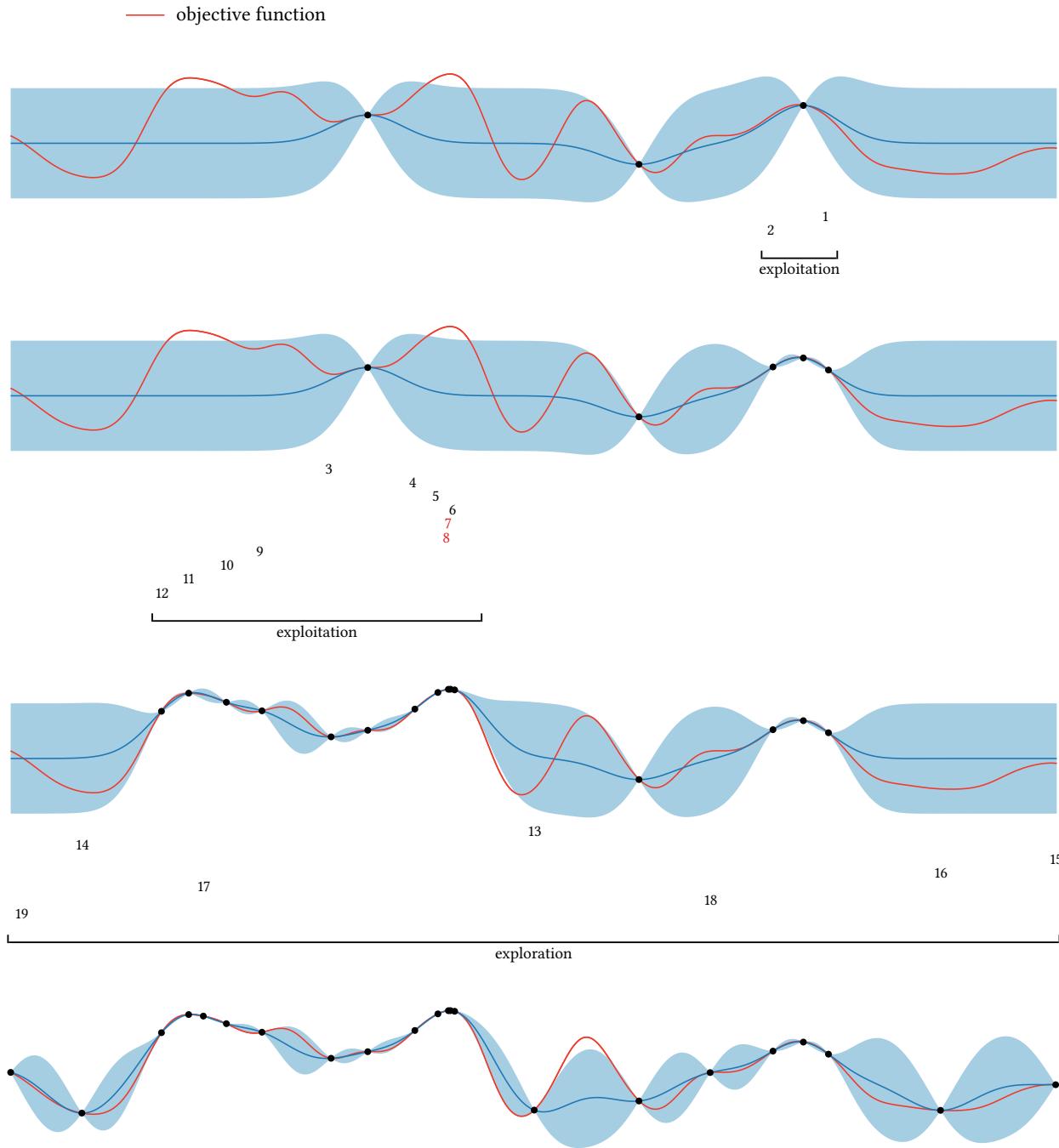


Figure 1.5: The posterior after the indicated number of steps of an example Bayesian optimization policy, starting from the posterior in figure 1.4. The marks show the points chosen by the policy, progressing from top to bottom. Observations sufficiently close to the optimum are marked in red; the optimum was located on iteration 7.

The “magic” of Bayesian optimization is that the intuitive behavior of this optimization policy is not the result of ad hoc design, but rather emerges *automatically* through the machinery of Gaussian processes and Bayesian decision theory that we will develop over the coming chapters. In this framework, building an optimization policy boils down to:

- choosing a model of the objective function,
- deciding what sort of data we seek to obtain, and
- systematically transforming these beliefs and preferences into an optimization policy.

Over the following chapters, we will develop tools for achieving each of these goals: Gaussian processes (chapters 2–4) for expressing what we believe about the objective function, utility functions (chapter 6) for expressing what we value in data, and Bayesian decision theory (chapter 5) for building optimization policies aware of the uncertainty encoded in the model and guided by the preferences encoded in the utility function. In chapter 7 we will combine these fundamental components to realize complete Bayesian optimization policies, at which point we will be equipped to replicate this example from first principles.

2

GAUSSIAN PROCESSES

The central object in optimization is an objective function $f: \mathcal{X} \rightarrow \mathbb{R}$, and the primary challenge in algorithm design is inherent *uncertainty* about this function: most importantly, where is the function maximized and what is its maximal value? Prior to optimization, we may very well have no idea. Optimization affords us the opportunity to acquire information about the objective – through observations of our own design – to shed light on these questions. However, this process is itself fraught with uncertainty, as we cannot know the outcomes and implications of these observations at the time of their design. Notably, we face this uncertainty even when we have a closed-form expression for the objective function, a favorable position as many objectives act as “black boxes.”

Reflecting on this situation, DIACONIS posed an intriguing question:¹ “what does it mean to ‘know’ a function?” The answer is unclear when an analytic expression, which might at first glance seem to encapsulate the essence of the function, is insufficient to determine features of interest. However, DIACONIS argued that although we may not know *everything* about a function, we often have *some* prior knowledge that can facilitate a numerical procedure such as optimization. For example, we may expect an objective function to be smooth (or rough), or to assume values in a given range, or to feature a relatively simple underlying trend, or to depend on some hidden low-dimensional representation we hope to uncover.² All of this knowledge could be instrumental in accelerating optimization if it could be systematically captured and exploited.

Having identifiable information about an objective function prior to optimization motivates the Bayesian approach we will explore throughout this book. We will address uncertainty in the objective function through the unifying framework of Bayesian inference, treating f – as well as ancillary quantities such as x^* and f^* (1.1) – as random variables to be inferred from observations revealed during optimization.

To pursue this approach, we must first determine how to build useful prior distributions for objective functions and how to compute a posterior belief given observations. If the system under investigation is well understood, we may be able to identify an appropriate parametric form $f(x; \theta)$ and infer the parameters θ directly. This approach is likely the best course of action when possible;³ however, many objective functions have no obvious parametric form, and most models used in Bayesian optimization are thus nonparametric to avoid undue assumptions.⁴

In this chapter we will introduce *Gaussian processes* (GPS), a convenient class of nonparametric regression models widely used in Bayesian optimization. We will begin by defining Gaussian processes and deriving some basic properties, then demonstrate how to perform inference from observations. In the case of exact observation and additive Gaussian noise, we can perform this inference *exactly*, resulting in an updated posterior Gaussian process. We will continue by considering some theoretical properties of Gaussian processes relevant to optimization and inference with non-Gaussian observation models.

¹ P. DIACONIS (1988). Bayesian Numerical Analysis. In: *Statistical Decision Theory and Related Topics IV*.

² We will explore all of these possibilities in the next chapter, p. 45.

Bayesian inference of the objective function:
§ 1.2, p. 8

³ V. DALIBARD et al. (2017). BOAT: Building Auto-Tuners with Structured Bayesian Optimization. www.2017.

⁴ The term “nonparametric” is something of a misnomer. A nonparametric objective function model has parameters but their dimension is infinite – we effectively parametrize the objective by its value at every point.

- 5 C. E. RASMUSSEN and C. K. I. WILLIAMS (2006). *Gaussian Processes for Machine Learning*. MIT Press.

multivariate normal distribution: appendix A,
p. 291

chapter 3: modeling with Gaussian processes,
p. 45

- 6 P. HENNIG et al. (2015). Probabilistic numerics and uncertainty in computations. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 471(2179):20150142.

- 7 P. HENNIG et al. (2022). *Probabilistic Numerics: Computation as Machine Learning*. Cambridge University Press.

- 8 If \mathcal{X} is finite, there is no distinction between a Gaussian process and a multivariate normal distribution, so only the infinite case is interesting for this discussion.

- 9 B. ØKSENDAL (2013). *Stochastic Differential Equations: An Introduction with Applications*. Springer-Verlag. [§ 2.1]

- 10 Writing the process as if it were a function-valued probability density function is an abuse of notation, but a useful and harmless one.

mean function, μ
covariance function (kernel), K
value of objective at x , ϕ

The literature on Gaussian processes is vast, and we do not intend this chapter to serve as a standalone introduction but rather as companion to the existing literature. Although our discussion will be comprehensive, our focus on optimization will sometimes bias its scope. For a broad overview, the interested reader may consult RASMUSSEN and WILLIAMS's classic monograph.⁵

2.1 DEFINITION AND BASIC PROPERTIES

A *Gaussian process* is an extension of the familiar multivariate normal distribution suitable for modeling functions on infinite domains. Gaussian processes inherit the convenient mathematical properties of the multivariate normal distribution without sacrificing computational tractability. Further, by modifying the structure of a GP, we can model functions with a rich variety of behavior; we will explore this capability in the next chapter. This combination of mathematical elegance and flexibility in modeling has established Gaussian processes as the workhorse of Bayesian approaches to numerical tasks, including optimization.^{6,7}

Definition

Consider an objective function $f: \mathcal{X} \rightarrow \mathbb{R}$ of interest over an arbitrary infinite domain \mathcal{X} .⁸ We will take a nonparametric approach and reason about the function as an infinite collection of random variables, one corresponding to the function value at every point in the domain. Mutual dependence between these random variables will then determine the statistical properties of the function's shape.

It is perhaps not immediately clear how we can specify a useful distribution over infinitely many random variables, a construction known as a *stochastic process*. However, a result known as the *Kolmogorov extension theorem* allows us to construct a stochastic process by defining only the distribution of arbitrary *finite* sets of function values, subject to natural consistency constraints.⁹ For a Gaussian process, these finite-dimensional distributions are all multivariate Gaussian, hence its name.

In this light, we build a Gaussian process by replacing the parameters in the finite-dimensional case – a mean vector and a positive semidefinite covariance matrix – by analogous *functions* over the domain. We specify a Gaussian process on f :¹⁰

$$p(f) = \mathcal{GP}(f; \mu, K)$$

by a *mean function* $\mu: \mathcal{X} \rightarrow \mathbb{R}$ and a positive semidefinite *covariance function* (or *kernel*) $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. The mean function determines the expected function value $\phi = f(x)$ at any location x :

$$\mu(x) = \mathbb{E}[\phi | x],$$

thus serving as a location parameter representing the function's central tendency. The covariance function determines how deviations from the

mean are structured, encoding expected properties of the function's behavior. Defining $\phi' = f(x')$, we have:

$$K(x, x') = \text{cov}[\phi, \phi' | x, x']. \quad (2.1)$$

The mean and covariance functions of the process allow us to compute any finite-dimensional marginal distribution on demand. Let $\mathbf{x} \subset \mathcal{X}$ be finite and let $\phi = f(\mathbf{x})$ be the corresponding function values, a vector-valued random variable. For the Gaussian process (2.1), the distribution of ϕ is multivariate normal with parameters determined by the mean and covariance functions:

$$p(\phi | \mathbf{x}) = \mathcal{N}(\phi; \mu, \Sigma), \quad (2.2)$$

where

$$\mu = \mathbb{E}[\phi | \mathbf{x}] = \mu(\mathbf{x}); \quad \Sigma = \text{cov}[\phi | \mathbf{x}] = K(\mathbf{x}, \mathbf{x}). \quad (2.3)$$

Here $K(\mathbf{x}, \mathbf{x})$ is the matrix formed by evaluating the covariance function for each pair of points: $\Sigma_{ij} = K(x_i, x_j)$, also called the *Gram matrix* of \mathbf{x} .

In many ways, Gaussian processes behave like “really big” Gaussian distributions, and one can intuit many of their properties from this heuristic alone. For example, the Gaussian marginal property in (2.2–2.3) corresponds precisely with the analogous formula in the finite-dimensional case (A.13). Further, this property automatically ensures global consistency in the following sense.¹¹ If \mathbf{x} is an arbitrary set of points and $\mathbf{x}' \supset \mathbf{x}$ is a superset, then we arrive at the same belief about ϕ whether we compute it directly from (2.2–2.3) or indirectly by first computing $p(\phi' | \mathbf{x}')$ then marginalizing (A.13).

Example and basic properties

Let us construct and explore an explicit Gaussian process for a function on the interval $\mathcal{X} = [0, 30]$. For the mean function we take the zero function $\mu \equiv 0$, indicating a constant central tendency. For the covariance function, we take the prototypical *squared exponential* covariance:

$$K(x, x') = \exp\left(-\frac{1}{2}|x - x'|^2\right). \quad (2.4)$$

Let us pause to consider the implications of this choice. First, note that $\text{var}[\phi | \mathbf{x}] = K(\mathbf{x}, \mathbf{x}) = 1$ at every point $x \in \mathcal{X}$, and thus the covariance function (2.4) also measures the *correlation* between the function values ϕ and ϕ' . This correlation decreases with the distance between x and x' , falling from unity to zero as these points become increasingly separated; see the illustration in the margin. We can loosely interpret this as a statistical consequence of continuity: function values at nearby locations are highly correlated, whereas function values at distant locations are effectively independent. This assumption also implies that observing the function at some point x provides nontrivial information about the function at sufficiently nearby locations (roughly when $|x - x'| < 3$). We will explore this implication further shortly.

value of objective at x', ϕ'

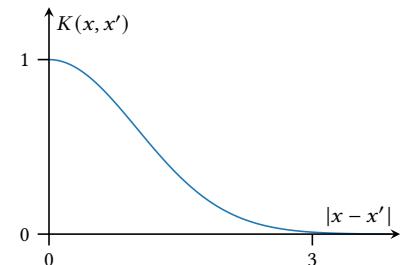
values of objective at $\mathbf{x}, \phi = f(\mathbf{x})$

Gram matrix of $\mathbf{x}, \Sigma = K(\mathbf{x}, \mathbf{x})$

¹¹ In fact, this is precisely the consistency required by the Kolmogorov extension theorem mentioned on the facing page.

marginalizing multivariate normal distributions, § A.2, p. 295

squared exponential covariance: § 3.3, p. 51



The squared exponential covariance (2.4) as a function of the distance between inputs.

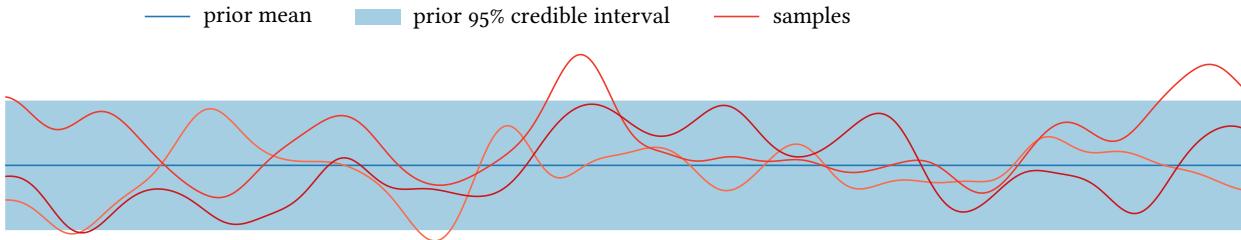


Figure 2.1: Our example Gaussian process on the domain $\mathcal{X} = [0, 30]$. We illustrate the marginal belief at every location with its mean and a 95% credible interval and also show three example functions sampled from the process.

predictive credible intervals

For a Gaussian process, the marginal distribution of any *single* function value is univariate normal (2.2):

$$p(\phi | x) = \mathcal{N}(\phi; \mu, \sigma^2); \quad \mu = \mu(x); \quad \sigma^2 = K(x, x), \quad (2.5)$$

where we have abused notation by overloading the symbol μ . This result allows us to derive pointwise credible intervals; for example, the familiar $\mu \pm 1.96\sigma$ is a 95% credible interval for ϕ . Examining our example GP, the marginal distribution of every function value is in fact *standard* normal. We provide a rough visual summary of the process via its mean function and pointwise 95% predictive credible intervals in figure 2.1. There is nothing terribly exciting we can glean from these marginal distributions alone, and no interesting structure in the process is yet apparent.

Sampling

sampling from a multivariate normal distribution: § A.2, p. 295

We may gain more insight by inspecting samples drawn from our example process reflecting the *joint* distribution of function values. Although it is impossible to represent an arbitrary function on \mathcal{X} in finite memory, we can approximate the sampling process by taking a dense grid $\mathbf{x} \subset \mathcal{X}$ and sampling the corresponding function values from their joint multivariate normal distribution (2.2). Plotting the sampled vectors against the chosen grid reveals curves approximating draws from the Gaussian process. Figure 2.1 illustrates this procedure for our example using a grid of 1000 equally spaced points. Each sample is smooth and has several local optima distributed throughout the domain – for some applications, this might be a reasonable model for an objective function on \mathcal{X} .

2.2 INFERENCE WITH EXACT AND NOISY OBSERVATIONS

We now turn to our attention to *inference*: given a Gaussian process prior on an objective function, how can we condition this initial belief on observations obtained during optimization?

Let us look at an example to build intuition before diving into the details. Figure 2.2 shows the effect of conditioning our example GP from the previous section on three exact measurements of the function. The

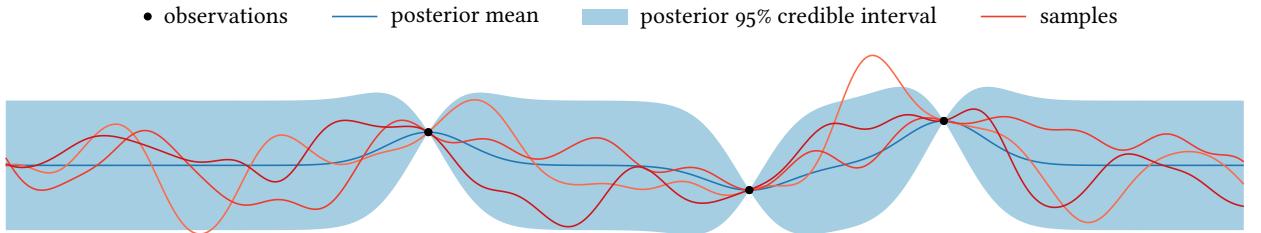


Figure 2.2: The posterior for our example scenario in figure 2.1 conditioned on three exact observations.

updated belief reflects both our prior assumptions and the information contained in the data, the hallmark of Bayesian inference. To elaborate, the posterior mean smoothly interpolates through the observed values, agreeing with both the measured values and the smoothness encoded in the prior covariance function. The posterior credible intervals are reduced in the neighborhood of the measured locations – where the prior covariance function encodes nontrivial dependence on at least one observed value – and vanish where the function value has been exactly determined. On the other hand, our marginal belief remains effectively unchanged from the prior in regions sufficiently isolated from the data, where the prior covariance function encodes effectively no correlation.

Conveniently, inference is straightforward for the pervasive observation models of exact measurement and additive Gaussian noise, where the self-conjugacy of the normal distribution yields a *Gaussian process* posterior with updated parameters we can compute in closed form. The reasoning underlying inference for both observation models is identical and is subsumed by a flexible general argument we will present first.

Inference from arbitrary jointly Gaussian observations

We may exactly condition a Gaussian process $p(f) = \mathcal{GP}(f; \mu, K)$ on the observation of *any* vector y sharing a joint Gaussian distribution with f :

$$p(f, y) = \mathcal{GP}\left(\begin{bmatrix} f \\ y \end{bmatrix}; \begin{bmatrix} \mu \\ \mathbf{m} \end{bmatrix}, \begin{bmatrix} K & \kappa^\top \\ \kappa & C \end{bmatrix}\right). \quad (2.6)$$

This notation, analogous to (A.12), extends the Gaussian process on f to include the entries of y ; that is, we assume the distribution of any finite subset of function and/or observed values is multivariate normal. We specify the joint distribution via the marginal distribution of y .¹²

$$p(y) = \mathcal{N}(y; \mathbf{m}, C) \quad (2.7)$$

and the cross-covariance function between y and f :

$$\kappa(x) = \text{cov}[y, \phi \mid x]. \quad (2.8)$$

vector of observed values, y

¹² We assume C is positive definite; if it were only positive semidefinite, there would be wasteful linear dependence among observations.

observation mean and covariance, \mathbf{m}, C

cross-covariance between observations and function values, κ

inference from exact observations: § 2.2, p. 22
affine transformations: § A.2, p. 294
derivatives and expectations: § 2.6, p. 30

conditioning a multivariate normal distribution: § A.2, p. 295

posterior mean and covariance, $\mu_{\mathcal{D}}, K_{\mathcal{D}}$

¹³ This is a useful exercise! The result will be a stochastic process with multivariate normal finite-dimensional distributions, a Gaussian process by definition (2.5).

noisy observation of y, z
vector of random errors, ϵ
noise covariance matrix, N

sums of normal vectors: § A.2, p. 296

Although it may seem absurd that we could identify and observe a vector satisfying such strong restrictions on its distribution, we can already deduce several examples from first principles, including:

- any vector of function values (2.2),
- any affine transformation of function values (A.10), and
- limits of such quantities, such as partial derivatives or expectations.

Further, we may condition on any of the above even if corrupted by independent additive Gaussian noise, as we will shortly demonstrate.

We may condition the joint distribution (2.6) on y analogously to the finite-dimensional case (A.14), resulting in a Gaussian process posterior on f . Writing $\mathcal{D} = y$ for the observed data, we have:

$$p(f \mid \mathcal{D}) = \mathcal{GP}(f; \mu_{\mathcal{D}}, K_{\mathcal{D}}), \quad (2.9)$$

where

$$\begin{aligned} \mu_{\mathcal{D}}(x) &= \mu(x) + \kappa(x)^{\top} C^{-1}(y - m); \\ K_{\mathcal{D}}(x, x') &= K(x, x') - \kappa(x)^{\top} C^{-1} \kappa(x'). \end{aligned} \quad (2.10)$$

This can be verified by computing the joint distribution of an arbitrary finite set of function values and y and conditioning on the latter (A.14).¹³

The above result provides a simple procedure for GP posterior inference from any vector of observations satisfying (2.6):

1. compute the marginal distribution of y (2.7),
2. derive the cross-covariance function κ (2.8), and
3. find the posterior distribution of f via (2.9–2.10).

We will realize this procedure for several special cases below. However, we will first demonstrate how we may seamlessly handle measurements corrupted by additive Gaussian noise and build intuition for the posterior distribution by dissecting its moments in terms of the statistics of the observations and the correlation structure of the prior.

Corruption by additive Gaussian noise

We pause to make one observation of immense practical importance: any vector satisfying (2.6) would continue to suffice even if corrupted by independent additive Gaussian noise, and thus we can use the above result to condition a Gaussian process on *noisy* observations as well.

Suppose that rather than observing y exactly, our measurement mechanism only allowed observing $z = y + \epsilon$ instead, where ϵ is a vector of random errors independent of y . If the errors are normally distributed with mean zero and known (arbitrary) covariance N :

$$p(\epsilon \mid N) = \mathcal{N}(\epsilon; \mathbf{0}, N), \quad (2.11)$$

then we have

$$p(z \mid N) = \mathcal{N}(z; m, C + N); \quad \text{cov}[z, \phi \mid x] = \text{cov}[y, \phi \mid x] = \kappa(x).$$

Thus we can condition on an observation of the corrupted vector \mathbf{z} by simply replacing \mathbf{C} with $\mathbf{C} + \mathbf{N}$ in the prior (2.6) and posterior (2.10).¹⁴ Note that the posterior converges to that from a direct observation of \mathbf{y} if we take the noise covariance $\mathbf{N} \rightarrow \mathbf{0}$ in the positive semidefinite cone, a reassuring result.

Interpretation of posterior moments

The moments of the posterior Gaussian process (2.10) contain update terms adjusting the prior moments in light of the data. These updates have intuitive interpretations in terms of the nature of the prior process and the observed values, which we may unravel with some care.

We can gain some initial insight by considering the case where we observe a *single* value with y distribution $\mathcal{N}(y; m, s^2)$ and breaking down its impact on our belief. Consider an arbitrary function value ϕ with prior distribution $\mathcal{N}(\phi; \mu, \sigma^2)$ (2.5) and define

$$z = \frac{y - m}{s}$$

to be the z -score of the observed value y and

$$\rho = \text{corr}[y, \phi | x] = \frac{\kappa(x)}{\sigma s}$$

to be the correlation between y and ϕ . Then the posterior mean and standard deviation of ϕ are, respectively:

$$\mu + \sigma\rho z; \quad \sigma\sqrt{1 - \rho^2}. \quad (2.12)$$

The z -score of the posterior mean, with respect to the prior distribution of ϕ , is ρz . An independent measurement with $\rho = 0$ thus leaves the prior mean unchanged, whereas a perfectly dependent measurement with $|\rho| = 1$ shifts the mean up or down by z standard deviations (depending on the sign of the correlation) to match the magnitude of the measurement's z -score. Measurements with partial dependence result in outcomes between these extremes. Further, *surprising* measurements – that is, those with large $|z|$ – yield larger shifts in the mean, whereas an entirely expected measurement with $y = m$ leaves the mean unchanged.

Turning to the posterior standard deviation, the measurement reduces our uncertainty in ϕ by a factor depending on the correlation ρ , but *not* on the value observed. An independent measurement again leaves the prior intact, whereas a perfectly dependent measurement collapses the posterior standard deviation to zero as the value of ϕ would be completely determined. The relative reduction in posterior uncertainty as a function of the absolute correlation is illustrated in the margin.

In the case of vector-valued observations, we can interpret similar structure in the posterior, although dependence between entries of \mathbf{y} must also be accounted for. We may factor the observation covariance matrix as

$$\mathbf{C} = \mathbf{S}\mathbf{P}\mathbf{S}, \quad (2.13)$$

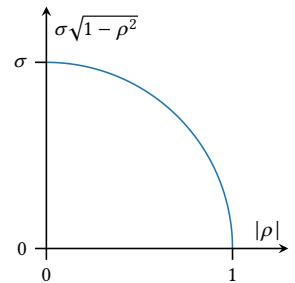
¹⁴ Assuming zero-mean errors is not strictly necessary but is overwhelmingly common in practice. A nonzero mean $\mathbb{E}[\epsilon] = \mathbf{n}$ is possible by further replacing $(\mathbf{y} - \boldsymbol{\mu})$ with $(\mathbf{y} - [\boldsymbol{\mu} + \mathbf{n}])$ in (2.10).

z-score of measurement y, z

correlation between measurement y and function value ϕ, ρ

posterior moments of ϕ from a scalar observation

interpretation of moments



The posterior standard deviation of ϕ as a function of the strength of relationship with $y, |\rho|$.

where \mathbf{S} is diagonal with $S_{ii} = \sqrt{C_{ii}} = \text{std}[y_i]$ and $\mathbf{P} = \text{corr}[\mathbf{y}]$ is the observation correlation matrix. We may then rewrite the posterior mean of ϕ as

$$\mu + \sigma \boldsymbol{\rho}^\top \mathbf{P}^{-1} \mathbf{z},$$

where \mathbf{z} and $\boldsymbol{\rho}$ represent the vectors of measurement z -scores and the cross-correlation between ϕ and \mathbf{y} , respectively:

$$z_i = \frac{y_i - m_i}{s_i}; \quad \rho_i = \frac{[\kappa(\mathbf{x})]_i}{\sigma s_i}.$$

¹⁵ It can be instructive to contrast the behavior of the posterior when conditioning on two highly correlated values versus two independent ones. In the former case, the posterior does not change much as a result of the second measurement, as correlation reduces the effective number of measurements.

¹⁶ \mathbf{P} is congruent to \mathbf{C} (2.13) and is thus positive definite from Sylvester's law of inertia.

¹⁷ For positive semidefinite \mathbf{A}, \mathbf{B} , $|\mathbf{A}| \leq |\mathbf{A} + \mathbf{B}|$.

¹⁸ The Löwner order is the partial order induced by the convex cone of positive-semidefinite matrices. For symmetric \mathbf{A}, \mathbf{B} , we define $\mathbf{A} \prec \mathbf{B}$ if and only if $\mathbf{B} - \mathbf{A}$ is positive definite:

K. LÖWNER (1934). Über monotone Matrixfunktionen. *Mathematische Zeitschrift* 38:177–216.

observed data, $\mathcal{D} = (\mathbf{x}, \phi)$

The posterior mean is now in the same form as the scalar case (2.12), with the introduction of the observation correlation matrix moderating the z -scores to account for dependence between the observed values.¹⁵

The posterior standard deviation of ϕ in the vector-valued case is

$$\sigma \sqrt{1 - \boldsymbol{\rho}^\top \mathbf{P}^{-1} \boldsymbol{\rho}},$$

again analogous to (2.12). Noting that the inverse correlation matrix \mathbf{P}^{-1} is positive definite,¹⁶ the posterior covariance again reflects a global reduction in the marginal uncertainty of every function value. In fact, the *joint* distribution of any set of function values has reduced uncertainty in the posterior in terms of the differential entropy (A.16), as¹⁷

$$|K(\mathbf{x}, \mathbf{x}) - \kappa(\mathbf{x})^\top \mathbf{C}^{-1} \kappa(\mathbf{x})| \leq |K(\mathbf{x}, \mathbf{x})|.$$

The reduction of uncertainty again depends on the strength of dependence between function values and the observed data, with independence ($\boldsymbol{\rho} = \mathbf{0}$) resulting in no change. The reduction also depends on the precision of the measurements: all other things held equal, observations with greater precision in terms of the Löwner order¹⁸ on the precision matrix \mathbf{C}^{-1} provide a globally better informed posterior. In particular, as $(\mathbf{C} + \mathbf{N})^{-1} \prec \mathbf{C}^{-1}$ for any noise covariance \mathbf{N} , noisy measurements (2.11) categorically provide *less* information about the function than direct observations, as one might hope.

Inference with exact function evaluations

We will now explicitly demonstrate the general process of Gaussian process inference for important special cases, beginning with the simplest possible observation mechanism: exact observation.

Suppose we have observed f at some set of locations \mathbf{x} , revealing the corresponding function values $\phi = f(\mathbf{x})$, and let $\mathcal{D} = (\mathbf{x}, \phi)$ denote this dataset. The observed vector shares a joint Gaussian distribution with any other set of function values by the GP assumption on f (2.2), so we may follow the above procedure to derive the posterior. The marginal distribution of ϕ is Gaussian (2.3):

$$p(\phi \mid \mathbf{x}) = \mathcal{N}(\phi; \boldsymbol{\mu}, \Sigma),$$

and the cross-covariance between an arbitrary function value and ϕ is by definition given by the covariance function:

$$\kappa(x) = \text{cov}[\phi, \phi | x, x] = K(x, x).$$

Appealing to (2.9–2.10) we have:

$$p(f | \mathcal{D}) = \mathcal{GP}(f; \mu_{\mathcal{D}}, K_{\mathcal{D}}),$$

where

$$\begin{aligned}\mu_{\mathcal{D}}(x) &= \mu(x) + K(x, x)\Sigma^{-1}(\phi - \mu); \\ K_{\mathcal{D}}(x, x') &= K(x, x') - K(x, x)\Sigma^{-1}K(x, x').\end{aligned}\tag{2.14}$$

Our previous figure 2.2 illustrates the posterior resulting from conditioning our GP prior in figure 2.1 on three exact measurements, with high-level analysis of its behavior in the accompanying text.

example and discussion

Inference with function evaluations corrupted by additive Gaussian noise

With the notable exception of optimizing the output of a deterministic computer program or simulation, observations of an objective function are typically corrupted by noise due to measurement limitations or statistical approximation; we must be able to handle such noisy observations to maximize utility. Fortunately, in the important case of additive Gaussian noise, we may perform exact inference following the general procedure described above. In fact, the derivation below follows directly from our previous discussion on arbitrary additive Gaussian noise, but the case of additive Gaussian noise in function evaluations is important enough to merit its own discussion.

Suppose we make observations of f at locations x , revealing corrupted values $y = \phi + \epsilon$. Suppose the measurement errors ϵ are independent of ϕ and normally distributed with mean zero and covariance N , which may optionally depend on x :

$$p(\epsilon | x, N) = \mathcal{N}(\epsilon; 0, N).\tag{2.15}$$

As before we aggregate the observations into a dataset $\mathcal{D} = (x, y)$.

The observation noise covariance can in principle be arbitrary;¹⁹ however, the most common models in practice are independent homoskedastic noise with scale σ_n :

$$N = \sigma_n^2 I,\tag{2.16}$$

and independent heteroskedastic noise with scale depending on location according to a function $\sigma_n: \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$:

$$N = \text{diag}(\sigma_n^2(x)).\tag{2.17}$$

For a given observation location x , we will simply write σ_n for the associated noise scale, leaving any dependence on x implicit.

arbitrary additive Gaussian noise: § 2.2, p. 20

¹⁹ Allowing nondiagonal N departs from our typical convention of assuming conditional independence between observations (1.3), but doing so does not complicate inference, so there is no harm in this generality.

special case: independent homoskedastic and heteroskedastic noise

observation noise scale, σ_n

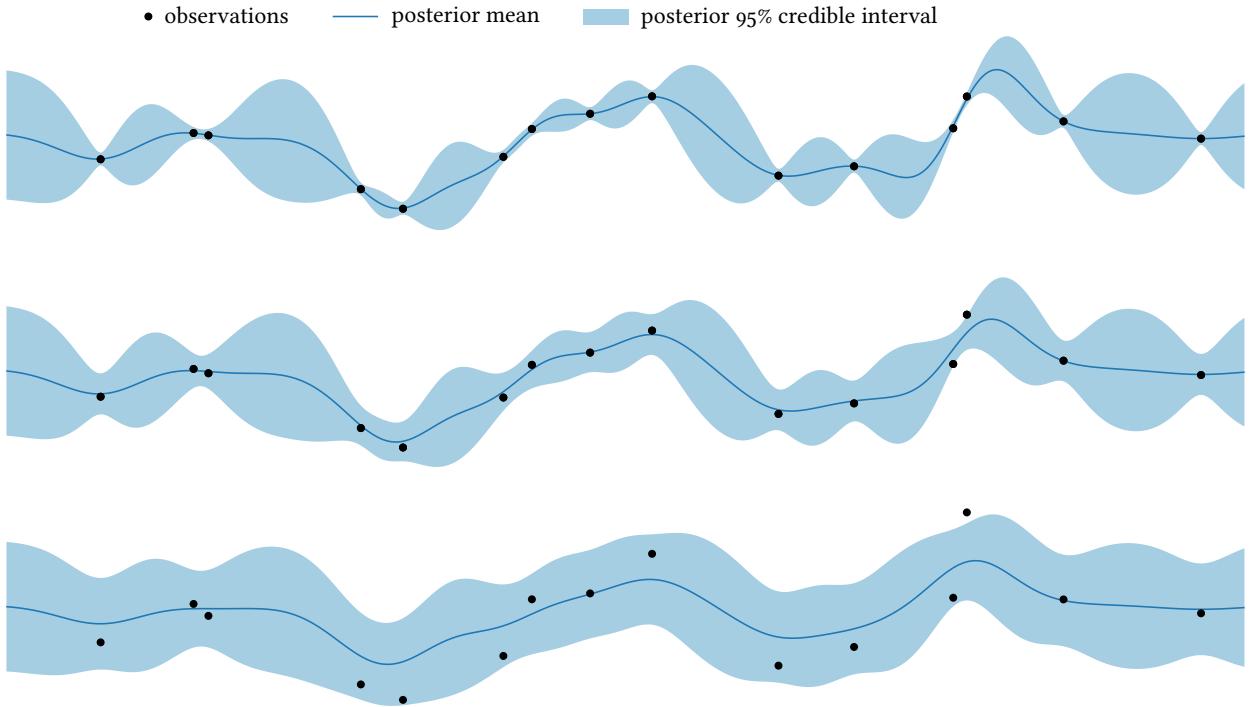


Figure 2.3: Posteriors for our example GP from figure 2.1 conditioned on 15 noisy observations with independent homoskedastic noise (2.16). The signal-to-noise ratio is 10 for the top example, 3 for the middle example, and 1 for the bottom example.

The prior distribution of the observations is now multivariate normal (2.3, A.15):

$$p(\mathbf{y} \mid \mathbf{x}, \mathbf{N}) = \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}, \Sigma + \mathbf{N}). \quad (2.18)$$

Due to independence of the noise, the cross-covariance remains the same as in the exact observation case:

$$\kappa(x) = \text{cov}[\mathbf{y}, \phi \mid \mathbf{x}, x] = K(\mathbf{x}, x).$$

Conditioning on the observed value now yields a GP posterior with

$$\begin{aligned} \mu_{\mathcal{D}}(x) &= \mu(x) + K(x, \mathbf{x})(\Sigma + \mathbf{N})^{-1}(\mathbf{y} - \boldsymbol{\mu}); \\ K_{\mathcal{D}}(x, x') &= K(x, x') - K(x, \mathbf{x})(\Sigma + \mathbf{N})^{-1}K(\mathbf{x}, x'). \end{aligned} \quad (2.19)$$

homoskedastic example and discussion

Figure 2.3 shows a sequence of posterior distributions resulting from conditioning our example GP on data corrupted by increasing levels of homoskedastic noise (2.16). As the noise level increases, the observations have diminishing influence on our belief, with some extreme values eventually being partially explained away as outliers. As measurements are assumed to be inexact, the posterior mean is not compelled to interpolate perfectly through the observations, as in the exact case (figure

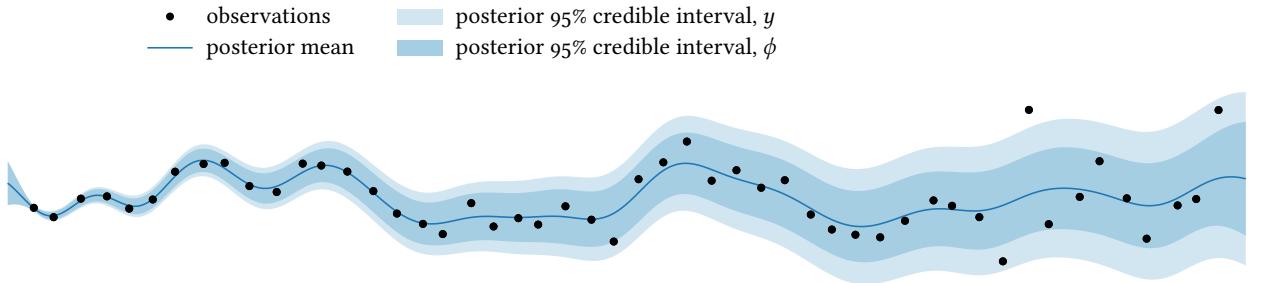


Figure 2.4: The posterior distribution for our example GP from figure 2.1 conditioned on 50 observations with heteroskedastic observation noise (2.17). We show predictive credible intervals for both the latent objective function and noisy observations; the standard deviation of the observation noise increases linearly from left to right.

2.2). Further, with increasing levels of noise, our posterior belief reflects significant residual uncertainty in the function, even in regions with multiple nearby observations.

We illustrate an example of Gaussian process inference with heteroskedastic noise (2.17) in figure 2.4, where the signal-to-noise ratio decreases smoothly from left-to-right over the domain. Although the observations provide relatively even coverage, our posterior uncertainty is minimal on the left-hand side of the domain – where the measurements provide maximal information – and increases as our observations become more noisy and less informative.

We will often require the posterior predictive distribution for a noisy measurement y that would result from observing at a given location x . The posterior distribution on f (2.19) provides the posterior predictive distribution for the latent function value $\phi = f(x)$ (2.5):

$$p(\phi | x, \mathcal{D}) = \mathcal{N}(\phi; \mu, \sigma^2); \quad \mu = \mu_{\mathcal{D}}(x); \quad \sigma^2 = K_{\mathcal{D}}(x, x),$$

but does not account for the effect of observation noise. In the case of independent additive Gaussian noise (2.16–2.17), deriving the posterior predictive distribution is trivial; we have (A.15):

$$p(y | x, \mathcal{D}, \sigma_n) = \mathcal{N}(y; \mu, \sigma^2 + \sigma_n^2). \quad (2.20)$$

This predictive distribution is illustrated in figure 2.4; the credible intervals for noisy measurements reflect inflation of the credible intervals for the underlying function value commensurate with the scale of the noise.

If the noise contains nondiagonal correlation structure, we must account for dependence between training and test errors in the predictive distribution. The easiest way to proceed is to recognize that the noisy observation process $y = \phi + \varepsilon$, as a function of x , is itself a Gaussian process with mean function μ and covariance function

$$C(x, x') = \text{cov}[y, y' | x, x'] = K(x, x') + N(x, x'),$$

heteroskedastic example and discussion

posterior predictive distribution for noisy observations

predictive distribution with correlated noise

covariance function for noisy measurements,
 C

covariance function for observation noise, N where N is the noise covariance: $N(x, x') = \text{cov}[\varepsilon, \varepsilon' | x, x']$. The posterior of the observation process is then a GP with

$$\begin{aligned}\mathbb{E}[y | x, \mathcal{D}] &= \mu(x) + C(x, \mathbf{x})(\Sigma + N)^{-1}(y - \boldsymbol{\mu}); \\ \text{cov}[y, y' | x, x', \mathcal{D}] &= C(x, x') - C(x, \mathbf{x})(\Sigma + N)^{-1}C(\mathbf{x}, x'),\end{aligned}\tag{2.21}$$

from which we can derive predictive distributions via (2.2).

2.3 OVERVIEW OF REMAINDER OF CHAPTER

In the remainder of this chapter we will cover some additional, somewhat niche and/or technical aspects of Gaussian processes that see occasional use in Bayesian optimization. Modulo mathematical nuances irrelevant in practical settings, an *intuitive* (but not entirely accurate!) summary follows:

- a *joint Gaussian process* (discussed below) allows us to model *multiple* related functions simultaneously, which is critical for some scenarios such as multifidelity and multiobjective optimization;
- GP sample paths are continuous if the mean function is continuous and the covariance function is continuous along the “diagonal” $x = x'$;
- GP sample paths are differentiable if the mean function is differentiable and the covariance function is differentiable along the “diagonal” $x = x'$;
- a function with a sufficiently smooth GP distribution shares a joint GP distribution with its gradient; among other things, this allows us to condition on (potentially noisy) derivative observations via exact inference;
- GP sample paths attain a maximum when sample paths are continuous and the domain is compact,
- GP sample paths attain a *unique* maximum under the additional condition that no two unique function values are perfectly correlated, and
- several methods are available for approximating the posterior process of a GP conditioned on information incompatible with exact inference.

If satisfied with the above summary, the reader may safely skip this material for now and move on with the next chapter. For those who wish to see the gritty details, dive in below!

2.4 JOINT GAUSSIAN PROCESSES

In some settings, we may wish to reason *jointly* about two-or-more related functions, such as an objective function and its gradient or an expensive objective function and a cheaper surrogate. To this end we can extend Gaussian processes to yield a joint distribution over the values assumed by multiple functions. The key to the construction is to “paste together” a collection of functions into a single function on a larger domain, then construct a standard GP on this combined function.

Definition

To elaborate, consider a set of functions $\{f_i: \mathcal{X}_i \rightarrow \mathbb{R}\}$ we wish to model.²⁰ We define the *disjoint union* of these functions $\sqcup f$ – defined on the disjoint union²¹ of their domains $\mathcal{X} = \bigsqcup \mathcal{X}_i$ – by insisting its restriction to each domain be compatible with the corresponding function:

$$\sqcup f: \mathcal{X} \rightarrow \mathbb{R}; \quad \sqcup f|_{\mathcal{X}_i} \equiv f_i.$$

We now can define a GP on $\sqcup f$ by choosing mean and covariance functions on \mathcal{X} as desired:

$$p(\sqcup f) = \mathcal{GP}(\sqcup f; \mu, K). \quad (2.22)$$

We will call this construction a *joint Gaussian process* on $\{f_i\}$.

It is often convenient to decompose the moments of a joint GP into their restrictions on relevant subspaces. For example, consider a joint GP (2.22) on $f: \mathcal{F} \rightarrow \mathbb{R}$ and $g: \mathcal{G} \rightarrow \mathbb{R}$. After defining

$$\begin{aligned} \mu_f &\equiv \mu|_{\mathcal{F}}; & \mu_g &\equiv \mu|_{\mathcal{G}}; \\ K_f &\equiv K|_{\mathcal{F} \times \mathcal{F}}; & K_g &\equiv K|_{\mathcal{G} \times \mathcal{G}}; & K_{fg} &\equiv K|_{\mathcal{F} \times \mathcal{G}}; & K_{gf} &\equiv K|_{\mathcal{G} \times \mathcal{F}}, \end{aligned}$$

we can see that f and g in fact have marginal GP distributions:²²

$$p(f) = \mathcal{GP}(f; \mu_f, K_f); \quad p(g) = \mathcal{GP}(g; \mu_g, K_g), \quad (2.23)$$

that are coupled by the *cross-covariance functions* K_{fg} and K_{gf} . Given vectors $\mathbf{x} \subset \mathcal{F}$ and $\mathbf{x}' \subset \mathcal{G}$, these compute the covariance between the corresponding function values $\boldsymbol{\phi} = f(\mathbf{x})$ and $\boldsymbol{\gamma} = g(\mathbf{x}')$:

$$\begin{aligned} K_{fg}(\mathbf{x}, \mathbf{x}') &= \text{cov}[\boldsymbol{\phi}, \boldsymbol{\gamma} \mid \mathbf{x}, \mathbf{x}']; \\ K_{gf}(\mathbf{x}, \mathbf{x}') &= \text{cov}[\boldsymbol{\gamma}, \boldsymbol{\phi} \mid \mathbf{x}, \mathbf{x}'] = K_{fg}(\mathbf{x}, \mathbf{x})^\top. \end{aligned} \quad (2.24)$$

When convenient we will notate a joint GP in terms of these decomposed functions, here writing:²³

$$p(f, g) = \mathcal{GP}\left(\begin{bmatrix} f \\ g \end{bmatrix}; \begin{bmatrix} \mu_f \\ \mu_g \end{bmatrix}, \begin{bmatrix} K_f & K_{fg} \\ K_{gf} & K_g \end{bmatrix}\right). \quad (2.25)$$

With this notation, the marginal GP property (2.23) is perfectly analogous to the marginal property of the multivariate Gaussian distribution (A.13).

We can also use this construction to define a GP on a *vector-valued* function $\mathbf{f}: \mathcal{X} \rightarrow \mathbb{R}^d$ by defining a joint Gaussian process on its d coordinate functions $\{f_i\}: \mathcal{X} \rightarrow \mathbb{R}$. In this case we typically write the resulting model using the standard notation $\mathcal{GP}(\mathbf{f}; \mu, K)$, where the mean and covariance functions are now understood to map to \mathbb{R}^d and $\mathbb{R}^{d \times d}$.

Example

We can demonstrate the behavior of a joint Gaussian process by extending our running example GP on $f: [0, 30] \rightarrow \mathbb{R}$. Recall the prior on f has

disjoint union of $\{f_i\}$, $\sqcup f$
disjoint union of $\{\mathcal{X}_i\}$, \mathcal{X}

²⁰ The domains need not be equal, but they often are in practice.

²¹ A disjoint union represents a point $x \in \mathcal{X}_i$ by the pair (x, i) , thereby combining the domains while retaining their identities.

joint Gaussian process

²² In fact, *any* restriction of a GP-distributed function has a GP (or multivariate normal) distribution.

²³ We also used this notation in (2.6), where the “domain” of the vector \mathbf{y} can be taken to be some finite index set of appropriate size.

extension to vector-valued functions

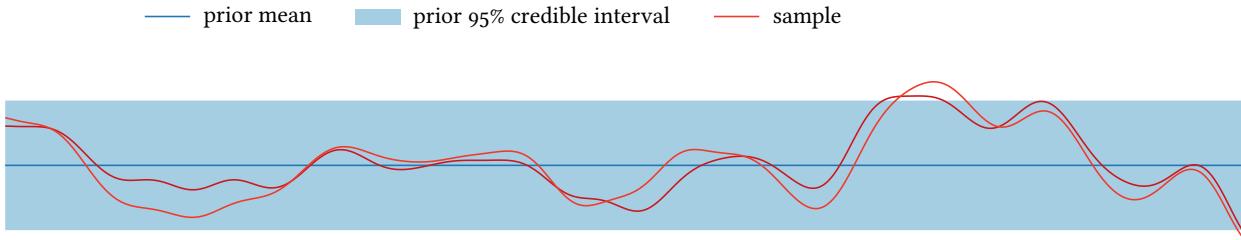


Figure 2.5: A joint Gaussian process over two functions on the shared domain $\mathcal{X} = [0, 30]$. The marginal belief over both functions is the same as our example GP from figure 2.1, but the cross-covariance (2.26) between the functions strongly couples their behavior. We also show a sample from the joint distribution illustrating the strong correlation induced by the joint prior.

zero mean function $\mu \equiv 0$ and squared exponential covariance function (2.4). We augment our original function with a companion function g , defined on the same domain, that has exactly the same marginal GP distribution. However, we couple the distribution of f and g by defining a nontrivial cross-covariance function K_{fg} (2.24):

$$K_{fg}(x, x') = 0.9K(x, x'), \quad (2.26)$$

where K is the marginal covariance function of f and g . A consequence of this choice is that for any given point $x \in \mathcal{X}$, the correlation of the corresponding function values $\phi = f(x)$ and $\gamma = g(x)$ is quite strong:

$$\text{corr}[\phi, \gamma | x] = 0.9. \quad (2.27)$$

We illustrate the resulting joint GP in figure 2.5. The marginal credible intervals for f (and now g) have not changed from our original example in figure 2.1. However, drawing a sample of the functions from their joint distribution reveals the strong coupling encoded in the prior (2.26–2.27).

Inference for joint Gaussian processes

inference from jointly Gaussian distributed observations: § 2.2, p. 18

The construction in (2.22) allows us to reason about a joint Gaussian process as if it were a single GP. This allows us to condition a joint GP on observations of jointly Gaussian distributed values following the procedure outlined previously. In figure 2.6, we condition the joint GP prior from figure 2.5 on ten observations: five exact observations of f on the left-hand side of the domain and five exact observations of g on the right-hand side. Due to the strong correlation between the two functions, an observation of either function strongly informs our belief about the other, even in regions where there are no direct observations.

2.5 CONTINUITY

In this and the following sections we will establish some important properties of Gaussian processes determined by the properties of their

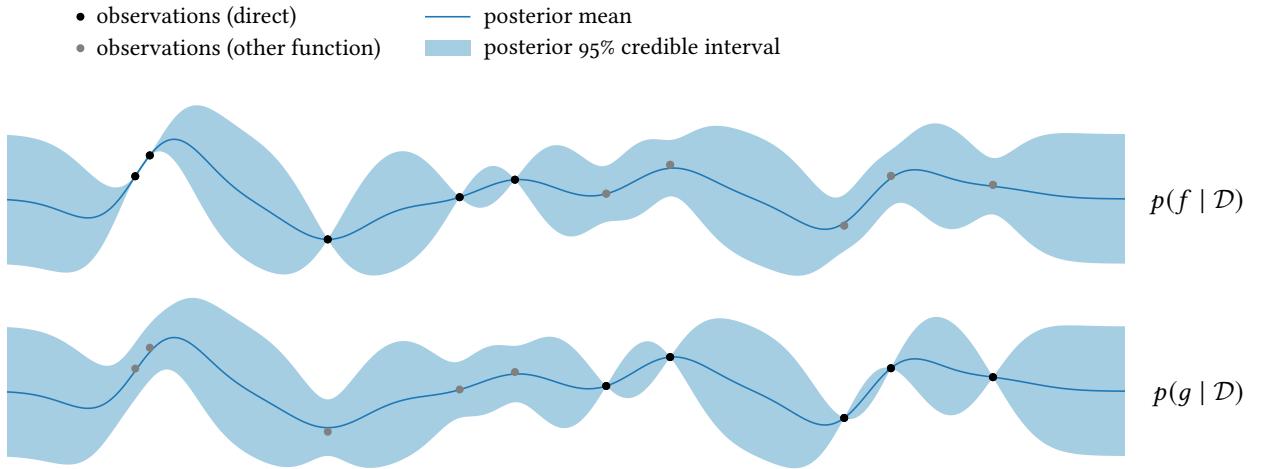


Figure 2.6: The joint posterior for our example joint GP prior in figure 2.5 conditioned on five exact observations of each function.

moments. As a GP is completely specified by its mean and covariance functions, it should not be surprising that the nature of these functions have far-reaching implications regarding properties of the function being modeled. A good familiarity with these implications can help guide model design in practice – the focus of the next two chapters.

To begin, a fundamental question regarding Gaussian processes is whether sample paths are almost surely continuous, and if so how many times differentiable they may be. This is obviously an important consideration for modeling and is also critical to ensure that global optimization is a well-posed problem, as we will discuss later in this chapter. Fortunately, continuity of Gaussian processes is a well-understood property that can be guaranteed almost surely under simple conditions on the mean and covariance functions.

existence of global maxima: § 2.7, p. 34

Suppose $f: \mathcal{X} \rightarrow \mathbb{R}$ has distribution $\mathcal{GP}(f; \mu, K)$. Recall that f is continuous at x if $f(x) - f(x') = \phi - \phi' \rightarrow 0$ when $x' \rightarrow x$. Continuity is thus a limiting property of differences in function values. But under the Gaussian process assumption, this difference is Gaussian distributed (2.5, A.9)! We have

$$p(\phi - \phi' | x, x') = \mathcal{N}(\phi - \phi'; m, s^2),$$

where

$$m = \mu(x) - \mu(x'); \quad s^2 = K(x, x) - 2K(x, x') + K(x', x').$$

Now if μ is continuous at x and K is continuous at $x = x'$, then both $m \rightarrow 0$ and $s^2 \rightarrow 0$ as $x \rightarrow x'$, and thus $\phi - \phi'$ converges in probability to 0. This intuitive condition of continuous moments is known as *continuity in mean square* at x ; if μ and K are both continuous over the entire domain (the latter along the “diagonal” $x = x'$), then we say the entire process is continuous in mean square.

continuity in mean square

sample path continuity

²⁴ R. J. ADLER and J. E. TAYLOR (2007). *Random Fields and Geometry*. Springer-Verlag. [§§ 1.3–1.4]

²⁵ Hölder continuity is a generalization of Lipschitz continuity. Effectively, the covariance function must, in some sense, be “predictably” continuous.

²⁶ W. RUDIN (1976). *Principles of Mathematical Analysis*. McGraw-Hill. [theorem 2.4]

²⁷ Following the discussion in the next section, they in fact are *infinitely* differentiable.

sequences of normal RVs: § A.2, p. 296

differentiability in mean square

joint GP between function and gradient

It turns out that continuity in mean square is not quite sufficient to guarantee that f is simultaneously continuous at every $x \in \mathcal{X}$ with probability one, a property known as *sample path continuity*. However, very slightly stronger conditions on the moments of a GP are sufficient to guarantee sample path continuity.²⁴ The following result is adequate for most settings arising in practice and may be proven as a corollary to the slightly weaker (and slightly more complicated) conditions assumed in ADLER and TAYLOR’s theorem 1.4.1.

Theorem. Suppose $\mathcal{X} \subset \mathbb{R}^d$ is compact and $f: \mathcal{X} \rightarrow \mathbb{R}$ has Gaussian process distribution $\mathcal{GP}(f; \mu, K)$, where μ is continuous and K is Hölder continuous.²⁵ Then f is almost surely continuous on \mathcal{X} .

The condition that $\mathcal{X} \subset \mathbb{R}^d$ be compact is equivalent to the domain being closed and bounded, by the Heine–Borel theorem.²⁶ Applying this result to our example GP in figure 2.1, we conclude that samples from the process are continuous with probability one as the domain $\mathcal{X} = [0, 30]$ is compact and the squared exponential covariance function (2.4) is Hölder continuous. Indeed, the generated samples are very smooth.²⁷

Sample path continuity can also be guaranteed on non-Euclidean domains under similar smoothness conditions.²⁴

2.6 DIFFERENTIABILITY

We can approach the question of differentiability by again reasoning about the limiting behavior of linear transformations of function values. Suppose $f: \mathcal{X} \rightarrow \mathbb{R}$ with $\mathcal{X} \subset \mathbb{R}^d$ has distribution $\mathcal{GP}(f; \mu, K)$, and consider the i th partial derivative of f at \mathbf{x} , if it exists:

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h},$$

where \mathbf{e}_i is the i th standard basis vector. For $h > 0$, the value in the limit is Gaussian distributed as a linear transformation of Gaussian-distributed random variables (A.9). Assuming the corresponding partial derivative of the mean exists at \mathbf{x} and the corresponding partial derivative with respect to each input of the covariance function exists at $\mathbf{x} = \mathbf{x}'$, then as $h \rightarrow 0$ the partial derivative converges in distribution to a Gaussian:

$$p\left(\frac{\partial f}{\partial x_i}(\mathbf{x}) \mid \mathbf{x}\right) = \mathcal{N}\left(\frac{\partial f}{\partial x_i}(\mathbf{x}); \frac{\partial \mu}{\partial x_i}(\mathbf{x}), \frac{\partial^2 K}{\partial x_i \partial x'_i}(\mathbf{x}, \mathbf{x})\right).$$

If this property holds for each coordinate $1 \leq i \leq d$, then f is said to be *differentiable in mean square* at \mathbf{x} .

If f is differentiable in mean square everywhere in the domain, the process itself is called differentiable in mean square, and we have the remarkable result that the function and its gradient have a *joint* Gaussian process distribution:

$$p(f, \nabla f) = \mathcal{GP}\left(\begin{bmatrix} f \\ \nabla f \end{bmatrix}; \begin{bmatrix} \mu \\ \nabla \mu \end{bmatrix}, \begin{bmatrix} K & K\nabla^\top \\ \nabla K & \nabla K \nabla^\top \end{bmatrix}\right). \quad (2.28)$$

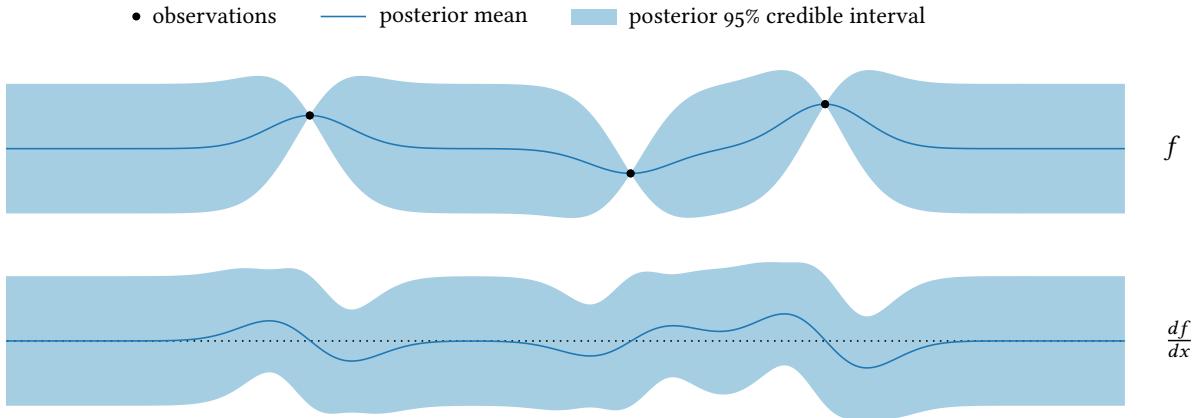


Figure 2.7: The joint posterior of the function and its derivative for our example Gaussian process from figure 2.2. The dashed line in the lower plot corresponds to a derivative of zero.

Here by writing the gradient operator ∇ on the left-hand side of K we mean the result of taking the gradient with respect to its *first* input, and by writing ∇^\top on the right-hand side of K we mean taking the gradient with respect to its *second* input and transposing the result. Thus $\nabla K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^d$ maps pairs of points to column vectors:

$$[\nabla K(\mathbf{x}, \mathbf{x}')]_i = \text{cov}\left[\frac{\partial f}{\partial x_i}(\mathbf{x}), f(\mathbf{x}') \mid \mathbf{x}, \mathbf{x}'\right] = \frac{\partial K}{\partial x_i}(\mathbf{x}, \mathbf{x}'),$$

and $K\nabla^\top: \mathcal{X} \times \mathcal{X} \rightarrow (\mathbb{R}^d)^*$ maps pairs of points to row vectors:

$$K\nabla^\top(\mathbf{x}, \mathbf{x}') = [\nabla K(\mathbf{x}', \mathbf{x})]^\top.$$

Finally, the function $\nabla K\nabla^\top: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^{d \times d}$ represents the result of applying both operations, mapping a pair of points to the covariance matrix between the entries of the corresponding gradients:

$$[\nabla K\nabla^\top(\mathbf{x}, \mathbf{x}')]_{ij} = \text{cov}\left[\frac{\partial f}{\partial x_i}(\mathbf{x}), \frac{\partial f}{\partial x'_j}(\mathbf{x}') \mid \mathbf{x}, \mathbf{x}'\right] = \frac{\partial^2 K}{\partial x_i \partial x'_j}(\mathbf{x}, \mathbf{x}').$$

As the gradient of f has a Gaussian process marginal distribution (2.28), we can reduce the question of *continuous* differentiability to sample path continuity of the gradient process following the discussion above.

Figure 2.7 shows the posterior distribution for the derivative of our example Gaussian process alongside the posterior for the function itself. We can observe a clear correspondence between the two distributions; for example, the posterior mean of the derivative vanishes at critical points of the posterior mean of the function. Notably, we have a great deal of residual uncertainty about the derivative, even at the observed locations. That is because the relatively high spacing between the existing observations limits our ability to accurately estimate the derivative

covariance between $\nabla f(\mathbf{x})$ and $f(\mathbf{x}')$, ∇K

transpose of covariance between $f(\mathbf{x})$ and $\nabla f(\mathbf{x}')$, $K\nabla^\top$

covariance between $\nabla f(\mathbf{x})$ and $\nabla f(\mathbf{x}')$, $\nabla K\nabla^\top$

continuous differentiability

example and discussion

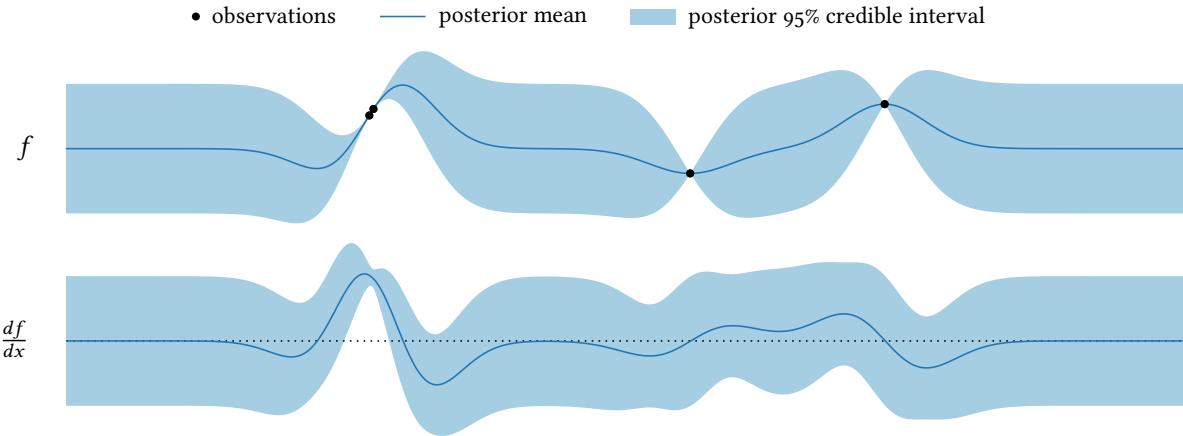


Figure 2.8: The joint posterior of the derivative of our example Gaussian process after adding a new observation nearby another suggesting a large positive slope. The dashed line in the lower plot corresponds to a derivative of zero.

anywhere. Adding an observation immediately next to a previous one significantly reduces the uncertainty in the derivative in that region by effectively providing a finite-difference approximation; see figure 2.8.

Conditioning on derivative observations

However, we can be more direct in specifying derivatives than finite differencing. We can instead condition the joint GP (2.28) *directly* on a derivative observation, as described previously. Figure 2.9 shows the joint posterior after conditioning on an exact observation of the derivative at to the left-most observation location, where the uncertainty in the derivative now vanishes entirely. This capability allows the seamless incorporation of derivative information into an objective function model. Notably, we can even condition a Gaussian process on *noisy* derivative observations as well, as we might obtain in stochastic gradient descent.

inference from jointly Gaussian distributed observations: § 2.2, p. 18

²⁸ For K we again only need to consider the “diagonal” $\mathbf{x} = \mathbf{x}'$.

²⁹ Recall the Hessian is symmetric (assuming the second partial derivatives are continuous) and thus redundant. The *half-vectorization* operator $\text{vech } \mathbf{A}$ maps the upper triangular part of a square, symmetric matrix \mathbf{A} to a vector.

We can reason about derivatives past the first recursively. For example, if μ and K are *twice* differentiable,²⁸ then the (e.g., half-vectorized²⁹) Hessian of f will also have a joint GP distribution with f and its gradient. Defining \mathbf{h} to be the operator mapping a function to its half-vectorized Hessian:

$$\mathbf{h}f = \text{vech } \nabla \nabla^T f,$$

for a Gaussian process with suitably differentiable moments, we have

$$p(\mathbf{h}f) = \mathcal{GP}(\mathbf{h}f; \mathbf{h}\mu, \mathbf{h}K\mathbf{h}^\top), \quad (2.29)$$

where we have used the same notational convention for the transpose. Further, f , ∇f , and $\mathbf{h}f$ will have a joint Gaussian process distribution given by augmenting (2.28) with the marginal in (2.29) and the cross-covariance functions

$$\text{cov}[\mathbf{h}f, f] = \mathbf{h}K; \quad \text{cov}[\mathbf{h}f, \nabla f] = \mathbf{h}K\nabla^\top$$

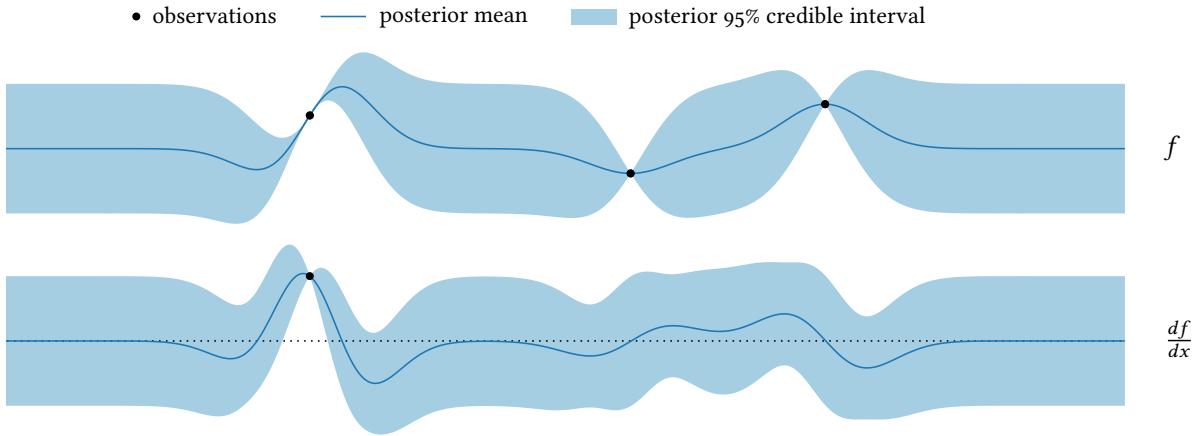


Figure 2.9: The joint posterior of the derivative of our example Gaussian process after adding an exact observation of the derivative at the indicated location. The dashed line in the lower plot corresponds to a derivative of zero.

We can continue further in this vein if needed; however, we rarely reason about derivatives of third-or-higher order in Bayesian optimization.³⁰

³⁰ This is true in classical optimization as well!

Other linear transformations

The joint GP distribution between a suitably smooth GP-distributed function and its gradient (2.28) is simply an infinite-dimensional analog of the general result that Gaussian random variables are jointly Gaussian distributed with arbitrary linear transformations (A.10), after noting that differentiation is a linear operator. We can extend this result to reason about other linear transformations of GP-distributed functions. DIACONIS's original motivation for studying Bayesian numerical methods was *quadrature*, the numerical estimation of intractable integrals.³¹ It turns out that Gaussian processes are a rather convenient model for this task: if $p(f) = \mathcal{GP}(f; \mu, K)$ and we want to reason about the expectation

$$Z = \int f(x) p(x) dx,$$

then (under mild conditions) we again have a joint Gaussian process distribution over f and Z .³² This enables both inference about Z and conditioning on noisy observations of integrals, such as a Monte Carlo estimate of an expectation. The former is the basis for *Bayesian quadrature*, an analog of Bayesian optimization bringing Bayesian experimental design to bear on numerical integration.^{31,33,34}

³¹ P. DIACONIS (1988). Bayesian Numerical Analysis. In: *Statistical Decision Theory and Related Topics IV*.

³² This can be shown, for example, by considering the limiting distribution of Riemann sums.

³³ A. O'HAGAN (1991). Bayes–Hermite quadrature. *Journal of Statistical Planning and Inference* 29(3):245–260

³⁴ C. E. RASMUSSEN and Z. GHAHRAMANI (2002). Bayesian Monte Carlo. *NeurIPS 2002*

2.7 EXISTENCE AND UNIQUENESS OF GLOBAL MAXIMA

The primary use of GPS in Bayesian optimization is to inform optimization decisions, which will be our focus for the majority of this book. Before continuing down this path, we pause to consider whether global

optimization of a GP-distributed function is a well-posed problem – in particular, whether the model guarantees the existence of a global maximum at all.

Consider a function $f: \mathcal{X} \rightarrow \mathbb{R}$ with distribution $\mathcal{GP}(f; \mu, K)$, and consider the location and value of its global optimum, if one exists:

$$x^* = \arg \max_{x \in \mathcal{X}} f(x); \quad f^* = \max_{x \in \mathcal{X}} f(x) = f(x^*).$$

mutual information and entropy search: § 7.6,
p. 135

As f is unknown, these quantities are random variables. Many Bayesian optimization algorithms operate by reasoning about the distributions of (and uncertainties in) these quantities induced by our belief on f .

There are two technical issues we must address. The first is whether we can be certain that a globally optimal value f^* exists when the objective function is random. If existence is not guaranteed, then its distribution is meaningless. The second issue is one of uniqueness: assuming the objective does attain a maximal value, can we be certain the optimum is unique? In general x^* is a *set*-valued random variable, and thus its distribution might have support over arbitrary subsets of the domain, rendering it complicated to reason about. However, if we could ensure the uniqueness of x^* , its distribution would have support on \mathcal{X} rather than its power set, allowing more straightforward inference.

Both the existence of f^* and uniqueness of x^* are tacitly assumed throughout the Bayesian optimization literature when building algorithms based on distributions of these quantities, but these properties are not guaranteed for arbitrary Gaussian processes. However, we can ensure these properties hold almost surely under mild conditions.

Existence of global maxima

To begin, guaranteeing the existence of an optimal value is straightforward if we suppose the domain \mathcal{X} is compact, a pervasive assumption in optimization. This is no coincidence! In this case, if f is continuous then it achieves a global optimum by the extreme value theorem.³⁵ Thus sample path continuity of f and compactness of \mathcal{X} is sufficient to ensure that f^* exists almost surely. Both conditions can be readily established: sample path continuity by following our previous discussion, and compactness of the domain by standard arguments (for example, ensuring that $\mathcal{X} \subset \mathbb{R}^d$ be closed and bounded).

Uniqueness of global maxima

We now turn to the question of uniqueness of x^* , which obviously only becomes a meaningful question after presupposing that f^* exists. Again, this condition is easy to ensure almost surely under simple conditions on the covariance function of a Gaussian process.

KIM and POLLARD considered this issue and provided straightforward conditions under which the uniqueness of x^* is guaranteed for a centered Gaussian process.^{36,37} Namely, no two unique points in the domain can

³⁵ W. RUDIN (1976). *Principles of Mathematical Analysis*. McGraw–Hill. [theorem 4.16]

sample path continuity: § 2.5, p. 28

³⁶ A centered Gaussian process has identically zero mean function $\mu \equiv 0$.

³⁷ J. KIM and D. POLLARD (1990). Cube Root Asymptotics. *The Annals of Statistics* 18(1):191–219. [lemma 2.6]

have perfectly correlated function values, a natural condition that can be easily verified.

Theorem (KIM and POLLARD, 1990). *Let \mathcal{X} be a compact metric space.³⁸ Suppose $f: \mathcal{X} \rightarrow \mathbb{R}$ has distribution $\mathcal{GP}(f; \mu \equiv 0, K)$, and that f is sample path continuous. If for all $x, x' \in \mathcal{X}$ with $x \neq x'$ we have*

$$\text{var}[\phi - \phi' | x, x'] = K(x, x) - 2K(x, x') + K(x', x') \neq 0,$$

then f almost surely has a unique maximum on \mathcal{X} .

ARCONES provided slightly weaker conditions for uniqueness of the supremum, avoiding the requirement of sample path continuity.³⁹

Counterexamples

Although the above conditions for ensuring existence of f^* and uniqueness of x^* are fairly mild, it is easy to construct counterexamples.

Consider a function on the closed unit interval, which we note is compact: $f: [0, 1] \rightarrow \mathbb{R}$. We endow f with a “white noise”⁴⁰ Gaussian process with

$$\mu(x) \equiv 0; \quad K(x, x') = [x = x'].$$

Now f almost surely does not have a maximum. Roughly, because the value of f at every point in the domain is independent of every other, there will almost always be a point with value exceeding any putative maximum.⁴¹ However, the conditions of sample path continuity were violated as the covariance is discontinuous at $x = x'$.

We may also construct a Gaussian process that almost surely achieves a maximum that is not unique. Consider a random function f defined on the (compact) interval $[0, 4\pi]$ defined by the parametric model

$$f(x) = \alpha \cos x + \beta \sin x,$$

where α and β are independent standard normal random variables. Then f has a Gaussian process distribution with

$$\mu(x) \equiv 0; \quad K(x, x') = \cos(x - x'). \quad (2.30)$$

Here μ is continuous and K is Hölder continuous, and thus f is sample path continuous and almost surely achieves a global maximum. However, f is also periodic with period 2π with probability one and will thus almost surely achieve its maximum *twice*. Note that the covariance function does not satisfy the conditions outlined in the above theorem, as any input locations separated by 2π have perfectly correlated function values.

2.8 INFERENCE WITH NON-GAUSSIAN OBSERVATIONS AND CONSTRAINTS

Gaussian process inference is tractable when the observed values are jointly Gaussian distributed with the function of interest (2.6). However, this may not always hold for all relevant information we may receive.

³⁸ Although unlikely to matter in practice, KIM and POLLARD allow \mathcal{X} to be σ -compact and show that the supremum (rather than the maximum) is unique under the same conditions.

³⁹ M. A. ARCONES (1992). On the arg max of a Gaussian process. *Statistics & Probability Letters* 15(5):373–374.

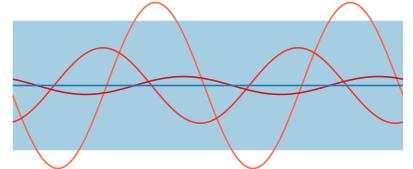
⁴⁰ It turns out this naïve model of white noise has horrible mathematical properties, but it is sufficient for this counterexample.

⁴¹ Let $Q = \mathbb{Q} \cap [0, 1] = \{q_i\}$ be the rationals in the domain and let f^* be a putative maximum. Defining $\phi_i = f(q_i)$, we must have $\phi_i \leq f^*$ for every i ; call this event A .

Define the event A_k by f^* exceeding the first k elements of Q . From independence,

$$\Pr(A_k) = \prod_{i=1}^k \Pr(\phi_i \leq f^*) = \Phi(f^*)^k,$$

so $\Pr(A_k) \rightarrow 0$ as $k \rightarrow \infty$. But $\{A_k\} \nearrow A$, so $\Pr(A) = 0$, and f^* is almost surely not the maximum.



Our counterexample GP without a unique maximum. Every sample achieves its maximum twice.

inference from jointly Gaussian distributed observations: § 2.2, p. 18

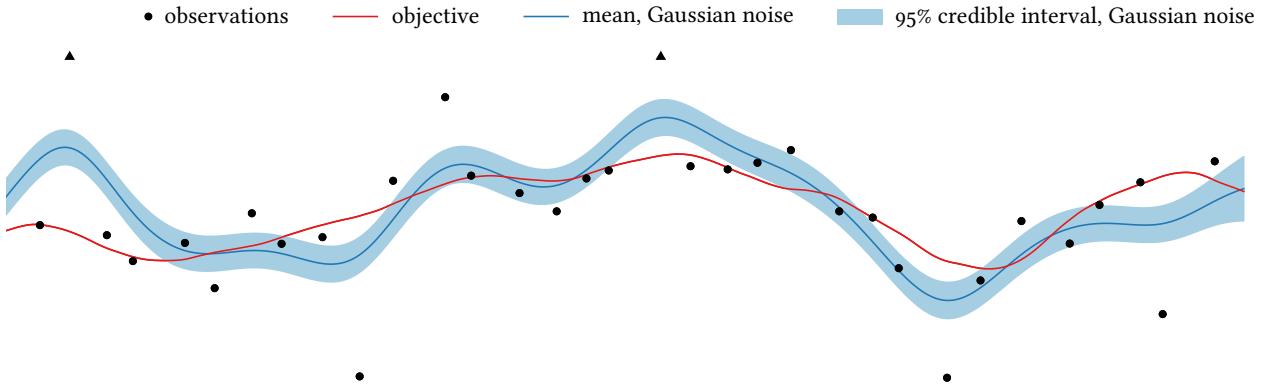
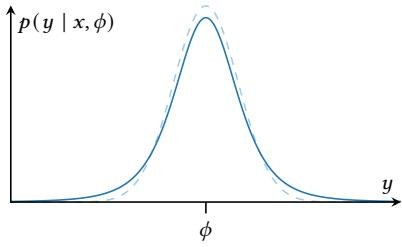


Figure 2.10: Regression with observations corrupted with heavy-tailed noise. The triangular marks indicate observations lying beyond the plotted range. Shown is the posterior distribution of an objective function (ground truth plotted in red) modeling the errors as Gaussian. The posterior is heavily affected by the outliers.



A Student- t error model (solid) with a Gaussian error model (dashed) for reference. The heavier tails of the Student- t model can better explain large outliers.

⁴² K. L. LANGE et al. (1989). Robust Statistical Modeling Using the t Distribution. *Journal of the American Statistical Association* 84(408):881–896.

differentiability, derivative observations: § 2.6,
p. 30

One obvious limitation is an incompatibility with naturally non-Gaussian observations. A scenario particularly relevant to optimization is heavy-tailed noise. Consider the data shown in figure 2.10, where some observations represent extreme outliers. These errors are poorly modeled as Gaussian, and attempting to infer the underlying objective function with the additive Gaussian noise model leads to overfitting and poor predictive performance. A Student- t error model with $v \approx 4$ degrees of freedom provides a robust alternative:⁴²

$$p(y | x, \phi) = T(y; \phi, \sigma_n^2, v).$$

The heavier tails of this model can better explain large outliers; unfortunately, the non-Gaussian nature of this model also renders exact inference impossible. We will demonstrate how to overcome this impasse.

Constraints on an objective function, such as bounds on given function values, can also provide valuable information during optimization, but many natural constraints cannot be reduced to observations that can be handled in closed form. Several Bayesian optimization policies impose hypothetical constraints on the objective function when designing each observation, requiring inference from intractable constraints even when the observations themselves pose no difficulties.

To see how constraints might arise in optimization, consider a Gaussian process belief on a one-dimensional objective f , and suppose we wish to condition on f on having a *local* maximum at a given location x . Assuming the function is twice differentiable, we can invoke the second-derivative test to encode this information in two constraints:

$$f'(x) = 0; \quad f''(x) < 0. \quad (2.31)$$

We can condition a GP on the first of these conditions by following our previous discussion. However, no GP is compatible with the second

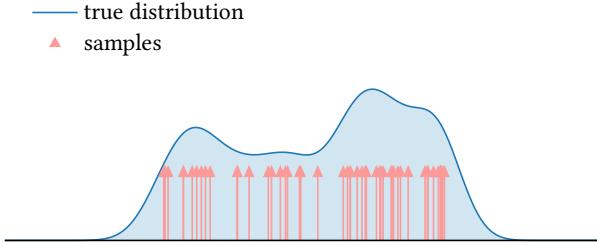


Figure 2.11: The probability density function of an example distribution along with 50 samples drawn independently from the distribution. In Monte Carlo approaches, the distribution is effectively approximated by a mixture of Dirac delta distributions at the sample locations.

condition as $f''(x)$ would necessarily have a Gaussian distribution with unbounded support (2.29). We need some other means to proceed.

Non-Gaussian observations: general case

We can address both non-Gaussian observations and constraints with the following general case, which is flexible enough to handle a large range of information. As in our discussion on exact inference, suppose there is some vector \mathbf{y} sharing a joint Gaussian process distribution with a function of interest f (2.6):

$$p(f, \mathbf{y}) = \mathcal{GP}\left(\begin{bmatrix} f \\ \mathbf{y} \end{bmatrix}; \begin{bmatrix} \mu \\ \mathbf{m} \end{bmatrix}, \begin{bmatrix} K & \kappa^\top \\ \kappa & C \end{bmatrix}\right).$$

Suppose we receive some information about \mathbf{y} in the form of information \mathcal{D} inducing a non-Gaussian posterior on \mathbf{y} . Here, it is convenient to adopt the language of factor graphs⁴³ and write the resulting posterior as proportional to the prior weighted by a function $t(\mathbf{y})$ encoding the available information, which may factorize:

$$p(\mathbf{y} | \mathcal{D}) \propto p(\mathbf{y}) t(\mathbf{y}) = \mathcal{N}(\mathbf{y}; \mathbf{m}, \mathbf{C}) \prod_i t_i(\mathbf{y}). \quad (2.32)$$

The functions $\{t_i\}$ are called *factors* or *local functions* that may comprise a likelihood augmented by any desired (hard or soft) constraints. The term “local functions” arises because each factor often depends only on a low-dimensional subspace of \mathbf{y} , often a single entry.⁴⁴

The posterior on \mathbf{y} (2.32) in turn induces a posterior on f :

$$p(f | \mathcal{D}) = \int p(f | \mathbf{y}) p(\mathbf{y} | \mathcal{D}) d\mathbf{y}. \quad (2.33)$$

At first glance, we may hope to resolve this posterior easily as $p(f | \mathbf{y})$ is a Gaussian process (2.9–2.10). Unfortunately, the non-Gaussian posterior on \mathbf{y} usually renders the posterior on f intractable.

Monte Carlo sampling

A Monte Carlo approach to approximating the f posterior (2.33) begins by drawing samples from the \mathbf{y} posterior (2.32):

$$\{\mathbf{y}_i\}_{i=1}^s \sim p(\mathbf{y} | \mathcal{D}).$$

⁴³ F. R. KSCHISCHANG et al. (2001). Factor Graphs and the Sum–Product Algorithm. *IEEE Transactions on Information Theory* 47(2):498–519.

factors, local functions, $\{t_i\}$

⁴⁴ For example, when observations are conditionally independent given the corresponding function values, the likelihood factorizes into a product of one-dimensional factors (1.3):

$$p(\mathbf{y} | \mathbf{x}, \boldsymbol{\phi}) = \prod_i p(y_i | x_i, \phi_i).$$

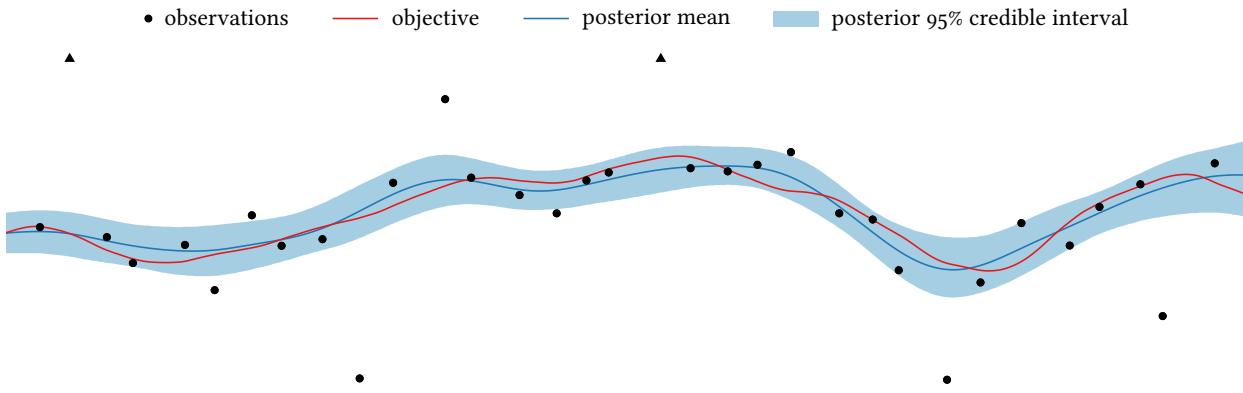


Figure 2.12: Regression with observations corrupted with heavy-tailed noise. The triangular marks indicate observations lying beyond the plotted range. Shown is the posterior distribution of an objective function (ground truth plotted in red) modeling the errors as Student- t distributed with $v = 4$ degrees of freedom. The posterior was approximated from 100 000 Monte Carlo samples. Comparing with the additive Gaussian noise model from figure 2.10, this model effectively ignores the outliers and the fit is excellent.

⁴⁵ *Handbook of Markov Chain Monte Carlo* (2011). Chapman & Hall.

⁴⁶ I. MURRAY et al. (2010). Elliptical slice sampling. *AISTATS 2010*.

We may generate these by appealing to one of numerous Markov chain Monte Carlo (MCMC) routines.⁴⁵ One natural choice would be *elliptical slice sampling*,⁴⁶ which is specifically tailored for latent Gaussian models of this form. Samples from a one-dimensional toy example distribution are shown in figure 2.11.

Given posterior samples of \mathbf{y} , we may then approximate (2.33) via the standard Monte Carlo estimator

$$p(f \mid \mathcal{D}) \approx \frac{1}{s} \sum_{i=1}^s p(f \mid \mathbf{y}_i) = \frac{1}{s} \sum_{i=1}^s \mathcal{GP}(f; \mu_{\mathcal{D}_i}, K_{\mathcal{D}}). \quad (2.34)$$

This is a mixture of Gaussian processes, each of the form in (2.9–2.10). The posterior mean functions depend on the corresponding \mathbf{y} samples, whereas the posterior covariance functions are identical as there is no dependence on the observed values. In this approximation, the marginal belief about any function value is then a mixture of univariate Gaussians:

$$p(\phi \mid x, \mathcal{D}) \approx \frac{1}{s} \sum_{i=1}^s \mathcal{N}(\phi; \mu_i, \sigma^2); \quad \mu_i = \mu_{\mathcal{D}_i}(x); \quad \sigma^2 = K_{\mathcal{D}}(x, x). \quad (2.35)$$

Although slightly more complex than the Gaussian marginals of a Gaussian process, this is often convenient enough for most needs.

A Monte Carlo approximation to the posterior for the heavy-tailed dataset from figure 2.10 is shown in figure 2.12. The observations were modeled as corrupted by Student- t errors with $v = 4$ degrees of freedom. The posterior was approximated using a truly excessive number of samples (100 000, with a burn-in of 10 000) from the \mathbf{y} posterior drawn using elliptical slice sampling.⁴⁶ The outliers in the data are ignored and the predictive performance is excellent.

example: Student- t observation model

Gaussian approximate inference

An alternative to sampling is *approximate inference*, where we make a parametric approximation to the \mathbf{y} posterior that yields a tractable posterior on f . In particular, if the posterior (2.32) were actually *normal*, it would induce a Gaussian process posterior on f . This insight is the basis for most approximations.

In this vein, we proceed by first – somehow – approximating the true posterior over \mathbf{y} with a multivariate Gaussian distribution:

$$p(\mathbf{y} \mid \mathcal{D}) \approx q(\mathbf{y} \mid \mathcal{D}) = \mathcal{N}(\mathbf{y}; \tilde{\mathbf{m}}, \tilde{\mathbf{C}}). \quad (2.36)$$

We are free to design this approximation as we see fit. There are several general-purpose approaches available, distinguished by how they approach maximizing the fidelity of fitting the true posterior (2.32). These include the Laplace approximation, Gaussian expectation propagation, and variational Bayesian inference. The first two of these methods are covered in appendix B, and NICKISCH and RASMUSSEN provide an extensive survey of these and other approaches in the context of Gaussian process binary classification.⁴⁷

Regardless of the details of the approximation scheme, the high-level result is the same – the normal approximation (2.36) in turn induces an approximate Gaussian process posterior on f . To demonstrate this, we consider the posterior on f that would arise from a direct observation of \mathbf{y} (2.9–2.10) and integrate against the approximate posterior (2.36):

$$p(f \mid \mathcal{D}) \approx \int p(f \mid \mathbf{y}) q(\mathbf{y} \mid \mathcal{D}) d\mathbf{y} = \mathcal{GP}(f; \mu_{\mathcal{D}}, K_{\mathcal{D}}), \quad (2.37)$$

where

$$\begin{aligned} \mu_{\mathcal{D}}(x) &= \mu(x) + \kappa(x)^T \mathbf{C}^{-1} (\tilde{\mathbf{m}} - \mathbf{m}); \\ K_{\mathcal{D}}(x, x') &= K(x, x') - \kappa(x)^T \mathbf{C}^{-1} (\mathbf{C} - \tilde{\mathbf{C}}) \mathbf{C}^{-1} \kappa(x'). \end{aligned} \quad (2.38)$$

For most approximation schemes, the posterior covariance on f simplifies to a nicer, more familiar form. Most approximations to the \mathbf{y} posterior (2.36) yield an approximate posterior covariance of the form

$$\tilde{\mathbf{C}} = \mathbf{C} - \mathbf{C}(\mathbf{C} + \mathbf{N})^{-1} \mathbf{C}, \quad (2.39)$$

where \mathbf{N} is positive definite. Although this might appear mysterious, it is actually a natural form: it is the posterior covariance that would result from observing \mathbf{y} corrupted by additive Gaussian noise with covariance \mathbf{N} (2.19), except we are now free to design the noise covariance to maximize the fit. For approximations of this form (2.39), the approximate posterior covariance function on f simplifies to

$$K_{\mathcal{D}}(x, x') = K(x, x') - \kappa(x)^T (\mathbf{C} + \mathbf{N})^{-1} \kappa(x'). \quad (2.40)$$

To demonstrate the power of approximate inference, we return to our motivating scenario of conditioning a one-dimensional process on having a local maximum at an identified point x , which we can achieve by

Laplace approximation: § B.1, p. 297

Gaussian expectation propagation: § B.2 p. 298

⁴⁷ H. NICKISCH and C. E. RASMUSSEN (2008). Approximations for Binary Gaussian Process Classification. *Journal of Machine Learning Research* 9(Oct):2035–2078.

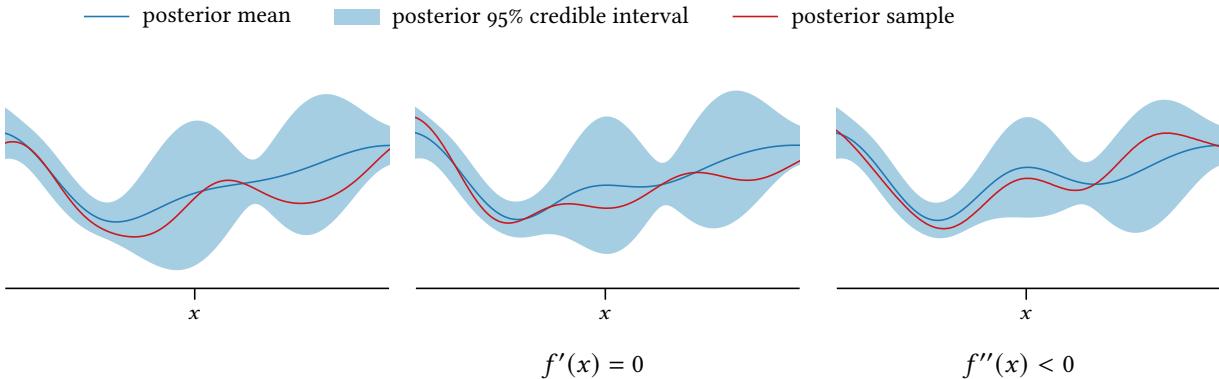


Figure 2.13: Approximately conditioning a Gaussian process to have a local maximum at the marked point x . We show each stage of the conditioning process with a sample drawn from the corresponding posterior. We begin with the unconstrained process (left), which we condition on the first derivative being zero at x using exact inference (middle). Finally we use Gaussian expectation propagation to approximately condition on the second derivative being negative at x .

derivative observations: § 2.6, p. 32

conditioning the first derivative to be zero and constraining the second derivative to be negative at x (2.31). We illustrate an approximation to the resulting posterior step-by-step in figure 2.13, beginning with the example Gaussian process in the left-most panel. We first condition the process on the first derivative observation $f'(x) = 0$ using *exact* inference; the result is shown in the middle panel. Both the updated posterior mean and the sample reflect this information; however, the sample displays a local *minimum* at x , as the second-derivative constraint has not yet been addressed.

To incorporate the second-derivative constraint, we begin with this updated GP and consider the second derivative $h = f''(x)$, which is Gaussian distributed prior to the constraint (2.29):

$$p(h) = \mathcal{N}(h; m, s^2).$$

The negativity constraint induces a posterior on h incorporating the factor $[h < 0]$ (2.32); see figure 2.14:

$$p(h | \mathcal{D}) \propto p(h) [h < 0].$$

The result is a truncated normal posterior on h . We may use Gaussian expectation propagation, which is especially convenient for handling bound constraints of this form, to produce a Gaussian approximation:

$$p(h | \mathcal{D}) \approx q(h | \mathcal{D}) = \mathcal{N}(h; \tilde{m}, \tilde{s}^2).$$

Incorporating the updated belief on h into the Gaussian process (2.38) yields the approximate posterior in the right-most panel of figure 2.13. Although there is still some residual probability that the second derivative is positive at x in the approximate posterior (approximately 8%; see figure 2.14), the belief reflects the desired information reasonably faithfully.

48 P. McCULLAGH and J. A. NELDER (1989). *Generalized Linear Models*. Chapman & Hall.

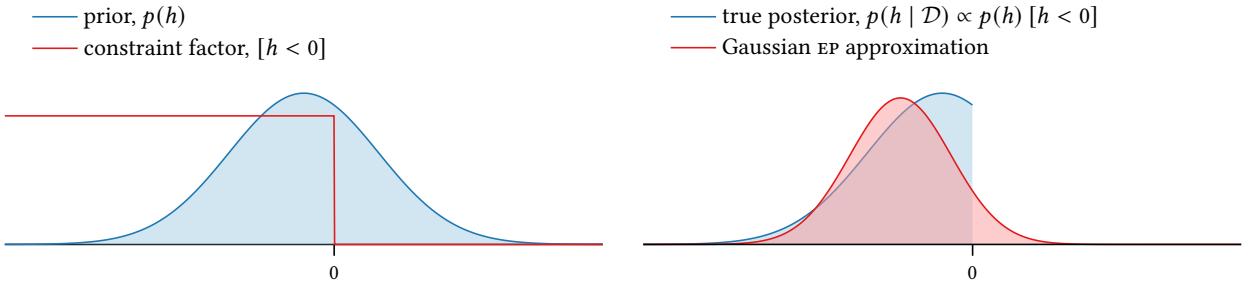


Figure 2.14: A demonstration of Gaussian expectation propagation. On the left we have a Gaussian belief on the second derivative, $p(h)$. We wish to constrain this value to be negative, introducing a step-function factor encoding the constraint, $[h < 0]$. The resulting distribution is non-Gaussian (right), but we can approximate it with a Gaussian, which induces an updated GP posterior on the function approximately incorporating the constraint.

Going beyond this example, we may use the approach outlined above to realize a general framework for Bayesian nonlinear regression by combining a GP prior on a latent function with an observation model appropriate for the task at hand, then approximating the posterior as desired. The convenience and modeling flexibility offered by Gaussian processes can easily justify any extra effort required for approximating the posterior. This can be seen as a nonlinear extension of the well-known family of *generalized linear models*.⁴⁸

This approach is quite popular and been realized countless times. Notable examples include binary classification using a logistic or probit observation model,⁴⁹ modeling point processes as a nonhomogenous Poisson process with unknown intensity,^{50,51} and robust regression with heavy-tailed additive noise such as Laplace⁵² or Student- t ^{53,54} distributed errors. With regard to the latter and our previous heavy-tailed noise example, a Laplace approximation to the posterior for the data in figures 2.10–2.12 with the Student- t observation model produces an approximate posterior in excellent agreement with the Monte Carlo approximation in figure 2.12; see figure 2.15. The cost of approximate inference in this case was dramatically (several orders of magnitude) cheaper than Monte Carlo sampling.

SUMMARY OF MAJOR IDEAS

Gaussian processes have been studied – in one form or another – for over 100 years.⁵⁵ Although we have covered a lot of ground in this chapter, we have only scratched the surface of an expansive body of literature. A good entry point to that literature is RASMUSSEN and WILLIAMS’s monograph, which focuses on machine learning applications of Gaussian processes but also covers their theoretical underpinnings and properties in depth.⁵⁶ A good companion to this work is the book of ADLER and TAYLOR, which takes a deep dive into the properties and geometry of sample paths, including statistical properties of their maxima.⁵⁷

- 49 H. NICKISCH and C. E. RASMUSSEN (2008). Approximations for Binary Gaussian Process Classification. *Journal of Machine Learning Research* 9(Oct):2035–2078.
- 50 J. MØLLER et al. (1998). Log Gaussian Cox Processes. *Scandinavian Journal of Statistics* 25(3): 451–482.
- 51 R. P. ADAMS et al. (2009). Tractable Nonparametric Bayesian Inference in Poisson Processes with Gaussian Process Intensities. *ICML* 2009.
- 52 M. KUSS (2006). Gaussian Process Models for Robust Regression, Classification, and Reinforcement Learning. PhD thesis. Technische Universität Darmstadt.[§ 5.4]
- 53 R. M. NEAL (1997). *Monte Carlo Implementation of Gaussian Process Models for Bayesian Regression and Classification*. Technical report (9702). Department of Statistics, University of Toronto.
- 54 P. JYLÄNKI et al. (2011). Robust Gaussian Process Regression with a Student- t Likelihood. *Journal of Machine Learning Research* 12(99): 3227–3257.
- 55 DIACONIS identified an early application of GPS by POINCARÉ for nonlinear regression:
- P. DIACONIS (1988). Bayesian Numerical Analysis. In: *Statistical Decision Theory and Related Topics IV*.
- H. POINCARÉ (1912). *Calcul des probabilités*. Gauthier-Villars.
- C. E. RASMUSSEN and C. K. I. WILLIAMS (2006). *Gaussian Processes for Machine Learning*. MIT Press.
- R. J. ADLER and J. E. TAYLOR (2007). *Random Fields and Geometry*. Springer-Verlag.

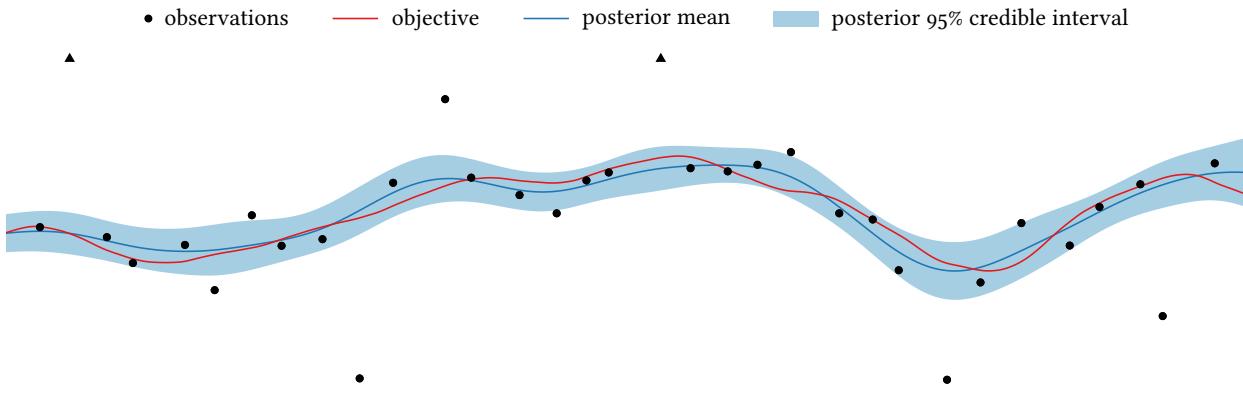


Figure 2.15: A Laplace approximation to the posterior from figure 2.12.

Fortunately, the basic definitions and properties covered in § 2.1 and exact inference procedure covered in § 2.2 already provide a sufficient foundation for the majority of practical applications of Bayesian optimization. This material also provides sufficient background knowledge for the majority of the remainder of the book. However, we wish to underscore the major results from this chapter at a high level.

- Gaussian processes extend the multivariate normal distribution to model functions on infinite domains. As in the finite-dimensional case, Gaussian processes are specified by their first two moments – a mean function and a positive-definite covariance function – which endow any finite set of function values with a multivariate normal distribution (2.2–2.3).
- Conditioning a Gaussian process on function observations that are either exact or corrupted by additive Gaussian noise yields a Gaussian process posterior with updated moments reflecting the assumptions in the prior and the information in the observations (2.9–2.10).
- In fact, we may condition a Gaussian process on the observation of *any* observations sharing a joint Gaussian distribution with the function of interest, including (potentially noisy) *derivative* observations. This is ultimately a result of the closure of Gaussian distributions under linear transformations and the linearity of differentiation.
- In the case of exact inference, the posterior moments of a Gaussian process can be rewritten in terms of correlations among function values and *z*-scores of the observed values in a manner that may be more intuitive than the standard formulas.
- We may extend Gaussian processes to jointly model multiple correlated functions via careful bookkeeping, a construction known as a *joint Gaussian process*. Joint GPS are widely used in optimization settings involving multiple objectives and/or cheaper surrogates for an expensive objective.
- Continuity and differentiability of Gaussian process sample paths can be guaranteed under mild assumptions on the mean and covariance

inference from arbitrary joint Gaussian observations: § 2.2, p. 22

derivative observations: § 2.6, p. 32

interpretation of posterior moments: § 2.2,
p. 21

joint Gaussian processes: § 2.4, p. 26

extensions and other settings: chapter 11,
p. 243

continuity: § 2.5, p. 28
differentiability: § 2.6, p. 30

functions. When these functions are sufficiently differentiable, a GP-distributed function shares a joint GP distribution with its gradient (2.28).

- The existence and uniqueness of global maxima for Gaussian process sample paths can be guaranteed under mild assumptions on the mean and covariance functions. Establishing these properties ensures that the location x^* and value f^* of the global maximum are well-founded random variables, which will be critical for some optimization methods introduced later in the book.⁵⁸

Looking forward, the focus of this chapter has been on theoretical rather than practical properties of Gaussian processes. A huge outstanding question is how to actually *design* a Gaussian process to model a given system. This will be our focus for the next two chapters. In the next chapter, we will explore model *construction*, and in the following chapter we will consider model *assessment* in light of available data.

Finally, we have not yet discussed any computational issues inherent to Gaussian process inference, including, most importantly, how the cost of computing the posterior grows with respect to the number of observations. We will discuss implementation details and scaling in a dedicated chapter later in the book.

existence and uniqueness of global maxima:
§ 2.7, p. 33

⁵⁸ In particular, policies grounded in information theory under the umbrella of “entropy search.” See § 7.6, p. 135 for more.

implementation and scaling of Gaussian process inference: § 9.1, p. 201

3

MODELING WITH GAUSSIAN PROCESSES

Bayesian optimization relies on a faithful model of the system of interest to make well-informed decisions. In fact, even more so than the details of the optimization policy, the fidelity of the underlying model of the objective function is the most decisive factor determining optimization performance. This has been long acknowledged, with MOCKUS for example commenting in his seminal work that:¹

The development of some system of a priori distributions suitable for different classes of the function f is probably the most important problem in the application of [the] Bayesian approach to... global optimization.

The importance of careful modeling has not waned in the intervening years, but our capacity for building sophisticated models has improved.

Recall our approach to modeling observations obtained during optimization combines a prior process for a (perhaps not directly observable) objective function (1.8) and an observation model linking the values of the objective to measured values (1.2). Both distributions must be specified before we can derive a posterior belief about the objective function (1.10) and predictive distribution for proposed observations (1.7), which together serve as the key enablers of Bayesian optimization policies.

In practice, the choice of observation model is often noncontroversial,² and our running prototypes of exact observation and additive Gaussian noise suffice for many systems. The bulk of modeling effort is thus spent crafting the prior process. Although specifying a *Gaussian* process is seemingly as simple as choosing a mean and covariance function, it can be difficult to intuit appropriate choices without a great deal of knowledge about the system of interest. As an alternative to prior knowledge, we may appeal to a data-driven approach, where we establish a space of candidate models and search through this space for those offering the best explanation of available data. Almost all Gaussian process models used in practice are designed in this manner, and we will lay the groundwork for this approach in this chapter and the next.

As a Gaussian process is specified by its first two moments, data-driven model design boils down to searching for the prior mean and covariance functions most harmonious with our observations. This can be a daunting task as the space of possibilities is limitless. However, we do not need to begin from scratch: there are mean and covariance functions available off-the-shelf for modeling a variety behavioral archetypes, and by systematically combining these components we may model functions with a rich variety of behavior. We will explore the world of possibilities in this chapter, while addressing details important to optimization.

Once we have established a space of candidate models, we will require some mechanism to differentiate possible choices based on their merits, a process known as *model assessment* that we will explore at length in the next chapter. We will begin the present discussion by revisiting the topic

¹ J. MOCKUS (1974). On Bayesian Methods for Seeking the Extremum. *Optimization Techniques IFIP Technical Conference*.

Bayesian inference of the objective function:
§ 1.2, p. 8

² However, we may not be certain about some details, such as the scale of observation noise, an issue we will address in the next chapter.

chapter 4: model assessment, selection, and averaging, p. 67

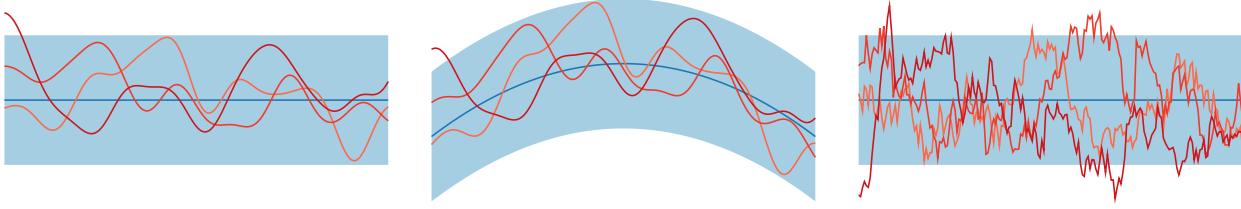


Figure 3.1: The importance of the prior mean function in determining sample path behavior. The models in the first two panels differ in their mean function but share the same covariance function. Sample path behavior is identical up to translation. The model in the third panel features the same mean function as the first panel but a different covariance function. Samples exhibit dramatically different behavior.

of prior mean and covariance functions with an eye toward practical utility.

3.1 THE PRIOR MEAN FUNCTION

Recall the mean function of a Gaussian process specifies the expected value of an arbitrary function value $\phi = f(x)$:

$$\mu(x) = \mathbb{E}[\phi | x].$$

Although this is obviously a fundamental concern, the choice of prior mean function has received relatively little consideration in the Bayesian optimization literature.

There are several reasons for this. To begin, it is actually the *covariance* function rather than the mean function that largely determines the behavior of sample paths. This should not be surprising: the mean function only affects the *marginal* distribution of function values, whereas the covariance function can further modify the *joint* distribution of function values. To elaborate, consider an arbitrary Gaussian process $\mathcal{GP}(f; \mu, K)$. Its sample paths are distributed identically to those from the corresponding centered process $f - \mu$, after shifting pointwise by μ . Therefore the sample paths of *any* Gaussian process with the same covariance function are effectively the same up to translation, and it is the covariance function determining their behavior otherwise; see the demonstration in figure 3.1.

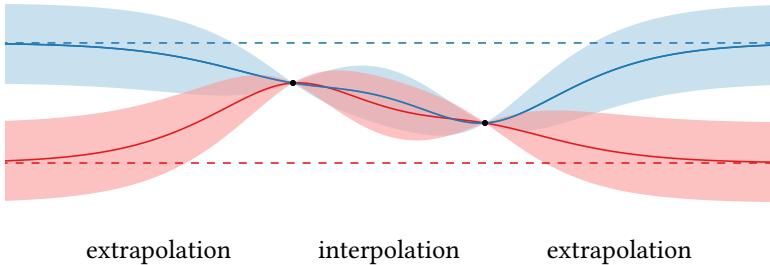
It is also important to understand the role of the prior mean function in the posterior process. Suppose we condition a Gaussian process $\mathcal{GP}(f; \mu, K)$ on the observation of a vector y with marginal distribution (2.7) and cross covariance function (2.24)

$$p(y) = \mathcal{N}(y; m, C); \quad \kappa(x) = \text{cov}[y, \phi | x].$$

The prior mean influences the posterior process *only* through the posterior mean (2.10):

$$\mu_D(x) = \mu(x) + \kappa(x)^\top C^{-1}(y - m).$$

3.1. THE PRIOR MEAN FUNCTION



We can roughly understand the behavior of the posterior mean by identifying two regimes determined by the strength of correlation between a given function value and the observations. In “interpolatory” regions, where function values have significant correlation with one-or-more observed value, the posterior mean is mostly determined by the data rather than the prior mean. On the other hand, in “extrapolatory” regions, where $\kappa(x) \approx 0$, the data have little influence and the posterior mean effectively equals the prior mean. Figure 3.2 demonstrates this effect.

Constant mean function

The primary impact of the prior mean on our predictions – and on an optimization policy informed by these predictions – is in the extrapolatory regime. However, extrapolation without strong prior knowledge can be a dangerous business. In Bayesian optimization, the prior mean is almost always taken to be a constant:

$$\mu(x; c) \equiv c, \quad (3.1)$$

in order to avoid any unwanted bias on our decisions caused by spurious structure in the prior process. This simple choice is supported empirically by a study comparing optimization performance across a range of problems as a function of the choice of prior mean.³

When adopting a constant mean, the value of the constant c is usually treated as a parameter to be estimated or (approximately) marginalized, as we will discuss in the next chapter. However, we can actually do better in some cases. Consider a parametric Gaussian process prior with constant mean (3.1) and arbitrary covariance function:

$$p(f | c) = \mathcal{GP}(f; \mu \equiv c, K),$$

and suppose we place a normal prior on c :

$$p(c) = \mathcal{N}(c; a, b^2). \quad (3.2)$$

Then we can marginalize the unknown constant mean *exactly* to derive the marginal Gaussian process

$$p(f) = \int p(f | c) p(c) dc = \mathcal{GP}(f; \mu \equiv a, K + b^2), \quad (3.3)$$

Figure 3.2: The influence of the prior mean on the posterior mean. We show two Gaussian process posteriors differing only in their prior mean functions, shown as dashed lines. In the “interpolatory” region between the observations, the posterior means are mostly determined by the data, but devolve to the respective prior means when extrapolating outside this region.

behavior in interpolatory and extrapolatory regions

³ G. DE ATH et al. (2020). What do you Mean? The Role of the Mean Function in Bayesian Optimization. *GECCO 2020*.

marginalizing constant prior mean

model selection and averaging: §§ 4.3–4.4,
p. 73

⁴ Noting that c and f form a joint Gaussian process, we may perform inference as described in § 2.4, p. 26 to reveal their joint posterior.

⁵ The basis functions can be arbitrarily complex, such as the output layer of a deep neural network:

J. SNOEK et al. (2015). Scalable Bayesian Optimization Using Deep Neural Networks. *ICML 2015*.

basis functions, ψ
weight vector, β

⁶ A. o'HAGAN (1978). Curve Fitting and Optimal Design for Prediction. *Journal of the Royal Statistical Society Series B (Methodological)* 40(1): 1–42.

⁷ C. E. RASMUSSEN and C. K. I. WILLIAMS (2006). *Gaussian Processes for Machine Learning*. MIT Press. [§ 2.7]

⁸ This mean function was proposed in the context of Bayesian optimization (with diagonal A) by

J. SNOEK et al. (2015). Scalable Bayesian Optimization Using Deep Neural Networks. *ICML 2015*,

who also proposed appropriate priors for A and b . The mean was also proposed in the related context of Bayesian quadrature (see § 2.6, p. 33) by

L. ACERBI (2018). Variational Bayesian Monte Carlo. *NeurIPS 2018*.

were the uncertainty in the mean has been absorbed into the prior covariance function. We may now use this prior directly, avoiding any estimation of c . The unknown mean will be automatically marginalized in both the prior and posterior process, and we may additionally derive the posterior belief over c given data if it is of interest.⁴

Linear combination of basis functions

We may extend the above result to marginalize the weights of an *arbitrary* linear combination of basis functions under a normal prior, making this a particularly convenient class of mean functions. Namely, consider a parametric mean function of the form

$$\mu(x; \beta) = \beta^\top \psi(x), \quad (3.4)$$

where the vector-valued function $\psi: \mathcal{X} \rightarrow \mathbb{R}^n$ defines the basis functions and β is a vector of weights.⁵

Now consider a parametric Gaussian process prior with a mean function of this form (3.4) and arbitrary covariance function K . Placing a multivariate normal prior on β ,

$$p(\beta) = \mathcal{N}(\beta; a, B), \quad (3.5)$$

and marginalizing yields the marginal prior,^{6,7}

$$p(f) = \mathcal{GP}(f; m, C),$$

where

$$m(x) = a^\top \psi(x); \quad C(x, x') = K(x, x') + \psi(x)^\top B \psi(x'). \quad (3.6)$$

We may recover the constant mean case above by taking $\psi(x) \equiv 1$.

Other options

Of course, we stress that a constant or linear mean function is by no means necessary, and when a system is understood sufficiently well to suggest a plausible alternative – perhaps the output of a relatively simple predictive model – it should be strongly considered.

One option that might be a reasonable choice in some optimization contexts is a concave quadratic mean:

$$\mu(x; A, b, c) = (x - b)^\top A^{-1} (x - b) + c, \quad (3.7)$$

where $A < 0$.⁸ This mean encodes that values near b (according to the Mahalanobis distance (A.8)) are expected to be higher than those farther away and could reasonably model an objective function expected to be “bowl-shaped” to a first approximation. The middle panel of figure 3.1 incorporates a mean of this form; note that the maxima of sample paths are of course not constrained to agree with that of the prior mean.

3.2 THE PRIOR COVARIANCE FUNCTION

The prior covariance function determines the covariance between the function values corresponding to a pair of input locations x and x' :

$$K(x, x') = \text{cov}[\phi, \phi' | x, x']. \quad (3.8)$$

The covariance function determines fundamental properties of sample path behavior, including continuity, differentiability, and aspects of the global optima, as we have already seen. More so than the mean function, careful design of the covariance function is critical to ensure fidelity in modeling. We will devote considerable discussion to this topic, beginning with some important and moving on to useful examples and mechanisms for systematically modifying and composing multiple covariance functions together to model complex behavior.

After appropriate normalization, a covariance function K may be loosely interpreted as a measure of similarity between pairs of points in the domain. Namely, given $x, x' \in \mathcal{X}$, the correlation between the corresponding function values is

$$\rho = \text{corr}[\phi, \phi' | x, x'] = \frac{K(x, x')}{\sqrt{K(x, x) K(x', x')}}, \quad (3.9)$$

and we may interpret the strength of this dependence as a measure of similarity between the input locations. This intuition can be useful, but some caveats are in order. To begin, note that correlation may be *negative*, which might be interpreted as indicating *dis*-similarity as the function values react to information with opposite sign.

Further, for a proposed covariance function K to be admissible, it must satisfy two global consistency properties ensuring that the collection of random variables comprising f are able to satisfy the purported relationships. First, we can immediately deduce from its definition (3.8) that K must be *symmetric* in its inputs. Second, the covariance function must be *positive semidefinite*; that is, given any finite set of points $\mathbf{x} \subset \mathcal{X}$, the Gram matrix $K(\mathbf{x}, \mathbf{x})$ must have only nonnegative eigenvalues.⁹

As an illustration of how positive semidefiniteness ensures statistical validity, note that a direct consequence is that $K(x, x) = \text{var}[\phi | x] \geq 0$, and thus marginal variance is always nonnegative. On a slightly less trivial level, consider a pair of points $\mathbf{x} = (x, x')$ and normalize the corresponding Gram matrix $\Sigma = K(\mathbf{x}, \mathbf{x})$ to yield the correlation matrix:

$$\mathbf{P} = \text{corr}[\phi | \mathbf{x}] = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix},$$

where ρ is given by (3.9). For this matrix to be valid, we must have $\rho \in [-1, 1]$. This happens precisely when \mathbf{P} is positive semidefinite, as its eigenvalues are $1 \pm \rho$. Finally, noting that \mathbf{P} is congruent to Σ ,¹⁰ we conclude the implied correlations are consistent if and only if Σ is positive semidefinite. With more than two points, the positive semidefiniteness of K ensures similar consistency at higher orders.

sample path continuity: § 2.5, p. 28

sample path differentiability: § 2.6, p. 30

existence and uniqueness of global maxima:

§ 2.7, p. 33

correlation between function values, ρ

symmetry

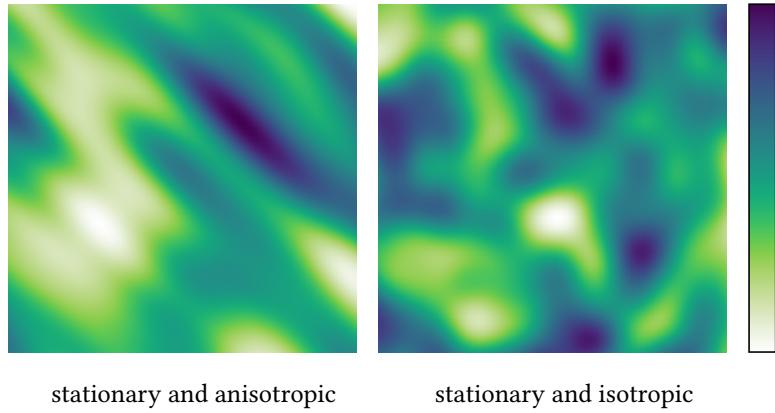
positive semidefiniteness

⁹ Symmetry guarantees the eigenvalues are real.

consequences of positive semidefiniteness

¹⁰ We have $\Sigma = S P S$, where S is diagonal with $S_{ii} = \sqrt{\Sigma_{ii}}$.

Figure 3.3: Left: a sample from a stationary Gaussian process in two dimensions. The joint distribution of function values is translation- but not rotation-invariant, as the function tends to vary faster in some directions than others. Right: a sample from an isotropic process. The joint distribution of function values is both translation- and rotation-invariant.



stationary covariance function, $K(x - x')$

¹¹ Of course this definition requires $x - x'$ to be well defined. This is trivial in Euclidean spaces; a fairly general treatment for more exotic spaces would assume an abelian group structure on \mathcal{X} with binary operation $+$ and inverse $-$ and define $x - x' = x + (-x')$.

isotropic covariance function, $K(d)$

¹² S. BOCHNER (1933). Monotone Funktionen, Stieltjessche Integrale und harmonische Analyse. *Mathematische Annalen* 108:378–410.

¹³ We do not quote the most general version of the theorem here; the result can be extended to complex-valued covariance functions on arbitrary locally compact abelian groups if necessary. It is remarkably universal.

Stationarity, isotropy, and Bochner's theorem

Some covariance functions exhibit structure giving rise to certain computational benefits. Namely, a covariance function $K(x, x')$ that only depends on the difference $x - x'$ is called *stationary*.¹¹ When convenient, we will abuse notation and write a stationary covariance function in terms of a single input, writing $K(x - x')$ for $K(x, x') = K(x - x', 0)$. If a GP has a stationary covariance function and *constant* mean function (3.1), then the process itself is also called stationary. A consequence of stationarity is that the distribution of any set of function values is invariant under translation; that is, the function “acts the same” everywhere from a statistical viewpoint. The left panel of figure 3.3 shows a sample from a 2d stationary GP, demonstrating this translation-invariant behavior.

Stationarity is a convenient assumption when modeling, as defining the *local* behavior around a single point suffices to specify the *global* behavior of an entire function. Many common covariance functions have this property as a result. However, this may not always be a valid assumption in the context of optimization, as an objective function may for example exhibit markedly different behavior near the optimum than elsewhere. We will shortly see some general approaches for addressing nonstationarity when appropriate.

If $\mathcal{X} \subset \mathbb{R}^n$, a covariance function $K(x, x')$ only depending on the Euclidean distance $d = |x - x'|$ is called *isotropic*. Again, when convenient, we will notate such a covariance with $K(d)$. Isotropy is a more restrictive assumption than stationarity – indeed it trivially implies stationarity – as it implies the covariance is invariant to both translation *and* rotation, and thus the function has identical behavior in every direction. An example sample from a 2d isotropic GP is shown in the right panel of figure 3.3. Many of the standard covariance functions we will define shortly will be isotropic on first definition, but we will again develop generic mechanisms to modify them in order to induce anisotropic behavior when desired.

BOCHNER's theorem is a landmark result characterizing stationary covariance functions in terms of their Fourier transforms:^{12, 13}

Theorem (BOCHNER, 1933). A continuous function $K: \mathbb{R}^n \rightarrow \mathbb{R}$ is positive semidefinite (that is, represents a stationary covariance function) if and only if we have

$$K(\mathbf{x}) = \int \exp(2\pi i \mathbf{x}^\top \boldsymbol{\xi}) d\nu,$$

where ν is a finite, positive Borel measure on \mathbb{R}^n . Further, this measure is symmetric around the origin; that is, $\nu(A) = \nu(-A)$ for any Borel set $A \subseteq \mathbb{R}^n$ where $-A$ is the “negation” of A : $-A = \{-a \mid a \in A\}$.

To summarize, BOCHNER’s theorem states that the Fourier transform of any stationary covariance function on \mathbb{R}^n is proportional to a probability measure and vice versa; the constant of proportionality is $K(\mathbf{0})$. The measure ν corresponding to K is called the *spectral measure* of K . When a corresponding density function κ exists, it is called the *spectral density* of K and forms a Fourier pair with K :

$$K(\mathbf{x}) = \int \exp(2\pi i \mathbf{x}^\top \boldsymbol{\xi}) \kappa(\boldsymbol{\xi}) d\boldsymbol{\xi}; \quad \kappa(\boldsymbol{\xi}) = \int \exp(-2\pi i \mathbf{x}^\top \boldsymbol{\xi}) K(\mathbf{x}) d\mathbf{x}. \quad (3.10)$$

The symmetry of the spectral measure implies a similar symmetry in the spectral density: $\kappa(\boldsymbol{\xi}) = \kappa(-\boldsymbol{\xi})$ for all $\boldsymbol{\xi} \in \mathbb{R}^n$.

BOCHNER’s theorem is surprisingly useful in practice, allowing us to approximate an arbitrary stationary covariance function by approximating (e.g., by modeling or sampling from) its spectral density. This is the basis of the *spectral mixture covariance* described in the next section, as well as the *sparse spectrum approximation* scheme, which facilitates the computation of some Bayesian optimization policies.

spectral measure, ν
spectral density, κ

symmetry of spectral density

sparse spectrum approximation: § 8.7, p. 178

3.3 NOTABLE COVARIANCE FUNCTIONS

It can be difficult to define a new covariance function for a given scenario *de novo*, as the positive-semidefinite criterion can be nontrivial to guarantee for what might otherwise be an intuitive notion of similarity. In practice, it is common to instead construct covariance functions by combining and transforming established “building blocks” modeling various types of atomic behavior while following rules guaranteeing the result will be valid. We describe several useful examples below.¹⁴

Our presentation will depart from most in that several of the covariance functions below will initially be defined without parameters that some readers may be expecting. We will shortly demonstrate how coupling these covariance functions with particular transformations of the function domain and output gives rise to common covariance function parameters such as *characteristic length scales* and *output scales*.

¹⁴ For a more complete survey, see

C. E. RASMUSSEN and C. K. I. WILLIAMS (2006). *Gaussian Processes for Machine Learning*. MIT Press. [chapter 4]

The Matérn family and squared exponential covariance

If there is one class of covariance functions to be familiar with, it is the *Matérn family*. This is a versatile family of covariance functions for modeling isotropic behavior on Euclidean domains $\mathcal{X} \subset \mathbb{R}^n$ of any

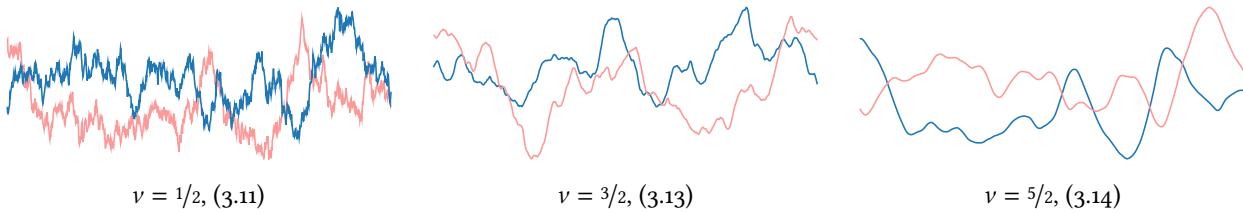


Figure 3.4: Samples from centered Gaussian processes with the Matérn covariance function with different values of the smoothness parameter v . Sample paths with $v = 1/2$ are continuous but not differentiable; incrementing this parameter by one unit increases the number of continuous derivatives by one.

sample path differentiability: § 2.6, p. 30

desired degree of smoothness, in terms of the differentiability of sample paths. The Matérn covariance $K_{M(v)}$ depends on a parameter $v \in \mathbb{R}_{>0}$ determining this smoothness; sample paths from a centered Gaussian process with this covariance are $\lceil v \rceil - 1$ times continuously differentiable. In practice v is almost always taken to be a half-integer, in which case the expression for the covariance assumes a simple form as a function of the Euclidean distance $d = |x - x'|$.

To begin with the extremes, the case $v = 1/2$ yields the so-called *exponential covariance*:

$$K_{M^{1/2}}(x, x') = \exp(-d). \quad (3.11)$$

Sample paths from a centered Gaussian process with exponential covariance are continuous but nowhere differentiable, which is perhaps too rough to be interesting in most optimization contexts. However, this covariance is often encountered in historical literature. In the one-dimensional case $\mathcal{X} \subset \mathbb{R}$, a Gaussian process with this covariance is known as a *Ornstein–Uhlenbeck (ou) process* and satisfies a continuous-time Markov property that renders its posterior moments particularly convenient.

Taking the limit of increasing smoothness $v \rightarrow \infty$ yields the *squared exponential covariance* from the previous chapter:

$$K_{SE}(x, x') = \exp\left(-\frac{1}{2}d^2\right). \quad (3.12)$$

Note that the squared exponential covariance does not belong to the Matérn family – which only models functions with finite smoothness – but instead represents an important limiting case. The squared exponential covariance is without a doubt the most prevalent covariance function in the statistical and machine learning literature. However, it may not always be a good choice in practice. Sample paths from a centered Gaussian process with squared exponential covariance are *infinitely* differentiable, which has been ridiculed as an absurd assumption for most physical processes.¹⁵ STEIN does not mince words on this, starting off a three-sentence “summary of practical suggestions” with “use the Matérn model” and devoting significant effort to discouraging the use of the squared exponential in the context of geostatistics.

Ornstein–Uhlenbeck (ou) process

Samples from a centered Gaussian process with squared exponential covariance K_{SE} .

¹⁵ M. L. STEIN (1999). *Interpolation of Spatial Data: Some Theory for Kriging*. Springer-Verlag. [§ 1.7]

Between these extremes are the cases $\nu = 3/2$ and $\nu = 5/2$, which respectively model once- and twice-differentiable functions:

$$K_{M^{3/2}}(x, x') = (1 + \sqrt{3}d) \exp(-\sqrt{3}d); \quad (3.13)$$

$$K_{M^{5/2}}(x, x') = (1 + \sqrt{5}d + \frac{5}{3}d^2) \exp(-\sqrt{5}d). \quad (3.14)$$

Figure 3.4 illustrates samples from centered Gaussian processes with different values of the smoothness parameters ν . The $\nu = 5/2$ case in particular has been singled out as a prudent off-the-shelf choice for Bayesian optimization when no better alternative is obvious.¹⁶

The spectral mixture covariance

Members of the Matérn family and the squared exponential covariance function express fairly simple correlation structure, with the covariance dropping monotonically to zero as the distance $d = |x - x'|$ increases. All differences in sample path behavior such as differentiability, etc. are expressed entirely through nuances in the tail behavior of the covariance functions; see the figure in the margin.

The Fourier transforms of these covariances are also broadly comparable: all are proportional to unimodal distributions centered on the origin. However, BOCHNER's theorem indicates that there is a vast world of stationary covariance functions indexed by the *entire* space of symmetric spectral measures, which may have considerably more complex structure. Several authors have sought to exploit this characterization to build stationary covariance functions with virtually unlimited flexibility.

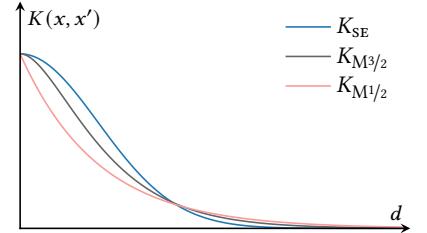
A notable contribution in this direction is the *spectral mixture* covariance function proposed by WILSON and ADAMS.¹⁷ The idea is simple but powerful: we parameterize a space of stationary covariance functions by some suitable family of mixture distributions in the Fourier domain representing their spectral density. The parameters of this spectral mixture distribution specify a covariance function via the correspondence in (3.10), and we can make the resulting family as rich as desired by adjusting the number of components in the mixture. WILSON and ADAMS proposed Gaussian mixtures for the spectral density, which are universal approximators and have a convenient Fourier transform. We define a Gaussian mixture spectral density κ as

$$k(\xi) = \sum_i w_i \mathcal{N}(\xi; \mu_i, \Sigma_i); \quad \kappa(\xi) = \frac{1}{2} [k(\xi) + k(-\xi)],$$

where the indirect construction via k ensures the required symmetry. Note that the weights $\{w_i\}$ must be positive but need not sum to unity. Taking the inverse Fourier transform (3.10), the corresponding covariance function is

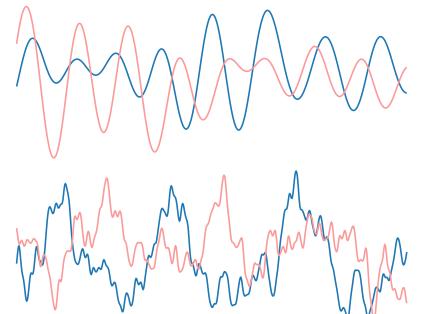
$$\begin{aligned} K_{SM}(\mathbf{x}, \mathbf{x}'; \{w_i\}, \{\mu_i\}, \{\Sigma_i\}) = \\ \sum_i w_i \exp(-2\pi^2(\mathbf{x} - \mathbf{x}')^\top \Sigma_i (\mathbf{x} - \mathbf{x}')) \cos(2\pi(\mathbf{x} - \mathbf{x}')^\top \mu_i). \end{aligned} \quad (3.15)$$

¹⁶ J. SNOEK et al. (2012). Practical Bayesian Optimization of Machine Learning Algorithms. *NeurIPS 2012*.



Some members of the Matérn family and the squared exponential covariance as a function of the distance between inputs. All decay to zero correlation as distance increases.

¹⁷ A. G. WILSON and R. P. ADAMS (2013). Gaussian Process Kernels for Pattern Discovery and Extrapolation. *ICML 2013*.



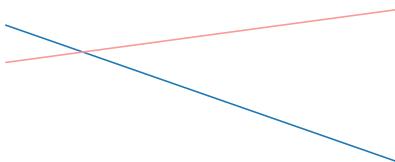
Samples from centered Gaussian processes with two realizations of a Gaussian spectral mixture covariance function, offering a glimpse into the flexibility of this class.

Inspecting this expression, we can see that every covariance function induced by a Gaussian mixture spectral density is infinitely differentiable, and one might object to this choice on the grounds of overly smooth sample paths. This can be mitigated by using enough mixture components to induce sufficiently complex structure in the covariance (on the order of ~ 5 is common). Another option would be to use a different family of spectral distributions; for example, a mixture of Cauchy distributions would induce a family of continuous but nondifferentiable covariance functions analogous to the exponential covariance (3.11), but this idea has not been explored.

Linear covariance function

¹⁸ Independence is usual but not necessary; an arbitrary joint prior would add a term of $2\mathbf{b}^\top \mathbf{x}$ to (3.16), where $\mathbf{b} = \text{cov}[\boldsymbol{\beta}, \boldsymbol{\beta}]$.

linear basis functions: § 3.1, p. 48



Samples from a centered Gaussian process with linear covariance K_{LIN} .

Another useful covariance function arises from a Bayesian realization of linear regression. Let the domain be Euclidean, $\mathcal{X} \subset \mathbb{R}^n$, and consider the model

$$f(\mathbf{x}) = \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{x},$$

where we have abused notation slightly to distinguish the constant term from the remaining coefficients. Following our discussion on linear basis functions, if we take independent¹⁸ normal priors on $\boldsymbol{\beta}$ and $\boldsymbol{\beta}$:

$$p(\boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta}; \mathbf{a}, b^2); \quad p(\boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta}; \mathbf{a}, \mathbf{B}),$$

we arrive at the so-called *linear covariance*:

$$K_{\text{LIN}}(\mathbf{x}, \mathbf{x}'; b, \mathbf{B}) = b^2 + \mathbf{x}^\top \mathbf{B} \mathbf{x}. \quad (3.16)$$

Although this covariance is unlikely to be of any direct use in Bayesian optimization (linear programming is much simpler!), it can be a useful component of more complex composite covariance structures.

3.4 MODIFYING AND COMBINING COVARIANCE FUNCTIONS

With the notable exception of the spectral mixture covariance, which can approximate any stationary covariance function, several of the covariances introduced in the last section are still too rigid to be useful.

In particular, consider any member of the Matérn family or the squared exponential covariance function (3.11–3.14), each of which encodes several explicit and possibly dubious assumptions about the function of interest. To begin, each prescribes unit variance for every function value:

$$\text{var}[\phi | x] = K(x, x) = 1, \quad (3.17)$$

which is an arbitrary, possibly inappropriate choice of scale. Further, each of these covariance functions fixes an isotropic *characteristic length scale* of correlation of approximately one unit:¹⁹ at a separation of $|x - x'| = 1$, the correlation between the corresponding function values drops to roughly

$$\text{corr}[\phi, \phi' | x, x'] \approx 0.5, \quad (3.18)$$

¹⁹ Although an important concept, there is no clear-cut definition of characteristic length scale. It is simply a convenient separation distance for which correlation remains appreciable, but beyond which correlation begins to noticeably decay.

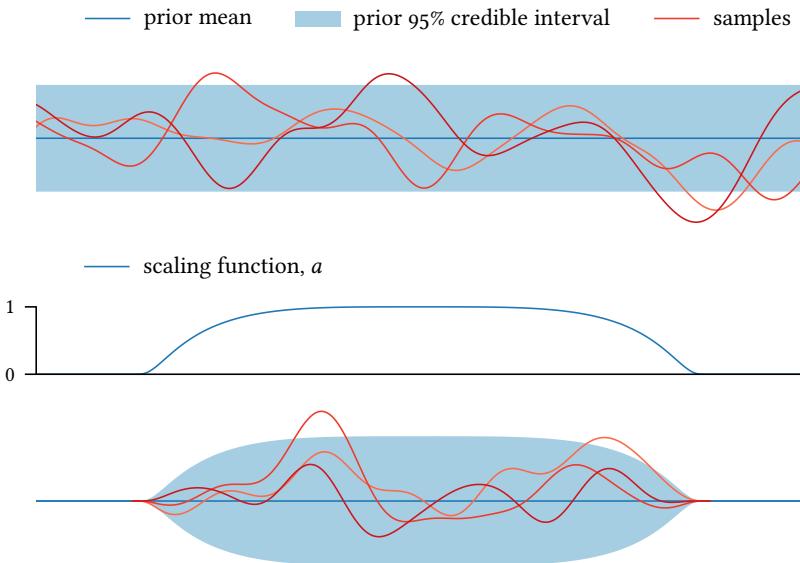


Figure 3.5: Scaling a stationary covariance by a nonconstant function (here, a smooth bump function of compact support) to yield a nonstationary covariance.

and this correlation continues to drop effectively to zero at a separation of approximately five units. Again, this choice of scale is arbitrary, and the assumption of isotropy is particularly restrictive.

A Gaussian process encodes strong assumptions regarding the joint distribution of function values (2.5), which may not be compatible with a given function “out of the box.” However, we can often improve model fit by appropriate transformations of the objective. In fact, *linear* transformations of function inputs and outputs are almost universally considered, although only implicitly by introducing parameters conveying the effects of these transformations. We will show how both linear and nonlinear transformations of function input and output give rise to common model parameters and lead to more expressive classes of models.

Scaling function outputs

We first address the issue of scale in function output (3.17) by considering the statistical effects of arbitrary scaling. Consider a random function $f: \mathcal{X} \rightarrow \mathbb{R}$ with covariance function K and let $a: \mathcal{X} \rightarrow \mathbb{R}$ be a known scaling function.²⁰ Then the pointwise product $af: x \mapsto a(x)f(x)$ has covariance function

$$\text{cov}[af | a] = a(x)K(x, x')a(x'), \quad (3.19)$$

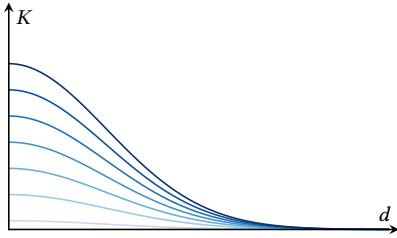
by the bilinearity of covariance. If the scaling function is *constant*, $a \equiv \lambda$, then we have

$$\text{cov}[\lambda f | \lambda] = \lambda^2 K. \quad (3.20)$$

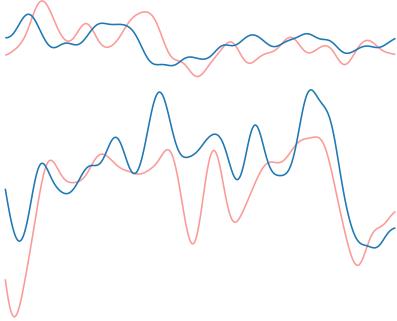
This simple result allows us to extend a “base” covariance K with fixed scale, as in (3.17), to a parametric family with arbitrary scale:

$$K'(x, x'; \lambda) = \lambda^2 K(x, x').$$

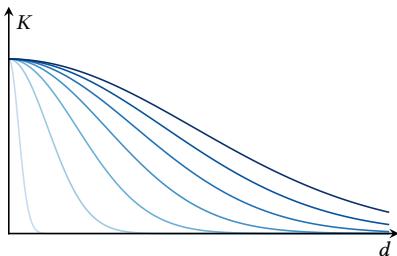
²⁰ For this result f need not have a GP distribution.



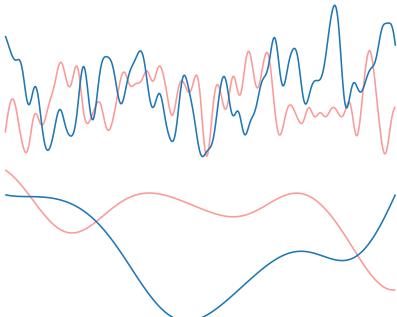
The squared exponential covariance K_{SE} scaled by a range of output scales λ (3.20).



Sample paths from centered GPS with smaller (top) and larger (bottom) output scales.



The squared exponential covariance K_{SE} dilated by a range of length scales ℓ (3.22).



Sample paths from centered GPS with shorter (top) and longer (bottom) characteristic length scales.

In this context the parameter λ is known as an *output scale*, or when the base covariance is stationary with $K(x, x) = 1$, the *signal variance*, as it determines the variance of any function value: $\text{var}[\phi | x, \lambda] = \lambda^2$. The illustration in the margin shows the effect of scaling the squared exponential covariance function by a series of increasing output scales.

We can also of course consider *nonlinear* transformations of the function output as well. This can be useful for modeling constraints – such as nonnegativity or boundedness – that are not compatible with the Gaussian assumption. However, a nonlinear transformation of a Gaussian process is no longer Gaussian, so it often more convenient to model the transformed function after “removing the constraint.”

We may use the general form of this scaling result (3.19) to transform a stationary covariance into a nonstationary one, as any nonconstant scaling is sufficient to break translation invariance. We show an example of such a transformation in figure 3.5, where we have scaled a stationary covariance by a bump function to create a prior on smooth functions with compact support.

Transforming the domain and length scale parameters

We now address the issue of the scaling of correlation as a function of distance (3.18) by introducing a powerful tool: transforming the domain of the function of interest into a more convenient space for modeling.

Namely, suppose we wish to reason about a function $f: \mathcal{X} \rightarrow \mathbb{R}$, and let $g: \mathcal{X} \rightarrow \mathcal{Z}$ be a map from the domain to some arbitrary space \mathcal{Z} , which might also be \mathcal{X} . If $K_{\mathcal{Z}}$ is a covariance function on \mathcal{Z} , then the composition

$$K_{\mathcal{X}}(x, x') = K_{\mathcal{Z}}(g(x), g(x')) \quad (3.21)$$

is trivially a covariance function on \mathcal{X} . This allows us to define a covariance for f indirectly by jointly designing a map g to another space and a corresponding covariance $K_{\mathcal{Z}}$ (and mean $\mu_{\mathcal{Z}}$) on that space. This approach offers a lot of flexibility, as we are free to design these components as we see fit to impose any desired structure.

We will spend some time exploring this idea, beginning with the relatively simple but immensely useful case of combining a *linear* transformation on a Euclidean domain $\mathcal{X} \subset \mathbb{R}^n$ with an *isotropic* covariance on the output. Perhaps the simplest example is the dilation $x \mapsto x/\ell$, which simply scales distance by ℓ^{-1} . Incorporating this transformation into an isotropic base covariance $K(d)$ on \mathcal{X} yields a parametric family of dilated versions:

$$K'(x, x'; \ell) = K(d/\ell). \quad (3.22)$$

If the base covariance has a characteristic length scale of one unit, the length scale of the dilated version will be ℓ ; for this reason, this parameter is simply called *the characteristic length scale* of the parameterized covariance (3.22). Adjusting the length scale allows us to model functions with a range of “wiggliness,” where shorter length scale implies more wiggly behavior; see the margin for examples.

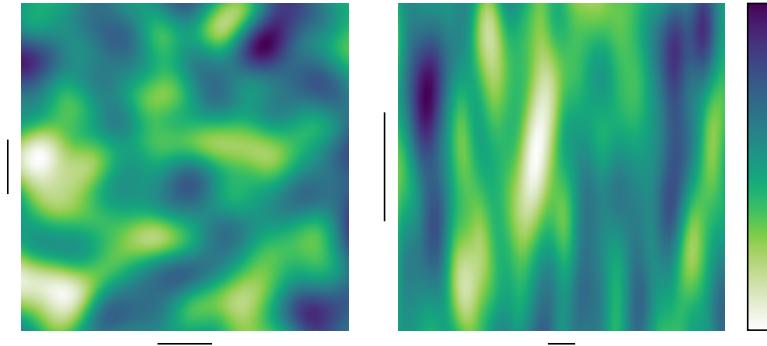


Figure 3.6: Left: a sample from a centered Gaussian process in two dimensions with isotropic squared exponential covariance. Right: a sample from a centered Gaussian process with an ARD squared exponential covariance. The length of the lines on each axis are proportional to the length scale along that axis.

Taking this one step further, we may consider dilating each axis by a separate factor:

$$x_i \mapsto x_i/\ell_i; \quad \mathbf{x} \mapsto [\text{diag } \boldsymbol{\ell}]^{-1}\mathbf{x}, \quad (3.23)$$

which induces the weighted Euclidean distance

$$d_{\boldsymbol{\ell}} = \sqrt{\sum_i \frac{(x_i - x'_i)^2}{\ell_i}}. \quad (3.24)$$

Geometrically, the effect of this map is to transform surfaces of equal distance around each point – which represent curves of constant covariance for an isotropic covariance – from spheres into axis-aligned ellipsoids; see the figure in the margin. Incorporating into an isotropic base covariance $K(d)$ produces a parametric family of *anisotropic* covariances with different characteristic length scales along each axis, corresponding to the parameters $\boldsymbol{\ell}$:

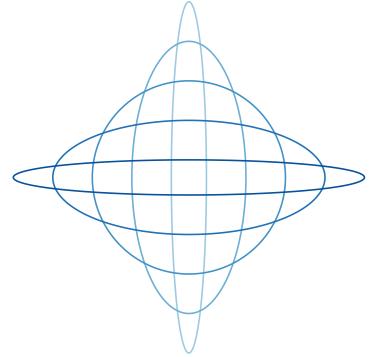
$$K'(\mathbf{x}, \mathbf{x}'; \boldsymbol{\ell}) = K(d_{\boldsymbol{\ell}}). \quad (3.25)$$

When the length scale parameters are inferred from data, this construction is known as *automatic relevance determination* (ARD). The motivation for the name is that if the function has only weak dependence on some mostly irrelevant dimension of the input, we could hope to infer a very long length scale for that dimension. The contribution to the weighted distance (3.24) for that dimension would then be effectively nullified, and the resulting covariance would effectively “ignore” that dimension.

Figure 3.6 shows samples from $2d$ centered Gaussian processes, comparing behavior with an isotropic covariance and an ARD modified version that contracts the horizontal and expands the vertical axis (see curves of constant covariance in the margin). The result is anisotropic behavior with a longer characteristic length scale in the vertical direction than in the horizontal direction, but with the behavior of local features remaining aligned with the axes overall.

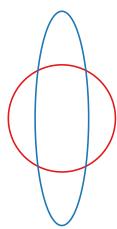
Finally, we may also consider an arbitrary linear transformation $g: \mathbf{x} \mapsto \mathbf{Ax}$, which induces the Mahalanobis distance (A.8)

$$d_{\mathbf{A}} = |\mathbf{Ax} - \mathbf{Ax}'|.$$

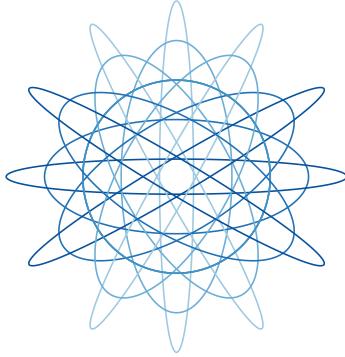


Possible surfaces of equal covariance with the center when combining separate dilation of each axis with an isotropic covariance.

automatic relevance determination, ARD



Surfaces of equal covariance with the center for the examples in figure 3.6: the isotropic covariance in the left panel (red), and the ARD covariance in the right panel (blue).



Possible surfaces of equal covariance with the center when combining an arbitrary linear transformation with an isotropic covariance.

high-dimensional domains: § 3.5, p. 61

As before, we may incorporate this map into an isotropic base covariance K to realize a family of anisotropic covariance functions:

$$K'(x, x'; A) = K(d_A). \quad (3.26)$$

Geometrically, an arbitrary linear map can transform surfaces of constant covariance from spheres into arbitrary ellipsoids; see the figure in the margin. The sample from the left-hand side of figure 3.3 was generated by composing an isotropic covariance with a map inducing both anisotropic scaling and rotation. The effect of the underlying transformation can be seen in the shapes of local features, which are not aligned with the axes.

Due to the inherent number of parameters required to specify a general transformation, this construction is perhaps most useful when the map is to a much lower-dimensional space: $\mathbb{R}^n \rightarrow \mathbb{R}^k$, $k \ll n$. This has been promoted as one strategy for modeling functions on high-dimensional domains suspected of having hidden low-dimensional structure – if this low-dimensional structure is along a linear subspace of the domain, we could capture it by an appropriately designed projection A . We will discuss this idea further in the next section.

Nonlinear warping

When using a covariance function with an inherent length scale, such as a Matérn or squared exponential covariance, *some* linear transformation of the domain is almost always considered, whether it be simple dilation (3.22), anisotropic scaling (3.25), or a general transformation (3.26). However, *nonlinear* transformations can also be useful for imposing structure on the domain, a process commonly referred to as *warping*.

To provide an example that may not often be useful in optimization but is illustrative nonetheless, suppose we wish to model a function $f: \mathbb{R} \rightarrow \mathbb{R}$ that we believe to be smooth and *periodic* with period p . None of the covariance functions introduced thus far would be able to induce the periodic correlations that this assumption would entail.²¹ A construction due to MACKAY is to compose a map onto a circle of radius $r = p/(2\pi)$:²²

$$x \mapsto \begin{bmatrix} r \cos x \\ r \sin x \end{bmatrix} \quad (3.27)$$

with a covariance function on that space reflecting any desired properties of f .²³ As this map identifies points separated by any multiple of the period, the corresponding function values are perfectly correlated, as desired. A sample from a Gaussian process employing this construction with a Matérn covariance after warping is shown in the margin.

A compelling use of warping is to build nonstationary models by composing a nonlinear map with a stationary covariance, an idea SNOEK et al. explored in the context of Bayesian optimization.²⁴ Many objective functions exhibit different behavior depending on the proximity to the optimum, suggesting that nonstationary models may sometimes be worth exploring. SNOEK et al. proposed a flexible family of warping functions for optimization problems with box-bounded constraints, where

²¹ We did see a periodic GP in the previous chapter (2.30); however, that model only had support on *perfectly* sinusoidal functions.

²² D. J. C. MACKAY (1998). Introduction to Gaussian Processes. *Neural Networks and Machine Learning*. Vol. 168. Springer-Verlag. [§ 5.2]

²³ The covariance on the circle is usually inherited from a covariance on \mathbb{R}^2 . The result of composing with the squared exponential covariance in particular is often called “the” periodic covariance, but we stress that any other covariance on \mathbb{R}^2 could be used instead.



A sample path of a centered GP with Matérn covariance with $v = 5/2$ (3.14) after applying the periodic warping function (3.27).

²⁴ J. SNOEK et al. (2014). Input Warping for Bayesian Optimization of Non-Stationary Functions. *ICML 2014*.

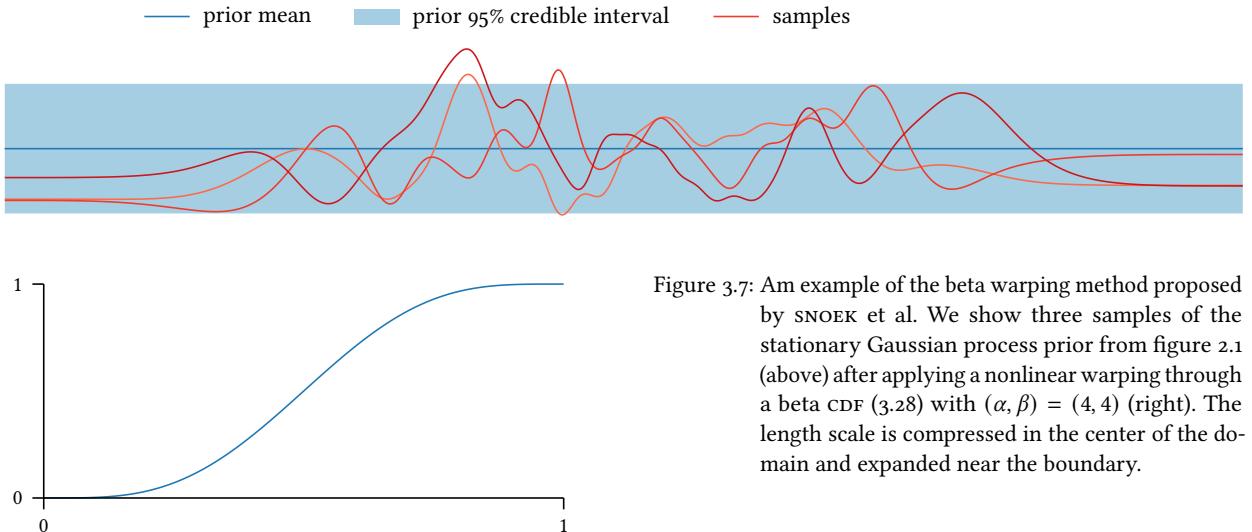


Figure 3.7: An example of the beta warping method proposed by SNOEK et al. We show three samples of the stationary Gaussian process prior from figure 2.1 (above) after applying a nonlinear warping through a beta CDF (3.28) with $(\alpha, \beta) = (4, 4)$ (right). The length scale is compressed in the center of the domain and expanded near the boundary.

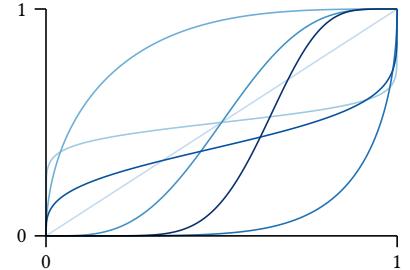
we may take the domain to be the unit cube by scaling and translating as necessary: $\mathcal{X} = [0, 1]^n$. The idea is to warp each coordinate of the input via the cumulative distribution function of a beta distribution:

$$x_i \mapsto I(x_i; \alpha_i, \beta_i), \quad (3.28)$$

where (α_i, β_i) are shape parameters and I is the regularized beta function. This represents a monotonic bijection on the unit interval that can assume several shapes; see the marginal figure for examples. The map may contract portions of the domain and expand others, effectively decreasing and increasing the length scale in those regions. Finally, taking $\alpha = \beta = 1$ recovers the identity map, allowing us to degrade gracefully to the unwarped case if desired.

In figure 3.7 we combine a beta warping on a one-dimensional domain with a stationary covariance on the output. The chosen warping shortens the length scale near the center of the domain and extends it near the boundary, which might be reasonable for an objective function expected to exhibit the most “interesting” behavior on the interior of its domain.

A recent innovation is to use sophisticated artificial neural networks as warping maps for modeling functions of high-dimensional data with complex structure. Notable examples of this approach include the families of *manifold Gaussian processes* introduced by CALANDRA et al.²⁵ and *deep kernels* introduced contemporaneously by WILSON et al.²⁶ Here the warping function was taken to be an arbitrary neural network, the output layer of which was fed into a suitable stationary covariance function. This gives a highly parameterized covariance function where the parameters of the base covariance and the neural map become parameters of the resulting model. In the context of Bayesian optimization, this can be especially useful when there is sufficient data to learn a useful representation of the domain via unsupervised methods.



Some examples of beta CDF warping functions (3.28).

²⁵ R. CALANDRA et al. (2016). Manifold Gaussian Processes for Regression. *IJCNN 2016*.

²⁶ A. G. WILSON et al. (2016). Deep Kernel Learning. *AISTATS 2016*.

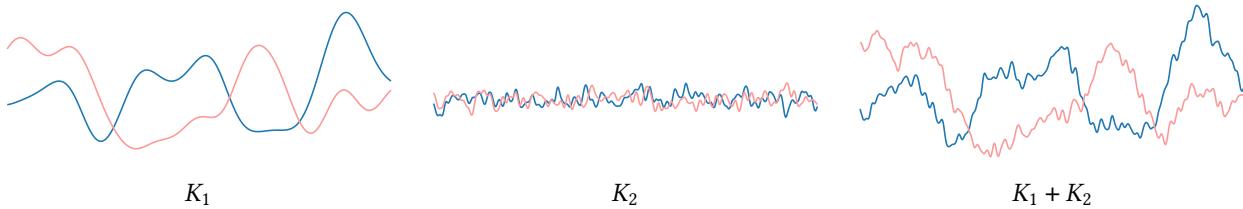


Figure 3.8: Samples from centered Gaussian processes with different covariance functions: (left) a squared exponential covariance, (middle) a squared exponential covariance with smaller output scale and shorter length scale, and (right) the sum of the two. Samples from the process with the sum covariance show smooth variation on two different scales.

Combining covariance functions

In addition to modifying covariance functions via scaling the output and/or transforming the domain, we may also combine multiple covariance functions together to model functions influenced by multiple random processes.

Let $f, g: \mathcal{X} \rightarrow \mathbb{R}$ be two centered, independent (not necessarily Gaussian) random functions with covariance functions K_f and K_g , respectively. By the properties of covariance, the sum and pointwise product of these functions have covariance functions with the same structure:

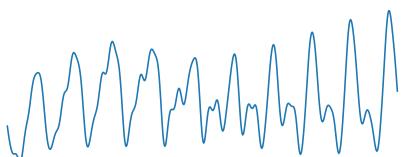
$$\text{cov}[f + g] = K_f + K_g; \quad \text{cov}[fg] = K_f K_g, \quad (3.29)$$

²⁷ The assumption of the processes being centered is needed for the product result only; otherwise, there would be additional terms involving scaled versions of each individual covariance as in (3.19). The sum result does not depend on any assumptions regarding the mean functions.

and thus covariance functions are closed under addition and pointwise multiplication.²⁷ Combining this result with (3.20), we have that *any* polynomial of covariance functions with nonnegative coefficients forms a valid covariance. This enables us to construct infinite families of increasingly complex covariance functions from simple components.

We may use a sum of covariance functions to model a function with independent additive contributions, such as random behavior on several length scales. Precisely such a construction is illustrated in figure 3.8. If the covariance functions are nonnegative and have roughly the same scale, the effect of addition is roughly one of logical disjunction: the sum will assume nontrivial values whenever any one of its constituents does.

Meanwhile, a product of covariance functions can loosely be interpreted in terms of logical conjunction, with function values having appreciable covariance only when every individual covariance function does. A prototypical example of this effect is a covariance function modeling functions that are “almost periodic,” formed by the product of a bump-shaped isotropic covariance function such as a squared exponential (3.12) with a warped version modeling perfectly periodic functions (3.27). The former moderates the influence of the latter by driving the correlation between function values to zero for inputs that are sufficiently separated, regardless of their positions in the periodic cycle. We show a sample from such a covariance in the margin, where the length scale of the modulation term is three times the period.



A sample from a centered Gaussian process with an “almost periodic” covariance function.

3.5 MODELING FUNCTIONS ON HIGH-DIMENSIONAL DOMAINS

Optimization on a high-dimensional domain can be challenging, as we can succumb to the *curse of dimensionality* if we are not careful. As an example, consider optimizing an objective function on the unit cube $[0, 1]^n$. Suppose we model this function with an isotropic covariance from the Matérn family, taking the length scale to be $\ell = 1/10$ so that ten length scales span the domain along each axis.²⁸ This choice implies that function values on the corners of the domain would be effectively independent, as $\exp(-10) < 10^{-4}$ (3.11) and $\exp(-50)$ is smaller still (3.12). If we were to demand even a modicum of confidence in these regions at termination, say by having a measurement within one length scale of every corner, we would need 2^n observations! This exponential growth in the number of observations required to cover the domain is the tyrannical curse of dimensionality.

However, compelling objectives do not tend to have so many degrees of freedom; if they did, we should perhaps give up on the idea of global optimization altogether. Rather, many authors have noted a tendency toward low *intrinsic dimensionality* in real-world problems: that is, most of the variation in the objective is confined to a low-dimensional subspace of the domain. This phenomenon has been noted for example in hyperparameter optimization²⁹ and optimizing the parameters of neural networks.³⁰ LEVINA and BICKEL suggested that “hidden” low-dimensional structure is actually a *universal* requirement for success on any task:³¹

There is a consensus in the high-dimensional data analysis community that the only reason any methods work in very high dimensions is that, in fact, the data are not truly high dimensional.

The global optimization community shares a similar consensus: typical high-dimensional objectives are not “truly” high dimensional. This intuition presents us with an opportunity: if we could only identify inherent low-dimensional structure during optimization, we could sidestep the curse of dimensionality by restricting our search accordingly.

Several strategies are available for capturing low intrinsic dimension with Gaussian process models. The general approach closely follows our discussion from the previous section: we identify some appropriate mapping from the high-dimensional domain to a lower-dimensional space, then model the objective function after composing with this embedding (3.21). This is one realization of the general class of manifold Gaussian processes,³² where the sought-after manifold is low dimensional. Adopting this approach then raises the issue of identifying *useful* families of mappings that can suitably reduce dimension while preserving enough structure of the objective to keep optimization feasible.

Neural embeddings

Given the success of deep learning in designing feature representations for complex, high-dimensional objects, *neural embeddings* – as used

curse of dimensionality

28 This is far from excessive: the domain for the marginal sampling examples in this chapter spans 15 length scales and there’s just enough room for interesting behavior to emerge.

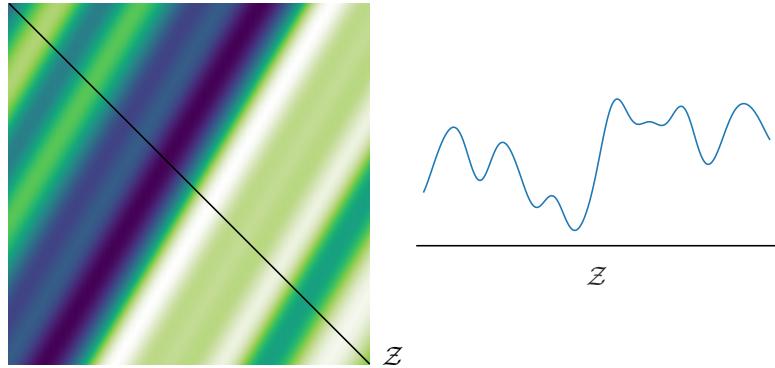
29 J. BERGSTRA and Y. BENGIO (2012). Random Search for Hyper-Parameter Optimization. *Journal of Machine Learning Research* 13:281–305.

30 C. LI et al. (2018a). Measuring the Intrinsic Dimension of Objective Landscapes. *ICLR* 2018.

31 E. LEVINA and P. J. BICKEL (2004). Maximum Likelihood Estimation of Intrinsic Dimension. *NeurIPS* 2004.

32 R. CALANDRA et al. (2016). Manifold Gaussian Processes for Regression. *IJCNN* 2016.

Figure 3.9: An objective function on a two-dimensional domain (left) with intrinsic dimension 1. The entire variation of the objective is determined on the one-dimensional linear subspace \mathcal{Z} corresponding to the diagonal black line, which we can model in its inherent dimension (right).



³³ A. G. WILSON et al. (2016). Deep Kernel Learning. *AISTATS 2016*.

³⁴ R. GÓMEZ-BOMBARELLI et al. (2018). Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules. *ACS Central Science* 4(2):268–276.

³⁵ J. SNOEK et al. (2015). Scalable Bayesian Optimization Using Deep Neural Networks. *ICML 2015*.

cost of Gaussian process inference: § 9.1, p. 201

in the family of deep kernels³³ – present a tantalizing option. Neural embeddings have shown some success in Bayesian optimization, where they can facilitate optimization over complex structured objects such as molecules by providing a nice continuous latent space to work in.³⁴

SNOEK et al. demonstrated excellent performance on hyperparameter tuning tasks by interpreting the output layer of a deep neural network as a set of custom nonlinear basis functions for Bayesian linear regression, as in (3.6).³⁵ An advantage of this particular construction is that Gaussian process inference and prediction is accelerated dramatically by adopting the linear covariance (3.6) – the cost of inference scales linearly with the number of observations, rather than cubically as in the general case.

Linear embeddings

Another line of attack is to search for a low-dimensional *linear* subspace of the domain encompassing the relevant variation in inputs and model the function after projection onto that space. For an objective f on a high-dimensional domain $\mathcal{X} \subset \mathbb{R}^n$, we consider models of the form

$$f(\mathbf{x}) = g(\mathbf{A}\mathbf{x}); \quad \mathbf{A} \in \mathbb{R}^{k \times n} \quad (3.30)$$

where $g: \mathbb{R}^k \rightarrow \mathbb{R}$ is a ($k \ll n$)-dimensional surrogate for f .

The simplest such approach is automatic relevance determination (3.25), where we learn separate length scales along each dimension.³⁶ Although the corresponding linear transformation (3.23) does not reduce dimension, axes with sufficiently long length scales are *effectively* eliminated, as they do not have strong influence on the covariance. This can be effective when some dimensions are likely to be irrelevant, but limits us to axis-aligned subspaces only.

A more flexible option is to consider *arbitrary* linear transformations in the model (3.26, 3.30), an idea that has seen significant attention for Gaussian process modeling in general³⁷ and for Bayesian optimization in particular.³⁸ Figure 3.9 illustrates a simple example where a one-dimensional objective function is embedded in two dimensions in a non-axis-aligned manner. Both axes would appear important for explaining the function when using ARD, but a one-dimensional subspace suffices

³⁶ J. BERGSTRA and Y. BENGIO (2012). Random Search for Hyper-Parameter Optimization. *Journal of Machine Learning Research* 13:281–305.

³⁷ F. VIVARELLI and C. K. I. WILLIAMS (1998). Discovering hidden features with Gaussian process regression. *NeurIPS 1998*.

³⁸ Z. WANG et al. (2016b). Bayesian Optimization in a Billion Dimensions via Random Embeddings. *Journal of Artificial Intelligence Research* 55:361–387

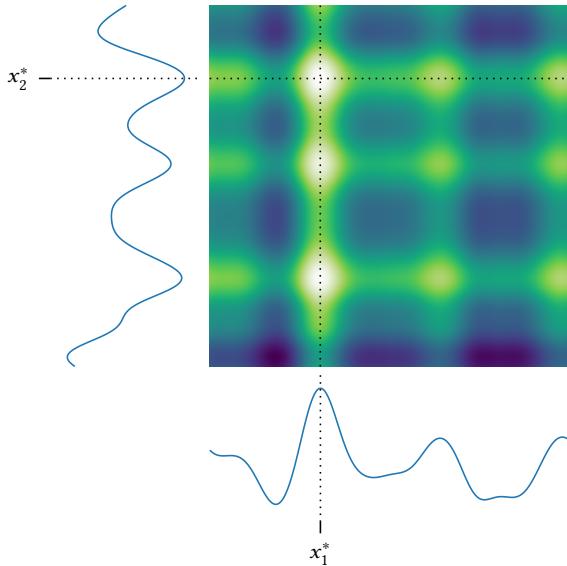


Figure 3.10: A sample from a GP in two dimensions with the decomposition $f(\mathbf{x}) = g_1(x_1) + g_2(x_2)$. Here a nominally two-dimensional function is actually the sum of two one-dimensional components defined along each axis with no interaction. The maximum of the function is achieved at the point corresponding to the maxima of the individual components.

if chosen carefully. This approach offers considerably more modeling flexibility than ARD at the expense of a k -fold increase in the number of parameters that must be specified. However, several algorithms have been proposed for efficiently identifying a suitable map \mathbf{A} ,^{39,40} and WANG et al. demonstrated success in optimizing objectives in extremely high dimension by simply searching along a *random* low-dimensional subspace. The authors also provided theoretical guarantees regarding the recoverability of the global optimum with this approach, assuming the hypothesis of low intrinsic dimensionality holds.

If more flexibility is desired, we may represent an objective function as a sum of contributions on multiple relevant linear subspaces:

$$f(\mathbf{x}) = \sum_i g_i(\mathbf{A}_i \mathbf{x}). \quad (3.31)$$

This decomposition is similar in spirit to the classical family of *generalized additive models*,⁴¹ where the linear maps can be arbitrary and of variable dimension. If we assume the additive components in (3.31) are independent, each with Gaussian process prior $\mathcal{GP}(g_i; \mu_i, K_i)$, then the resulting model for f is a Gaussian process with additive moments (3.29):

$$\mu(\mathbf{x}) = \sum_i \mu_i(\mathbf{A}_i \mathbf{x}); \quad K(\mathbf{x}, \mathbf{x}') = \sum_i K_i(\mathbf{A}_i \mathbf{x}, \mathbf{A}_i \mathbf{x}').$$

Several specific schemes have been proposed for building such decompositions. One convenient approach is to partition the coordinates of the input into disjoint groups and add a contribution defined on each subset.^{42,43} Figure 3.10 shows an example, where a two-dimensional objective is the sum of independent axis-aligned components. We might use such a model when every feature of the input is likely to be relevant but only through interaction with a limited number of additional

- 39 J. DJOLONGA et al. (2013). High-Dimensional Gaussian Process Bandits. *NeurIPS 2013*.
- 40 R. GARNETT et al. (2014). Active Learning of Linear Embeddings for Gaussian Processes. *UAI 2014*.

41 T. HASTIE and R. TIBSHIRANI (1986). Generalized Additive Models. *Statistical Science* 1(3): 297–318.

42 K. KANDASAMY et al. (2015). High Dimensional Bayesian Optimisation and Bandits via Additive Models. *ICML 2015*.

43 J. R. GARDNER et al. (2017). Discovering and Exploiting Additive Structure for Bayesian Optimization. *AISTATS 2017*.

- ⁴⁴ P. ROLLAND et al. (2018). High-Dimensional Bayesian Optimization via Additive Models with Overlapping Groups. *AISTATS 2018*.
- ⁴⁵ M. MUTNÝ and A. KRAUSE (2018). Efficient High Dimensional Bayesian Optimization with Additivity and Quadrature Fourier Features. *NeurIPS 2018*.
- ⁴⁶ E. GILBOA et al. (2013). Scaling Multidimensional Gaussian Processes using Projected Additive Approximations. *ICML 2013*.
- ⁴⁷ C.-L. LI et al. (2016). High Dimensional Bayesian Optimization via Restricted Projection Pursuit Models. *AISTATS 2016*.

variables. An advantage of a disjoint partition is that we may reduce optimization of the high-dimensional objective to separate optimization of each of its lower-dimensional components (3.31). Several other additive schemes have been proposed as well, including partitions with overlapping groups^{44,45} and decompositions of the general form (3.31) with arbitrary projection matrices.^{46,47}

SUMMARY OF MAJOR IDEAS

Specifying a Gaussian process entails choosing a mean and covariance function for the function of interest. As we saw in the previous chapter, the structure of these functions has important implications regarding sample path behavior, and as we will see in the next chapter, important implications regarding its ability to explain a given set of data.

In practice, the design of a Gaussian process model is usually data-driven: we establish some space of candidate models to consider, then search this space for the models providing the best explanation of available data. In this chapter we offered some guidance for the construction of models – or parametric spaces of models – as possible explanations of a given system. We will continue the discussion in the next chapter by taking up the question of assessing model quality in light of data. Below we summarize the important ideas arising in the present discussion.

- prior mean function: § 3.1, p. 46
- impact on sample path behavior: figure 3.1, p.46 and surrounding discussion
- impact on extrapolation: figure 3.2, p. 47 and surrounding discussion

- The mean function of a Gaussian process determines the expected value of function values. Although an important concern, the mean function can only affect sample path behavior through pointwise translation, and most interesting properties are determined by the covariance function instead.
- Nonetheless, the mean function has important implications for prediction, namely, in *extrapolation*. When making predictions in locations poorly explained by available data – that is, locations where function value are not strongly correlated with any observation – the prior mean function effectively determines the posterior predictive mean.
- There are no restrictions on the mean function of a Gaussian process, and we are free to use any sensible choice in a given scenario. In practice, the mean function is usually taken to have some relatively simple parametric form, such as a constant (3.1) or a low-order polynomial (3.7). Such choices are both simple and unlikely to cause grossly undesirable extrapolatory behavior.
- When the mean function includes a *linear* combination of basis functions, we may exactly marginalize the coefficients under a multivariate normal prior (3.5). The result is a marginal Gaussian process where uncertainty in the linear terms of the mean is absorbed into the covariance function (3.6). As an important special case, we may marginalize the value of a constant mean (3.3) under a normal prior (3.2).
- The covariance function of a Gaussian process is critical to determining the behavior of its sample paths. To be valid, a covariance function must

prior covariance function: § 3.2, p. 49

be symmetric and positive semidefinite. The latter condition can be difficult to guarantee for arbitrary “similarity measures,” but covariance functions are closed under several natural operations, allowing us to build complex covariance functions from simple building blocks.

- In particular, sums and pointwise products of covariance functions are valid covariance functions, and by extension any polynomial expression of covariance functions with positive coefficients.
- Many common covariance functions are invariant to translation of their inputs, a property known as *stationarity*. An important result known as *BOCHNER’s theorem* provides a useful representation for the space of stationary covariance functions: their Fourier transforms are symmetric, finite measures, and vice versa. This result has important implications for modeling and computation, as the Fourier representation can be much easier to work with than the covariance function itself.
- Numerous useful covariance functions are available “off-the-shelf.” The family of *Matérn covariances* – and its limiting case the *squared exponential covariance* – can model functions with any desired degree of smoothness (3.11–3.14). A notable special case is the Matérn covariance with $\nu = 5/2$ (3.14), which has been promoted as a reasonable default.
- The *spectral mixture covariance* (3.15) appeals to BOCHNER’s theorem to provide a parametric family of covariance functions able to approximate *any* stationary covariance.
- Covariance functions can be modified by arbitrary scaling of function outputs (3.19) and/or arbitrary transformation of function inputs (3.21). This ability allows us to create *parametric* families of covariance functions with tunable behavior.
- Considering arbitrary *constant* scaling of function outputs gives rise to parameters known as *output scales* (3.20).
- Considering arbitrary *dilations* of function inputs gives rise to parameters known as *characteristic length scales* (3.22). Taking the dilation to be anisotropic introduces a characteristic length scale for each input dimension, a construction known as *automatic relevance determination* (ARD). With an ARD covariance, setting a given dimension’s length scale very high effectively “turns off” its influence on the model.
- *Nonlinear* warping of function inputs is also possible. This enables us to easily build custom *nonstationary* covariance functions by combining a nonlinear warping with a stationary base covariance.
- Optimization can be especially challenging in high dimensions due to the curse of dimensionality. However, if an objective function has intrinsic low-dimensional structure, we can avoid some of the challenges by finding a structure-preserving mapping to a lower-dimensional space and modeling the function on the “smaller” space. This idea has repeatedly proven successful, and several general-purpose constructions are available.

sums and products of covariance functions:
§ 3.4, p. 55

stationarity: § 3.2, p. 50

BOCHNER’s theorem: § 3.2, p. 51

the Matérn family and squared exponential covariance: § 3.3, p. 51

spectral mixture covariance: § 3.3, p. 53

scaling function outputs: § 3.4, p. 55
transforming function inputs: § 3.4, p. 56

nonlinear warping: figure 3.7, p. 59 and surrounding discussion

modeling functions on high-dimensional domains: § 3.5, p. 61

4

MODEL ASSESSMENT, SELECTION, AND AVERAGING

The previous chapter offered a glimpse into the flexibility of Gaussian processes, which can evidently model functions with a wide range of behavior. However, a critical question remains: how we can identify *which* models are appropriate in a given situation?

The difficulty of this question is compounded by several factors. To begin, the number of possible choices is staggering. *Any* function can serve as a mean function for a Gaussian process, and we may construct arbitrary complex covariance functions through a variety of mechanisms. Even if we fix the general form of the moment functions, introducing natural parameters such as output and length scales yields an infinite spectrum of possible models.

Further, many systems to be optimized act as “black boxes,” about which we may have little prior knowledge. Before optimization, we may have only a vague notion of which models might be reasonable for a given objective function or how any parameters of these models should be set. We might even be uncertain about aspects of the observation process, such as the nature or precise scale of observation noise. Therefore, we may find ourselves in the unfavorable position of having infinitely many possible models to choose from and no idea how to choose!

Acquiring *data*, however, provides a way out of this conundrum. After obtaining some observations of the system, we may determine which models are the most compatible with the data and thereby establish preferences over possible choices, a process known as *model assessment*. Model assessment is a surprisingly complex and nuanced subject – even if we limit the scope to Bayesian methods – and no method can rightfully be called “the” Bayesian approach.¹ In this chapter we will present one convenient framework for model assessment via Bayesian inference over models, which are evaluated based on their ability to explain observed data and our prior beliefs.

We will begin our presentation by carefully defining the models we will be assessing and discussing how we may build useful spaces of models for consideration. With Gaussian processes, these spaces will most often be built from what we will call *model structures*, comprising a parametric mean function, covariance function, and observation model; in the context of model assessment, the parameters of these model components are known as *hyperparameters*. We will then show how to perform Bayesian inference over the hyperparameters of a model structure from observations, resulting in a *model posterior* enabling model assessment and other tasks. We will later extend this process to multiple model structures and show how we can even *automatically* search for better model structures.

Central to this approach is a fundamental measure of model fit known as the *marginal likelihood* of the data or *model evidence*. Gaussian process models are routinely selected by maximizing this score, which can produce excellent results when sufficient data are available to unambiguously determine the best-fitting models. However, model construction

prior mean function: § 3.1, p. 46

prior covariance function: § 3.2, p. 49

¹ The interested reader can find an overview of this rich subject in:

A. VEHTARI and J. OJANEN (2012). A survey of Bayesian predictive methods for model assessment, selection and comparison. *Statistics Surveys* 6:142–228.

models and model structures: § 4.1, p. 68

Bayesian inference over parametric model spaces: § 4.2, p. 70

multiple model structures: § 4.5, p. 78

automating model structure search: § 4.6, p. 81

marginal likelihood, model evidence: § 4.2, p. 71

model selection via MAP inference: § 4.3, p. 73

model averaging: § 4.4, p. 74

model, $p(y | x)$

model induced by prior process and observation model

model structure

in the context of Bayesian optimization is unusual as the expense of gathering observations relegates us to the realm of *small data*. Effective modeling with small datasets requires careful consideration of model uncertainty: models explaining the data equally well may disagree drastically in their predictions, and committing to a single model may yield biased predictions with poorly calibrated uncertainty – and disappointing optimization performance as a result. *Model averaging* is one solution that has proven effective in Bayesian optimization, where the predictions of multiple models are combined in the interest of robustness.

4.1 MODELS AND MODEL STRUCTURES

In *model assessment*, we seek to evaluate a space of models according to their ability to explain a set of observations $\mathcal{D} = (\mathbf{x}, \mathbf{y})$. Before taking up this problem in earnest, let us establish exactly what we mean by “model” in this context, which is a model for the given *observations*, rather than of a latent function alone as in the previous chapter.

For this discussion we will define a *model* to be a prior probability distribution over the measured values \mathbf{y} that would result from observing at a set of locations \mathbf{x} : $p(\mathbf{y} | \mathbf{x})$. In the overarching approach we have adopted for this book, a model is specified *indirectly* via a prior process on a latent function f and an observation model linking this function to the observed values:

$$[p(f), p(\mathbf{y} | \mathbf{x}, \phi)]. \quad (4.1)$$

Given explicit choices for these components, we may form the desired distribution by marginalizing the latent function values $\phi = f(\mathbf{x})$ through the observation model:

$$p(\mathbf{y} | \mathbf{x}) = \int p(\mathbf{y} | \mathbf{x}, \phi) p(\phi | \mathbf{x}) d\phi. \quad (4.2)$$

All models we will consider below will be of this composite form (4.1), but the assessment framework we will describe will accommodate arbitrary models.

Spaces of candidate models

² Although defining a space of candidate models may seem natural and innocuous, this is actually a major point of contention between different approaches to Bayesian model assessment. If we subscribe to the maxim “all models are wrong,” we might conclude that the *true* model will *never* be contained in any space we define, no matter how expansive. However, some are likely “more wrong” than others, and we can still reasonably establish preferences over the given space.

To proceed, we must establish some space of candidate models we wish to consider as possible explanations of the observed data.² Although this space can in principle be arbitrary, with Gaussian process models it is convenient to consider *parametric* collections of models defined by parametric forms for the observation model and the prior mean and covariance functions of the latent function. We invested significant effort in the last chapter laying the groundwork to enable this approach: a running theme was the introduction of flexible parametric mean and covariance functions that can assume a wide range of different shapes – perfect building blocks for expressive model spaces.

We will call a particular combination of observation model, prior mean function μ , and prior covariance function K a *model structure*.

Corresponding to each model structure is a natural model space formed by exhaustively traversing the joint parameter space:

$$\mathcal{M} = \left\{ [p(f | \theta), p(y | x, \phi, \theta)] \mid \theta \in \Theta \right\}, \quad (4.3)$$

where

$$p(f | \theta) = \mathcal{GP}(f; \mu(x; \theta), K(x, x'; \theta)).$$

We have indexed the space by a vector θ , the entries of which jointly specify any necessary parameters from their joint range Θ . The entries of θ are known as *hyperparameters* of the model structure, as they parameterize the prior distribution for the observations, $p(y | x, \theta)$ (4.2).

In many cases we may be happy with a single suitably flexible model structure for the data, in which case we can proceed with the corresponding space (4.3) as the set of candidate models. We may also consider multiple model structures for the data by taking a discrete union of such spaces, an idea we will return to later in this chapter.

Example

Let us momentarily take a step back from abstraction and create an explicit model space for optimization on the interval $\mathcal{X} = [a, b]$.³ Suppose our initial beliefs are that the objective will exhibit stationary behavior with a constant trend near zero, and that our observations will be corrupted by additive noise with unknown signal-to-noise ratio.

For the observation model, we take homoskedastic additive Gaussian noise, a reasonable choice when there is no obvious alternative:

$$p(y | \phi, \sigma_n) = \mathcal{N}(y; \phi, \sigma_n^2), \quad (4.4)$$

and leave the scale of the observation noise σ_n as a parameter. Turning to the prior process, we assume a constant mean function (3.1) with a zero-mean normal prior on the unknown constant:

$$\mu(x; c) \equiv c; \quad p(c) = \mathcal{N}(c; 0, b^2),$$

and select the Matérn covariance function with $\nu = 5/2$ (3.14) with unknown output scale λ (3.20) and unknown length scale ℓ (3.22):

$$K(x, x'; \lambda, \ell) = \lambda^2 K_{M^{5/2}}(d/\ell).$$

Following our discussion in the last chapter, we may eliminate one of the parameters above by marginalize the unknown constant mean under its assumed prior,⁴ leaving us with the identically zero mean function and an additive contribution to the covariance function (3.3):

$$\mu(x) \equiv 0; \quad K(x, x'; \lambda, \ell) = b^2 + \lambda^2 K_{M^{5/2}}(d/\ell). \quad (4.5)$$

This, combined with (4.4), completes the specification of a model structure with three hyperparameters: $\theta = [\sigma_n, \lambda, \ell]^\top$. Figure 4.1 illustrates

model space, \mathcal{M}

vector of hyperparameters, θ
range of hyperparameter values, Θ

multiple model structures: § 4.5, p. 78

³ The interval can be arbitrary; our discussion will be purely qualitative.

observation model: additive Gaussian noise with unknown scale

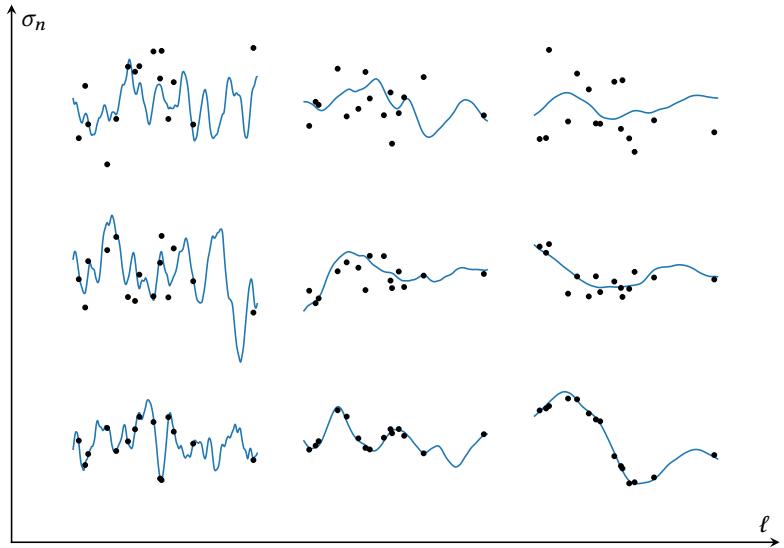
prior mean function: constant mean with unknown value

prior covariance function: Matérn $\nu = 5/2$ with unknown output and length scales

eliminating mean parameter via
marginalization: § 3.1, p. 47

⁴ We would ideally marginalize the other parameters as well, but it would not result in a Gaussian process, as we will discuss shortly.

Figure 4.1: Samples from our example model space for a range of the hyperparameters: σ_n , the observation noise scale, and ℓ , the characteristic length scale. The output scale λ is fixed for each example. Each example demonstrates a sample of the latent function and observations resulting from measurements at a fixed set of 15 locations \mathbf{x} . Elements of the model space can model functions with short- or long-scale correlations that are observed with a range of fidelity from virtually exact observation to extreme noise.



samples from the joint prior over the objective function and the observed values \mathbf{y} that would result from measurements at 15 locations \mathbf{x} (4.2) for a range of these hyperparameters. Even this simple model space is quite flexible, offering degrees of freedom for the variation in the objective function and the precision of our measurements.

4.2 BAYESIAN INFERENCE OVER PARAMETRIC MODEL SPACES

Given a space of candidate models, we now turn to the question of assessing the quality of these models in light of data. There are multiple paths forward,⁵ but Bayesian inference offers one effective solution. By accepting that we can never be absolutely certain regarding which model is the most faithful representation of a given system, we can – as with anything unknown in the Bayesian approach – treat that “best model” as a random variable to be inferred from data and prior beliefs.

We will limit this initial discussion to parametric model spaces built from a single model structure (4.1), which will simplify notation and allow us to conflate models and their corresponding hyperparameters $\boldsymbol{\theta}$ as convenient. We will consider more complex spaces comprising multiple alternative model structures presently.

Model prior

We first endow the model space with a prior encoding which models are more plausible a priori, $p(\boldsymbol{\theta})$.⁶ For convenience, it is common to design the model hyperparameters such that the uninformative (and possibly improper) “uniform prior”

$$p(\boldsymbol{\theta}) \propto 1 \quad (4.6)$$

⁵ A. VEHTARI and J. OJANEN (2012). A survey of Bayesian predictive methods for model assessment, selection and comparison. *Statistics Surveys* 6:142–228.

⁶ As it is most likely that *no* model among the candidates actually *generated* the data, some authors have suggested that any choice of prior is dubious. If this bothers the reader, it can help to frame the inference as being over the model “closest to the truth” rather than over the “true model” itself.

model prior, $p(\boldsymbol{\theta})$



Figure 4.2: The dataset for our model assessment example, generated using a hidden model from the space on the facing page.

may be used, in which case the model prior may not be explicitly acknowledged at all. However, it can be helpful to express at least weakly informative prior beliefs – especially when working with small datasets – as it can offer gentle regularization away from patently absurd choices. This should be possible for most hyperparameters in practice. For example, when modeling a physical system, it would be unlikely that interaction length scales of say one nanometer and one kilometer would be equally plausible a priori; we might capture this intuition with a wide prior on the logarithm of the length scale.

Model posterior

Given a set of observations $\mathcal{D} = (\mathbf{x}, \mathbf{y})$, we may appeal to Bayes' theorem to derive the posterior distribution over the candidate models:

$$p(\theta | \mathcal{D}) \propto p(\theta) p(\mathbf{y} | \mathbf{x}, \theta). \quad (4.7)$$

The model posterior provides support to the models most consistent with our prior beliefs and the observed data. Consistency with the data is encapsulated by the $p(\mathbf{y} | \mathbf{x}, \theta)$ term, the prior PDF over observations evaluated on the actual data.⁷ This value is known as the *model evidence* or the *marginal likelihood* of the data, as it serves as a likelihood in Bayes' theorem (4.7) and, in our class of latent function models, is computed by marginalizing the latent function values at the observed locations (4.2).

model posterior, $p(\theta | \mathcal{D})$

⁷ Recall that this distribution is precisely what a model defines: § 4.1, p. 68.

model evidence, marginal likelihood,
 $p(\mathbf{y} | \mathbf{x}, \theta)$

Marginal likelihood and Bayesian Occam's razor

Model assessment becomes trivial in light of the model posterior if we simply establish preferences over models according to their posterior probability. When using the uniform model prior (4.6) (perhaps implicitly), the model posterior is proportional to the marginal likelihood alone, which can be then used directly for model assessment.

It is commonly argued that the model evidence encodes automatic penalization for model complexity, a phenomenon known as *Bayesian Occam's razor*.⁸ MACKAY outlines a simple argument for this effect by noting that a model $p(\mathbf{y} | \mathbf{x})$ must integrate to unity over all possible measurements \mathbf{y} . Thus if a “simpler” model wishes to become more “complex” by putting support over a wider range of possible observations, it can only do so by reducing the support for the datasets that are already well explained; see the illustration in the margin.

The marginal likelihood of a given dataset can be conveniently computed in closed form for Gaussian process models with additive Gaussian

⁸ D. J. C. MACKAY (2003). *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press. [chapter 28]

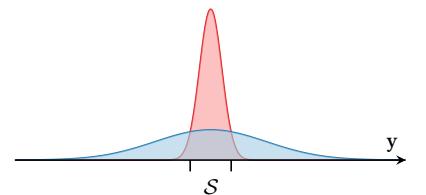
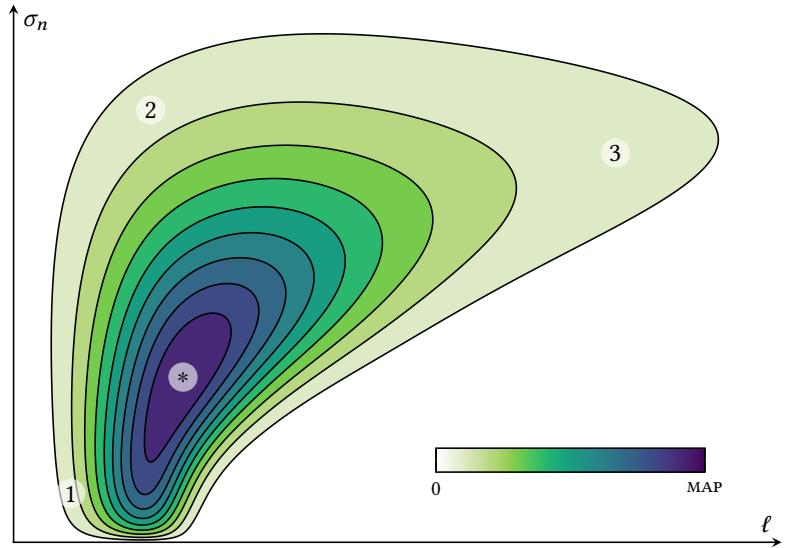


Figure 4.3: The posterior distribution over the model space from figure 4.1 (the range of the axes are compatible with that figure) conditioned on the dataset in figure 4.2. The output scale is fixed (to its true value) for the purposes of illustration. Significant uncertainty remains in the exact values of the hyperparameters, but the model posterior favors models featuring either short length scales with low noise or long length scales with high noise. The points marked 1–3 are referenced in figure 4.4; the point marked * is the MAP (figure 4.5).



noise or exact observation. In this case, we have (2.18):

$$p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}, \Sigma + \mathbf{N}),$$

where $\boldsymbol{\mu}$ and Σ are the prior mean and covariance of the latent objective function values $\boldsymbol{\phi}$ (2.3), and \mathbf{N} is the observation noise covariance matrix (the zero matrix for exact observation) – all of which may depend on $\boldsymbol{\theta}$. As this value can be exceptionally small and have high dynamic range, the logarithm of the marginal likelihood is usually preferred for computational purposes (A.6–A.7):

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) = & \\ -\frac{1}{2} & [(\mathbf{y} - \boldsymbol{\mu})^\top (\Sigma + \mathbf{N})^{-1} (\mathbf{y} - \boldsymbol{\mu}) + \log |\Sigma + \mathbf{N}| + n \log 2\pi]. \end{aligned} \quad (4.8)$$

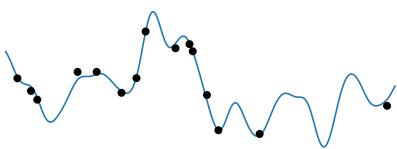
The first term of this expression is the sum of the squared Mahalanobis norms (A.8) of the observations under the prior and represents a measure of data fit. The second term serves as a complexity penalty: the volume of any confidence ellipsoid under the prior is proportional to $|\Sigma + \mathbf{N}|$, and thus this term scales according to the volume of the model’s support in observation space. The third term simply ensures normalization.

Return to example

Let us return to our example scenario and model space. We invite the reader to consider the hypothetical set of 15 observations in figure 4.2 from our example system of interest and contemplate which models from our space of candidates in figure 4.1 might be the most compatible with these observations.⁹

We illustrate the model posterior given this data in figure 4.3, where, in the interest of visualization, we have fixed the covariance output

⁹ The dataset was realized using a moderate length scale (30 length scales spanning the domain) and a small amount of additive noise, shown below. But this is *impossible* to know from inspection of the data alone, and many alternative explanations are just as plausible according to the model posterior!



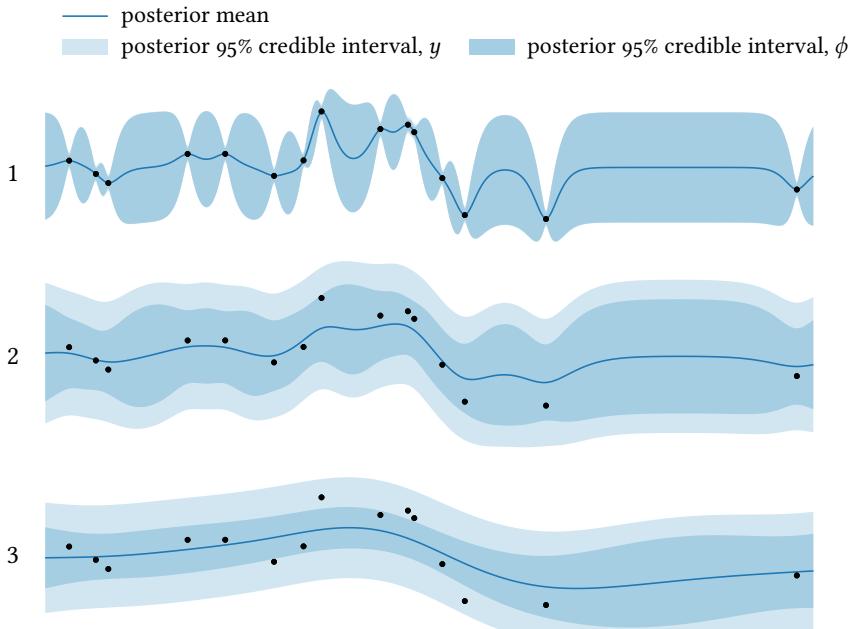


Figure 4.4: Posterior distributions given the observed data corresponding to the three settings of the model hyperparameters marked in figure 4.3. Although remarkably different in their interpretations, each model represents an equally plausible explanation in the model posterior. Model 1 favors near-exact observations with a short length scale, and models 2–3 favor large observation noise with a range of length scales.

scale to its true value and set the range of the axes to be compatible with the samples from figure 4.1. The model prior was designed to be weakly informative regarding the expected order of magnitude of the hyperparameters by taking independent, wide Gaussian priors on the logarithm of the observation noise and covariance length scale.¹⁰

The first observation we can make regarding the model posterior is that it is remarkably *broad*, with many settings of the model hyperparameters remaining plausible after observing the data. However, the model posterior does express a preference for models with either low noise and short length scale or high noise combined with a range of compatible length scales. Figure 4.4 provides examples of objective function and observation posteriors corresponding to the hyperparameters indicated in figure 4.3. Although each is equally plausible in the posterior,¹¹ their explanations of the data are diverse.

4.3 MODEL SELECTION VIA POSTERIOR MAXIMIZATION

Winnowing down a space of candidate models to a *single* model for use in inference and prediction is known as *model selection*. Model selection becomes straightforward if we agree to rank candidates according to the model posterior, as we may then select the maximum a posteriori (MAP) (4.7) model:¹²

$$\hat{\theta} = \arg \max_{\theta} p(\theta) p(y | x, \theta). \quad (4.9)$$

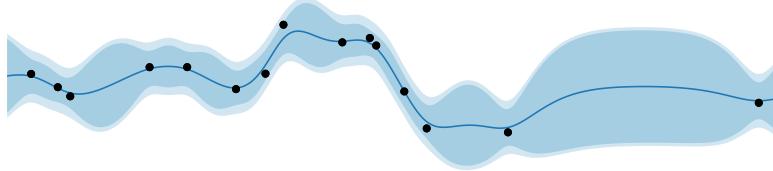
When the model prior is flat (4.6), the MAP model corresponds to the maximum likelihood estimate (MLE) of the model hyperparameters. Fig-

¹⁰ Both parameters are nonnegative, so the prior has support on the entire parameter range.

¹¹ The posterior probability density of these points is approximately 10% of the maximum.

¹² If we only wish to find the maximum, there is no benefit to normalizing the posterior.

Figure 4.5: The predictions of the maximum a posteriori (MAP) model from the example data in figure 4.2.



acceleration via gradient-based optimization

gradient of log marginal likelihood with respect to θ : § C.1, p. 303

model-marginal objective posterior, $p(f | \mathcal{D})$
model-marginal predictive distribution,
 $p(y | x, \mathcal{D})$

ure 4.5 shows the predictions made by the MAP model for our running example; in this case, the MAP hyperparameters are in fact a reasonable match to the parameters used to generate the example dataset.

When the model space is defined over a continuous space of hyperparameters, computation of the MAP model can be significantly accelerated via gradient-based optimization. Here it is advisable to work in the log domain, where the objective becomes the unnormalized log posterior:

$$\log p(\theta) + \log p(y | x, \theta). \quad (4.10)$$

The log marginal likelihood is given in (4.8), noting that μ , Σ , and N are all implicitly functions of the hyperparameters θ . This objective (4.10) is differentiable with respect to θ assuming the Gaussian process prior moments, the noise covariance, and the model prior are as well, in which case we may appeal to off-the-shelf gradient methods for solving (4.9). However, a word of warning is in order: the model posterior is not guaranteed to be concave and may have multiple local maxima, so multistart optimization is prudent.

4.4 MODEL AVERAGING

Reliance on a single model is questionable when the model posterior is not well determined by the data. For example, in our running example, a diverse range of models are consistent with the data (figures 4.3–4.4). Committing to a single model in this case may systematically bias our predictions and underestimate predictive uncertainty – note how the diversity in predictions from figure 4.4 is lost in the MAP model (4.5).

An alternative is to *marginalize* the model with respect to the model posterior, a process known as *model averaging*:

$$p(f | \mathcal{D}) = \int p(f | \mathcal{D}, \theta) p(\theta | \mathcal{D}) d\theta; \quad (4.11)$$

$$p(y | x, \mathcal{D}) = \iint p(y | x, \phi, \theta) p(\phi | x, \mathcal{D}, \theta) p(\theta | \mathcal{D}) d\phi d\theta, \quad (4.12)$$

where we have marginalized the hyperparameters of both the objective and observation models. Model averaging is more consistent with the ideal Bayesian convention of marginalizing nuisance parameters when possible¹³ and promises robustness to model misspecification, at least in the chosen model space.

Unfortunately, neither of these model-marginal distributions (4.11–4.12) can be computed exactly for Gaussian process models except in

¹³ Although it may be unusual to consider the choice of model a “nuisance!”

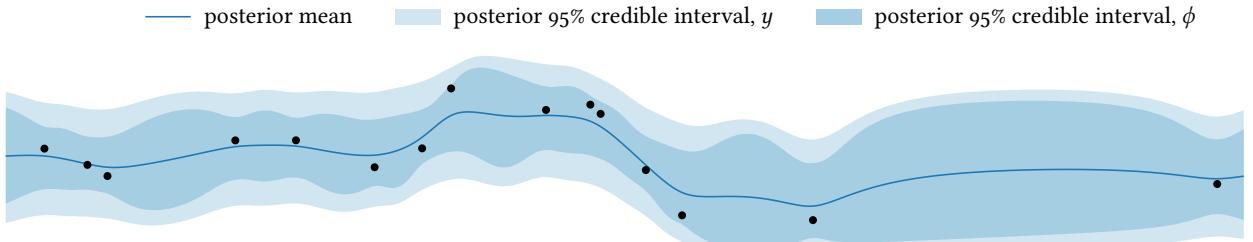


Figure 4.6: A Monte Carlo estimate to the model-marginal predictive distribution (4.11) for our example scenario using 100 samples drawn from the model posterior in figure 4.3 (4.14–4.15); see illustration in margin. Samples from the objective function posterior display a variety of behavior due to being associated with different hyperparameters.

some special cases,¹⁴ so we must resort to approximation if we wish to pursue this approach. In fact, maximum a posteriori estimation can be interpreted as one rather crude approximation scheme where the model posterior is replaced by a Dirac delta distribution at the MAP hyperparameters:

$$p(\boldsymbol{\theta} \mid \mathcal{D}) \approx \delta(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}).$$

This can be defensible when the dataset is large compared to the number of hyperparameters, in which case the model posterior is often unimodal with little residual uncertainty. However, large datasets are the exception rather than the rule in Bayesian optimization, and more sophisticated approximations can pay off when model uncertainty is significant.

¹⁴ A notable example is marginalizing the coefficients of a linear prior mean against a Gaussian prior: § 3.1, p. 47.

Monte Carlo approximation

Monte Carlo approximation is one straightforward path forward. Drawing a set of hyperparameter samples from the model posterior,

$$\{\boldsymbol{\theta}_i\}_{i=1}^s \sim p(\boldsymbol{\theta} \mid \mathcal{D}), \quad (4.13)$$

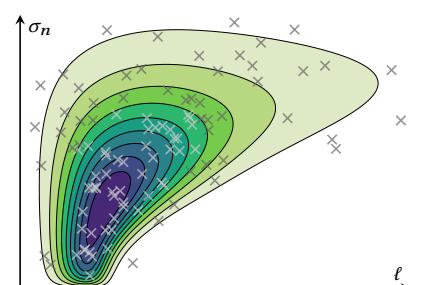
yields the following simple Monte Carlo estimates:

$$p(f \mid \mathcal{D}) \approx \frac{1}{s} \sum_{i=1}^s \mathcal{GP}(f; \mu_{\mathcal{D}}(\boldsymbol{\theta}_i), K_{\mathcal{D}}(\boldsymbol{\theta}_i)); \quad (4.14)$$

$$p(y \mid x, \mathcal{D}) \approx \frac{1}{s} \sum_{i=1}^s \int p(y \mid x, \phi, \boldsymbol{\theta}_i) p(\phi \mid x, \mathcal{D}, \boldsymbol{\theta}_i) d\phi. \quad (4.15)$$

The objective function posterior is approximated by a *mixture* of Gaussian processes corresponding to the sampled hyperparameters, and the posterior predictive distribution for observations is then derived by integrating a Gaussian mixture (2.35) against the observation model.

Any Markov chain Monte Carlo procedure could be used to generate the hyperparameter samples (4.13); a variation on Hamiltonian Monte



The 100 hyperparameter samples used to produce figure 4.6.

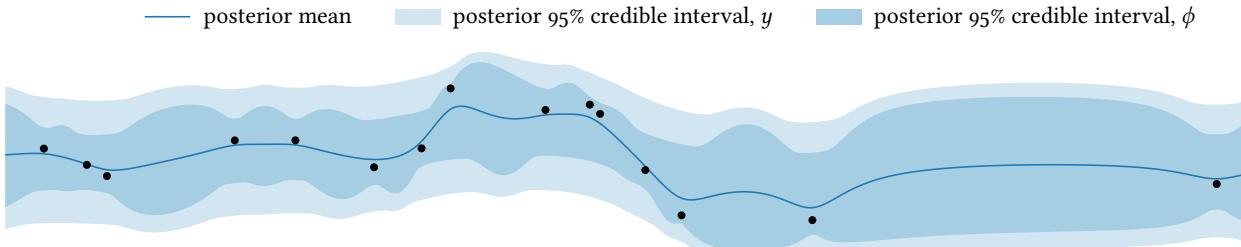


Figure 4.7: An approximation to the model-marginal posterior (4.11) using the central composite design approach proposed by RUE et al. A total of nine hyperparameter samples are used for the approximation, illustrated in the margin below.

¹⁵ M. D. HOFFMAN and A. GELMAN (2014). The No-U-turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research* 15(4):1593–1623.

Carlo (HMC) such as the no u-turn sampler (NUTS) would be a reasonable choice when the gradient of the log posterior (4.10) is available, as it can exploit this information to accelerate mixing.¹⁵

Figure 4.6 demonstrates a Monte Carlo approximation to the model-marginal posterior (4.11–4.12) for our running example. Comparing with the MAP approximation in figure 4.5, the predictive uncertainty of both objective function values and observations has increased considerably due to accounting for model uncertainty in the predictive distributions.

Deterministic approximation schemes

The downside of Monte Carlo approximation is relatively inefficient use of the hyperparameter samples – the price of random sampling rather than careful design. This inefficiency in turn leads to an increased computational burden for inference and prediction from having to derive a GP posterior for each sample. Several more efficient (but less accurate) alternative approximations for hyperparameter marginalization have also been proposed. A common simplifying tactic taken by these cheaper procedures is to approximate the hyperparameter posterior with a multivariate normal via a Laplace approximation:

$$p(\boldsymbol{\theta} | \mathcal{D}) \approx \mathcal{N}(\boldsymbol{\theta}; \hat{\boldsymbol{\theta}}, \mathbf{C}), \quad (4.16)$$

where $\hat{\boldsymbol{\theta}}$ is the MAP (4.9). Integrating this approximation into (4.11) gives

$$p(f | \mathcal{D}) \approx \int \mathcal{GP}(f; \mu_{\mathcal{D}}(\boldsymbol{\theta}), K_{\mathcal{D}}(\boldsymbol{\theta})) \mathcal{N}(\boldsymbol{\theta}; \hat{\boldsymbol{\theta}}, \mathbf{C}) d\boldsymbol{\theta}. \quad (4.17)$$

Unfortunately this integral remains intractable due to the nonlinear dependence of the posterior moments on the hyperparameters, but reducing to this common form allows us to derive *deterministic* approximations against a single assumed posterior.

¹⁶ H. RUE et al. (2009). Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society Series B (Methodological)* 71(2):319–392.

¹⁷ G. E. P. BOX and K. B. WILSON (1951). On the Experimental Attainment of Optimum Conditions. *Journal of the Royal Statistical Society Series B (Methodological)* 13(1):1–45.

RUE et al. introduced several approximation schemes representing different tradeoffs between efficiency and fidelity.¹⁶ Notable among these is a simple, sample-efficient procedure grounded in classical experimental design. Here a central composite design¹⁷ in hyperparameter

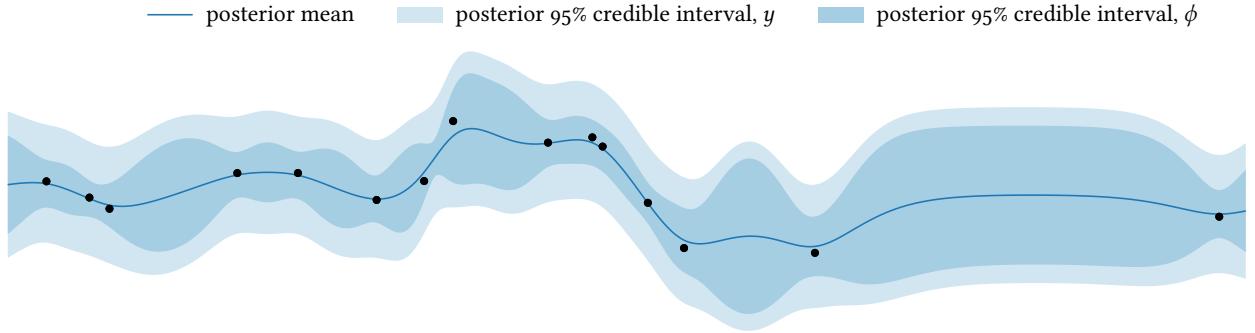


Figure 4.8: The approximation to the model-marginal posterior (4.11) for our running example using the approach proposed by OSBORNE et al.

space is transformed to agree with the moments of (4.16), then used as nodes in a numerical quadrature approximation to (4.17). The resulting approximation again takes the form of a (now weighted) mixture of Gaussian processes (4.14): the MAP model augmented by a small number of additional models designed to reflect the important variation in the hyperparameter posterior. The number of hyperparameter samples required by this scheme grows relatively slowly with the dimension of the hyperparameter space: less than 100 for $|\theta| \leq 8$ and less than 1000 for $|\theta| \leq 21$.¹⁸ The nine samples required for our running example are shown in the marginal figure. Figure 4.7 shows the resulting approximate posterior; comparing with the gold-standard Monte Carlo approximation from figure 4.6, the agreement is excellent.

An even more lightweight approximation was proposed by OSBORNE et al., which despite its crudeness is arguably still preferable to MAP estimation and can be used as a drop-in replacement.¹⁹ This approach again relies on a Laplace approximation to the hyperparameter posterior (4.16–4.17). The key observation is that under the admittedly strong assumption that the posterior mean were in fact *linear* in θ and the posterior covariance *independent* of θ , we could resolve (4.17) in closed form. We proceed by taking the best linear approximation to the posterior mean around the MAP:²⁰

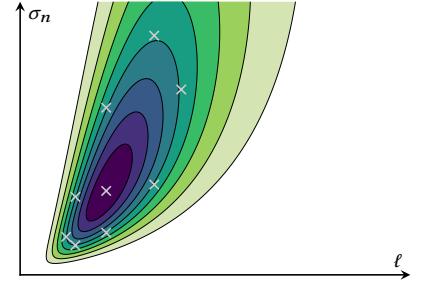
$$\mu_{\mathcal{D}}(x; \theta) \approx \mu_{\mathcal{D}}(x; \hat{\theta}) + g(x)^T(\theta - \hat{\theta}); \quad g(x) = \frac{\partial \mu_{\mathcal{D}}(x; \theta)}{\partial \theta}(\hat{\theta}),$$

and assuming the MAP posterior covariance is universal: $K_{\mathcal{D}}(\theta) \approx K_{\mathcal{D}}(\hat{\theta})$. The result is a *single* Gaussian process approximation to the posterior:

$$p(f | \mathcal{D}) \approx \mathcal{GP}(f; \hat{\mu}_{\mathcal{D}}, \hat{K}_{\mathcal{D}}), \quad (4.18)$$

where

$$\hat{\mu}_{\mathcal{D}}(x) = \mu_{\mathcal{D}}(x; \hat{\theta}); \quad \hat{K}_{\mathcal{D}}(x, x') = K_{\mathcal{D}}(x, x'; \hat{\theta}) + g(x)^T C g(x').$$



A Laplace approximation to the model posterior (performed in the log domain) and hyperparameter settings corresponding to the central composite design proposed by RUE et al. The samples do a good job covering the support of the true posterior.

¹⁸ S. M. SANCHEZ and P. J. SANCHEZ (2005). Very Large Fractional Factorial and Central Composite Designs. *ACM Transactions on Modeling and Computer Simulation* 15(4):362–377.

¹⁹ M. A. OSBORNE et al. (2012). Active Learning of Model Evidence Using Bayesian Quadrature. *NeurIPS 2012*.

²⁰ This is analogous to the linearization step in the extended Kalman filter, whereas the central composite design approach is closer to the unscented Kalman filter in pushing samples through the nonlinear transformation.

uncertainty in additive noise scale σ_n

²¹ This term vanishes if the hyperparameters are completely determined by the data, in which case the approximation regresses gracefully to the MAP estimate.

²² In general we have

$$p(y | x, \mathcal{D}) = \iint p(y | x, \phi, \sigma_n) p(\phi, \sigma_n | x, \mathcal{D}) d\phi d\sigma_n,$$

and we have resolved the integral on ϕ using the single-GP approximation.

This is the MAP model with covariance inflated by a term determined by the dependence of the posterior mean on the hyperparameters, \mathbf{g} , and the uncertainty in the hyperparameters, \mathbf{C} .²¹

OSBORNE et al. did not address how to account for uncertainty in observation model parameters when approximating $p(y | x, \mathcal{D})$, but we can derive a natural approach for independent additive Gaussian noise with unknown scale σ_n . Given x , let $p(\phi | x, \mathcal{D}) \approx \mathcal{N}(\phi; \mu, \sigma^2)$ as in (4.18). We must approximate²²

$$p(y | x, \mathcal{D}) \approx \int \mathcal{N}(y; \mu, \sigma^2 + \sigma_n^2) p(\sigma_n | x, \mathcal{D}) d\sigma_n.$$

A moment-matched approximation $p(y | x, \mathcal{D}) \approx \mathcal{N}(y; m, s^2)$ is possible by appealing to the law of total variance:

$$m = \mathbb{E}[y | x, \mathcal{D}] \approx \mu; \quad s^2 = \text{var}[y | x, \mathcal{D}] \approx \sigma^2 + \mathbb{E}[\sigma_n^2 | x, \mathcal{D}].$$

If the noise scale is parameterized by its logarithm, then the Laplace approximation (4.16) in particular yields

$$p(\log \sigma_n | x, \mathcal{D}) \approx \mathcal{N}(\log \hat{\sigma}_n; \log \hat{\sigma}_n, s^2); \quad \mathbb{E}[\sigma_n^2 | x, \mathcal{D}] \approx \hat{\sigma}_n^2 \exp(2s^2).$$

Thus we predict with the MAP estimate $\hat{\sigma}_n$ inflated by a factor commensurate with the residual uncertainty in the noise contribution.

Figure 4.8 shows the resulting approximation for our running example. Although not perfect, the predictive uncertainty in the observations is more faithful than the MAP model from figure 4.5, which severely underestimates the most likely scale of observation noise.

4.5 MULTIPLE MODEL STRUCTURES

We have now covered model inference, selection, and averaging with a *single* parametric model space (4.1). With a bit of extra bookkeeping, we may extend this framework to handle multiple model structures comprising different combinations of parametric prior moments and observation models.

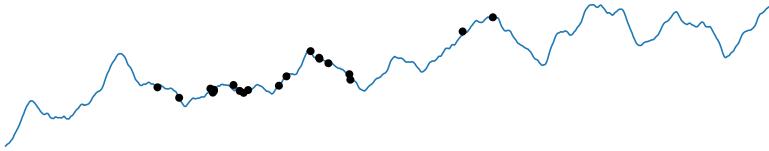
To begin, we may build a space of candidate models by taking a discrete union of parametric spaces as in (4.1), with one built from each desired model structure: $\{\mathcal{M}_i\}$. It is natural to index this space by (θ, \mathcal{M}) , where θ is understood to be a vector of hyperparameters associated with the specified model structure; the size and interpretation of this vector may differ across structures. All that remains is to derive our previous results while managing this compound structure–hyperparameter index.

We may define a model prior over this compound space by combining a prior over the chosen model structures with priors over the hyperparameters of each:

$$p(\theta, \mathcal{M}) = \Pr(\mathcal{M}) p(\theta | \mathcal{M}). \quad (4.19)$$

Given data, the model posterior has a similar form as before (4.7):

$$p(\theta, \mathcal{M} | \mathcal{D}) = \Pr(\mathcal{M} | \mathcal{D}) p(\theta | \mathcal{D}, \mathcal{M}). \quad (4.20)$$



The structure-conditional hyperparameter posterior $p(\theta | \mathcal{D}, \mathcal{M})$ is as in (4.7) and may be reasoned about following our previous discussion. The model structure posterior is then given by

$$\Pr(\mathcal{M} | \mathcal{D}) \propto \Pr(\mathcal{M}) p(y | x, \mathcal{M}); \quad (4.21)$$

$$p(y | x, \mathcal{M}) = \int p(y | x, \theta, \mathcal{M}) p(\theta | \mathcal{M}) d\theta. \quad (4.22)$$

The expression in (4.22) is the normalizing constant of the structure-conditional hyperparameter posterior (4.7), which we could ignore when there was only a single model structure. This integral is in general intractable, but several approximations are feasible. One effective choice is the Laplace approximation (4.16), which provides an approximation to the integral as a side effect (B.2). The classical *Bayesian information criterion* (BIC) may be seen as an approximation to this approximation.²³

Model selection may now be pursued by maximizing the model posterior over the model space as before, although we may no longer appeal to gradient methods as the model space is not continuous with multiple model structures. A simple approach would be to find the MAP hyperparameters for each of the model structures separately, then use these MAP points to approximate (4.22) for each structure via the Laplace approximation or BIC. This would be sufficient to estimate (4.20–4.21) and maximize over the MAP models.

Turning to model averaging, the model-marginal posterior to the objective function is:

$$p(f | \mathcal{D}) = \sum_i \Pr(\mathcal{M}_i | \mathcal{D}) p(f | \mathcal{D}, \mathcal{M}_i). \quad (4.23)$$

The structure-conditional, hyperparameter-marginal distribution on each space $p(f | \mathcal{D}, \mathcal{M})$ is as before (4.11) and may be approximated following our previous discussion. These are now combined in a mixture distribution weighted by the model structure posterior (4.21).

Multiple structure example

We now present an example of model inference, selection, and averaging over multiple model structures using the dataset in figure 4.9.²⁴ The data were sampled from a Gaussian process with linear prior mean (a linear trend with positive slope is evident) and Matérn $\nu = 3/2$ prior covariance (3.13), with a small amount of additive Gaussian noise. We also show a sample from the objective function posterior corresponding to the true model generating the data for reference.

Figure 4.9: The objective and dataset for our multiple-model example.

model structure posterior, $\Pr(\mathcal{M} | \mathcal{D})$

Laplace approximation: § B.1, p. 297

²³ S. KONISHI and G. KITAGAWA (2008). *Information Criteria and Statistical Modeling*. Springer-Verlag. [chapter 9]

model selection

model averaging

²⁴ The data are used as a demo in the code released with:

C. E. RASMUSSEN and C. K. I. WILLIAMS (2006). *Gaussian Processes for Machine Learning*. MIT Press.

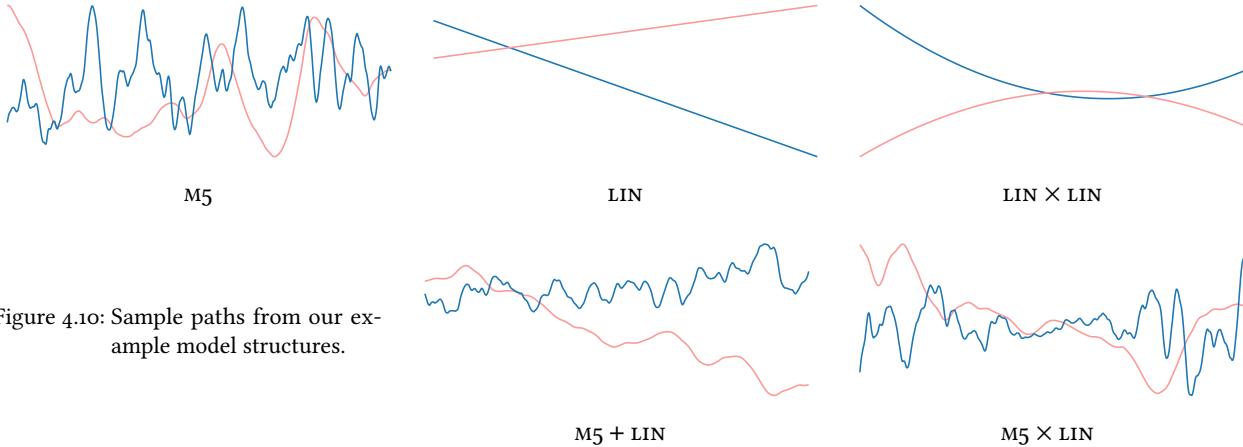


Figure 4.10: Sample paths from our example model structures.

initial model structure: p. 69

We build a model space comprising several model structures by augmenting our previous space with structures incorporating additional covariance functions. The treatment of the prior mean (unknown constant marginalized against a Gaussian prior) and observation model (additive Gaussian noise with unknown scale) will remain the same for all. The model structures reflect a variety of hypotheses positing potential linear or quadratic behavior:

- M5: the Matérn $\nu = 5/2$ covariance (3.14) from our previous example;
- LIN: the linear covariance (3.16), where the prior on the slope is vague and centered at zero and the prior on the intercept agrees with the M5 model;
- LIN × LIN: the product of two linear covariances designed as above, modeling a latent quadratic function with unknown coefficients;
- M5 + LIN: the sum of a Matérn $\nu = 5/2$ and linear covariance designed as in the corresponding individual model structures; and
- M5 × LIN: the product of a Matérn $\nu = 5/2$ and linear covariance designed as in the corresponding individual model structures.

Objective function samples from models in each of these structures are shown in figure 4.10. Among these, the model structure closest to the truth is arguably M5 + LIN.

Following the above discussion, we find the MAP hyperparameters for each of these model structures separately and use a Laplace approximation (4.16) to approximate the hyperparameter posterior on each space, along with the normalizing constant (4.22). Normalizing over the structures provides an approximate model structure posterior:

$$\begin{aligned} \Pr(M5 | \mathcal{D}) &\approx 10.8\%; \\ \Pr(M5 + LIN | \mathcal{D}) &\approx 71.8\%; \\ \Pr(M5 \times LIN | \mathcal{D}) &\approx 17.0\%, \end{aligned}$$

with the remaining model structures (LIN and LIN × LIN) sharing the re-

approximation to model structure posterior

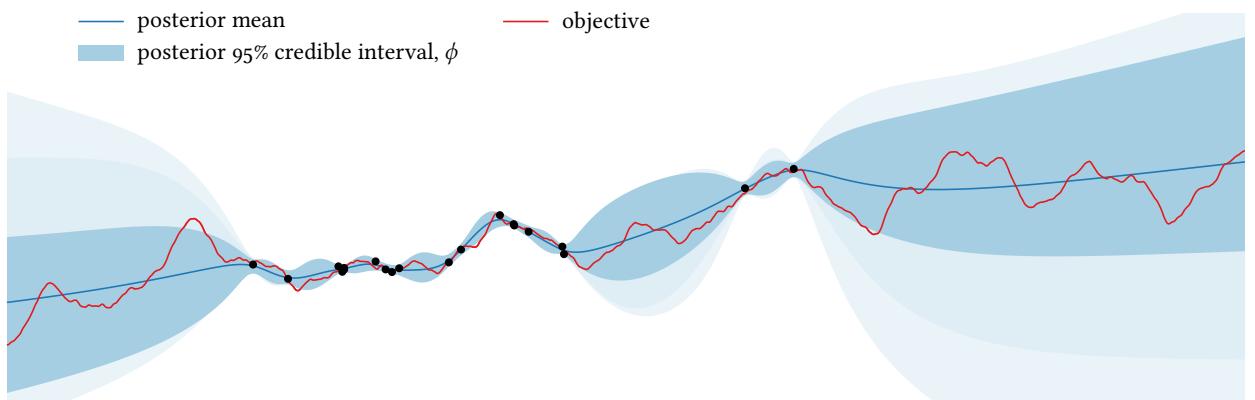


Figure 4.11: An approximation to the model-marginal posterior (4.23) for our multiple-model example. The posterior on each model structure is approximated separately as a mixture of Gaussian processes following RUE et al. (see figure 4.7); these are then combined by weighting by an approximation of the model structure posterior (4.21). We show the result with three superimposed, transparent credible intervals, which are shaded with respect to their weight in contributing to the final approximation.

maining 0.4%. The `M5 + LIN` model structure is the clear winner, and there is strong evidence that the purely polynomial models are insufficient for explaining the data alone.

Figure 4.11 illustrates an approximation to the model-marginal posterior (4.23), approximated by applying the central composite design approach of RUE et al. to each of the model structures separately, then combining these into a large Gaussian process mixture by weighting by the approximate model structure posterior. The highly asymmetric credible intervals reflect the diversity in explanations for the data offered by the chosen model structures, and the combined model makes reasonable predictions of our example objective function sampled from the true model.

For this example, averaging over the model structure has important implications regarding the behavior of the resulting optimization policy. Figure 4.12 illustrates a common acquisition function²⁵ built from the off-the-shelf `M5` model, as well as from the structure-marginal model. The former chooses to exploit near what it believes is a local optimum, but the latter has a strong belief in an underlying linear trend and chooses to explore the right-hand side of the domain instead. For our example objective function sample, this would in fact reveal the global optimum with the next observation.

approximation to marginal predictive distribution

averaging over a space of Gaussian processes in policy computation: § 8.10, p. 192

²⁵ specifically, expected improvement: § 7.3 p. 127

4.6 AUTOMATING MODEL STRUCTURE SEARCH

We now have a comprehensive framework for reasoning about model uncertainty, including methods for model assessment, selection, and averaging across one or multiple model structures. However, it is still not clear how we should determine which model structures to consider for a

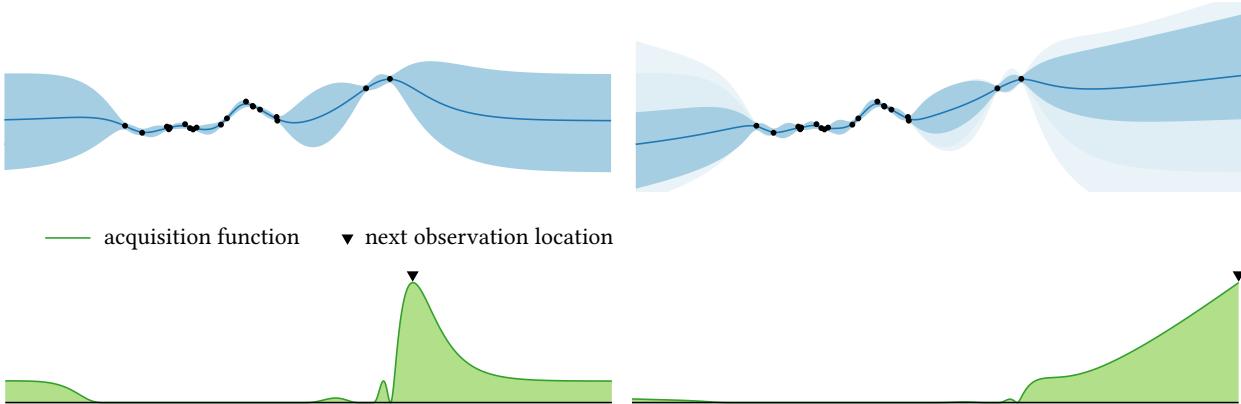


Figure 4.12: Optimization policies built from the MAP M5 model (left) and the structure-marginal posterior (right). The M5 model chooses to exploit near the local optimum, but the structure-marginal model is aware of the underlying linear trend and chooses to explore the right-hand side as a result.

given system. This is critical as our model inference procedure requires the space of candidate models to be predefined; equivalently, the model prior (4.19) is implicitly set to zero for every model outside this space. Ideally, we would simply enumerate every possible model structure and average over all of them, but even a naïve approximation of this ideal would entail overwhelming computational effort.

However, the set of model structures we consider for a given dataset can be *adaptively* tailored as we gather data. One powerful idea is to appeal to metaheuristics such as local search: by establishing a suitable space of candidate model structures, we can dynamically explore this space for the best explanations of available data.

To enable this approach, we must first establish a sufficiently rich space of candidate model structures. DUVENAUD et al. proposed one convenient mechanism for defining such a space via a simple productive grammar.²⁶ The idea is to appeal to the closure of covariance functions under addition and pointwise multiplication to systematically build up families of increasingly complex models from simple components. We begin by choosing a set of so-called *base kernels*, \mathcal{B} , modeling relatively simple behavior, then extend this set to an infinite family of compositions via the following context-free grammar:

$$\begin{aligned} K &\rightarrow B \\ K &\rightarrow K + K \\ K &\rightarrow KK \\ K &\rightarrow (K). \end{aligned}$$

The symbol B in the first rule represents any desired base kernel. The five model structures considered in our multiple-structure example above in fact represent five members of the language generated by this grammar with the base kernels $\mathcal{B} = \{K_{M5/2}, K_{LIN}\}$ or simply $\mathcal{B} = \{M5, LIN\}$. The

²⁶ D. DUVENAUD et al. (2013). Structure Discovery in Nonparametric Regression through Compositional Kernel Search. *ICML* 2013.

addition and multiplication of covariance functions: § 3.4, p. 60
base kernels, \mathcal{B}

grammar however also generates arbitrarily more complicated expressions such as

$$(M_5 + (M_5 + \text{LIN})M_5)(M_5 + M_5) + \text{LIN}. \quad (4.24)$$

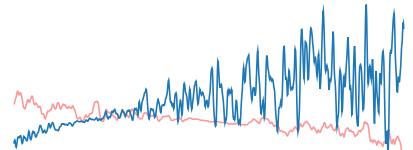
We are free to design the base kernels in this construction to capture any potential atomic behavior in the objective function. For example, if the domain is high dimensional and we suspect that the objective may depend only on interactions of small groups of mostly independent variables, we might design the base kernels to model variations in single variables at a time, then rely on the grammar to generate an array of possible interaction structures.

Other spaces of model structures have also been proposed for automated structure search. With an eye toward high-dimensional domains, GARDNER et al. for example considered spaces of additive model structures indexed by every possible partition of the input variables.²⁷ This is an expressive class of model structures, but the number of partitions grows so rapidly that exhaustive search is not feasible.

Once a space of candidate model structures has been established, we may develop a search procedure seeking the most promising structures to explain a given dataset. Several approaches have been proposed for this search with a range of complexity, all of which frame the problem in terms of optimizing some figure of merit over the space. Although any score could be used in this context, a natural choice is an approximation to the (unnormalized) model structure posterior (4.21) such as the Laplace approximation or the Bayesian information criterion, and every method we will describe uses one of these two scores.

DUVENAUD et al. suggested a greedy search procedure for spaces generated by their grammar, wherein the base kernels are first evaluated, and the best among them is subjected to productive grammar rules generating similar structures to search next.²⁶ We continue in this manner as desired, alternating between evaluating the newly proposed structures, then using the grammar to expand around the best-seen structure to generate new proposals. This simple procedure is easy to implement and offers a strong baseline.

MALKOMES et al. refined this approach by replacing greedy search with *Bayesian optimization* over the space of model structures.²⁸ As in the DUVENAUD et al. procedure, the authors pose the problem in terms of maximizing a score over model structures: a Laplace approximation of the (log) unnormalized structure posterior (4.21). This objective function was then modeled using a Gaussian process, which informed a sequential Bayesian optimization procedure seeking to effectively manage the exploration–exploitation tradeoff in the space of candidate structures. The Gaussian process in model space requires a covariance function over model structures, and the authors proposed an exotic “kernel kernel” evaluating the similarity of proposed structures in terms of the overlap between their hyperparameter-marginal priors for the given dataset. The resulting optimization procedure was found to rapidly locate promising models across a range of regression tasks.



Samples from objective function models incorporating the example covariance structure (4.24).

additive decompositions, § 3.5, p. 61

²⁷ J. R. GARDNER et al. (2017). Discovering and Exploiting Additive Structure for Bayesian Optimization. *AISTATS 2017*.

²⁸ G. MALKOMES et al. (2016). Bayesian optimization for automated model selection. *NeurIPS 2016*.

²⁹ G. MALKOMES and R. GARNETT (2018). Automating Bayesian optimization with Bayesian optimization. *NeurIPS 2018*.

³⁰ J. R. GARDNER et al. (2017). Discovering and Exploiting Additive Structure for Bayesian Optimization. *AISTATS 2017*.

³¹ J. SNOEK et al. (2012). Practical Bayesian Optimization of Machine Learning Algorithms. *NeurIPS 2012*.

models and model structures: § 4.1, p. 68

Bayesian inference over (parametric) model spaces: § 4.2, p. 70

Follow-on work demonstrated a *completely automated* Bayesian optimization system built on this structure search procedure avoiding any manual modeling at all.²⁹ The key idea was to dynamically maintain a set of plausible model structures throughout optimization. Predictions are made via model averaging over this set, offering robustness to model misspecification when computing the outer optimization policy. Every time a new observation is obtained, the set of model structures is then updated via a continual Bayesian optimization in model space given the new data. This interleaving of Bayesian optimization in data space and model space offered promising performance.

Finally, GARDNER et al. offered an alternative to *optimization* over model structures by constructing a Markov chain Monte Carlo routine to *sample* model structures from their posterior (4.21).³⁰ The proposed sampler was a realization of the Metropolis–Hastings algorithm with a custom proposal distribution making minor modifications to the incumbent structure. In the case of the additive decompositions considered in that work, this step consisted of applying random atomic operations such as merging or splitting components of the existing decomposition. Despite the absolutely enormous number of possible additive decompositions, this MCMC routine was able to quickly locate promising structures, and averaging over the sampled structures for prediction resulted in superior optimization performance as well.

SUMMARY OF MAJOR IDEAS

We have presented a convenient framework for model assessment, selection, and averaging grounded in Bayesian inference; this is the predominant approach with Gaussian process models. In the context of Bayesian optimization, perhaps the most important development was the notion of *model averaging*, which has proven beneficial to empirical performance³¹ and has become standard practice.

- Model assessment entails deriving preferences over a space of candidate models of a given system in light of available data.
- In its purest form, a *model* in this context is a prior distribution over observed values y arising from observations at a given set of locations x , $p(y | x)$. A convenient mechanism for specifying a model is via a prior process for a latent function, $p(f)$, and an observation model conditioned on this function, $p(y | x, \phi)$ (4.1–4.2).
- With Gaussian process models, it is convenient to work with combinations of *parametric* forms for the prior mean function, prior covariance function, and observation model, a construct we call a *model structure*. A model structure defines a space of corresponding models by traversing its parameter space (4.3), allowing us to build expressive model spaces.
- Once we delineate a space of candidate models, model assessment becomes straightforward if we make the – perhaps dubious but nonetheless practical – assumption that the mechanism generating our data is con-

tained within this space. This allows us to treat that true model as a random variable and proceed via Bayesian inference over the chosen model space.

- This inference proceeds as normal. We first define a *model prior* capturing any initial beliefs over the model space. Then, given a set of observations $\mathcal{D} = (\mathbf{x}, \mathbf{y})$, the model posterior is proportional to the model prior and a measure of model fit known as the *marginal likelihood* or *model evidence*, the probability (density) of the observed data under the model.
- In addition to quantifying fit, the model evidence encodes an automatic penalty for model complexity, an effect known as *Bayesian Occam's razor*.
- Model evidence can be computed in closed form for Gaussian process models with additive Gaussian observation noise (4.8).
- Model inference is especially convenient when the model space is a single model structure, but can be extended to spaces built from multiple model structures with a bit of extra bookkeeping.
- The model posterior provides a simple means of model assessment by establishing preferences according to posterior probability. If we must commit to a *single* model to explain the data – a task known as *model selection* – we then select the maximum a posteriori (MAP) model.
- Model selection may not be prudent when the model posterior is very flat, which is common when observations are scarce. In this case many models may be compatible with the data but incompatible in their predictions, which should be accounted for in the interest of robustness. *Model averaging* is a natural solution, where we marginalize the unknown model when making predictions according to the model posterior.
- Model averaging cannot in general be performed in closed form for Gaussian process models; however, we may proceed via MCMC sampling (4.14–4.15) or by appealing to more lightweight approximation schemes.
- Appealing to metaheuristics allows us to *automatically* search a space of candidate model *structures* to explain a given dataset. Once sufficiently mature, such schemes may some day enable fully automated Bayesian optimization pipelines that sidestep explicit modeling altogether.

The next chapter marks a major departure from our discussion thus far, which has focused on modeling and making predictions from data. We will now shift our attention from inference to decision making, with the goal of building effective optimization policies informed by the models we have now fully developed. This endeavor will consume the bulk of the remainder of the book.

The first step will be to develop a framework for optimal decision making under uncertainty. Our work to this point will serve an essential component of this framework, as every such decision will be made with reference to a posterior belief about what might happen as a result. In the context of optimization, this belief will take the form of a posterior predictive distribution for proposed observations given data, and our investment in building faithful models will pay off in spades.

Bayesian Occam's razor, § 4.2, p. 71

multiple model structures: § 4.5, p. 78

model selection via MAP inference: § 4.3, p. 73

model averaging: § 4.4, p. 74

approximations to model-marginal posterior:
figures 4.6–4.8 and surrounding text

automating model structure search: § 4.6, p. 81

5

DECISION THEORY FOR OPTIMIZATION

Optimization entails a series of decisions. Most obviously, we must repeatedly decide where to make each successive observation guided by the available data. Some settings also demand we decide when to terminate optimization, weighing the potential benefit from continuing against any costs that may be incurred. It is not obvious how we should make these decisions, especially in the face of incomplete and constantly evolving knowledge about the objective function that is only refined via the outcomes of our own actions.

In the previous four chapters, we established Bayesian inference as a framework for reasoning about uncertainty that offers partial guidance. The primary obstacle to decision making during optimization is uncertainty about the objective function, and, by extension, the outcomes of proposed observations. Bayesian inference allows us to reason about an unknown objective function with a probability distribution over plausible functions that we may seamlessly update as we gather new information. This belief over the objective function in turn enables prediction of proposed observations via the posterior predictive distribution.

How can we use these beliefs to guide our decisions? Bayesian inference offers no direct answer, but in this chapter we will bridge this gap. We will develop *Bayesian decision theory* as a principled means of decision making under uncertainty and apply this approach in the context of optimization, demonstrating how to use a probabilistic belief about an objective function to inform intelligent optimization policies.

Recall our model of sequential optimization outlined in algorithm 1.1, repeated for convenience on the following page. We begin with an arbitrary set of data, which we build upon through a sequence of observations of our own design. The core of the procedure is an *optimization policy*, which examines any already gathered data and makes the fundamental decision of where to make the next observation. With a policy in hand, optimization proceeds by repeating a straightforward pattern: the policy selects the next observation location, then we acquire the requested measurement and update our data accordingly. We repeat this process until satisfied, at which point we return the collected data.

Barring the question of termination, the behavior of this procedure is entirely determined by the policy, and constructing optimization policies will be our primary concern in this and the following chapters. We will begin with sheer audacity: we will derive the *optimal* policy – in terms of maximizing the expected quality of the returned data – in a generic setting. The reader may wonder why this book is so long if the optimal policy is apparently so simple. As it turns out, this theoretically optimal procedure is usually impossible to compute and rarely of practical value. However, our careful derivation will shed light on how we might derive effective approximations. This is a common theme in Bayesian optimization and will be our focus in chapters 7 and 8.

The question of when to terminate optimization also represents a decision that can be of critical importance in some applications. A

Bayesian decision theory

formalization of optimization, § 1.1, p. 2

optimization policy

optimal optimization policies: § 5.2, p. 91

running time and approximation: § 5.3, p. 99

chapter 7: common Bayesian optimization policies, p. 123

chapter 8: computation of Bayesian optimization policies for Gaussian processes, p. 157

Algorithm 1.1: Sequential optimization.

input: initial dataset \mathcal{D}	► can be empty
repeat	
$x \leftarrow \text{POLICY}(\mathcal{D})$	► select the next observation location
$y \leftarrow \text{OBSERVE}(x)$	► observe at the chosen location
$\mathcal{D} \leftarrow \mathcal{D} \cup \{(x, y)\}$	► update dataset
until termination condition reached	► e.g., budget exhausted
return \mathcal{D}	

stopping rule

¹ A predominant example is a preallocated budget on the number of allowed observations, in which case we are compelled to stop after exhausting the budget regardless of progress.

optimal stopping rules: § 5.4, p. 103

utility functions for optimization: chapter 6,
p. 109

acquisition function, infill function, figure of
merit

acquisition function, $\alpha(x; \mathcal{D})$

procedure for inspecting an observed dataset and deciding whether to stop or continue optimization is called a *stopping rule*. In optimization, the stopping rule is often fixed and known before we begin, in which case we do not need to worry over its design.¹ However, in some scenarios, we may wish instead to consider our evolving understanding of the objective function and the expected cost of further observations to dynamically decide when to stop, requiring more subtle adaptive stopping rules. We will also address termination decisions in this chapter and will again begin by deriving the *optimal* – but intractable – stopping procedure, which will inspire efficient and effective approximations.

Practical optimization routines will return datasets that reflect significant progress on our global optimization problem (1.1) in some way. For example, we may wish to return datasets containing near-optimal values of the objective function. Alternatively, we may be satisfied returning datasets that indirectly reveal likely locations of the global optimum or achieve some other related goal. We will formalize this notion of a returned dataset’s *utility* shortly and use it to guide optimization. First we pause to introduce a useful and pervasive technique for implicitly defining an optimization policy by maximizing a score function over the domain.

Defining optimization policies via acquisition functions

A convenient mechanism for defining an optimization policy is by specifying an intermediate so-called *acquisition function* (also called an *infill function* or *figure of merit*) that provides a score to each potential observation location commensurate with its propensity for aiding the optimization task. We then define a policy by observing at a point judged most promising by the acquisition function. Nearly all Bayesian optimization policies are defined in this manner, and this relationship is so intimate that the phrase “acquisition function” is often used interchangeably with “policy” in the literature and conversation, with *maximization* of the acquisition function understood.

Specifically, an acquisition function $\alpha: \mathcal{X} \rightarrow \mathbb{R}$ assigns a score to each point in the domain reflecting our preferences over locations for the next observation. Of course, these preferences will presumably depend on the data we have already observed. To make this dependence explicit, we adopt the notation $\alpha(x; \mathcal{D})$ for a general acquisition function, where available data serve as parameters. In the Bayesian approach, acquisition

functions are defined by deriving the posterior belief of the objective function given the data, $p(f | \mathcal{D})$, then defining preferences with respect to this belief.

An acquisition function α encodes preferences over potential observation locations by inducing a total order over the domain: given data \mathcal{D} , observing at a point x is preferred to at another point x' whenever $\alpha(x; \mathcal{D}) > \alpha(x'; \mathcal{D})$. Thus a rational action in light of these preferences is one maximizing the acquisition function:²

$$x \in \arg \max_{x' \in \mathcal{X}} \alpha(x'; \mathcal{D}). \quad (5.1)$$

Solving (5.1) maps a set of observed data \mathcal{D} to a point $x \in \mathcal{X}$ to observe next, exactly the role of an optimization policy. At first this idea may sound absurd: we have proposed solving a global optimization problem (1.1) by repeatedly solving global optimization problems (5.1)! To resolve this apparent paradox, we note that acquisition functions in common use have properties rendering their optimization considerably more tractable than the problem we ultimately wish to solve. Typical acquisition functions are both cheap to evaluate and analytically differentiable, allowing the use of off-the-shelf optimizers when computing the policy (5.1). The objective function, on the other hand, is assumed to be expensive to evaluate, and its gradient is often unavailable. Therefore we can reduce a difficult, expensive problem to a series of simpler, inexpensive problems – a reasonable pursuit!

Numerous acquisition functions have been proposed for Bayesian optimization, and we will describe many popular choices in detail in chapter 7. The most prominent means to constructing acquisition functions is *Bayesian decision theory*, an approach to optimal decision making we will discuss over the remainder of the chapter.

encoding preferences with an acquisition function

² Ties may be broken arbitrarily.

the paradox of Bayesian optimization: global optimization via...global optimization?

common Bayesian optimization policies:
chapter 7, p. 123

5.1 INTRODUCTION TO BAYESIAN DECISION THEORY

Bayesian decision theory is a framework for decision making under uncertainty that is flexible enough to handle effectively any scenario. Instead of presenting the entire theory in complete abstraction, we will introduce the essential concepts with an eye to the context of optimization. For a more in-depth and theoretical treatment, the interested reader may refer to numerous comprehensive reviews of the subject.³ A good familiarity with this material can demystify some key ideas that are often glossed over in the Bayesian optimization literature, as it serves as the “hidden origin” of many common acquisition functions.

In this section we will introduce to the Bayesian approach to decision making and demonstrate how to make optimal decisions in the case of a single isolated decision. Ultimately, we will require a theory for making a *sequence* of decisions to reason over an entire optimization session. In the next section, we will extend the line of reasoning presented below to address sequential decision making and the construction of optimization policies.

³ The following would be excellent companion texts:

M. H. DEGROOT (1970). *Optimal Statistical Decisions*. McGraw-Hill.

J. O. BERGER (1985). *Statistical Decision Theory and Bayesian Analysis*. Springer-Verlag.

action space, \mathcal{A}
unknown variables affecting decision outcome, ψ
relevant observed data, \mathcal{D} posterior belief about ψ , $p(\psi \mathcal{D})$
utility function, $u(a, \psi, \mathcal{D})$
expected utility

⁴ Typical presentations of Bayesian decision theory omit the data from the utility function, but including it offers more generality, and this allowance will be important when we turn our attention to optimization policies.

Isolated decisions

A decision problem under uncertainty has two defining characteristics. The first is the *action space* \mathcal{A} , the set of all available decisions. Our task is to select an action from this space. For example, in sequential optimization, an optimization policy decision must select a point in the domain \mathcal{X} for observation, and so we have $\mathcal{A} = \mathcal{X}$.

The second critical feature is the presence of *uncertain* elements of the world influencing the outcomes of our actions, complicating our decision. Let ψ represent a random variable encompassing any relevant uncertain elements when making and evaluating a decision. Although we may lack perfect knowledge, Bayesian inference allows us to reason about ψ in light of data via the posterior distribution $p(\psi | \mathcal{D})$, and we will use this belief to inform our decision.

Suppose now we must select a decision from an action space \mathcal{A} under uncertainty in ψ , informed by a set of observed data \mathcal{D} . To guide our choice, we select a real-valued *utility function* $u(a, \psi, \mathcal{D})$. This function measures the quality of selecting the action a if the true state of the world were revealed to be ψ , with higher utilities indicating more favorable outcomes. The arguments to a utility function comprise everything required to judge the quality of a decision in hindsight: the proposed action a , what we know (the data \mathcal{D}), and what we don't know (the uncertain elements ψ).⁴

We cannot know the exact utility that would result from selecting any given action *a priori*, due to our incomplete knowledge of ψ . We can, however, compute the *expected* utility that would result from selecting an action a , according to our posterior belief:

$$\mathbb{E}[u(a, \psi, \mathcal{D}) | a, \mathcal{D}] = \int u(a, \psi, \mathcal{D}) p(\psi | \mathcal{D}) d\psi. \quad (5.2)$$

This expected utility maps each available action to a real value, inducing a total order and providing a straightforward mechanism for making our decision. We pick an action maximizing the expected utility:

$$a \in \arg \max_{a' \in \mathcal{A}} \mathbb{E}[u(a', \psi, \mathcal{D}) | a', \mathcal{D}]. \quad (5.3)$$

This decision is optimal in the sense that no other action results in greater expected utility. (By definition!) This procedure for acting optimally under uncertainty – computing expected utility with respect to relevant unknown variables and maximizing to select an action – is the central tenant of Bayesian decision making.⁵

Example: recommending a point for use after optimization

With this abstract decision-making framework established, let us analyze an example decision that might be faced in the context of optimization. Consider a scenario where the purpose of optimization is to identify a single point $x \in \mathcal{X}$ for perpetual use in a production system, preferring

J. VON NEUMANN and O. Morgenstern (1944). *Theory of Games and Economic Behavior*. Princeton University Press. [appendix A]

locations achieving higher values of the objective function. If we run an optimizer and it returns some dataset \mathcal{D} , which point should we select for our final recommendation?

We may model this choice as a decision problem with action space $\mathcal{A} = \mathcal{X}$, where we must reason under uncertainty about the objective function f . We first select a utility function quantifying the quality of a given recommendation x in hindsight. One natural choice would be

$$u(x, f) = f(x) = \phi,$$

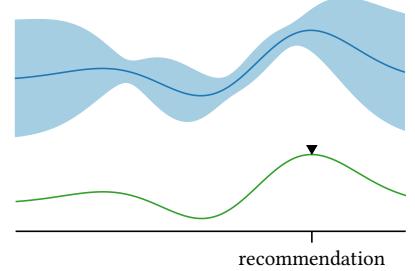
which rewards points for achieving high values of the objective function. Now if our optimization procedure returned a dataset \mathcal{D} , the expected utility from recommending a point x is simply the posterior mean of the corresponding function value:

$$\mathbb{E}[u(x, f) | x, \mathcal{D}] = \mathbb{E}[\phi | x, \mathcal{D}] = \mu_{\mathcal{D}}(x). \quad (5.4)$$

Therefore, an optimal recommendation maximizes the posterior mean:

$$x \in \arg \max_{x' \in \mathcal{X}} \mu_{\mathcal{D}}(x').$$

Of course, other considerations in a given scenario such as risk aversion might suggest some other utility function or action space would be more appropriate, in which case we are free to select any alternative as we see fit. We will discuss terminal recommendations at length in the next chapter, including alternative utility functions and action spaces.



Optimal terminal recommendation. Above: posterior belief about an objective function given the data returned by an optimizer, $p(f | \mathcal{D})$. Below: the expected utility for our example, the posterior mean $\mu_{\mathcal{D}}(x)$. The optimal recommendation maximizes the expected utility.

terminal recommendations: § 6.1, p. 109

5.2 SEQUENTIAL DECISIONS WITH A FIXED BUDGET

We have now introduced Bayesian decision theory as a framework for computing optimal decisions informed by data. The key idea is to measure the post hoc quality of a decision with an appropriately designed utility function, then choose actions maximizing expected utility according to our beliefs. We will now apply this idea to the construction of optimization policies. This setting is considerably more complicated because each decision we make over the course of optimization will shape the context of all future decisions.

Modeling policy decisions

To define an optimization routine, we must design a policy to adaptively design a sequence of observations seeking the optimum. Following our discussion in the previous section, we will model each of these choices as a decision problem under uncertainty. Some aspects of this modeling will be straightforward and others will take some care. To begin, the action space of each decision is the domain \mathcal{X} , and we must act under uncertainty about the objective function f , which induces uncertainty about the outcomes of proposed observations. Fortunately, we may make each decision guided by any data obtained from previous decisions.

Bayesian inference of the objective function:
 § 1.2, p. 8

optimization utility function, $u(\mathcal{D})$

utility functions for optimization: chapter 6,
 p. 109

To reason about uncertainty in the objective function, we follow the path laid out in the preceding chapters and maintain a probabilistic belief throughout optimization, $p(f \mid \mathcal{D})$. We make no assumptions regarding the nature of this distribution, and in particular it need not be a Gaussian process. Equipped with this belief, we may reason about the result of making an observation at some point x via the posterior predictive distribution $p(y \mid x, \mathcal{D})$ (1.7), which will play a key role below.

The ultimate purpose of optimization is to collect and return a dataset \mathcal{D} . Before we can reason about what data we should acquire, we must first clarify what data we *would like* to acquire. Following the previous section, we will accomplish this by defining a utility function $u(\mathcal{D})$ to evaluate the quality of data returned by an optimizer. This utility function will serve to establish preferences over optimization outcomes: all other things being equal, we would prefer to return a dataset with higher utility than any dataset with lower utility. As before, we will use this utility to guide the design of policies, by making observations that, in expectation, promise the biggest improvement in utility. We will define and motivate several utility functions used for optimization in the next chapter, and some readers may wish to jump ahead to that discussion for explicit examples before continuing. In the following, we will develop the general theory in terms of an arbitrary utility function.

Uncertainty faced during optimization

Suppose \mathcal{D} is a dataset of previous observations and that we must select the next observation location x . This is the core decision defining an optimization policy, and we will make all such decisions in the same manner: by maximizing the expected utility of the data we will return.

Although this sounds straightforward, let us consider the uncertainty faced when contemplating this decision in more detail. When evaluating a potential action x , uncertainty in the objective function induces uncertainty in the corresponding value y we will observe. Bayesian inference allows us to reason about this uncertain outcome via the posterior predictive distribution (1.7), and we may hope to be able to address this uncertainty without much trouble. However, we must also consider that evaluating at x would add the unknown observation (x, y) to our dataset, and that the contents of this updated dataset would be consulted for all future decisions. Thus we must reason not only about the outcome of the present observation but also its impact on the *entire remainder of optimization*. This requires special attention and distinguishes sequential decisions from the isolated decisions discussed in the last section.

Intuitively, we might suspect that decisions made closer to termination should be easier, as fewer future decisions depend on their outcomes. This is indeed the case, and it will be prudent to define optimization policies *in reverse*.⁶ We will first reason about the *final* decision – when we are freed from the burden of having to ponder any future observations – and proceed backwards to the choice of the first observation location, working out optimal behavior every step along the way.

⁶ In fact, we have already begun by analyzing a decision *after* optimization has completed!

In this section we will consider the construction of optimization policies assuming that we have a fixed and known budget on the number of observations we will make. This scenario is both common in practice and convenient for analysis, as we can for now ignore the question of when to terminate optimization. Note that this assumption effectively implies that every observation has a constant acquisition cost, which may not always be reasonable. We will address variable observation costs and the question of when to stop optimization later in this chapter.

Assuming a fixed observation budget allows us to reason about optimization policies in terms of the number of observations remaining to termination, which will always be known. The problem we will consider in this section then becomes the following: provided an arbitrary set of data, how should we design our next evaluation location when exactly τ observations remain before termination? In sequential decision making, this value is known as the decision *horizon*, as it indicates how far we must look ahead into the future when reasoning about the present.

To facilitate our discussion, we pause to define notation for future data that will be encountered during optimization relative to the present. When considering an observation at some point x , we will call the value resulting from an observation there y . We will then call the dataset available at the next stage of optimization $\mathcal{D}_1 = \mathcal{D} \cup \{(x, y)\}$, where the subscript indicates the number of future observations incorporated into the current data. We will write (x_2, y_2) for the following observation, which when acquired will form \mathcal{D}_2 , etc. Our final observation τ steps in the future will then be (x_τ, y_τ) , and the dataset returned by our optimization procedure will be \mathcal{D}_τ , with utility $u(\mathcal{D}_\tau)$.

This utility of the data we return is our ultimate concern and will serve as the utility function used to design every observation. Note we may write this utility in the same form we introduced in our general discussion:

$$u(\mathcal{D}_\tau) = u(\underbrace{\mathcal{D},}_{\text{known}} \underbrace{x,}_{\text{action}} \underbrace{y, x_2, y_2, \dots, x_\tau, y_\tau}_{\text{unknown}}),$$

which expresses the terminal utility in terms of a proposed current action x , the known data \mathcal{D} , and the unknown future data to be obtained: the not-yet observed value y , and the locations $\{x_2, \dots, x_\tau\}$ and values $\{y_2, \dots, y_\tau\}$ of any following observations.

Following our treatment of isolated decisions, we evaluate a potential observation location x via the expected utility at termination ultimately obtained if we observe at that point next:

$$\mathbb{E}[u(\mathcal{D}_\tau) | x, \mathcal{D}], \quad (5.5)$$

and define an optimization policy via maximization:

$$x \in \arg \max_{x' \in \mathcal{X}} \mathbb{E}[u(\mathcal{D}_\tau) | x', \mathcal{D}]. \quad (5.6)$$

On its surface, this proposal is relatively simple. However, we must now consider how to actually *compute* the expected terminal utility (5.5).

fixed, known budget

cost-aware optimization: § 5.4, p. 103

number of remaining observations (horizon), τ

putative next observation and dataset: (x, y) , \mathcal{D}_1

putative following observation and dataset: (x_2, y_2) , \mathcal{D}_2

putative final observation and dataset: (x_τ, y_τ) , \mathcal{D}_τ

expected terminal utility, $\mathbb{E}[u(\mathcal{D}_\tau) | x, \mathcal{D}]$

Explicitly writing out the expectation over the future data in (5.5) yields the following expression:

$$\int \cdots \int u(\mathcal{D}_\tau) p(y | x, \mathcal{D}) \prod_{i=2}^{\tau} p(x_i, y_i | \mathcal{D}_{i-1}) dy d\{(x_i, y_i)\} \quad (5.7)$$

This integral certainly appears unwieldy! In particular, it is unclear how to reason about uncertainty in our future actions, as we should hope that these actions are made to maximize our welfare rather than generated by a random process. We will show how to compute this expression under the bold but rational assumption that we *make all future decisions optimally*,⁷ and this analysis will reveal the optimal optimization policy.

⁷ This is known as BELLMAN's *principle of optimality*, and will be discussed further later in this section.

⁸ This procedure is often called "backward induction," where we consider the last decision first and work backward in time. Our approach of a forward induction on the horizon is equivalent.

defining a policy by maximizing an acquisition function: § 5, p. 88

isolated decisions: § 5.1, p. 89

We will proceed via induction on the number of evaluations remaining before termination, τ . We will first determine optimal behavior when only one observation remains and then inductively consider increasingly long horizons.⁸ For this analysis it will be useful to introduce notation for the expected *increase* in utility achieved when beginning from an arbitrary dataset \mathcal{D} , making an observation at x , and then continuing optimally until termination τ steps in the future. We will write

$$\alpha_\tau(x; \mathcal{D}) = \mathbb{E}[u(\mathcal{D}_\tau) | x, \mathcal{D}] - u(\mathcal{D})$$

for this quantity, which is simply the expected terminal utility (5.5) shifted by the utility of our existing data, $u(\mathcal{D})$. It is no coincidence this notation echoes our notation for acquisition functions! We will characterize the optimal optimization policy by a family of acquisition functions defined in this manner.

Fixed budget: one observation remaining

We first consider the case where only one observation remains before termination; that is, the horizon is $\tau = 1$. In this case the terminal dataset will be the current dataset augmented with a single additional observation. As there are no following decisions to consider, we may analyze the decision using the framework we have already developed for isolated decisions. The marginal gain in utility from a final evaluation at x is an expectation over the corresponding value y with respect to the posterior predictive distribution:

$$\alpha_1(x; \mathcal{D}) = \int u(\mathcal{D}_1) p(y | x, \mathcal{D}) dy - u(\mathcal{D}). \quad (5.8)$$

The optimal observation maximizes the expected marginal gain:

$$x \in \arg \max_{x' \in \mathcal{X}} \alpha_1(x'; \mathcal{D}), \quad (5.9)$$

and leads to our returning a dataset with expected utility

$$u(\mathcal{D}) + \alpha_1^*(\mathcal{D}); \quad \alpha_1^*(\mathcal{D}) = \max_{x' \in \mathcal{X}} \alpha_1(x'; \mathcal{D}). \quad (5.10)$$

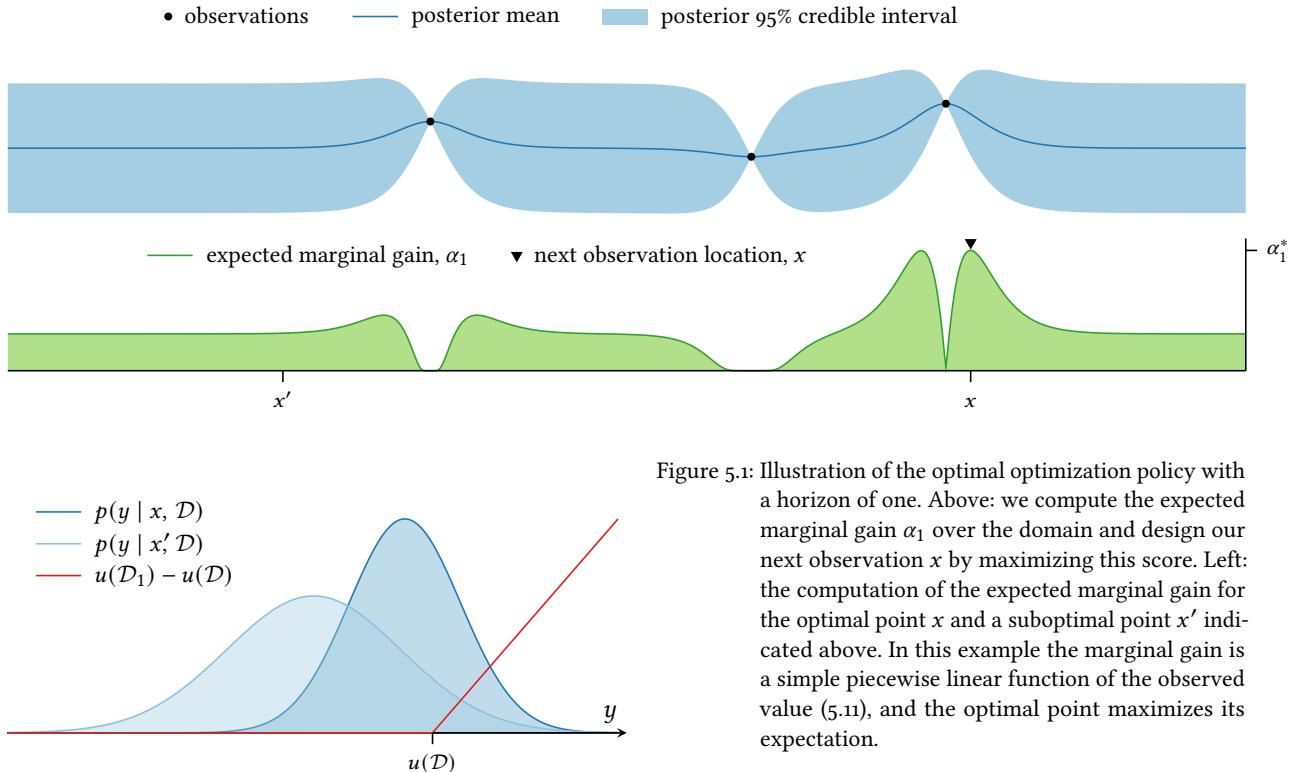


Figure 5.1: Illustration of the optimal optimization policy with a horizon of one. Above: we compute the expected marginal gain α_1 over the domain and design our next observation x by maximizing this score. Left: the computation of the expected marginal gain for the optimal point x and a suboptimal point x' indicated above. In this example the marginal gain is a simple piecewise linear function of the observed value (5.11), and the optimal point maximizes its expectation.

Here we have defined the symbol $\alpha^*(\mathcal{D})$ to represent the expected increase in utility when starting with \mathcal{D} and continuing *optimally* for τ additional observations. This is called the *value* of the dataset with a horizon of τ and will serve a central role below. We have now shown how to compute the value of any dataset with a horizon of $\tau = 1$ (5.10) and how to identify a corresponding optimal action (5.9). This completes the base case of our argument.

We illustrate the optimal optimization policy with one observation remaining in figure 5.1. In this scenario the belief over the objective function $p(f | \mathcal{D})$ is a Gaussian process, and for simplicity we assume our observations reveal exact values of the objective. We consider an intuitive utility function: the maximal objective value contained in the data, $u(\mathcal{D}) = \max f(\mathbf{x})$.⁹ The marginal gain in utility offered by a putative final observation (x, y) is then a piecewise linear function of the observed value:

$$u(\mathcal{D}_1) - u(\mathcal{D}) = \max\{y - u(\mathcal{D}), 0\}; \quad (5.11)$$

that is, the utility increases linearly if we exceed the previously best-seen value and otherwise remains constant. To design the optimal final observation, we compute the expectation of this quantity over the domain and choose the point maximizing it, as shown in the top panels. We also illustrate the computation of this expectation for the optimal choice and

value of \mathcal{D} with horizon τ , $\alpha^*(\mathcal{D})$

illustration of one-step optimal optimization policy

⁹ This is a special case of the *simple reward* utility function, which we discuss further in the next chapter (§ 6.1, p. 109). The corresponding expected marginal gain is the well-known *expected improvement* acquisition function (§ 7.3, p. 127).

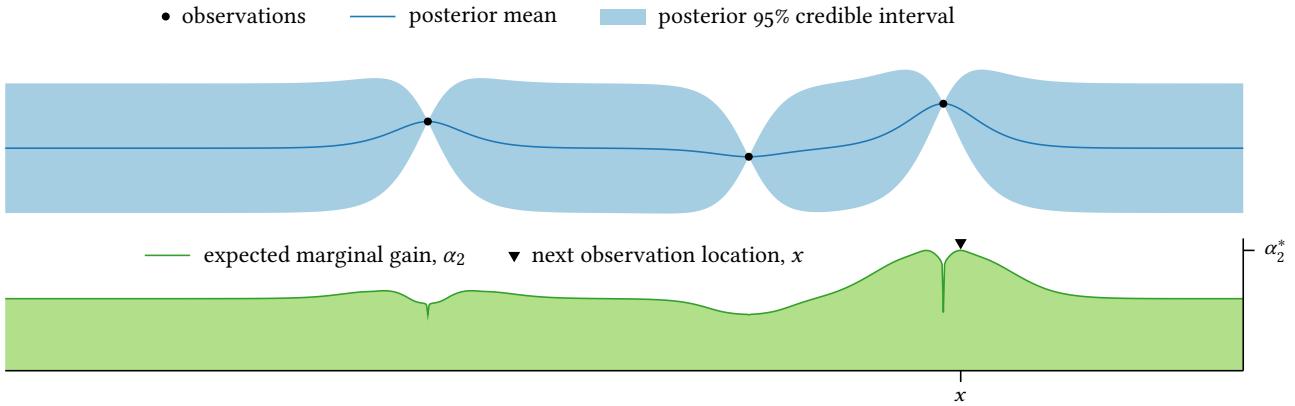
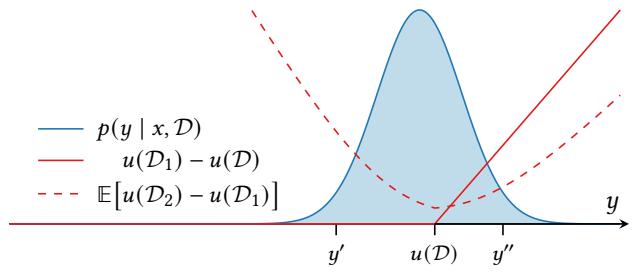


Figure 5.2: Illustration of the optimal optimization policy with a horizon of two. Above: the expected two-step marginal gain α_2 . Right: computation of α_2 for the optimal point x . The marginal gain is decomposed into two components (5.13): the immediate gain $u(\mathcal{D}_1) - u(\mathcal{D})$ and the expected future gain $\mathbb{E}[u(\mathcal{D}_2) - u(\mathcal{D}_1)]$. The chosen point offers a high expected future reward even if the immediate reward is zero; see the facing page for the scenarios resulting from the marked values.



a suboptimal alternative in the bottom panel. We expect an observation at the chosen location to improve utility by a greater amount than any alternative.

Fixed budget: two observations remaining

Rather than proceeding immediately to the inductive case, let us consider the specific case of two observations remaining: $\tau = 2$. Suppose we have obtained an arbitrary dataset \mathcal{D} and must decide where to make the *penultimate* observation x . The reasoning for this special case presents the inductive argument most clearly.

We again consider the expected increase in utility by termination, now after two observations:

$$\alpha_2(x; \mathcal{D}) = \mathbb{E}[u(\mathcal{D}_2) | x, \mathcal{D}] - u(\mathcal{D}). \quad (5.12)$$

Nominally this expectation requires marginalizing the observation y , as well as the final observation location x_2 and its value y_2 (5.7). However, if we assume optimal future behavior, we can simplify our treatment of the final decision x_2 . First we rewrite the two-step expected gain α_2 in terms of the one-step expected gain α_1 , a function for which we have already established a good understanding. We write the two-step difference in

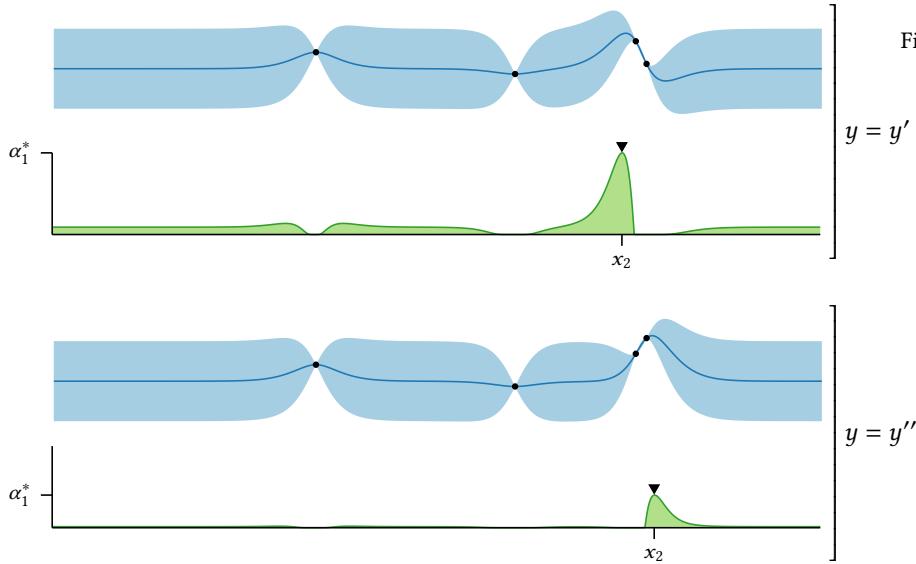


Figure 5.3: The posterior of the objective function given two possible observations resulting from the optimal two-step observation x illustrated on the facing page. The relatively low value y' offers no immediate reward, but reveals a new local optimum and the expected future reward from the optimal final decision x_2 is high. The relatively high value y'' offers a large immediate reward and respectable prospects from the optimal final decision as well.

utility as a telescoping sum:

$$u(\mathcal{D}_2) - u(\mathcal{D}) = [u(\mathcal{D}_1) - u(\mathcal{D})] + [u(\mathcal{D}_2) - u(\mathcal{D}_1)],$$

which yields

$$\alpha_2(x; \mathcal{D}) = \alpha_1(x; \mathcal{D}) + \mathbb{E}[\alpha_1(x_2; \mathcal{D}_1) | x, \mathcal{D}].$$

That is, the expected increase in utility after two observations can be decomposed as the expected increase after our first observation x – the expected *immediate gain* – plus the expected additional increase from the final observation x_2 – the expected *future gain*.

decomposition of expected marginal gain

It is still not clear how to address the second term in this expression. However, from our analysis of the base case, we can reason as follows. Given y (and thus knowledge of \mathcal{D}_1), the *optimal* final decision x_2 (5.9) results in an expected marginal gain of $\alpha_1^*(\mathcal{D}_1)$, a quantity we know how to compute (5.10). Therefore, assuming optimal future behavior, we have:

$$\alpha_2(x; \mathcal{D}) = \alpha_1(x; \mathcal{D}) + \mathbb{E}[\alpha_1^*(\mathcal{D}_1) | x, \mathcal{D}], \quad (5.13)$$

which expresses the desired quantity as an expectation with respect to the current observation y only – the future value α_1^* (5.10) does not depend on either x_2 (due to maximization) or y_2 (due to expectation). The optimal penultimate observation location maximizes the expected gain as usual:

$$x \in \arg \max_{x' \in \mathcal{X}} \alpha_2(x'; \mathcal{D}). \quad (5.14)$$

and provides an expected terminal utility of

$$u(\mathcal{D}) + \alpha_2^*(\mathcal{D}); \quad \alpha_2^*(\mathcal{D}) = \max_{x' \in \mathcal{X}} \alpha_2(x'; \mathcal{D}).$$

illustration of two-step optimal optimization policy

This demonstrates we can achieve optimal behavior for a horizon of $\tau = 2$ and compute the value of any dataset with this horizon.

The optimal policy with two observations remaining is illustrated in figures 5.2 and 5.3. The former shows the expected two-step marginal gain α_2 and the optimal action. This quantity depends both on the immediate gain from the next observation and the expected future gain from the optimal final action. The chosen observation appears quite promising: even if the result offers no immediate gain, it will likely provide information that can be exploited with the optimal final decision x_2 . We show the situation that would be faced in the final stage of optimization for two potential values in figure 5.3. The relatively low value y' offers no immediate gain but sets up an encouraging final decision, whereas the relatively high value y'' offers a significant immediate gain with some chance of further improvement.

Fixed budget: inductive case

We now present the general inductive argument, which closely follows the $\tau = 2$ analysis above. Let τ be an arbitrary decision horizon, and for the sake of induction assume we can compute the value of any dataset with a horizon of $\tau - 1$. Suppose we have an arbitrary dataset \mathcal{D} and must decide where to make the next observation. We will show how to do so optimally and how to compute its value with a horizon of τ .

Consider the τ -step expected gain in utility from observing at some point x :

$$\alpha_\tau(x; \mathcal{D}) = \mathbb{E}[u(\mathcal{D}_\tau) | x, \mathcal{D}] - u(\mathcal{D}),$$

¹⁰ Namely:

$$\begin{aligned} u(\mathcal{D}_\tau) - u(\mathcal{D}) &= \\ &[u(\mathcal{D}_1) - u(\mathcal{D})] + [u(\mathcal{D}_\tau) - u(\mathcal{D}_1)]. \end{aligned}$$

Now if we knew y (and thus \mathcal{D}_1), optimal continued behavior would provide an expected further gain of $\alpha_{\tau-1}^*(\mathcal{D}_1)$, a quantity we can compute via the inductive hypothesis. Therefore, assuming optimal behavior for all remaining decisions, we have:

$$\alpha_\tau(x; \mathcal{D}) = \alpha_1(x; \mathcal{D}) + \mathbb{E}[\alpha_{\tau-1}^*(\mathcal{D}_1) | x, \mathcal{D}], \quad (5.15)$$

which is an expectation with respect to y of a function we can compute. To find the optimal decision and the τ -step value of the data, we maximize:

$$x \in \arg \max_{x' \in \mathcal{X}} \alpha_\tau(x'; \mathcal{D}); \quad (5.16)$$

$$\alpha_\tau^*(\mathcal{D}) = \max_{x' \in \mathcal{X}} \alpha_\tau(x'; \mathcal{D}). \quad (5.17)$$

This demonstrates we can achieve optimal behavior for a horizon of τ given an arbitrary dataset and compute its corresponding value, establishing the inductive case and completing our analysis.

We pause to note that the value of any dataset with null horizon is $\alpha_0^*(\mathcal{D}) = 0$, and thus the expressions in (5.15–5.17) are valid for any horizon and compactly express the proposed policy. Further, we have actually shown that this policy is *optimal* in the sense of maximizing expected terminal utility over the space of all policies, at least with respect to our model of the objective function and observations. This follows from our induction: the base case is established in (5.9), and the inductive case by the sequential maximization in (5.16).¹¹

Bellman optimality and the Bellman equation

Substituting (5.15) into (5.17), we may derive the following recursive definition of the value in terms of the value of future data:

$$\alpha_t^*(\mathcal{D}) = \max_{x' \in \mathcal{X}} \left\{ \alpha_1(x'; \mathcal{D}) + \mathbb{E}[\alpha_{t-1}^*(\mathcal{D}_1) | x', \mathcal{D}] \right\}. \quad (5.18)$$

This is known as the *Bellman equation* and is a central result in the theory of optimal sequential decisions.¹² The treatment of future decisions in this equation – recursively assuming that we will always act to maximize expected terminal utility given the available data – reflects BELLMAN’s *principle of optimality*, which characterizes optimal sequential decision policies in terms of the optimality of subpolicies:¹³

An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision.

That is, to make a sequence of optimal decisions, we make the first decision optimally, then make all following decisions optimally given the outcome!

5.3 COST AND APPROXIMATION OF THE OPTIMAL POLICY

Although the framework presented in the previous section is conceptually simple and theoretically attractive, the optimal policy is unfortunately prohibitive to compute except for very short decision horizons.

To demonstrate the key computational barrier, consider the selection of the penultimate observation location. The expected two-step marginal gain to be maximized is (5.13):

$$\alpha_2(x; \mathcal{D}) = \alpha_1(x; \mathcal{D}) + \mathbb{E}[\alpha_1^*(\mathcal{D}_1) | x, \mathcal{D}].$$

The second term appears to be a straightforward expectation over the one-dimensional random variable y . However, evaluating the integrand in this expectation requires solving a nontrivial global optimization problem (5.10)! Even with only two evaluations remaining, we must solve a doubly nested global optimization problem, an onerous task.

Close inspection of the recursively defined optimal policy (5.15–5.16) reveals that when faced with a horizon of τ , we must solve τ nested

optimal policy: compact notation

optimality

¹¹ Since ties in (5.16) may be broken arbitrarily, this argument does not rule out the possibility of there being multiple, equally good optimal policies.

¹² R. BELLMAN (1952). On the Theory of Dynamic Programming. *Proceedings of the National Academy of Sciences* 38(8):716–719.

Bellman equation

BELLMAN’s principle of optimality

¹³ R. BELLMAN (1957). *Dynamic Programming*. Princeton University Press.

“unrolling” the optimal sequential policy

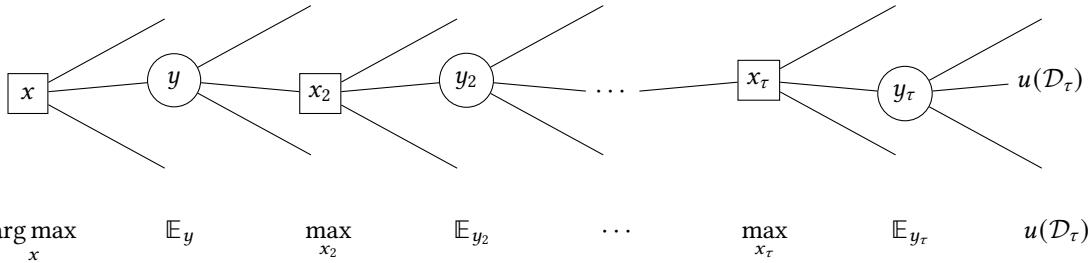


Figure 5.4: The optimal optimization policy as a decision tree. Squares indicate decisions (the choice of each observation), and circles represent expectations with respect to random variables (the outcomes of observations). Only one possible optimization path is shown; dangling edges lead to different futures, and all possibilities are always considered. We maximize the expected terminal utility $u(\mathcal{D}_\tau)$, recursively assuming optimal future behavior.

optimization problems to find the optimal decision. Temporarily adopting compact notation, we may “unroll” the optimal policy as follows:

$$\begin{aligned}
 x &\in \arg \max \alpha_\tau; \\
 \alpha_\tau &= \alpha_1 + \mathbb{E}[\alpha_{\tau-1}^*] \\
 &= \alpha_1 + \mathbb{E}[\max \alpha_{\tau-1}] \\
 &= \alpha_1 + \mathbb{E}[\max\{\alpha_1 + \mathbb{E}[\alpha_{\tau-2}^*]\}] \\
 &= \alpha_1 + \mathbb{E}\left[\max\left\{\alpha_1 + \mathbb{E}\left[\max\{\alpha_1 + \mathbb{E}[\max\{\alpha_1 + \dots\}]\}\right]\right\}\right].
 \end{aligned}$$

The design of each optimal decision requires repeated maximization over the domain and expectation over unknown observations until the horizon is reached. This computation is visualized as a decision tree in figure 5.4, where it is clear that each unknown quantity contributes a significant branching factor. Computing the expected utility at x exactly requires a complete traversal of this tree.

The cost of computing the optimal policy evidently grows with the horizon. Let us perform a careful running time analysis for a naïve implementation via exhaustive traversal of the decision tree in figure 5.4 with off-the-shelf procedures. Suppose we use an optimization routine for each maximization and a numerical quadrature routine for each expectation encountered in this computation. If we allow n evaluations of the objective for each call to the optimizer and q observations of the integrand for each call to the quadrature routine, then each decision along the horizon will contribute a multiplicative factor of $\mathcal{O}(nq)$ to the total running time. Computing the optimal decision with a horizon of τ thus requires $\mathcal{O}(n^\tau q^\tau)$ work, an exponential growth in running time with respect to the horizon.

Evidently, the computational effort required for realizing the optimal policy quickly becomes intractable, and we must find some alternative mechanism for designing effective optimization policies. General approximation schemes for the optimal policy have been studied in depth under the name *approximate dynamic programming*,¹⁴ and usually operate as

running time of optimal policy

evaluation budget for optimization, n
evaluation budget for quadrature, q

¹⁴ Detailed references are provided by:

W. B. POWELL (2011). *Approximate Dynamic Programming: Solving the Curses of Dimensionality*. John Wiley & Sons.

D. P. BERTSEKAS (2017). *Dynamic Programming and Optimal Control*. Vol. 1. Athena Scientific.

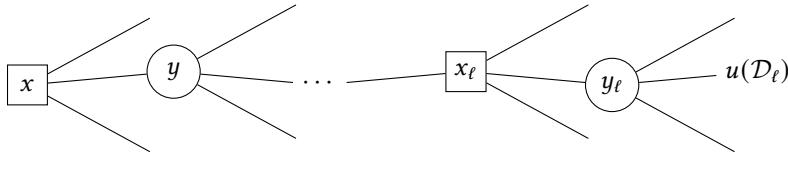


Figure 5.5: A lookahead approximation to the optimal optimization policy. We choose the optimal decision for a limited horizon $\ell \ll \tau$ decisions, ignoring any observations that would follow.

follows. We begin with the intractable optimal expected marginal gain (5.15):

$$\alpha_\tau(x; \mathcal{D}) = \alpha_1(x; \mathcal{D}) + \mathbb{E}[\alpha_{\tau-1}^*(\mathcal{D}_1) | x, \mathcal{D}],$$

and substitute a tractable approximation for the ‘‘hard’’ part of the expression: the recursively defined future value α^* (5.18). The result is an acquisition function inducing a suboptimal – but rationally guided – approximate policy. Two particular approximations schemes have proven useful in Bayesian optimization: *limited lookahead* and *rollout*.

Limited lookahead

One widespread and surprisingly effective approximation is to simply limit how many future observations we consider in each decision. This is practical as decisions closer to termination require substantially less computation than earlier decisions.

With this in mind, we can construct a natural family of approximations to the optimal policy defined by artificially limiting the horizon used throughout optimization to some computationally feasible maximum ℓ . When faced with an infeasible decision horizon τ , we make the crude approximation

$$\alpha_\tau(x; \mathcal{D}) \approx \alpha_\ell(x; \mathcal{D}),$$

and by maximizing this score, we act optimally under the incorrect but convenient assumption that only ℓ observations remain. This effectively assumes $u(\mathcal{D}_\tau) \approx u(\mathcal{D}_\ell)$.¹⁵ This may be reasonable if we expect decreasing marginal gains, implying a significant fraction of potential gains can be attained within the truncated horizon. This scheme is often described (sometimes disparagingly) as *myopic*, as we limit our sight to only the next few observations rather than looking ahead to the full horizon.

A policy that designs each observation to maximize the limited-horizon acquisition function $\alpha_{\min\{\ell, \tau\}}$ is called an ℓ -step lookahead policy.¹⁶ This is also called a *rolling horizon* strategy, as the fixed horizon ‘‘rolls along’’ with us as we go. By limiting the horizon, we bound the computational effort required for each decision to at-most $\mathcal{O}(n^\ell q^\ell)$ time with the implementation described above. This can be a considerable savings when the observation budget is much greater than the selected lookahead. A lookahead policy is illustrated as a decision tree in figure 5.5. Comparing to the optimal policy in figure 5.4, we simply ‘‘cut off’’ and ignore any portion of the tree lying deeper than ℓ steps in the future.

¹⁵ Equivalently, we approximate the true future value $\alpha_{\tau-1}^*$ with $\alpha_{\ell-1}^*$.

myopic approximations

ℓ -step lookahead
rolling horizon

¹⁶ We take the minimum to ensure we don’t look beyond the true horizon, which would be nonsense.

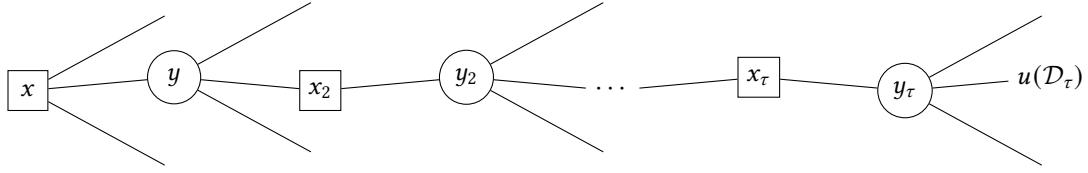


Figure 5.6: A decision tree representing a rollout policy. Comparing to the optimal policy in figure 5.4, we simulate future decisions starting with x_2 using an efficient but *suboptimal* heuristic policy, rather than the intractable optimal policy. We maximize the expected terminal utility $u(\mathcal{D}_\tau)$, assuming potentially suboptimal future behavior.

common Bayesian optimization policies:
chapter 7, p. 123

Particularly important in Bayesian optimization is the special case of *one-step lookahead*, which successively maximizes the expected marginal gain after acquiring a single additional observation, α_1 . One-step lookahead is the most efficient lookahead approximation (barring the absurdity that would be “zero-step” lookahead), and it is often possible to derive closed-form, analytically differentiable expressions for α_1 , enabling efficient implementation. Many well-known acquisition functions represent one-step lookahead approximations for some implicit choice of utility function, as we will see in chapter 7.

base policy, heuristic policy

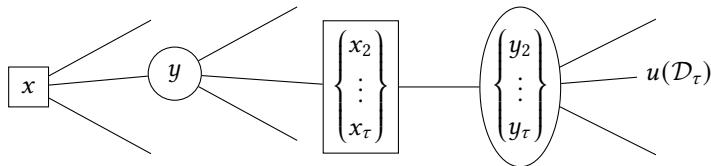
choice of base policy

Rollout

The optimal policy evaluates a potential observation location by simulating the entire remainder of optimization following that choice, recursively assuming we will use the *optimal* policy for every future decision. Although sensible, this is clearly intractable. *Rollout* is an approach to approximate policy design that emulates the structure of the optimal policy, but using a tractable *suboptimal* policy to simulate future decisions.

A rollout policy is illustrated as a decision tree in figure 5.6. Given a putative next observation (x, y) , we use an inexpensive so-called *base* or *heuristic* policy to simulate a plausible – but perhaps suboptimal – realization of the following decision x_2 . Note there is no branching in the tree corresponding to this decision, as it does not depend on the exhaustively enumerated subtree required by the optimal policy. We then take an expectation with respect to the unknown value y_2 as usual. Given a putative value of y_2 , we use the base policy to select x_3 and continue in this manner until reaching the decision horizon. We use the terminal utilities in the resulting pruned tree to estimate the expected marginal gain α_τ , which we maximize as a function of x .

There are no constraints on the design of the base policy used in rollout; however, for this approximation to be sensible, we must choose something relatively efficient. One common and often effective choice is to simulate future decisions with one-step lookahead. If we again use off-the-shelf optimization and quadrature routines to traverse the rollout decision tree in figure 5.6 with this particular choice, the running time of the policy with a horizon of τ is $\mathcal{O}(n^2 q^\tau)$, significantly faster



than the optimal policy. Although there is still exponential growth with respect to q , we typically have $q \ll n$,¹⁷ so we can usually entertain farther horizons with rollout than with limited lookahead with the same amount of computational effort.

Due to the flexibility in the design of the base policy, rollout is a remarkably flexible approximation scheme. For example, we can combine rollout with the idea of limiting the decision horizon to yield approximate policies with tunable running time. In fact, we can interpret ℓ -step lookahead as a special case of rollout, where the base policy designs the next $\ell - 1$ decisions optimally assuming a myopic horizon and then simply *terminates early*, discarding any remaining budget.

We may also adopt a base policy that designs all remaining observations *simultaneously*. Ignoring the dependence between these decisions can provide a computational advantage while retaining awareness of the evolving decision horizon, and such *batch rollout* schemes have proven useful in Bayesian optimization. A batch rollout policy is illustrated as a decision tree in figure 5.7. Although we account for the entire horizon, the tree depth is reduced dramatically compared to the optimal policy.

5.4 COST-AWARE OPTIMIZATION AND TERMINATION AS A DECISION

Thus far we have only considered the construction of optimization policies under a known budget on the total number of observations. Although this scenario is pervasive, it is not universal. In some situations, we might wish instead to use our evolving beliefs about the objective function to decide *dynamically* when termination is the best course of action.

Dynamic termination can be especially prudent when we want to reason explicitly about the cost of data acquisition during optimization. For example, if this cost were to *vary* across the domain, it would not be sensible to define a budget in terms of function evaluations. However, by accounting for observation costs in the utility function, we can reason about cost–benefit tradeoffs during optimization and seek to terminate whenever the expected cost of further observation outweighs any expected benefit it might provide.

Modeling termination decisions and the optimal policy

We consider a modification to the sequential decision problem we analyzed in the known-budget case, wherein we now allow ourselves to

Figure 5.7: A batch rollout policy as a decision tree. Given a putative value for the next evaluation (x, y) , we design all remaining decisions simultaneously using a batch base policy and take the expectation of the terminal utility with respect to their values.

¹⁷ For estimating a one-dimensional expectation we might take q on the order of roughly 10, but for optimizing a nonconvex acquisition function over the domain we might take n on the order of thousands or more.

limited lookahead as rollout

batch rollout

action space, \mathcal{A} termination option, \emptyset bound on total number of observations, τ_{\max}

¹⁸ It is possible to consider unbounded sequential decision problems, but this is probably not of practical interest in Bayesian optimization:

M. H. DEGROOT (1970). *Optimal Statistical Decisions*. McGraw-Hill. [§ 12.7]

¹⁹ This can be proven through various “information never hurts” (in expectation) results.

terminate optimization at any time of our choosing. Suppose we are at an arbitrary point of optimization and have already obtained data \mathcal{D} . We face the following decision: should we terminate optimization immediately and return \mathcal{D} ? If not, where should we make our next observation?

We model this scenario as a decision problem under uncertainty with an action space equal to the domain \mathcal{X} , representing potential observation locations if we decide to continue, augmented with a special additional action \emptyset representing immediate termination:

$$\mathcal{A} = \mathcal{X} \cup \{\emptyset\}. \quad (5.19)$$

For the sake of analysis, after the termination action has been selected, it is convenient to model the decision process as not actually terminating, but rather continuing with the collapsed action space $\mathcal{A} = \{\emptyset\}$ – once you terminate, there’s no going back.

As before, we may derive the optimal optimization policy in the adaptive termination case via induction on the decision horizon τ . However, we must address one technical issue: the base case of the induction, which analyzes the “final” decision, breaks down if we allow the possibility of a nonterminating sequence of decisions. To sidestep this issue, we assume there is a fixed and known upper bound τ_{\max} on the total number of observations we may make, at which point optimization is compelled to terminate regardless of any other concern. This is not an overly restrictive assumption in the context of Bayesian optimization. Because observations are assumed to be expensive, we can adopt some suitably absurd upper bound without issue; for example, $\tau_{\max} = 1\,000\,000$ would suffice for an overwhelming majority of plausible scenarios.¹⁸

After assuming the decision process is bounded, our previous inductive argument carries through after we demonstrate how to compute the value of the termination action. Fortunately, this is straightforward: termination does not augment our data, and once this action is taken, no other action will ever again be allowed. Therefore the expected marginal gain from termination is always zero:

$$\alpha_\tau(\emptyset; \mathcal{D}) = 0. \quad (5.20)$$

With this, substituting \mathcal{A} for \mathcal{X} in (5.15–5.17) now gives the optimal policy.

Intuitively, the result in (5.20) implies that termination is only the optimal decision if there is no observation offering positive expected gain in utility. For the utility functions described in the next chapter – all of which are agnostic to costs and measure optimization progress alone – reaching this state is actually *impossible*.¹⁹ However, explicitly accounting for observation costs in addition to optimization progress in the utility function resolves this issue, as we will demonstrate.

Example: cost-aware optimization

To illustrate the behavior of a policy allowing early termination, we return to our motivating scenario of accounting for observation costs.

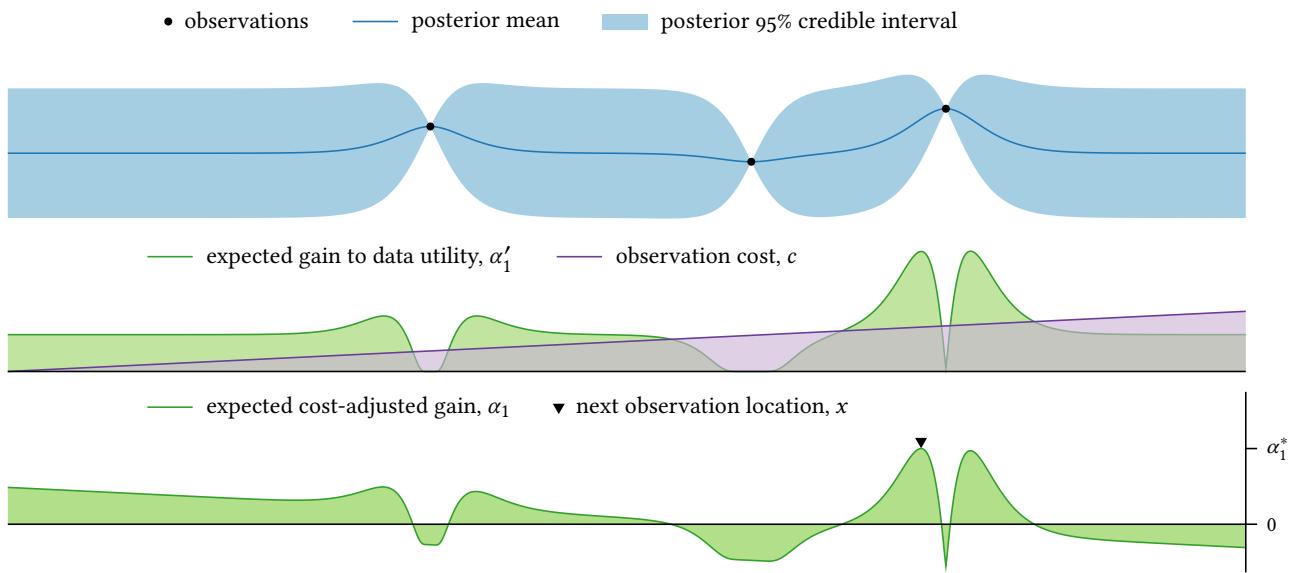


Figure 5.8: Illustration of one-step lookahead with the option to terminate. With a linear utility and additive costs, the expected marginal gain α_1 is the expected marginal gain to the data utility α'_1 adjusted for the cost of acquisition c . For some points, the cost-adjusted expected gain is negative, in which case we would prefer immediate termination to observing there. However, continuing with the chosen point is expected to increase the utility of the current data.

Consider the objective function belief in the top panel of figure 5.8 (which is identical to that from our running example from figures 5.1–5.3) and suppose that the cost of observation now depends on location according to a known cost function $c(x)$,²⁰ illustrated in the middle panel.

If we wish to reason about observation costs in the optimization policy, we must account for them somehow, and the most natural place to do so is in the utility function. Depending on the situation, there are many ways we could proceed;²¹ however, one natural approach is to first select a utility function measuring the quality of a returned dataset alone, ignoring any costs incurred to acquire it. We call this quantity the *data utility* and notate it with $u'(\mathcal{D})$. The data utility is akin to the cost-agnostic utility from the known-budget case, and any one of the options described in the next chapter could reasonably fill this role.

We now adjust the data utility to account for the cost of data acquisition. In many applications, these costs are additive, so that the total cost of gathering a dataset \mathcal{D} is simply

$$c(\mathcal{D}) = \sum_{x \in \mathcal{D}} c(x). \quad (5.21)$$

If the acquisition cost can be expressed in the same units as the data utility – for example, if both can be expressed in monetary terms²² – then we might reasonably evaluate a dataset \mathcal{D} by the cost-adjusted utility:

$$u(\mathcal{D}) = u'(\mathcal{D}) - c(\mathcal{D}). \quad (5.22)$$

²⁰ We will consider unknown and stochastic costs in § 11.1, p. 243.

observation cost function, $c(x)$

²¹ We wish to stress this point – there is considerable flexibility beyond the scheme we describe.

data utility, $u'(\mathcal{D})$
utility functions for optimization: chapter 6,
p. 109

observation costs, $c(\mathcal{D})$

²² Some additional discussion on this natural approach can be found in:

H. RAIFFA and R. SCHLAIFER (1961). *Applied Statistical Decision Theory*. Division of Research, Graduate School of Business Administration, Harvard University. [chapter 4]

Demonstration: one-step lookahead with cost-aware utility

Returning to the scenario in figure 5.8, let us adopt a cost-aware utility function of the above form (5.22) and consider the behavior of a one-step lookahead approximation to the optimal optimization policy.

For these choices, if we were to continue optimization by evaluating at a point x , the resulting one-step marginal gain in utility would be:

$$u(\mathcal{D}_1) - u(\mathcal{D}) = [u'(\mathcal{D}_1) - u'(\mathcal{D})] - c(x),$$

the cost-adjusted marginal gain in the data utility alone. Therefore the expected marginal gain in utility is:

$$\alpha_1(x; \mathcal{D}) = \alpha'_1(x; \mathcal{D}) - c(x),$$

where α'_1 is the one-step expected gain in the data utility (5.8). That is, we simply adjust what would have been the acquisition function in the cost-agnostic setting by subtracting the cost of data acquisition. To prefer evaluating at x to immediate termination, this quantity must have positive expected value (5.20).

The resulting policy is illustrated in figure 5.8. The middle panel shows the cost-agnostic acquisition function α'_1 (from figure 5.1), which is then adjusted for observation cost in the bottom panel. This renders the expected marginal gain negative in some locations, where observations are not expected to be worth their cost. However, in this case there are still regions where observation is favored to termination, and optimization continues at the selected location. Comparing with the cost-agnostic setting in figure 5.1, the optimal observation has shifted from the right-hand side to the left-hand side of the previously best-seen point, as an observation there is more cost effective.

SUMMARY OF MAJOR IDEAS

defining optimization policies via acquisition functions: p. 88

- Optimization policies can be conveniently defined via an *acquisition function* assigning a score to each potential observation location. We then design observations by maximizing the acquisition function (5.1).

- *Bayesian decision theory* is a general framework for optimal decision making under uncertainty, through which we can derive optimal optimization policies and stopping rules.

- The key elements of a decision problem under uncertainty are:

- an *action space* \mathcal{A} , from which we must choose an action a ,
- uncertainty in elements ψ relevant to the decision, represented by a posterior belief $p(\psi | \mathcal{D})$, and
- a *utility function* $u(a, \psi, \mathcal{D})$ quantifying the quality of the action a assuming a given realization of the uncertain elements ψ .

Given these, an optimal decision maximizes the expected utility (5.2–5.3).

- Optimization policy decisions may be cast in this framework by defining a utility function for the data returned by an optimizer, then designing

each observation location to maximize the expected utility with respect to all future data yet to be obtained (5.5–5.6).

- To ensure the optimality of a *sequence* of decisions, we must recursively assume the optimality of all future decisions. This is known as BELLMAN’s *principle of optimality*. Under this assumption, the optimal policy can be derived inductively and assumes a simple recursive form (5.15–5.17).
- The cost of computing the optimal policy grows exponentially with the decision horizon, but several techniques under the umbrella *approximate dynamic programming* provide tractable approximations. Two notable examples are *limited lookahead*, where the decision horizon is artificially limited, and *rollout*, where future decisions are simulated suboptimally.
- Through careful accounting, we may explicitly account for the (possibly nonuniform) cost of data acquisition in the utility function. Offering a termination option and computing the resulting optimal policy then allows us to adaptively terminate optimization when continuing optimization becomes a losing battle of cost versus expected gain.

BELLMAN’s principle of optimality: § 5.2, p. 99

computational burden and approximation of the optimal policy: § 5.3, p. 99

termination as a decision: § 5.4, p. 103

common Bayesian optimization policies:
chapter 7, p. 123

In the next chapter we will discuss several prominent utility functions for measuring the quality of a dataset returned by an optimization procedure. In the following chapter, we will demonstrate how many common acquisition functions for Bayesian optimization may be realized by performing one-step lookahead with these utility functions.

6

UTILITY FUNCTIONS FOR OPTIMIZATION

In the last chapter we introduced Bayesian decision theory, a framework for decision making under uncertainty through which can derive theoretically optimal optimization policies. Central to this approach is the notion of a *utility function* evaluating the quality of a dataset returned from an optimization routine. Given a model of the objective function, conveying our *beliefs* in the face of uncertainty, and a utility function, expressing our *preferences* over outcomes, computing the optimal policy is purely mechanical: we design every observation to maximize the expected utility of the returned dataset (5.15–5.17). Setting aside computational issues, adopting this approach entails only two major decisions: how to build an objective function model consistent with our beliefs and how to design a utility function consistent with our preferences.

Neither of these tasks is trivial! Beliefs and preferences are so innate to the human experience that distilling them down to mathematical symbols can be challenging. Fortunately, expressive and mathematically convenient options for both are readily available. We devoted significant attention to model building in the first part of this book, and we will address the construction of utility functions in this chapter. We will introduce a number of common utility functions designed for optimization, each carrying a different perspective on how optimization performance should be quantified. We hope that the underlying motivation for these utility functions may inspire the design of novel alternatives when called for. In the next chapter, we will demonstrate how approximating the optimal optimization policy corresponding to the utility functions described here yields many widespread Bayesian optimization algorithms.

Although we will be using Gaussian process models in our illustrations throughout the chapter, we will not assume the objective function model is a Gaussian process in our discussion. As in the previous chapters, we will use the notation $\mu_{\mathcal{D}}(x) = \mathbb{E}[\phi | x, \mathcal{D}]$ for the posterior mean of the objective function; this should not be interpreted as implying any particular model structure beyond admitting a posterior mean.

6.1 EXPECTED UTILITY OF TERMINAL RECOMMENDATION

The purpose of optimization is often to explore a space of possibilities in search of the single best alternative, and after investing in optimization, we commit to using some chosen point in a subsequent procedure. In this context, the only purpose of the data collected during optimization is to help select this final point. For example, in hyperparameter tuning, we may evaluate numerous hyperparameters during model development, only to use the apparently best settings found in a production system.

Selecting a point for permanent use represents a *decision*, which we may analyze using Bayesian decision theory. If the sole purpose of optimization is to inform a final decision, it is natural to design the policy to maximize the expected utility of the terminal decision directly, and several popular policies are defined in this manner.

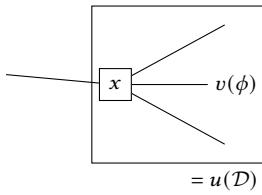
common Bayesian optimization policies: § 7,
p. 123

posterior mean function, $\mu_{\mathcal{D}}$

selecting a point for permanent use

Bayesian decision theory: chapter 5, p. 87

- ¹ Dependence on ϕ alone is not strictly necessary. For example, in the interest of robustness we might wish to ensure that function values are high in the *neighborhood* of our recommendation was well. This would be possible in the same framework by redefining the utility function as desired.



We may also interpret this class of utility functions as augmenting the decision tree in figure 5.4 with a final layer corresponding to the terminal decision. The utility of the data is then the expected utility of this subtree, assuming optimal behavior.

bound on uncertainty

² M. A. OSBORNE et al. (2009). Gaussian Processes for Global Optimization. *LION* 3.

Formalization of terminal recommendation decision

Suppose we have run an optimization routine, which returned a dataset $\mathcal{D} = (\mathbf{x}, \mathbf{y})$, and suppose we now wish to recommend a point $x \in \mathcal{X}$ for use in some task, with performance determined by the underlying objective function value $\phi = f(x)$.¹ This represents a decision under uncertainty about ϕ , informed by the predictive distribution, $p(\phi | x, \mathcal{D})$.

To completely specify the decision problem, we must identify an action space $\mathcal{A} \subseteq \mathcal{X}$ for our recommendation and a utility function $v(\phi)$ evaluating a recommendation in hindsight according to its objective value ϕ . Given these, a rational recommendation maximizes the expected utility:

$$x \in \arg \max_{x' \in \mathcal{A}} \mathbb{E}[v(\phi') | x', \mathcal{D}].$$

The expected utility of the optimal recommendation only depends on the data returned by the optimizer; it does not depend on the optimal recommendation x (due to maximization) nor its objective value ϕ (due to expectation). This suggests a natural utility for use in optimization: the expected quality of an optimal terminal recommendation given the data,

$$u(\mathcal{D}) = \max_{x' \in \mathcal{A}} \mathbb{E}[v(\phi') | x', \mathcal{D}]. \quad (6.1)$$

In the context of the sequential decision tree from figure 5.4, this utility function effectively “collapses” the expected utility of a final decision into a utility for the returned data; see the illustration in the margin. We are free to select the action space and utility function for the final recommendation as we see fit; we provide some advice below.

Choosing an action space

We begin with the action space $\mathcal{A} \subseteq \mathcal{X}$. One extreme option is to restrict our choice to only the visited points \mathbf{x} . This ensures at least some knowledge of the objective function at the recommended point, which may be prudent when the objective function model may be misspecified. The other extreme is the maximally permissive alternative: the entire domain \mathcal{X} , allowing us to recommend any point, including those arbitrarily far from our observations. The wisdom of recommending an unvisited point for perpetual use is ultimately a question of faith in the model’s beliefs.

Compromises between these extremes have also been occasionally suggested in the literature. OSBORNE et al. for example proposed restricting the choice of final recommendation to only those points where the objective function is known with acceptable tolerance.² Such a scheme can limit unwanted surprise from recommending points where the objective function value is not known with sufficient certainty. One might accomplish this in several ways; OSBORNE et al. adopted a parametric, data-dependent action space of the form

$$\mathcal{A}(\varepsilon; \mathcal{D}) = \{x \mid \text{std}[\phi | x, \mathcal{D}] \leq \varepsilon\},$$

where ε is a threshold specifying the largest acceptable uncertainty.

Choosing a utility function and risk tolerance

In addition to selecting an action space, we must also select a utility function $v(\phi)$ evaluating a recommendation at x in light of the corresponding function value ϕ . As our focus is on maximization (1.1), it is clear that the utility should be monotonically increasing in ϕ , but it is not necessarily clear what shape this function should assume. The answer depends on our *risk tolerance*, a concept demonstrated in the margin. When making our final recommendation, we may wish to consider not only the *expected* function value of a given point but also our *uncertainty* in this value, as points with greater uncertainty may result in more surprising and potentially disappointing results.

By controlling the shape of the utility function $v(\phi)$, we may induce different behavior with respect to risk. The simplest and most common option encountered in Bayesian optimization is a linear utility:

$$v(\phi) = \phi. \quad (6.2)$$

In this case, the expected utility from recommending x is simply the posterior mean of ϕ , as we have already seen (5.4):

$$\mathbb{E}[v(\phi) | x, \mathcal{D}] = \mu_{\mathcal{D}}(x),$$

and an optimal recommendation maximizes the posterior mean over the action space:

$$x = \arg \max_{x' \in \mathcal{A}} \mu_{\mathcal{D}}(x').$$

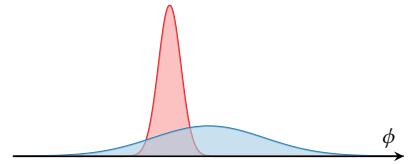
Uncertainty in the objective function is not considered in this decision at all! Rather, we are indifferent between points with equal expected value, regardless of their uncertainty – that is, we are *risk neutral*.

Risk neutrality is computationally convenient due to the simple form of the expected utility, but may not always reflect our true preferences. In the margin we show beliefs over the objective values for two potential recommendations with equal expected value but significantly different risk. In many scenarios we would have a clear preference between the two alternatives, but a risk-neutral utility induces complete indifference.

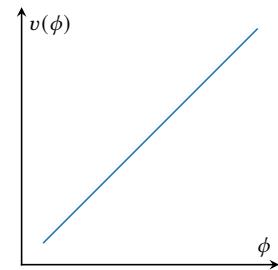
A useful concept when reasoning about risk preferences is the so-called *certainty equivalent*. Consider a risky potential recommendation x , that is, a point for which we do not know the objective value exactly. The certainty equivalent for x is the value of a hypothetical *risk-free* alternative for which our preferences would be indifferent. That is, the certainty equivalent for x corresponds to an objective function value ϕ' such that

$$v(\phi') = \mathbb{E}[v(\phi) | x, \mathcal{D}].$$

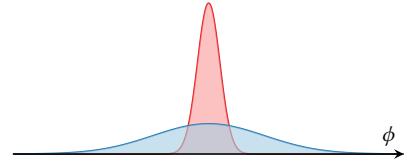
Under a risk-neutral utility function, the certainty equivalent of a point x is simply its expected value: $\phi' = \mu_{\mathcal{D}}(x)$. Thus we would abandon a potential recommendation for another only if it had greater expected value, independent of risk. However, we may encode risk-aware preferences with appropriately designed *nonlinear* utility functions.



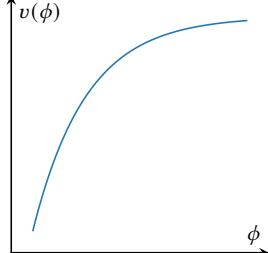
Consider the illustrated beliefs about the objective function value corresponding to two possible recommendations. The blue option has a higher expected value, but also greater uncertainty, and proposing it entails some risk. The red alternative has a lower expected value but is perhaps a safer option. A risk-averse agent might prefer the red point, whereas a risk-tolerant agent might prefer the blue point.



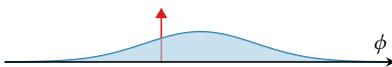
A risk-neutral (linear) utility function.



Beliefs over two recommendations with equal expected value. A risk-neutral agent would be indifferent between these alternatives, a risk-averse agent would prefer the red option, and a risk-seeking agent would prefer the blue option.



A risk-averse (concave) utility function.



A risk-averse agent may be indifferent between the risky recommendation in blue and its risk-free certainty equivalent with lower expected value in red.

If our preferences indicate *risk aversion*, we might be willing to recommend a point with lower expected value if it also entailed less risk. We may induce risk-averse preferences by adopting a utility function that is a concave function of the objective value. In this case, by Jensen's inequality we have

$$v(\phi') = \mathbb{E}[v(\phi) | x, \mathcal{D}] \leq v(\mathbb{E}[\phi | x, \mathcal{D}]) = v(\mu_{\mathcal{D}}(x)),$$

and thus the certainty equivalent of a risky recommendation is *less* than its expected value; see the example in the margin. Similarly, we may induce *risk-seeking* preferences with a convex utility function, in which case the certainty equivalent of a risky recommendation is *greater* than its expected value – our preferences encode an inclination toward gambling. Risk-averse and risk-seeking utilities are rarely encountered in the Bayesian optimization literature; however, they may be preferable in some practical settings, as risk neutrality is often questionable.

Numerous risk-averse utility functions have been proposed in the economics and decision theory literature,³ and a full discussion is beyond the scope of this book. However, one natural approach is to quantify the risk associated with recommending an uncertain value ϕ by its standard deviation:

$$\sigma = \text{std}[\phi | x, \mathcal{D}].$$

Now we may establish preferences over potential recommendations consistent with⁴ a weighted combination of a point x 's expected reward, $\mu = \mu_{\mathcal{D}}(x)$, and its risk, σ :⁵

$$\mu + \beta\sigma.$$

Here β serves as a tunable risk-tolerance parameter: values $\beta < 0$ penalize risk and induce risk-averse behavior, values $\beta > 0$ reward risk and induce risk-seeking behavior, and $\beta = 0$ induces risk neutrality (6.2).

Two particular utility functions from this general framework are widely encountered in Bayesian optimization, both representing the expected utility of a risk-neutral optimal terminal recommendation.

Simple reward

Suppose an optimization routine returned data $\mathcal{D} = (\mathbf{x}, \mathbf{y})$ to inform a terminal recommendation, and that we will make this decision using the risk-neutral utility function $v(\phi) = \phi$ (6.2). If we limit the action space of this recommendation to only the locations evaluated during optimization \mathbf{x} , the expected utility of the optimal recommendation is the so-called *simple reward*:^{6,7}

$$u(\mathcal{D}) = \max \mu_{\mathcal{D}}(\mathbf{x}). \quad (6.3)$$

In the special case of exact observations, where $\mathbf{y} = f(\mathbf{x}) = \phi$, the simple reward reduces to the maximal function value encountered during optimization:

$$u(\mathcal{D}) = \max \phi. \quad (6.4)$$

⁶ This name contrasts with the *cumulative reward*: § 6.2, p. 114.

⁷ One technical caveat is in order: when the dataset is empty, the maximum degenerates and we have $u(\emptyset) = -\infty$.

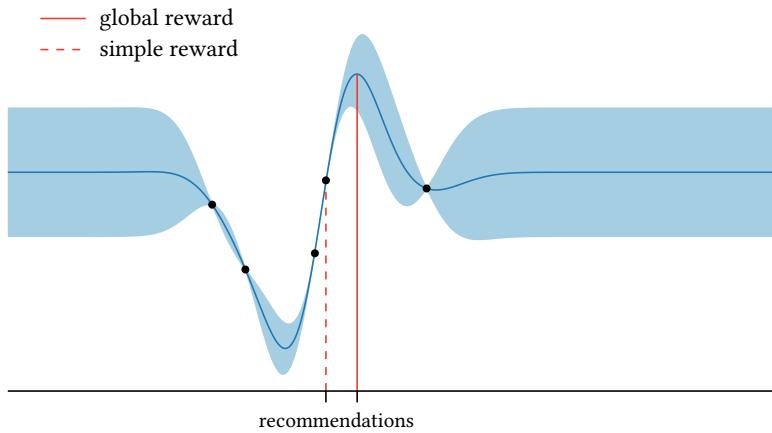


Figure 6.1: The terminal recommendations corresponding to the simple reward and global reward for an example dataset comprising five observations. The prior distribution for the objective for this demonstration is illustrated in figure 6.3.

One-step lookahead with the simple reward utility function produces a widely-used acquisition function known as *expected improvement*, which we will discuss in detail in the next two chapters.

expected improvement: § 7.3 p. 127

Global reward

Another prominent utility is the *global reward*.⁸ Here we again consider a risk-neutral terminal recommendation, but now expand the action space for this recommendation to the entire domain \mathcal{X} . The expected utility of this recommendation is the global maximum of the posterior mean:

$$u(\mathcal{D}) = \max_{x \in \mathcal{X}} \mu_{\mathcal{D}}(x). \quad (6.5)$$

An example dataset exhibiting a large discrepancy between the simple reward (6.3) and global reward (6.5) utilities is illustrated in figure 6.1. The larger action space underlying global reward leads to a markedly different and somewhat riskier recommendation.

⁸ “Global simple reward” would be a more accurate (but annoyingly bulky) name.

One-step lookahead with global reward (6.5) yields the *knowledge gradient* acquisition function, which we will also consider at length in the following chapters.

knowledge gradient: § 7.4 p. 129

A tempting, but nonsensical alternative

There is an alternative utility deceptively similar to the simple reward that is sometimes encountered in the Bayesian optimization literature, namely the maximum *noisy* observed value contained in the dataset:

$$u(\mathcal{D}) \stackrel{?}{=} \max \mathbf{y}. \quad (6.6)$$

In the case of *exact* observations of the objective function, this value coincides with the simple reward (6.4), which has a natural interpretation as the expected utility of a particular optimal terminal recommendation. However, this correspondence does *not* hold in the case of inexact or noisy observations, and the proposed utility is rendered absurd.

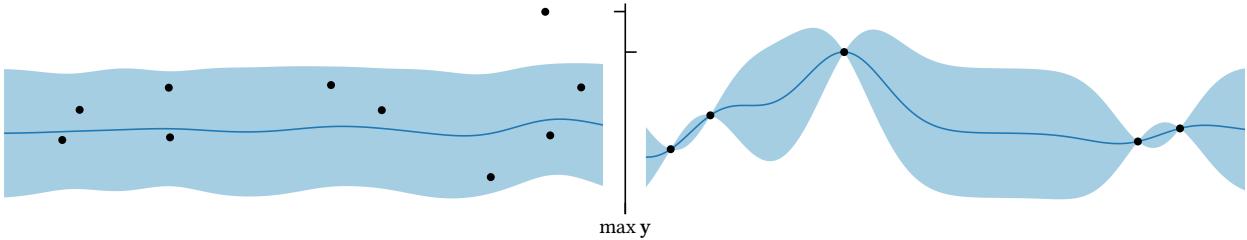


Figure 6.2: The utility $u(\mathcal{D}) = \max y$ would prefer the excessively noisy dataset on the left to the less-noisy dataset on the right with smaller maximum value. The data on the left reveal little information about the objective function, and the maximum observed value is very likely to be an outlier, whereas the data on the right indicate reasonable progress.

large-but-noisy observations are not necessarily preferable

example and discussion

approximation to the simple reward

expected improvement: § 7.3 p. 127

This is simple to demonstrate by contemplating the preferences over outcomes encoded in the utility, which may not align with intuition. This disparity is especially notable in situations with excessively noisy observations, where the maximum value observed will likely reflect spurious noise rather than actual optimization progress.

Figure 6.2 shows an extreme but illustrative example. We consider two optimization outcomes over the same domain, one with excessively noisy observations and the other with exact measurements. The noisy dataset contains a large observation on the right-hand side of the domain, but this is almost certainly the result of noise, as indicated by the objective function posterior. Although the other dataset has a lower maximal value, the observations are more trustworthy and represent a plainly better outcome. But the proposed utility (6.6) prefers the noisier dataset! On the other hand, both the simple and global reward utilities prefer the noiseless dataset, as the data produce a larger effect on the posterior mean – and thus yield more promising recommendations.

Of course, errors in noisy measurements are not always as extreme as in this example. When the signal-to-noise ratio is relatively high, the utility (6.6) can serve as a reasonable *approximation* to the simple reward. We will discuss this approximation scheme further in the context of expected improvement.

6.2 CUMULATIVE REWARD

Simple and global reward are motivated by supposing that the goal of optimization is to discover the best *single* point from a space of alternatives. To this end, we evaluate data according to the highest function value revealed and assume that the values of any suboptimal points encountered are irrelevant.

In other settings, the value of *every* individual observation might be significant, for example, if the optimization procedure is controlling a critical external system. If the consequences of these decisions are nontrivial, we might wish to discourage observing where we might encounter unexpectedly low objective function values.

Cumulative reward encourages obtaining observations with large *average* value. For a dataset $\mathcal{D} = (\mathbf{x}, \mathbf{y})$, its cumulative reward is simply the sum of the observed values:

$$u(\mathcal{D}) = \sum_i y_i. \quad (6.7)$$

One notable use of cumulative reward is in *active search*, a simple mathematical model of scientific discovery. Here, we successively select points for investigation seeking novel members of a rare, valuable class $\mathcal{V} \subset \mathcal{X}$. Observing at a point $x \in \mathcal{X}$ yields a *binary* observation indicating membership in the desired class: $y = [x \in \mathcal{V}]$. Most studies of active search seek to maximize the cumulative reward (6.7) of the gathered data, hoping to discover as many valuable items as possible.

active search: § 11.11, p. 277

6.3 INFORMATION GAIN

Simple, global, and cumulative reward judge optimization performance based solely on having found high objective function values, a natural and pragmatic concern. *Information theory*⁹ provides an alternative approach to measuring utility that has received significant use in Bayesian optimization. An information-theoretic approach to sequential decision making (including optimization) identifies some random variable that we wish to learn about through our observations. We then evaluate performance by quantifying the amount of information about this random variable revealed by data, favoring datasets containing more information. This line of reasoning gives rise to the notion of *information gain*.

Let ω be a random variable of interest that we wish to determine through the observation of data. The choice of ω is open-ended and should be guided by the application at hand. Natural choices aligned with optimization include the location of the global optimum, x^* , and the maximal value of the objective, f^* (1.1), each of which has been considered in depth in this context.

We may quantify our initial uncertainty about ω via the (differential) *entropy* of its prior distribution, $p(\omega)$:

$$H[\omega] = - \int p(\omega) \log p(\omega) d\omega.$$

The *information gain* offered by a dataset \mathcal{D} is then the reduction in entropy when moving from the prior to the posterior distribution:

$$u(\mathcal{D}) = H[\omega] - H[\omega | \mathcal{D}], \quad (6.8)$$

where $H[\omega | \mathcal{D}]$ is the differential entropy of the posterior:¹⁰

$$H[\omega | \mathcal{D}] = - \int p(\omega | \mathcal{D}) \log p(\omega | \mathcal{D}) d\omega.$$

Somewhat confusingly, some authors use an alternative definition of information gain – the *Kullback–Leibler (KL) divergence* between the

⁹ A broad introduction to information theory is provided by the classical text

T. M. COVER and J. A. THOMAS (2006). *Elements of Information Theory*. John Wiley & Sons,

and a treatment focusing on the connections to Bayesian inference can be found in

D. J. C. MACKAY (2003). *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press.

entropy

¹⁰ A caveat is in order regarding this notation, which is not standard. In information theory $H[\omega | \mathcal{D}]$ denotes the *conditional entropy* of ω given \mathcal{D} , which is the expectation of the given quantity over the observed values \mathbf{y} . For our purposes it will more useful for this to signify the differential entropy of the notationally parallel posterior $p(\omega | \mathcal{D})$. When needed, we will write conditional entropy with an explicit expectation: $\mathbb{E}[H[\omega | \mathcal{D}] | \mathbf{x}]$.

Kullback–Leibler (KL) divergence

¹¹ A simple example: suppose $\omega \in (0, 1)$ is the unknown bias of a coin, with prior

$$p(\omega) = \text{Beta}(\omega; 2, 1); \quad H \approx -0.193.$$

After flipping and observing “tails,” the posterior becomes

$$p(\omega | \mathcal{D}) = \text{Beta}(\omega; 2, 2); \quad H \approx -0.125.$$

The information “gained” was

$$H[\omega] - H[\omega | \mathcal{D}] \approx -0.068 < 0.$$

Of course, the most likely outcome of the flip a priori was “heads,” so the outcome was surprising. Indeed the *expected* information gain before the experiment was

$$H[\omega] - \mathbb{E}[H[\omega | \mathcal{D}]] \approx 0.137 > 0.$$

¹² See p. 138 for a proof.

mutual information, entropy search: § 7.6
p. 135

information-theoretic policies as the scientific method: § 7.6, p. 136

model averaging: §§ 4.4–4.5, p. 74
model-agnostic alternatives

¹³ The effect on simple and global reward is to maximize a model-marginal posterior mean, and the effect on information gain is to evaluate changes in model-marginal beliefs about ω .

posterior distribution and the prior distribution:

$$u(\mathcal{D}) = D_{\text{KL}}[p(\omega | \mathcal{D}) \| p(\omega)] = \int p(\omega | \mathcal{D}) \log \frac{p(\omega | \mathcal{D})}{p(\omega)} d\omega. \quad (6.9)$$

That is, we quantify the information contained in data by how much our belief in the ω changes as a result of collecting it. This definition has some convenient properties compared to the previous one (6.8); namely, the expression in (6.9) is invariant to reparametrization of ω and always nonnegative, whereas “surprising” observations may cause the information gain in (6.8) to become negative.¹¹ However, the previous definition as the direct reduction in entropy may be more intuitive.

Fortunately (and perhaps surprisingly!), there is a strong connection between these two “information gains” (6.8–6.9) in the context of sequential decision making. Namely, their expected values with respect to observed values are equal, and thus maximizing expected utility with either leads to identical decisions.¹² For this reason, the reader may simply choose whichever definition they find more intuitive.

One-step lookahead with (either) information gain yields an acquisition function known as *mutual information*. This is the basis for a family of related Bayesian optimization procedures sharing the moniker *entropy search*, which we will discuss further in the following chapters.

Unlike the other utility functions discussed thus far, information gain is not intimately linked to optimization, and may be adapted to a wide variety of tasks by selecting the random variable ω appropriately. Rather, this scheme of refining knowledge through experiment is effectively a mathematical formulation of scientific inquiry.

6.4 DEPENDENCE ON MODEL OF OBJECTIVE FUNCTION

One striking feature of most of the utility functions defined in this chapter is implicit dependence on an underlying model of the objective function. Both the simple and global reward are defined in terms of the posterior mean function $\mu_{\mathcal{D}}$, and information gain about the location or value of the optimum is defined in terms of the posterior belief about these values, $p(x^*, f^* | \mathcal{D})$; both of these quantities are byproducts of the objective function posterior.

One way to mitigate model dependence in the computation of utility is via model averaging (4.11, 4.23).¹³ We may also attempt to define purely model-agnostic utility functions in terms of the data alone, without reference to a model; however, the possibilities are somewhat limited if we wish the resulting utility to be sensible. Cumulative reward (6.2) is one example, as it depends only on the observed values y . The maximum function value observed is another possibility (6.6), but, as we have shown, it is dubious when observations are corrupted by noise. Other similarly defined alternatives may suffer the same fate – for additive noise with zero mean, the expected contribution from noise to the cumulative reward is zero; however, noise will bias many other natural measures such as order statistics (including the maximum) of the observations.

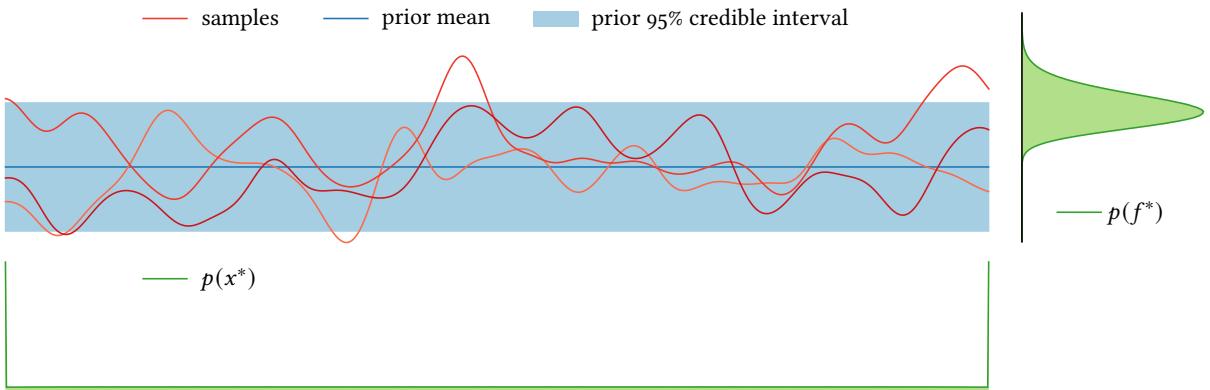


Figure 6.3: The objective function prior used throughout our utility function comparison. Marginal beliefs of function values are shown, as well as the induced beliefs over the location of the global optimum, $p(x^*)$, and the value of the global optimum, $p(f^*)$. Note that there is a significant probability that the global optimum is achieved on the boundary of the domain, reflected by large point masses.

6.5 COMPARISON OF UTILITY FUNCTIONS

We have now presented several utility functions for evaluating a dataset returned by an optimization routine. Each utility quantifies progress on our model optimization problem (1.1) in some way, but it may be difficult at this point to appreciate their, sometimes subtle, differences in approach. Here we will present and discuss example datasets for which different utility functions diverge in their opinion of quality.

We particularly wish to contrast the behavior of the simple reward (6.3) with other utility functions. Simple reward is probably the most prevalent utility in the Bayesian optimization literature (especially in applications), as it corresponds to the widespread expected improvement acquisition function. A distinguishing feature of simple reward is that it evaluates data based only on *local* properties of the objective function posterior. This locality is both computationally convenient and pragmatic. Simple reward is derived from the premise that we will be recommending one of the points observed during the course of optimization for permanent use, and thus it is sensible to judge performance based on the objective function values at the observed locations alone.

Several alternatives instead measure *global* properties of the objective function posterior. The global reward (6.5), for example, considers the entire posterior mean function, reflecting a willingness to recommend an unevaluated point after termination. Information gain (6.8) about the location or value of the optimum considers the posterior entropy of these quantities, again a global property. The consequences of reasoning about local or global properties of the posterior can sometimes lead to significant disagreement between the simple reward and other utilities.

In the following examples, we consider optimization on an interval with exact measurements. We model the objective function with a Gaussian process with constant mean (3.1) and squared exponential

local vs. global properties of posterior

expected improvement: § 7.3 p. 127

model of objective function

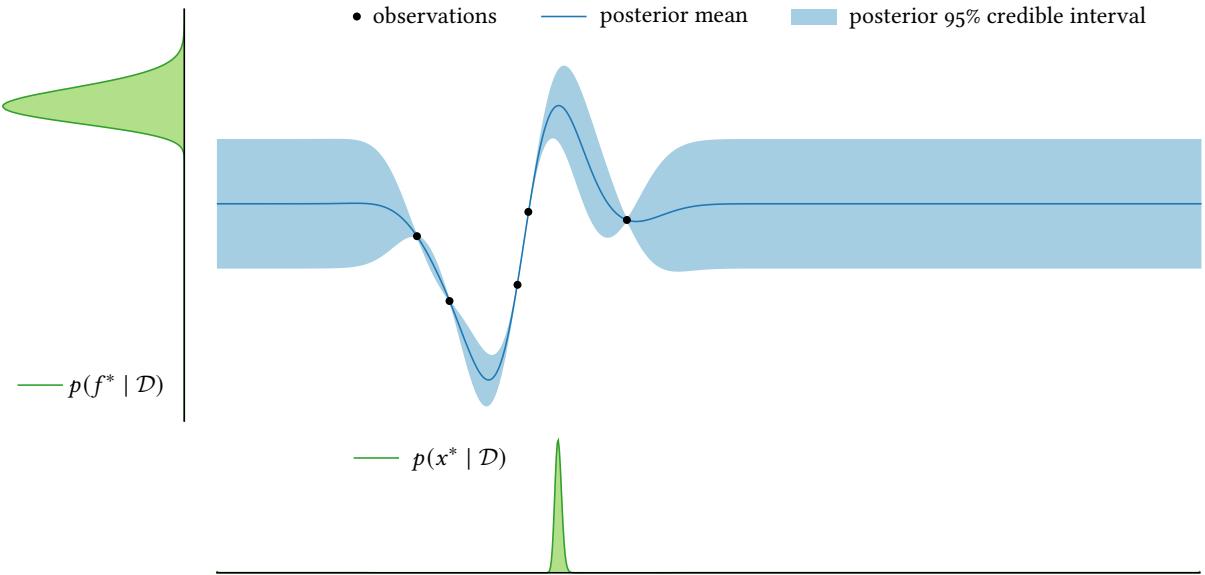


Figure 6.4: An example dataset of five observations and the resulting posterior belief of the objective function. This dataset exhibits relatively low simple reward (6.3) but relatively high global reward (6.5) and information gain (6.8) about the location x^* and value f^* of the optimum.

¹⁴ For this model, a unique optimum will exist with probability one; see § 2.7, p. 34 for more details.

covariance (3.12). This prior is illustrated in figure 6.3, along with the induced beliefs about the location x^* and value f^* of the global optimum. Both distributions reflect considerable uncertainty.¹⁴ We will examine two datasets that might be returned by an optimizer using this model and discuss how different utility functions would evaluate these outcomes.

Good global outcome but poor local outcome

Consider the dataset in figure 6.4 and the resulting posterior belief about the objective and its optimum. In this example, the simple reward is relatively low as the posterior mean at our observations is unremarkable. In fact, every observation was lower than the prior mean, a seemingly unlucky outcome. However, the global reward is relatively high: the data imply a steep derivative in one location, inducing high values of the posterior mean away from our data. This is a significant accomplishment from the point of view of the global reward, as the model expects a terminal recommendation in that region to be especially valuable.

Figure 6.1 shows the optimal final recommendations associated with these two utility functions. The simple reward recommendation prioritizes safety over reward, whereas the global reward recommendation reflects more risk tolerance. Neither is inherently better: although the global reward recommendation has a larger expected value, this expectation is computed using a model that might be mistaken. Further, comparing the posterior distribution in figure 6.4 with the prior in figure

low simple reward
high global reward
final recommendations
high information gain

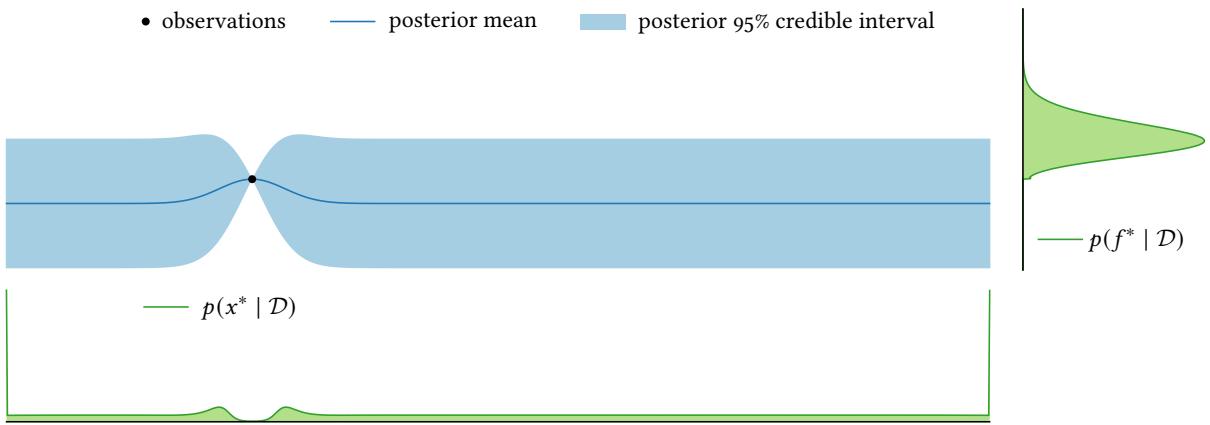


Figure 6.5: An example dataset containing a single observation and the resulting posterior belief of the objective function. This dataset exhibits relatively high simple reward (6.3) but relatively low global reward (6.5) and information gain (6.8) about the location x^* and value f^* of the optimum.

6.3, we see this example dataset also induces a significant reduction in our uncertainty about both the location and value of the global optimum, despite not containing any particularly notable values itself. Therefore, despite a somewhat low simple reward, observing this dataset results in relatively high information gain about these quantities.

Good local outcome but poor global outcome

We illustrate a different example dataset in figure 6.5. The dataset contains a single observation in the center of the domain with value somewhat higher than the prior mean. Although this dataset may not appear particularly impressive, its simple reward is higher than the previous dataset, as this observation exceeds every value seen in that scenario.

However, this dataset has lower value than the previous dataset when evaluated other utility functions. Its global reward is lower than in the first scenario, as the global maximum of the posterior mean is lower. This can be verified by visual inspection of figures 6.4 and 6.5, whose vertical axes are compatible. Further, the single observation in this scenario provides nearly no information regarding the location nor the value of the global maximum. The observation of a moderately high value provides only weak evidence that the global optimum may be located nearby, barely influencing our posterior belief about x^* . The observation does truncate our belief of the value of the global optimum f^* , but only rules out a relatively small portion of its lower tail.

high simple reward

low global reward

low information gain

SUMMARY OF MAJOR IDEAS

In Bayesian decision theory, preferences over outcomes are encoded by a *utility function*, which in the context of optimization policy design,

¹⁵ Just like human taste, there is no right or wrong when it comes to preferences, at least not over *certain* outcomes. The von Neumann–Morgenstern theorem mentioned on p. 90 entails rationality axioms, but these only apply to preferences over *uncertain* outcomes.

expected utility of terminal recommendation:
§ 6.1, p. 109

risk tolerance: § 6.1, p. 111

simple reward: § 6.1, p. 112
global reward: § 6.1, p. 113

cumulative reward: § 6.2, p. 114

information gain: § 6.3, p. 115

comparison of utility functions: § 6.5, p. 117

assesses the quality of data returned by an optimization routine, $u(\mathcal{D})$. The optimization policy then seeks to design observations to maximize the *expected* utility of the returned data. The general theory presented in the last chapter makes no assumptions regarding the utility function.¹⁵ However, in the context of optimization, some utility functions are particularly easy to motivate.

- In many cases there is a decision following optimization in which we must recommend a single point in the domain for perpetual use. In this case, it is sensible to define an optimization utility function in terms of the expected utility of the optimal terminal recommendation informed by the returned data. This requires fully specifying that terminal recommendation, including its action space and utility function, after which we may “pass through” the optimal expected utility (6.1).
- When designing a terminal recommendation – especially when we may recommend points with residual uncertainty in their underlying objective value – it may be prudent to consider our *risk tolerance*. Careful design of the terminal utility allows for us to tune our appetite for risk, in terms of trading off the a point’s expected value against its uncertainty. Most utilities encountered in Bayesian optimization are *risk neutral*, but this need not necessarily be the case.
- Two notable realizations of this scheme are *simple reward* (6.3) and *global reward* (6.5), both of which represent the expected utility of an optimal terminal recommendation with a risk-neutral utility. The action space for simple reward is the points visited during optimization, and the action space for global reward is the entire domain.
- The simple reward simplifies when observations are exact (6.4).
- An alternative to the simple reward is the *cumulative reward* (6.7), which evaluates a dataset based on the *average*, rather than maximum, value observed. This does not see too much direct use in policy design, but is an important concept for the analysis of algorithms.
- *Information gain* provides an information-theoretic approach to quantifying the value of data in terms of the information provided by the data regarding some quantity of interest. This can be quantified by either measuring the reduction in differential entropy moving from the prior to the posterior (6.8) or by the KL divergence between the posterior and prior (6.9) – either induces the same one-step lookahead policy.
- In the context of optimization, information gain regarding either the location x^* or value f^* of the global optimum (1.1) are judicious realizations of this general approach to utility design.
- An important feature distinguishing simple reward from most other utility functions is its dependence on the posterior belief at the observed locations *alone*, rather than the posterior belief over the entire objective function. Even in relatively simple examples, this may lead to disagreement between simple reward and other utility functions in judging the quality of a given dataset.

The utility functions presented in this chapter form the backbone of the most popular Bayesian optimization algorithms. In particular, many common policies are realized by maximizing the one-step expected marginal gain to one of these utilities, as we will show in the next chapter.

one-step lookahead: § 5.3, p. 101

7

COMMON BAYESIAN OPTIMIZATION POLICIES

The heart of an optimization routine is its policy, which sequentially designs each observation in light of available data.¹ In the Bayesian approach to optimization, policies are designed with reference to a probabilistic belief about the objective function, with this belief guiding the policy in making decisions likely to yield beneficial outcomes. Numerous Bayesian optimization policies have been proposed in the literature, many of which enjoy widespread use. In this chapter we will present an overview of popular Bayesian optimization policies and emphasize common themes in their construction. In the next chapter we will provide explicit computational details for implementing these policies with Gaussian process models of the objective function.

Nearly all Bayesian optimization algorithms result from one of two primary approaches to policy design. The most popular is *Bayesian decision theory*, the focus of the previous two chapters. In chapter 5 we introduced Bayesian decision theory as a general framework for deriving optimal, but computationally prohibitive, optimization policies. In this chapter, we will apply the ideas underlying these optimal procedures to realize computationally tractable and practically useful policies. We will see that a majority of popular Bayesian optimization algorithms can be interpreted in a uniform manner as performing one-step lookahead for some underlying utility function.

Another avenue for policy design is to adopt algorithms for *multi-armed bandits* to the optimization setting. A multi-armed bandit is a finite-dimensional model of sequential optimization with noisy observations. We consider an agent faced with a finite set of alternatives (“arms”), who is compelled to select a sequence of items from this set. Choosing a given item yields a stochastic reward drawn from an unknown distribution associated with that arm. We seek a sequential policy for selecting arms maximizing the expected cumulative reward (6.2).²

Multi-armed bandits have seen decades of sustained study, and some policies have strong theoretical guarantees on their performance, suggesting these policies may also be useful for optimization. To this end, we may model optimization as an *infinite-armed bandit*, where each point in the domain $x \in \mathcal{X}$ represents an arm with uncertain reward depending on the objective function value $\phi = f(x)$. Our belief about the objective function then provides a mechanism to reason about these rewards and derive a policy. This analogy has inspired several Bayesian optimization policies, many of which enjoy strong performance guarantees.

A central concern in bandit problems is the *exploration-exploitation dilemma*: we must repeatedly decide whether to allocate resources to an arm already known to yield high reward (“exploitation”) or to an arm with uncertain reward to learn about its reward distribution (“exploration”). Exploitation may yield a high instantaneous reward, but exploration may provide valuable information for improving future rewards. This tradeoff between instant payoff and learning for the future has been called “a conflict evident in all human action.”³ A similar choice is faced

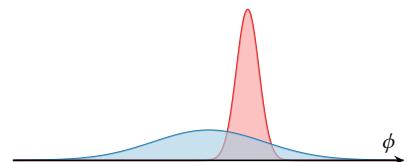
¹ The reader may wish to recall our model optimization procedure: Algorithm 1.1, p. 3.

Bayesian decision theory: chapter 5, p. 87

utility functions for optimization: chapter 6, p. 109

multi-armed bandits

² The name references a gambler contemplating how to allocate their bankroll among a wall of slot machines. Slot machines are known as “one-armed bandits” in American vernacular, as they eventually steal all your money.



Exploration vs. exploitation. We show reward distributions for two possible options. The blue option returns higher expected reward, but the green option reflects more uncertainty and may actually be superior. Which should we prefer?

³ P. WHITTLE (1982). *Optimization Over Time: Dynamic Programming and Stochastic Control*. Vol. 1. John Wiley & Sons.

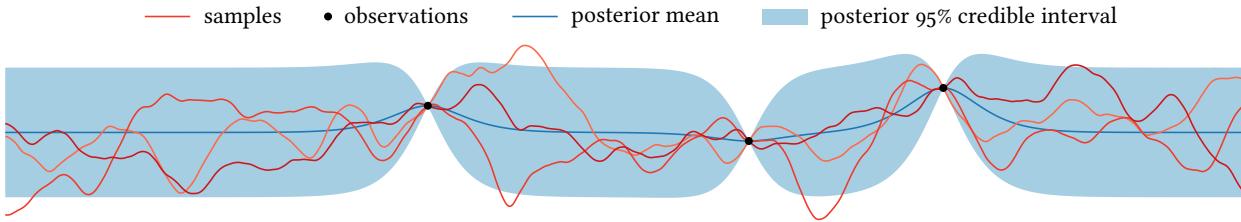


Figure 7.1: The scenario we will consider for illustrating optimization policies. The objective function prior is a Gaussian process with constant mean and Mátern covariance with $\nu = 5/2$ (3.14). We show the marginal predictive distributions and three samples from the posterior conditioned on the indicated observations.

throughout optimization, as we must continually decide whether to focus on a suspected local maximum (exploitation) or to explore unknown regions of the domain seeking new maxima (exploration). We will see that typical Bayesian optimization policies reflect consideration of this dilemma in some way, whether by explicit design or as a consequence of decision-theoretic reasoning.

Before diving into policy design, we pause to introduce a running example we will carry through the chapter and notation to facilitate our discussion. We will then derive a series of policies stemming from Bayesian decision theory, and finally consider bandit-inspired algorithms.

7.1 EXAMPLE OPTIMIZATION SCENARIO

⁴ Take note of the legend; it will not be repeated.

Gaussian processes: chapter 2, p. 15

objective function for simulation

utility functions for optimization: chapter 6,
p. 109

optimal policies: § 5.2, p. 91

Throughout this chapter we will demonstrate the behavior of optimization policies on an example scenario illustrated in figure 7.1.⁴ We consider a one-dimensional objective function observed without noise and adopt a Gaussian process prior belief about this function. The prior mean function is constant (3.1), and the prior covariance function is a Mátern covariance with $\nu = 5/2$ (3.14). The parameters are fixed so that the domain spans exactly 30 length scales. We condition this prior on three observations, inducing two local maxima in the posterior mean and a range of marginal predictive uncertainty.

We will illustrate the behavior of policies by simulating optimization to design a sequence of additional observations for this running example. The ground truth objective function we will use for these simulations is shown in figure 7.2, and was drawn from the corresponding model. The objective features numerous undiscovered local maxima and exhibits an unusually high global maximum on the left-hand side of the domain.

7.2 DECISION-THEORETIC POLICIES

Central to decision-theoretic optimization is a utility function $u(\mathcal{D})$ measuring the quality of a dataset returned by an optimizer. After selecting a utility function and a model of the objective function and our observations, we may design each observation to maximize the expected utility

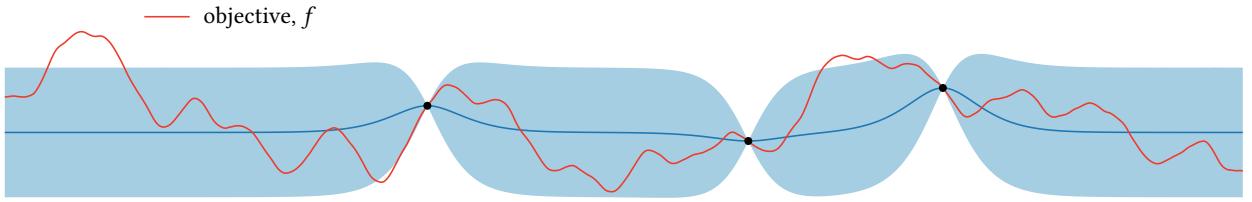


Figure 7.2: The true objective function used for simulating optimization policies.

of the returned data (5.16). This policy is optimal in the average case: it maximizes the expected utility of the returned dataset over the space of all possible policies.⁵ Unfortunately, optimality comes at a great cost. Computing the optimal policy requires recursive simulation of the entire remainder of optimization, a random process due to uncertainty in the outcomes of our observations. In general, the cost of computing the optimal policy grows exponentially with the horizon, the number of observations remaining before termination.

However, the structure of the optimal policy suggests a natural family of *lookahead* approximations based on fixing a computationally tractable maximum horizon throughout optimization. This line of reasoning has lead to many of the practical policies available for Bayesian optimization. In fact, most popular algorithms represent *one-step lookahead*, where in each iteration we greedily maximize the expected utility after obtaining only a single additional observation. Although these policies are maximally myopic, they are also maximally efficient among lookahead approximations and have delivered impressive empirical performance in a wide range of settings.

It may seem surprising that such dramatically myopic policies have any use at all. There is a huge difference between the scale of reasoning in one-step lookahead compared with the optimal procedure, which may consider hundreds of future decisions or more when designing an observation. However, the situation is somewhat more nuanced than it might appear. In a seminal paper, KUSHNER argued that myopic policies may in fact show *better* empirical performance than a theoretically optimal policy, and his argument remains convincing:⁶

Since a mathematical model of $[f]$ is available, it is theoretically possible, once a criterion of optimality is given, to determine the mathematically optimal sampling policy. However...determination of the optimum sampling policies is extremely difficult. Because of this, the development of our sampling laws has been guided primarily by heuristic considerations.⁷ There are some advantages to the approximate approach...[and] its use may yield better results than would a procedure that is optimum for the model. Although the model selected for $[f]$ is the best we have found for our purposes, it is sometimes too general...

⁵ To be precise, optimality is defined with respect to a model for the objective function $p(f)$, an observation model $p(y | x, \phi)$, a utility function $u(\mathcal{D})$, and an upper bound on the number of observations allowed τ . Bayesian decision theory provides a policy achieving the maximal expected utility at termination with respect to these choices.

running time of optimal policy and efficient approximations: § 5.3, p. 99

limited lookahead: § 5.3, p. 101

⁶ H. J. KUSHNER (1964). A New Method of Locating the Maximum Point of an Arbitrary Multi-peak Curve in the Presence of Noise. *Journal of Basic Engineering* 86(1):97–106.

⁷ Specifically, maximizing probability of improvement: § 7.5, p. 131.

What could possibly cause such a seemingly contradictory finding? As KUSHNER suggests, one possible reason could be model misspecification. The optimal policy is only defined with respect to a chosen model of the objective function and our observations, which is bound to be imperfect. By relying *less* on the model's belief, we may gain some robustness alongside considerable computational savings.

The intimate relationship between many Bayesian optimization methods and one-step lookahead is often glossed over, with a policy often introduced *ex nihilo* and the implied choice of utility function left unstated. This disconnect can sometimes lead to policies that are nonsensical from a decision-theoretic perspective or that incorporate implicit approximations that may not always be appropriate. We intend to clarify these connections here. We hope that our presentation can help guide practitioners in navigating the increasingly crowded space of available policies when presented with a novel scenario.

One-step lookahead

notation for one-step lookahead policies

proposed next point x with putative value y
updated dataset $\mathcal{D}' = \mathcal{D} \cup (x, y)$

expected marginal gain

acquisition functions: § 5, p. 88

value of sample information

Let us review the generic procedure for developing a one-step lookahead policy and adopt standard notation to facilitate their description. Suppose we have selected an arbitrary utility function $u(\mathcal{D})$ to evaluate a returned dataset. Suppose further that we have already gathered an arbitrary dataset $\mathcal{D} = (\mathbf{x}, \mathbf{y})$ and wish to select the next evaluation location. This is the fundamental role of an optimization policy.

If we were to choose some point x , we would observe a corresponding value y and update our dataset, forming $\mathcal{D}' = (\mathbf{x}', \mathbf{y}') = \mathcal{D} \cup \{(x, y)\}$. Note that in our discussion on decision theory in chapter 5, we notated this updated dataset with the symbol \mathcal{D}_1 , as we needed to be able to distinguish between datasets after the incorporation of a variable number of additional observations. As our focus in this chapter will be on one-step lookahead, we can simplify notation by dropping subscripts indicating time. Instead, we will systematically use the prime symbol to indicate future quantities after the acquisition of the next observation.

In one-step lookahead, we evaluate a proposed point x via the expected marginal gain in utility after incorporating an observation there (5.8):

$$\alpha(x; \mathcal{D}) = \mathbb{E}[u(\mathcal{D}') \mid x, \mathcal{D}] - u(\mathcal{D}),$$

which serves as an acquisition function inducing preferences over possible observation locations. We design each observation by maximizing this score:

$$x \in \arg \max_{x' \in \mathcal{X}} \alpha(x'; \mathcal{D}). \quad (7.1)$$

When the utility function $u(\mathcal{D})$ represents the expected utility of a decision informed by the data, such as a terminal recommendation following optimization, the expected marginal gain is also known as the *value of sample information* from observing at x . This term originates from the study of decision making in an economic context. Consider

utility function $u(\mathcal{D})$	expected one-step marginal gain
simple reward, (6.3)	expected improvement, § 7.3
global simple reward, (6.5)	knowledge gradient, § 7.4
unit for improving simple reward information gain, (6.8) or (6.9)	probability of improvement, § 7.5
cumulative reward, (6.7)	mutual information, § 7.6
	posterior mean, § 7.10

an agent who must make a decision under uncertainty, and suppose they have access to a third party who is willing to provide potentially insightful advice in exchange for a fee. By reasoning about the potential impact of this advice on the ultimate decision, we may quantify the expected value of the information,^{8,9} and determine whether the offered advice is worth the investment.

Due to its simplicity and inherent computational efficiency, one-step lookahead is a pervasive approximation scheme in Bayesian optimization. Table 7.1 provides a list of common acquisition functions, each representing the expected one-step marginal gain to a corresponding utility function. We will discuss each in detail below.

7.3 EXPECTED IMPROVEMENT

Adopting the simple reward utility function (6.3) and performing one-step lookahead defines the *expected improvement* acquisition function. Sequential maximization of expected improvement is perhaps the most widespread policy in all of Bayesian optimization.

Suppose that we wish to locate a single location in the domain with the highest possible objective value and ultimately wish to recommend one of the points investigated during optimization for permanent use. The simple reward utility function evaluates a dataset \mathcal{D} precisely by the expected value of an optimal final recommendation informed by the data, assuming risk neutrality:

$$u(\mathcal{D}) = \max \mu_{\mathcal{D}}(\mathbf{x}).$$

Suppose we have already gathered observations $\mathcal{D} = (\mathbf{x}, \mathbf{y})$ and wish to choose the next evaluation location. Expected improvement is derived by measuring the expected marginal gain in utility, or the instantaneous *improvement*, $u(\mathcal{D}') - u(\mathcal{D})$,¹⁰ offered by making the next observation at a proposed location x :¹¹

$$\alpha_{EI}(x; \mathcal{D}) = \int [\max \mu_{\mathcal{D}'}(\mathbf{x}')] p(y | x, \mathcal{D}) dy - \max \mu_{\mathcal{D}}(\mathbf{x}). \quad (7.2)$$

Expected improvement reduces to a particularly nice expression in the case of exact observations of the objective, where the utility takes a simpler form (6.4). Suppose that, when we elect to make an observation at a location x , we observe the exact objective value $\phi = f(x)$. Consider

Table 7.1: Summary of one-step lookahead optimization policies.

⁸ J. MARSCHAK and R. RADNER (1972). *Economic Theory of Teams*. Yale University Press. [§ 2.12]

⁹ H. RAIFFA and R. SCHLAIFER (1961). *Applied Statistical Decision Theory*. Division of Research, Graduate School of Business Administration, Harvard University. [§ 4.5]

simple reward: § 6.1, p. 109

risk neutrality: § 6.1, p. 109

¹⁰ This reasoning is the same for *all* one-step lookahead policies, which could all be described as maximizing “expected improvement.” But this name has been claimed for the simple reward utility alone.

¹¹ As mentioned in the last chapter, simple reward degenerates with an empty dataset; expected improvement does as well. In that case we can simply ignore the second term and compute the first, which for zero-mean additive noise becomes the mean function of the prior process.

expected improvement without noise

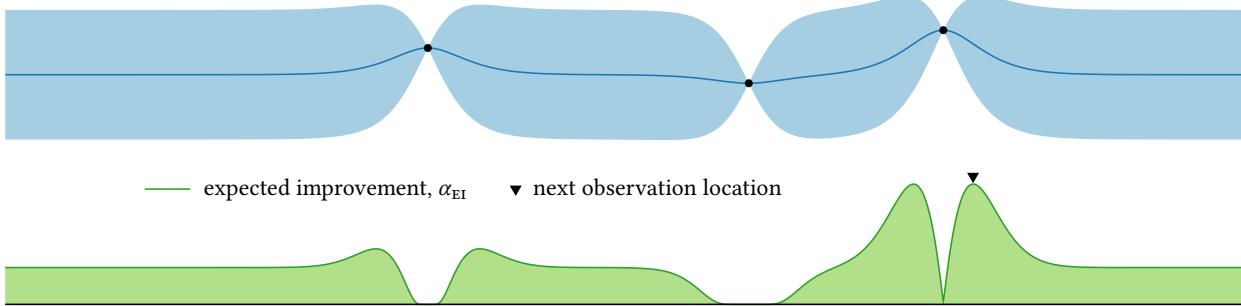


Figure 7.3: The expected improvement acquisition function (7.2) corresponding to our running example.

maximal value observed, incumbent ϕ^*

¹² The value ϕ^* is incumbent as it is currently “holding office” as our standing recommendation until it is deposed by a better candidate.

a dataset $\mathcal{D} = (\mathbf{x}, \phi)$, and define $\phi^* = \max \phi$ to be the so-called *incumbent*: the maximal objective value yet seen.¹² As a consequence of exact observation, we have

$$u(\mathcal{D}) = \phi^*; \quad u(\mathcal{D}') = \max(\phi^*, \phi);$$

and thus

$$u(\mathcal{D}') - u(\mathcal{D}) = \max(\phi - \phi^*, 0).$$

Substituting into (7.2), in the noiseless case we have

$$\alpha_{EI}(x; \mathcal{D}) = \int \max(\phi - \phi^*, 0) p(\phi | x, \mathcal{D}) d\phi. \quad (7.3)$$

example and interpretation

Expected improvement is illustrated for our running example in figure 7.3. In this case, maximizing expected improvement will select a point near the previous best point found, an example of exploitation. Notice that the expected improvement vanishes near regions where we have existing observations. Although these locations may be likely to yield values higher than ϕ^* due to relatively high expected value, the relatively narrow credible intervals suggest that the magnitude of any improvement is likely to be small. Expected improvement is thus considering the exploration-exploitation dilemma in the selection of the next observation location, and the tradeoff between these two concerns is considered automatically.

Figure 7.4 shows the posterior belief of the objective after sequentially maximizing expected improvement to gather 20 additional observations of our example objective function. The global optimum was efficiently located. The distribution of the sample locations, with more evaluations in the most promising regions, reflects consideration of the exploration-exploitation dilemma. However, there seems to have been a focus on exploitation throughout the entire process; the first ten observations for example never strayed from the initially known local optimum. This behavior is a reflection of the simple reward utility function underlying the policy, which only rewards the discovery of high objective function values at *observed* locations. As a result, one-step lookahead may

simulated optimization and interpretation

exploitative behavior resulting from myopia

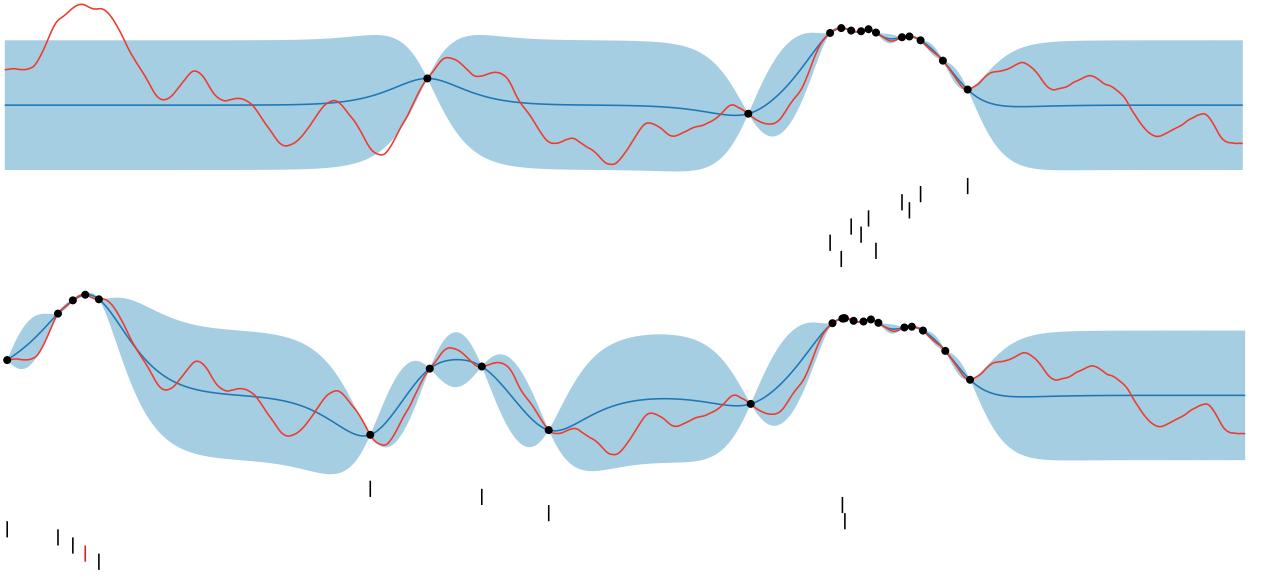


Figure 7.4: The posterior after 10 (top) and 20 (bottom) steps of the optimization policy induced by the expected improvement acquisition function (7.2) on our running example. The tick marks show the points chosen by the policy, progressing from top to bottom, during iterations 1–10 (top) and 11–20 (bottom). Observations within 0.2 length scales of the optimum are marked in red; the optimum was located on iteration 19.

rationally choose to make marginal improvements to the value of the best-seen point, even if the underlying function value is known with a fair amount of confidence.

7.4 KNOWLEDGE GRADIENT

Adopting the global reward utility (6.5) and performing one-step-lookahead yields an acquisition function known as the *knowledge gradient*.

global reward: § 6.1, p. 109

Assume that, just as in the situation leading to the derivation of expected improvement, we again wish to identify a single point in the domain maximizing the objective function. However, imagine that at termination we are willing to commit to a location possibly never evaluated during optimization. To this end, we adopt the global reward utility function to measure our progress:

$$u(\mathcal{D}) = \max_{x \in \mathcal{X}} \mu_{\mathcal{D}}(x),$$

which rewards data for increasing the posterior mean, irrespective of location. Computing the one-step marginal gain to this utility results in the knowledge gradient acquisition function:

$$\alpha_{KG}(x; \mathcal{D}) = \int \left[\max_{x' \in \mathcal{X}} \mu_{\mathcal{D}'}(x') \right] p(y | x, \mathcal{D}) dy - \max_{x' \in \mathcal{X}} \mu_{\mathcal{D}}(x'). \quad (7.4)$$

The knowledge gradient moniker was coined by FRAZIER and POWELL,¹³ who interpreted the global reward as the amount of “knowledge”

¹³ P. FRAZIER and W. POWELL (2007). The Knowledge Gradient Policy for Offline Learning with Independent Normal Rewards. *ADPRL* 2007.

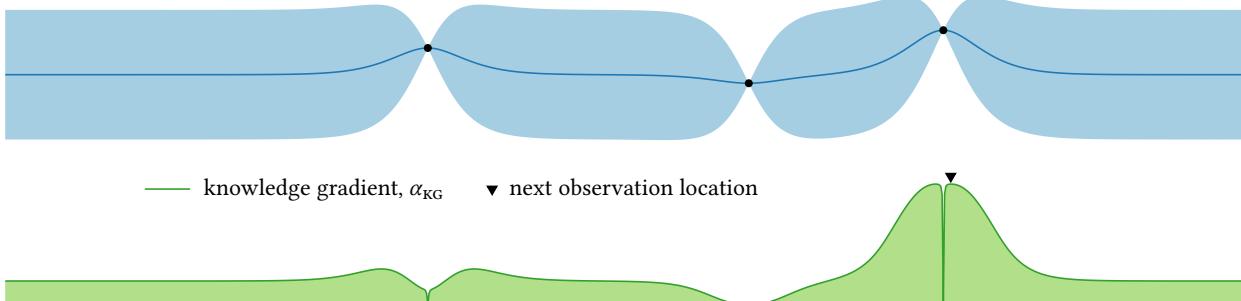
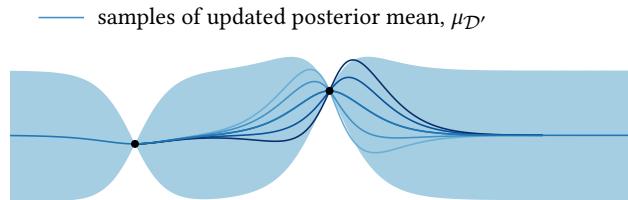


Figure 7.5: The knowledge gradient acquisition function (7.4) corresponding to our running example.

Figure 7.6: Samples of the updated posterior mean when evaluating at the location chosen by the knowledge gradient, illustrated in figure 7.5. Only the right-hand section of the domain is shown.



example and interpretation
reason for selected observation

about the global maximum offered by a dataset \mathcal{D} . The knowledge gradient $\alpha_{KG}(x; \mathcal{D})$ can then be interpreted as the expected (discrete-time) change in knowledge offered by a measurement at x .

The knowledge gradient is illustrated for our running example in figure 7.5. Perhaps surprisingly, the chosen observation location is remarkably close to the previously best-seen point. At first glance, this may seem wasteful, as we are already fairly confident about the value we might observe.

However, the knowledge gradient seeks to maximize the *global* maximum of the posterior mean, regardless of its location. With this in mind, we may reason as follows. There must be a local maximum of the objective function in the neighborhood of the best-seen point, but our current knowledge is insufficient to pinpoint its location. Further, as the relevant local maximum is probably not located precisely at this point, the objective function is either increasing or decreasing as it passes through. If we were to learn the *derivative* of the objective at this point, we would adjust our posterior belief to reflect that knowledge. Regardless of the sign or exact value of the derivative, our updated belief would reflect the discovery of a new, higher local maximum of the posterior mean in the indicated direction. By evaluating at the location selected by the knowledge gradient, we can effectively estimate the derivative of the objective; this is the principle behind finite differencing.

In figure 7.6, we show samples of the updated posterior mean function $\mu_{\mathcal{D}'}(x)$ derived from sampling from the predictive distribution at the chosen evaluation location and conditioning. Indeed, these samples exhibit newly located global maxima on either side of the selected point, depending on the sign of the implied derivative. Note that the locations

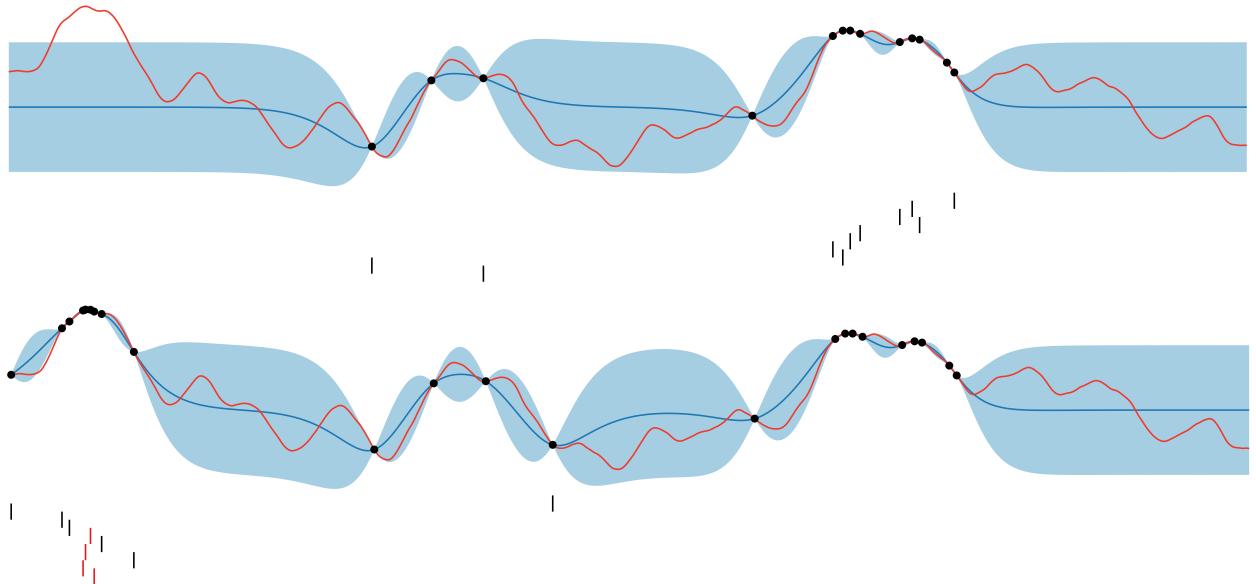


Figure 7.7: The posterior after 10 (top) and 20 (bottom) steps of the optimization policy induced by the knowledge gradient acquisition function (7.4) on our running example. The tick marks show the points chosen by the policy, progressing from top to bottom, during iterations 1–10 (top) and 11–20 (bottom). Observations within 0.2 length scales of the optimum are marked in red; the optimum was located on iteration 15.

of these new maxima coincide with local maxima of the expected improvement acquisition function; see figure 7.3 for comparison. This is not a coincidence! One way to interpret this relation is that, due to rewarding large values of the posterior mean at observed locations only, expected improvement must essentially guess on which side the hidden local optimum of the objective lies and hope to be correct. The knowledge gradient, on the other hand, considers identifying this maximum on either side a success, and guessing is not necessary.

Figure 7.7 illustrates the behavior of the knowledge gradient policy on our example optimization scenario. The global optimum was located efficiently. Comparing the decisions made by the knowledge gradient to those made by expected improvement (see figure 7.4), we can observe a somewhat more even exploration of the domain, including in local maxima. The knowledge gradient policy does not necessarily need to expend observations to verify a suspected maximum, instead putting more trust into the model to have correct beliefs in these regions.

simulated optimization and interpretation

more exploration than expected improvement from more-relaxed utility

7.5 PROBABILITY OF IMPROVEMENT

As its name suggests, the *probability of improvement* acquisition function computes the probability of an observed value to improve upon some chosen threshold, regardless of the magnitude of this improvement.

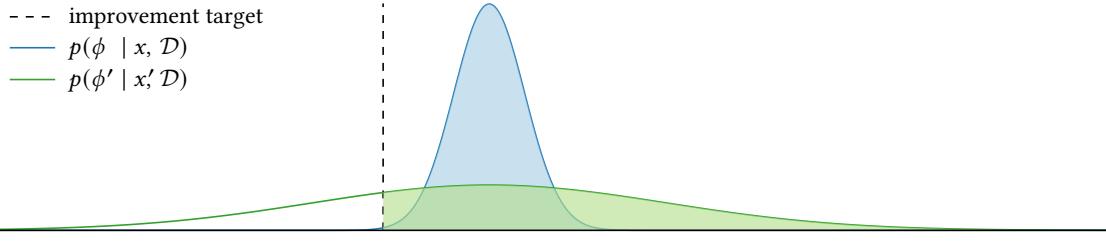


Figure 7.8: An illustrative example comparing the behavior of probability of improvement with expected improvement computed with respect to the dashed target. The predictive distributions for two points x and x' are shown. The distributions have equal mean but the distribution at x' has larger predictive standard deviation. The shaded regions represent the region of improvement. The relatively safe x is preferred by probability of improvement, whereas the more-risky x' is preferred by expected improvement.

simple reward: § 6.1, p. 109

desired margin of improvement, ε
desired improvement threshold, τ

utility formulation

noiseless case

comparison with expected improvement

Consider the simple reward of an already gathered dataset $\mathcal{D} = (\mathbf{x}, \mathbf{y})$:

$$u(\mathcal{D}) = \max \mu_{\mathcal{D}}(\mathbf{x}).$$

The probability of improvement acquisition function scores a proposed observation location x according to the probability that an observation there will improve this utility by at least some margin ε . Let us denote the desired utility threshold with $\tau = u(\mathcal{D}) + \varepsilon$; we will use both the absolute threshold τ and the marginal threshold ε in the following discussion as convenient. The probability of improvement is then the probability that the updated utility $u(\mathcal{D}')$ exceeds the chosen threshold:

$$\alpha_{\text{PI}}(x; \mathcal{D}, \tau) = \Pr(u(\mathcal{D}') > \tau \mid x, \mathcal{D}). \quad (7.5)$$

We may interpret probability of improvement in the Bayesian decision-theoretic framework as computing the expected one-step marginal gain in a peculiar choice of utility function: a utility offering unit reward for each observation increasing the simple reward by the desired amount.

In the case of exact observation, we have

$$u(\mathcal{D}) = \max f(\mathbf{x}) = \phi^*; \quad u(\mathcal{D}') = \max(\phi^*, \phi),$$

and we may write the probability of improvement in the somewhat simpler form

$$\alpha_{\text{PI}}(x; \mathcal{D}, \tau) = \Pr(\phi > \tau \mid x, \mathcal{D}). \quad (7.6)$$

In this case, the probability of improvement is simply the complementary cumulative distribution function of the predictive distribution evaluated at the improvement threshold τ . This form of probability of improvement is sometimes encountered in the literature, but our modification in terms of the simple reward allows for inexact observations as well.

It can be illustrative to compare the preferences over observation locations implied by the probability of improvement and expected improvement acquisition functions. In general, probability of improvement is somewhat more risk-averse than expected improvement, because probability of improvement would prefer a certain improvement of modest

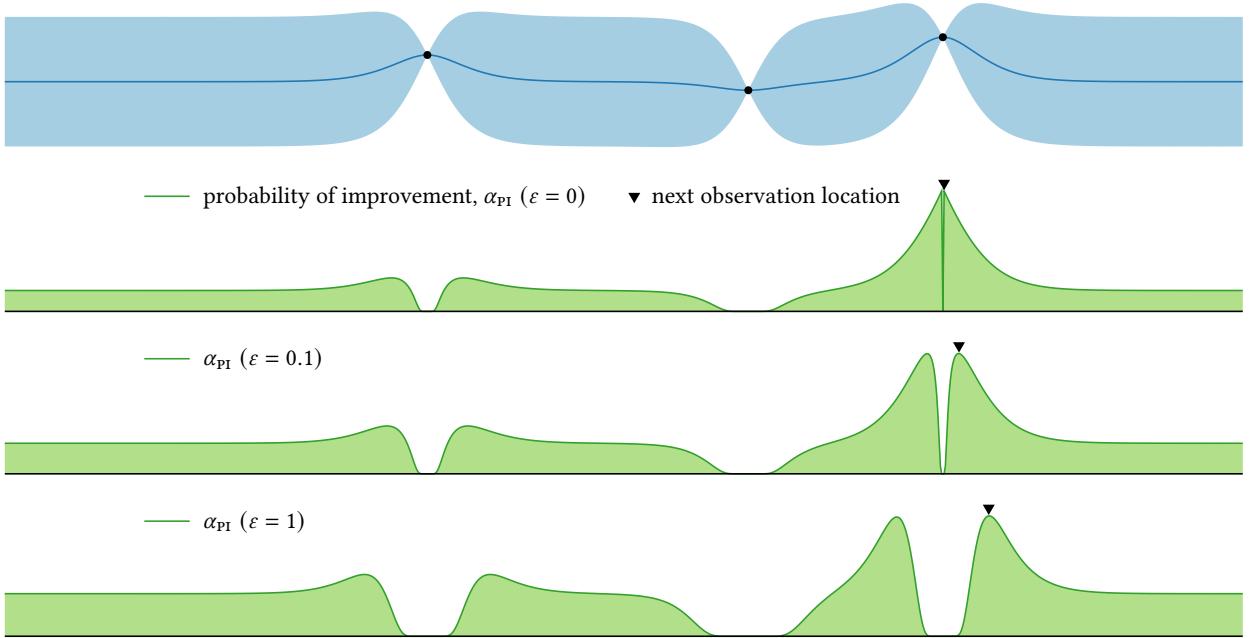


Figure 7.9: The probability of improvement acquisition function (7.5) corresponding to our running example for different values of the target improvement ε . The target is expressed as a fraction of the range of the posterior mean over the space. Increasing the target improvement leads to increasingly exploratory behavior.

magnitude to an uncertain improvement of potentially large magnitude. Figure 7.8 illustrates this phenomenon. Shown are the predictive distributions for the objective function values at two points x and x' . Both points have equal predictive means; however, x' has a significantly larger predictive standard deviation. We consider improvement with respect to the illustrated target. The shaded regions represent the regions of improvement; the probability mass of these regions equal the probabilities of improvement. Improvement is near certain at x ($\alpha_{PI} = 99.9\%$), whereas it is somewhat smaller at x' ($\alpha_{PI} = 72.6\%$), and thus probability of improvement would prefer to observe at x . The *expected* improvement at x , however, is small compared to x' with its longer tail:

$$\frac{\alpha_{EI}(x'; \mathcal{D})}{\alpha_{EI}(x; \mathcal{D})} = 1.28.$$

The expected improvement at x' is 28% larger than at x , indicating a preference for a less-certain but potentially larger payout.

The role of the improvement target

The magnitude of the required improvement plays a crucial role in shaping the behavior of probability of improvement policies. By adjusting this parameter, we may encourage exploration (with large ε) or exploitation (with small ε). Figure 7.9 shows the probability of improvement for

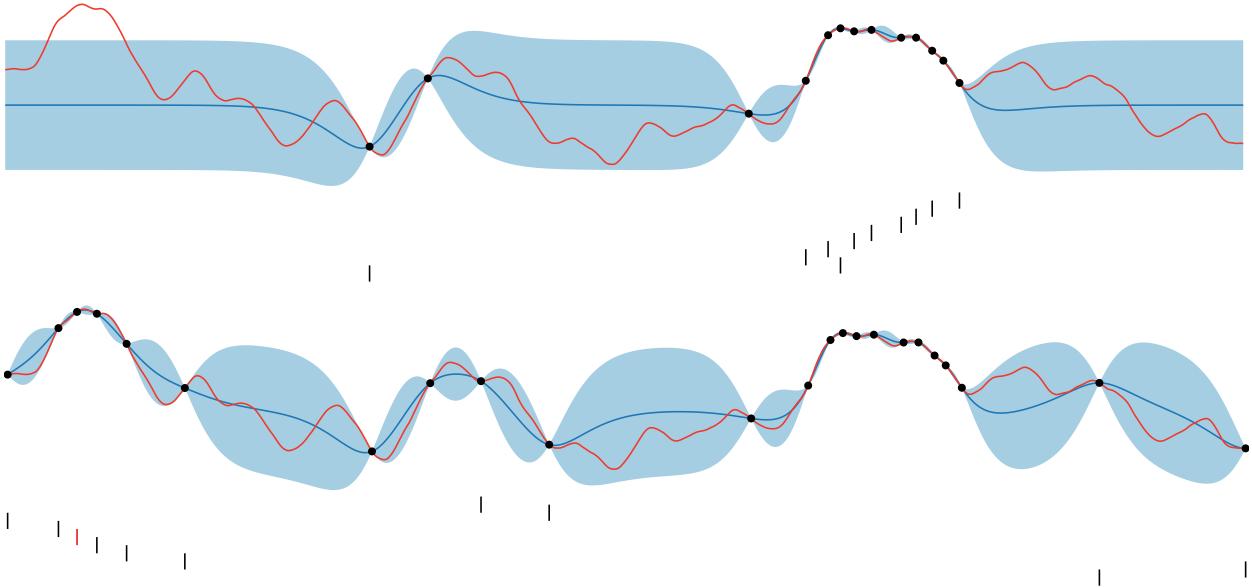


Figure 7.10: The posterior after 10 (top) and 20 (bottom) steps of the optimization policy induced by probability of improvement with $\varepsilon = 0.1[\max \mu_{\mathcal{D}}(x) - \min \mu_{\mathcal{D}}(x)]$ (7.5) on our running example. The tick marks show the points chosen by the policy, progressing from top to bottom, during iterations 1–10 (top) and 11–20 (bottom). Observations within 0.2 length scales of the optimum are marked in red; the optimum was located on iteration 15.

our example scenario with thresholds corresponding to infinitesimal improvement, a modest improvement, and a significant improvement. The shift towards exploratory behavior for larger improvement thresholds can be clearly seen.

In figure 7.10, we see 20 evaluations chosen by maximizing probability of improvement with the target dynamically set to 10% of the range of the posterior mean function. The global optimum was located, and the domain appears sufficiently explored. Although performance was quite reasonable here, the improvement threshold was set somewhat arbitrarily, and it is not always clear how one should set this parameter.

the $\varepsilon = 0$ case

On one extreme, some authors define a parameter-free (and perhaps too literal) version of probability of improvement by fixing the improvement target to $\varepsilon = 0$, rewarding even infinitesimal improvement to the current data. Intuitively, this low bar can induce overly exploitative behavior. Examining the probability of improvement with $\varepsilon = 0$ for our running example in figure 7.9, we see that the acquisition function is maximized directly next to the previously best-found point. This decision represents extreme exploitation and potentially undesirable behavior. The situation after applying probability of improvement with $\varepsilon = 0$ to select 20 additional observation locations, shown in figure 7.11, clearly demonstrates a drastic focus on exploitation. Notably, the global optimum was not identified.

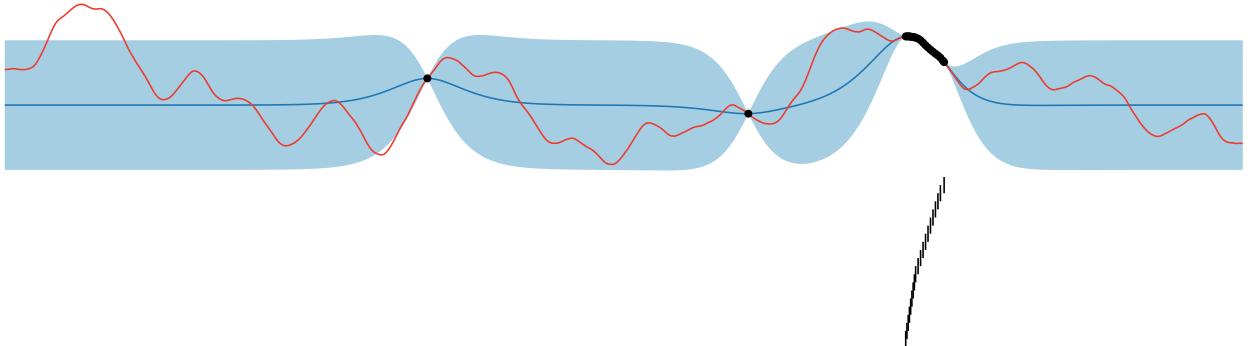


Figure 7.11: The posterior after 20 steps of the optimization policy induced by probability of improvement with $\varepsilon = 0$ (7.5) on our running example. The tick marks show the points chosen by the policy, progressing from top to bottom.

Evidently we must carefully select the desired improvement threshold to achieve ideal behavior. JONES provided some simple, data-driven advice for choosing improvement thresholds that remains sound.¹⁴ Define

$$\mu^* = \max_{x \in \mathcal{X}} \mu_D(x); \quad r = \max \mu_D(\mathbf{x}) - \min \mu_D(\mathbf{x});$$

to represent the global maximum of the posterior mean and the range of the posterior mean at the observed locations. JONES suggests considering targets of the form

$$\mu^* + \alpha r,$$

where $\alpha \geq 0$ controls the amount of desired improvement in terms of the range of observed data. He provides a table of 27 suggested values for α in the range $[0, 3]$ and remarks that the points optimizing the set of induced acquisition functions typically cluster together in a small number of locations, each representing a different tradeoff between exploration and exploitation.¹⁵ JONES continued to recommend selecting one point from each of these clusters to evaluate in parallel, defining a batch optimization policy. Although this may not always be possible, the recommended parameterization of the desired improvement is natural and would be appropriate for general use.

This proposal is illustrated for our running example in figure 7.12. We begin with the posterior after selecting 15 points in our previous demo (see figure 7.10), and indicate the points maximizing the probability of improvement for JONES's proposed improvement targets. The points cluster together in four regions reflecting varying exploration-exploitation tradeoffs.

¹⁴ D. R. JONES (2001). A Taxonomy of Global Optimization Methods Based on Response Surfaces. *Journal of Global Optimization* 21(4):345–383.

¹⁵ The proposed values for α given by JONES are compiled below.

α		
0	0.07	0.25
0.0001	0.08	0.3
0.001	0.09	0.4
0.01	0.1	0.5
0.02	0.11	0.75
0.03	0.12	1
0.04	0.13	1.5
0.05	0.15	2
0.06	0.2	3

7.6 MUTUAL INFORMATION AND ENTROPY SEARCH

A family of information-theoretic optimization policies have been proposed in recent years, most with variations on the name *entropy search*.

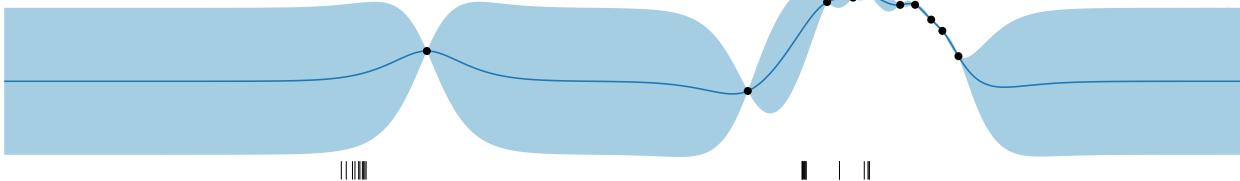


Figure 7.12: The points maximizing probability of improvement using the 27 improvement thresholds proposed by JONES, beginning with the posterior from figure 7.10 after 10 total observations have been obtained. The tick marks show the chosen points and cluster together in four regions representing different tradeoffs between exploration and exploitation.

¹⁶ T. M. COVER and J. A. THOMAS (2006). *Elements of Information Theory*. John Wiley & Sons.

¹⁷ D. J. C. MACKAY (2003). *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press.

information-theoretic decision making as a model of the scientific method

information gain: § 6.3, p. 115

The acquisition function in these methods is *mutual information*, a measure of dependence between random variables that is a central concept in information theory.^{16,17}

The reasoning underlying entropy search policies is somewhat different from and more general than the other acquisition functions we have considered thus far, all of which ultimately focus on maximizing the posterior mean function. Although this is a pragmatic concern, but one that is intimately linked to optimization. Information-theoretic experimental design is instead motivated by an abstract pursuit of *knowledge*, and may be interpreted as a mathematical formulation of the scientific method.

We begin by identifying some unknown feature of the world that we wish to learn about; in the context of Bayesian inference, this will be some random variable ω . We then view each observation we make as an opportunity to learn about this random variable, and seek to gather data that will, in aggregate, provide considerable information about ω . This process is analogous to a scientist designing a sequence of experiments to understand some natural phenomenon, where each experiment may be chosen to challenge or confirm constantly evolving beliefs.

The framework of information theory allows us to formalize this process. We may quantify the amount of information provided about a random variable ω by a dataset \mathcal{D} via the *information gain*, a concept for which we provided two definitions in the last chapter. Adopting either definition as a utility function and performing one-step lookahead yields mutual information as an acquisition function.

Information-theoretic optimization policies select ω such that its determination gives insight into our optimization problem (1.1). However, by selecting different choices for ω , we can generate radically different policies, each attempting to learn about a different aspect of the system of interest. Maximizing mutual information has long been promoted as a general framework for optimal experimental design,¹⁸ and this framework has been applied in numerous *active learning* settings.¹⁹

Before showing how mutual information arises in a decision-theoretic context, we pause to define the concept and derive some important properties.

¹⁸ D. V. LINDLEY (1956). On a Measure of the Information Provided by an Experiment. *The Annals of Mathematical Statistics* 27(4):986–1005.

¹⁹ B. SETTLES (2012). *Active Learning*. Morgan & Claypool.

Mutual information

Let ω and ψ be random variables with probability density functions $p(\omega)$ and $p(\psi)$. The mutual information between ω and ψ is

$$I(\omega; \psi) = \iint p(\omega, \psi) \log \frac{p(\omega, \psi)}{p(\omega) p(\psi)} d\omega d\psi. \quad (7.7)$$

This expression may be recognized as the Kullback–Leibler divergence between the joint distribution of the random variables and the product of their marginal distributions:

$$I(\omega; \psi) = D_{\text{KL}}[p(\omega, \psi) \parallel p(\omega) p(\psi)].$$

We may extend this definition to conditional probability distributions as well. Given an arbitrary set of observed data \mathcal{D} , we define the conditional mutual information between ω and ψ by:²⁰

$$I(\omega; \psi | \mathcal{D}) = \iint p(\omega, \psi | \mathcal{D}) \log \frac{p(\omega, \psi | \mathcal{D})}{p(\omega | \mathcal{D}) p(\psi | \mathcal{D})} d\omega d\psi.$$

Here we have simply conditioned all distributions on the data and applied the definition in (7.7) to the posterior beliefs.

Several properties of mutual information are immediately evident from its definition. First, mutual information is symmetric in its arguments:

$$I(\omega; \psi) = I(\psi; \omega). \quad (7.8)$$

We also have that if ω and ψ are independent, then $p(\omega, \psi) = p(\omega) p(\psi)$ and the mutual information is zero:

$$I(\omega; \psi) = \iint p(\omega, \psi) \log \frac{p(\omega) p(\psi)}{p(\omega) p(\psi)} d\omega d\psi = 0.$$

Further, recognition of mutual information as a Kullback–Leibler divergence implies several additional inherited properties, including nonnegativity. Thus mutual information attains its minimal value when ω and ψ are independent.

We may also manipulate (7.7) by twice applying the identity

$$p(\omega, \psi) = p(\psi) p(\omega | \psi)$$

to derive an equivalent expression for the mutual information:

$$\begin{aligned} I(\omega; \psi) &= \iint p(\omega, \psi) \log \frac{p(\omega, \psi)}{p(\omega) p(\psi)} d\omega d\psi \\ &= \iint p(\omega, \psi) \log p(\omega | \psi) d\omega d\psi - \int p(\omega) \log p(\omega) d\omega \\ &= \int p(\psi) \left[\int p(\omega | \psi) \log p(\omega | \psi) d\omega \right] d\psi + H[\omega] \\ &= H[\omega] - \mathbb{E}[H[\omega | \psi]]. \end{aligned} \quad (7.9)$$

definition

conditional mutual information, $I(\omega; \psi | \mathcal{D})$

²⁰ Some authors use the notation $I(\omega; \psi | \mathcal{D})$ to represent the *expectation* of the given quantity with respect to the dataset \mathcal{D} . In optimization, we will always have an explicit dataset in hand, in which case the provided definition is more useful.

symmetry

nonnegativity

expected reduction in entropy

²¹ It is important to note that this is true only in expectation. Consider two random variables x and y with the following joint distribution. x takes value 0 or 1 with probability $1/2$ each. If x is 0, y takes value 0 or 1 with probability $1/2$ each. If x is 1, y takes value 0 or -1 with probability $1/2$ each. The entropy of x is 1 bit and the entropy of y is 1.5 bits. Observing x always yields 0.5 bits about y . However, observing y produces either *no information* about x (0 bits), with probability $1/2$, or *complete information* about x (1 bit), with probability $1/2$. So the information gain about x from y and about y from x is actually *never* equal. However, the *expected* information gain is equal, $I(x; y) = 0.5$ bits.

²² Setting $u(\mathcal{D}) = D_{\text{KL}}[p(\omega | \mathcal{D}) \| p(\omega)]$, we have:

$$\begin{aligned} \mathbb{E}[u(\mathcal{D}') | x, \mathcal{D}] &= \mathbb{E}\left[\int p(\omega | \mathcal{D}') \log p(\omega | \mathcal{D}') d\omega | x, \mathcal{D}\right] \\ &\quad - \int p(\omega | \mathcal{D}) \log p(\omega) d\omega \\ &= -\mathbb{E}[H[\omega | x, \mathcal{D}'] | \mathcal{D}] \\ &\quad - \int p(\omega | \mathcal{D}) \log p(\omega) d\omega. \end{aligned}$$

Here the second term is known as the *cross entropy* between $p(\omega)$ and $p(\omega | \mathcal{D})$. We can also rewrite the utility in similar terms:

$$\begin{aligned} u(\mathcal{D}) &= \int p(\omega | \mathcal{D}) \frac{\log p(\omega | \mathcal{D})}{\log p(\omega)} d\omega \\ &= \int p(\omega | \mathcal{D}) \log p(\omega | \mathcal{D}) d\omega \\ &\quad - \int p(\omega | \mathcal{D}) \log p(\omega) d\omega \\ &= -H[\omega | \mathcal{D}] \\ &\quad - \int p(\omega | \mathcal{D}) \log p(\omega) d\omega. \end{aligned}$$

If we subtract, the cross-entropy terms cancel and we obtain mutual information:

$$\begin{aligned} \mathbb{E}[u(\mathcal{D}') | x, \mathcal{D}] - u(\mathcal{D}) &= \\ H[\omega | \mathcal{D}] - \mathbb{E}[H[\omega | \mathcal{D}'] | \mathcal{D}]. \end{aligned}$$

Thus the mutual information between ω and ψ is the expected decrease in the differential entropy of ω if we were to observe ψ . Due to symmetry (7.8), we may swap the roles of ω and ψ to derive an equivalent expression in the other direction:

$$I(\omega; \psi) = H[\omega] - \mathbb{E}_\psi[H[\omega | \psi]] = H[\psi] - \mathbb{E}_\omega[H[\psi | \omega]]. \quad (7.10)$$

Observing either ω or ψ will, in expectation, provide the same amount of information about the other: the mutual information $I(\omega; \psi)$.²¹

Maximizing mutual information as an optimization policy

Mutual information arises naturally in Bayesian sequential experimental design as the one-step expected information gain resulting from an observation. In the previous chapter, we introduced two different methods for quantifying this information gain. The first was the reduction in the differential entropy of ω from the prior to the posterior:

$$\begin{aligned} u(\mathcal{D}) &= H[\omega] - H[\omega | \mathcal{D}] \\ &= \int p(\omega | \mathcal{D}) \log p(\omega | \mathcal{D}) d\omega - \int p(\omega) \log p(\omega) d\omega. \end{aligned} \quad (7.11)$$

The second was the Kullback–Leibler divergence between the posterior and the prior:

$$u(\mathcal{D}) = D_{\text{KL}}[p(\omega | \mathcal{D}) \| p(\omega)] = \int p(\omega | \mathcal{D}) \log \frac{p(\omega | \mathcal{D})}{p(\omega)} d\omega. \quad (7.12)$$

Remarkably, performing one-step lookahead with either choice yields mutual information as an acquisition function.

Let us first compute the expected marginal gain in (7.11). In this case the marginal information gain is:

$$H[\omega | \mathcal{D}] - H[\omega | \mathcal{D}'],$$

and the expected marginal information gain is then:

$$\begin{aligned} \alpha_{\text{MI}}(x; \mathcal{D}) &= H[\omega | \mathcal{D}] - \mathbb{E}[H[\omega | \mathcal{D}'] | x, \mathcal{D}] \\ &= I(y; \omega | x, \mathcal{D}), \end{aligned} \quad (7.13)$$

where we have recognized the expected reduction in entropy in (7.13) as the mutual information between y and ω given the putative location x and the available data \mathcal{D} (7.9). It is simple to verify that the expected marginal improvement to the alternative information gain definition (7.12) gives the same expression; several terms cancel when computing the expectation, and those that remain are identical to those in (7.13).²²

Due to the symmetry of mutual information, we have several equivalent forms for this acquisition function (7.10):

$$\begin{aligned} \alpha_{\text{MI}}(x; \mathcal{D}) &= I(y; \omega | x, \mathcal{D}) \\ &= H[\omega | \mathcal{D}] - \mathbb{E}_y[H[\omega | \mathcal{D}'] | x, \mathcal{D}] \end{aligned} \quad (7.14)$$

$$= H[y | x, \mathcal{D}] - \mathbb{E}_\omega[H[y | \omega, x, \mathcal{D}] | x, \mathcal{D}]. \quad (7.15)$$

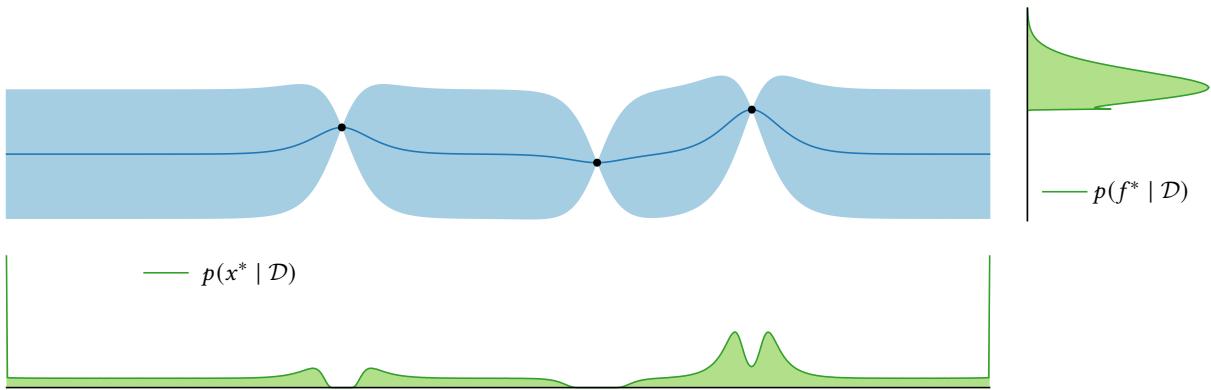


Figure 7.13: The posterior belief about the location of the global optimum, $p(x^* | \mathcal{D})$, and about the value of the global optimum, $p(f^* | \mathcal{D})$, for our running example. Note the significant probability mass associated with the optimum lying on the boundary.

Depending on the application, one of these two forms may be preferable, and maximizing either results in the same policy.

Adopting mutual information as an acquisition function for optimization requires that ω be selected to support the optimization task. Two natural options present themselves: the location of the global optimum, x^* , and the maximum value attained, $f^* = f(x^*)$ (1.1). Both have received extensive consideration, and we will discuss each in turn.

*Mutual information with x^**

Several authors have proposed mutual information with the location of the global optimum x^* as an acquisition function:²³

$$\alpha_{x^*}(x; \mathcal{D}) = I(y; x^* | x, \mathcal{D}). \quad (7.16)$$

The distribution of x^* is illustrated for our running example in figure 7.13. Even for this simple example, the distribution of the global optimum is nontrivial and multimodal. In fact, in this case, there is a significant probability that the global maximum occurs on the *boundary* of the domain, which has Lebesgue measure zero, so x^* does not even have a proper probability density function.²⁴ We will nonetheless use the notation $p(x^* | \mathcal{D})$ in our discussion below.

The middle panel of figure 7.14 shows the mutual information with x^* (7.16) for our running example. The next evaluation location will be chosen to investigate the neighborhood of the best-seen point. It is interesting to compare the behavior of the mutual information with other acquisition functions near the boundary of the domain. Although there is a significant probability that the maximum is achieved on the boundary, mutual information indicates that we cannot expect to reveal much information regarding x^* by measuring there. More information tends to be revealed by evaluating away from the boundary, as we can reduce our

²³ We tacitly assume the location of the global optimum is unique to simplify discussion. Technically x^* is a set-valued random variable. For Gaussian process models, the uniqueness of x^* can be guaranteed under mild assumptions (§ 2.7, p. 34).

²⁴ A proper treatment would separate the probability density on the interior of \mathcal{X} from the distribution restricted to the boundary, but in practice we will only ever be sampling from this distribution, as it is simply too complicated to work with directly.

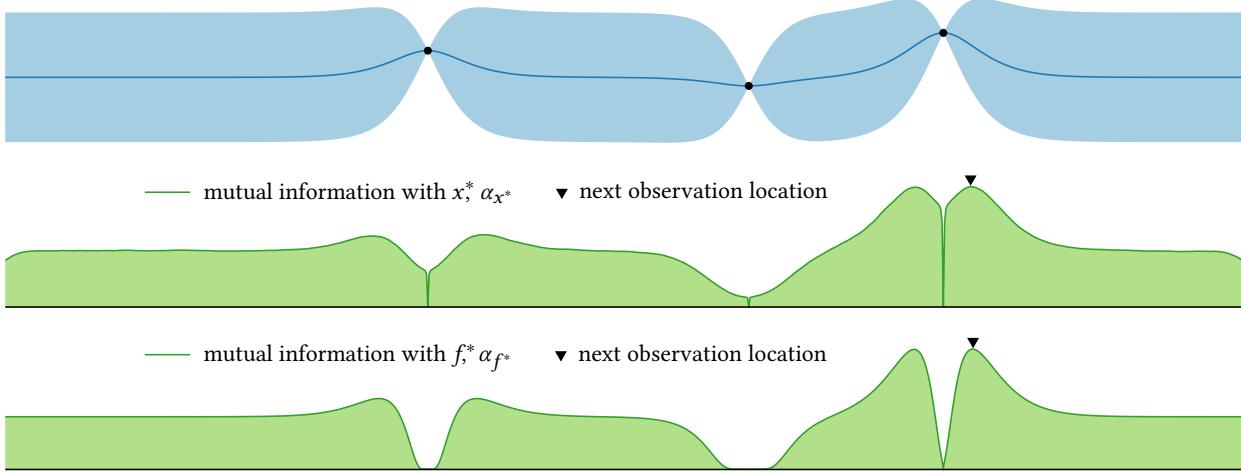


Figure 7.14: The mutual information between the observed value and the location of the global optimum, α_{x^*} (middle panel), and between the observed value and the value of the global optimum, α_{f^*} (bottom panel), for our running example.

uncertainty about the objective function over a larger volume. Expected improvement (figure 7.3) and probability of improvement (figure 7.9), on the other hand, are computed only from inspection of the marginal predictive distribution $p(y | x, \mathcal{D})$. As a result, they cannot differentiate observation locations based on their global impact on our belief.

In figure 7.15, we demonstrate 20 steps of optimization by maximizing the mutual information with x^* (7.16) for our scenario. We also show the posterior belief about the maximum location at termination, $p(x^* | \mathcal{D})$. The global optimum was discovered efficiently and with remarkable confidence. Further, the posterior mode matches the true optimal location.

*Mutual information with f^**

Mutual information with the value of the global optimum f^* has also been investigated as an acquisition function:²⁵

$$\alpha_{f^*}(x; \mathcal{D}) = I(y; f^* | x, \mathcal{D}). \quad (7.17)$$

The distribution of this quantity is illustrated for our running example in figure 7.13. There is a sharp mode in the distribution corresponding to the best-seen value in fact being near-optimal, and there is no mass below the best-seen value, as it serves as a lower bound on the maximum due to the assumption of exact observation.

The bottom panel of figure 7.14 shows the mutual information with f^* (7.17) for our running example. The next evaluation location will be chosen to investigate the neighborhood of the best-seen point. It is interesting to contrast this surface with that of the mutual information with x^* in figure (7.16). Mutual information with f^* is heavily penalized near

²⁵ Again we assume in this discussion that a global optimum f^* exists almost surely. This assumption can be guaranteed for Gaussian process models under mild assumptions (§ 2.7, p.34).

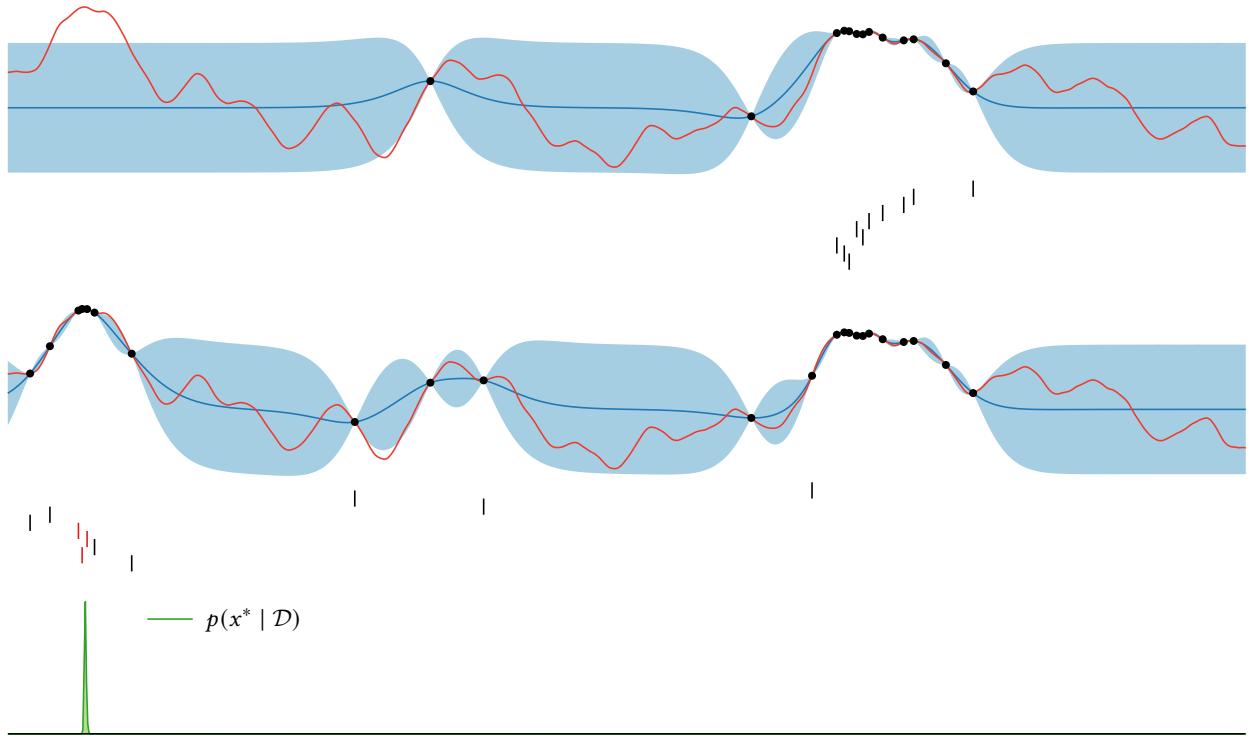


Figure 7.15: The posterior after 10 (top) and 20 (bottom) steps of maximizing the mutual information between the observed value y and the location of the global maximum x^* (7.16) on our running example. The tick marks show the points chosen by the policy, progressing from top to bottom. Observations within 0.2 length scales of the optimum were marked in red; the optimum was located on iteration 16.

existing observations, even those with relatively high values. Observing at these points would contribute relatively little information about the value of the optimum, as their predictive distributions already reflect a narrow range of possibilities and contribute little to the distribution of f^* . Further, points on the boundary were less favored when seeking to learn about x^* . However, these points are expected to provide just as much information about f^* as neighboring points.

Figure 7.16 illustrates 25 evaluations chosen by sequentially maximizing the mutual information with f^* (7.17) for our scenario, along with the posterior belief about the value of the maximum given these observations, $p(f^* | \mathcal{D})$. The global optimum was discovered after 23 iterations, somewhat slower than the alternatives described above. The value of the optimum is known with almost complete confidence at termination.

7.7 MULTI-ARMED BANDITS AND OPTIMIZATION

Several Bayesian optimization algorithms have been inspired by policies for *multi-armed bandits*, a model system for sequential decision making

multi-armed bandits

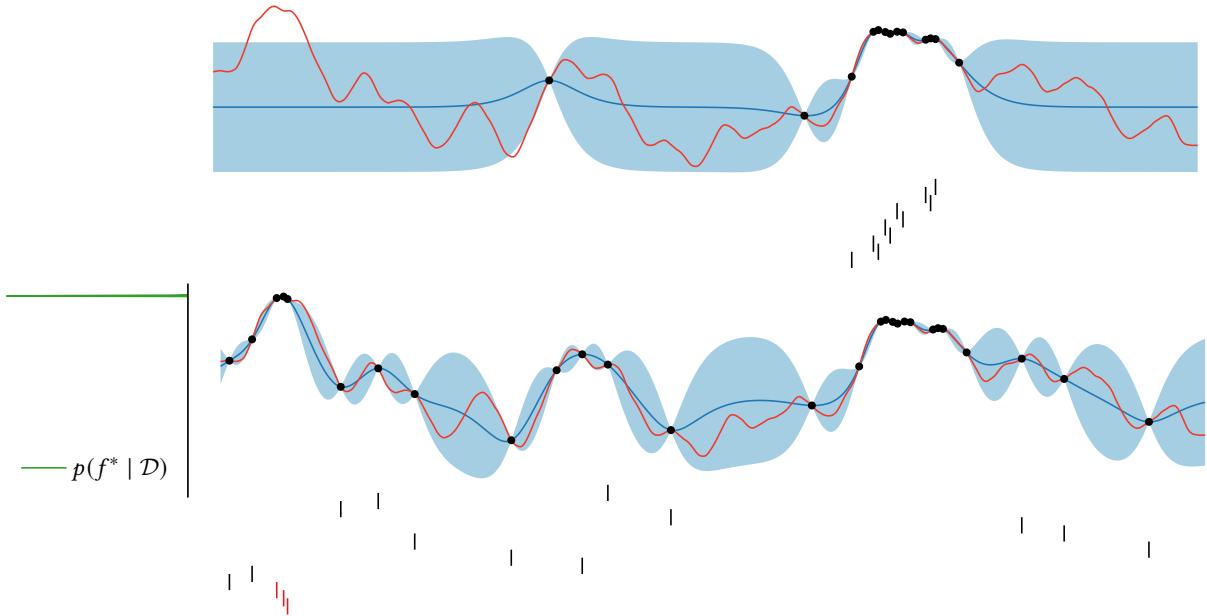


Figure 7.16: The posterior after 10 (top) and 25 (bottom) steps of maximizing the mutual information between the observed value y and the location of the global maximum f^* (7.17) on our running example. The tick marks show the points chosen by the policy, progressing from top to bottom. Observations within 0.2 length scales of the optimum were marked in red; the optimum was located on iteration 23.

under uncertainty. A multi-armed bandit problem can be interpreted as a particular finite-dimensional analog of sequential optimization, and effective algorithm design in both settings requires addressing many shared concerns.

The classical multi-armed bandit problem considers a finite set of “arms” \mathcal{X} and an agent who must select a sequence of items from this set. Selecting an arm x results in a stochastic reward y drawn from an unknown distribution $p(y | x)$ associated with that arm; these rewards are assumed to be independent of time and conditionally independent given the chosen arm. The goal of the agent is to select a sequence of arms $\{x_i\}$ to maximize the cumulative reward (6.2) received, $\sum y_i$.

Multi-armed bandits have been studied as a model of many sequential decision processes arising in practice. For example, consider a doctor caring for a series of patients with the same condition, who must determine which of two possible treatments is the best course of action. We could model the sequence of treatment decisions as a two-armed bandit, with patient outcomes determining the rewards. The Hippocratic oath compels the doctor to discover the optimal arm (the best treatment) as efficiently and confidently as possible to minimize patient harm, creating a dilemma of effective assignment. In fact, this scenario of sequential clinical trials was precisely the original motivation for studying multi-armed bandits.²⁶

²⁶ W.R. THOMPSON (1933). On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of Two Samples. *Biometrika* 25(3-4):285–294.

To facilitate the following discussion, for each arm $x \in \mathcal{X}$, we define $\phi = \mathbb{E}[y | x]$ to be its expected reward and will aggregate these into a vector \mathbf{f} when convenient. We also define

$$x^* \in \arg \max_{x \in \mathcal{X}} \mathbb{E}[y | x] = \arg \max \mathbf{f}; \quad f^* = \max_{x \in \mathcal{X}} \mathbb{E}[y | x] = \max \mathbf{f}$$

to be the index of an arm with maximal expected reward and the value of that optimal reward, respectively.²⁷ If the reward distributions associated with each arm were known a priori, the optimal policy would be trivial: we would always select the arm with the highest expected reward. This policy generates expected reward f^* in each iteration, and it is clear from linearity of expectation that this is optimal. Unfortunately, the reward distributions are *unknown* to the agent and must be learned from observations instead. This complicates policy design considerably.

The only way we can learn about the reward distributions is to allocate resources to each arm and observe the outcomes. If the reward distributions have considerable spread and/or overlap with each other, a large number of observations may be necessary before the agent can confidently conclude which arm is optimal. The agent thus faces an exploration–exploitation dilemma, constantly forced to decide whether to select an arm believed to have high expected reward (exploitation) or whether to sample an uncertain arm to better understand its reward distribution (exploration). Ideally, the agent would have a policy that *efficiently* explores the arms, so that in the limit of many decisions the agent would eventually allocate an overwhelming majority of resources to the best possible alternative.

Dozens of policies for the multi-armed bandit problem have been proposed and studied from both the Bayesian and frequentist perspectives. Numerous variations on the basic formulation outlined above have also received consideration in the literature, and the interested reader may refer to one of several available exhaustive surveys for more information.^{28,29,30}

The Bayesian optimal policy

A multi-armed bandit is fundamentally a sequential decision problem under uncertainty, and we may derive an optimal expected-case policy following our discussion in chapter 5. The selection of each arm is a decision with action space \mathcal{X} , and we must act under uncertainty about the expected reward vector \mathbf{f} . Over the course of τ decisions, we will gather a dataset $\mathcal{D}_\tau = (\mathbf{x}_\tau, \mathbf{y}_\tau)$, seeking to maximize the cumulative reward (6.2): $u(\mathcal{D}_\tau) = \sum y_i$.

The key to the Bayesian approach is maintaining a belief about the expected reward of each arm. We begin by choosing a prior over the expected rewards, $p(\mathbf{f})$, and an observation model for the observed rewards given the index of an arm and its expected reward, $p(y | x, \phi)$.³¹ Now given an arbitrary set of previous observations \mathcal{D} , we may derive a posterior belief about the expected rewards, $p(\mathbf{f} | \mathcal{D})$.

optimal policy with known rewards

²⁷ All of the notation throughout this section is chosen to align with that for optimization. In a multi-armed bandit, an arm x is associated with expected reward $\phi = \mathbb{E}[y | x]$. In optimization with zero-mean additive noise, a point x is associated with expected observed value $\phi = f(x) = \mathbb{E}[y | x]$.

challenges in policy design

²⁸ D. A. BERRY and B. FRISTEDT (1985). *Bandit Problems: Sequential Allocation of Experiments*. Chapman & Hall.

²⁹ S. BUBECK and N. CESA-BIANCHI (2012). Regret Analysis of Stochastic and Nonstochastic Multi-armed Bandit Problems. *Foundations and Trends in Machine Learning* 5(1):1–122.

³⁰ T. LATTIMORE and C. SZEPESVÁRI (2020). *Bandit Algorithms*. Cambridge University Press.

belief about expected rewards

³¹ This model is often conditionally independent of the arm given the expected reward, allowing the definition of a single observation model $p(y | \phi)$.

optimal policy: § 5.2, p. 91

³² See § 7.10 for an analogous result in optimization.

running time of optimal policy: § 5.3, p. 99

³³ R. AGRAWAL (1995). The Continuum-Armed Bandit Problem. *SIAM Journal on Control and Optimization* 33(6):1926–1951.

correlation of rewards

³⁴ In the multi-armed bandit literature, bandits with correlated rewards are known as *restless bandits*, as our belief about arms may change even when they are not selected (left alone):

P. WHITTLE (1988). Restless Bandits: Activity Allocation in a Changing World. *Journal of Applied Probability* 25(A):287–298

cumulative vs. simple reward

The optimal policy may now be derived following our previous analysis. We make each decision by maximizing the expected reward by termination, recursively assuming optimal future behavior (5.15–5.17). Notably, the optimal decision for the *last* round is the arm maximizing the posterior mean reward, reflecting pure exploitation.³²

$$x_\tau \in \arg \max \mathbb{E}[f \mid \mathcal{D}_{\tau-1}].$$

More exploratory behavior begins with the penultimate decision and increases with the decision horizon.

Unfortunately, the cost of computing the optimal policy increases exponentially with the horizon. We must therefore find some mechanism to design computationally efficient but empirically effective policies for use in practice. This is precisely the same situation we face in optimization!

Optimization as an infinite-armed bandit

We may model continuous optimization as an *infinite-armed bandit* problem,³³ and this analogy has proven fruitful. Suppose we seek to optimize an objective function $f: \mathcal{X} \rightarrow \mathbb{R}$, where the domain \mathcal{X} is now infinite. We assume as usual that we can observe this function at any point x of our choosing, revealing an observation y with distribution $p(y \mid x, \phi); \phi = f(x)$. For the bandit analogy to be maximally appropriate, we will further assume that the expected value of this observation is ϕ : $\mathbb{E}[y \mid \phi] = \phi$; this is not unduly restrictive and is satisfied for example by zero-mean additive noise.

Assuming an evaluation budget of τ observations, we may formulate this scenario to a multi-armed bandit. We interpret each point $x \in \mathcal{X}$ as one of infinitely many arms, with each arm returning an expected reward determined by the underlying objective function value ϕ . With some care, we may now adapt a multi-armed bandit policy to this setting.

In the traditional multi-armed bandit problem, we assume that rewards are conditionally independent given the chosen arm index. As a consequence, the reward from any selected arm provides no information about the rewards of other arms. However, this independence would render an infinite-armed bandit hopeless, as we would never be able to determine the best arm with a finite budget. Instead, we must assume that the rewards are *correlated* over the domain, so that each observation can potentially inform us about the rewards of every other arm.³⁴

This assumption is natural in optimization; the objective function must reflect some nontrivial structure, or optimization would also be hopeless. In the Bayesian framework, we formalize our assumptions regarding the structure of correlations between function values by choosing an appropriate prior distribution, which we may condition on available data to form the posterior belief, $p(f \mid \mathcal{D})$. In our bandit analogy, this distribution encapsulates beliefs about the expected rewards of each arm that can be used to derive effective policies.

Why should we reduce optimization to the multi-armed bandit at all? Notably, in optimization we are usually concerned with identifying

a *single* point in the domain maximizing the function, and variations on the simple reward directly measure progress toward this end. It may seem odd to focus on maximizing the *cumulative* reward, which judges a dataset based on the *average* value observed rather than the maximum.

These aims are not necessarily incompatible. In the limit of many observations, we hope to guarantee the best arm will be eventually identified so that we can guarantee convergence to optimal behavior. Bandit algorithms are typically analyzed in terms of their *cumulative regret*, the difference between the cumulative reward received and that expected from the optimal policy. If this quantity decreases sufficiently quickly, we may conclude that the optimal arm is eventually identified and selected. In our infinite-armed case, this implies the global optimum of the objective will eventually be located,³⁵ suggesting the multi-armed bandit reduction is indeed reasonable.

cumulative regret and convergence: § 10.1, p. 214

³⁵ In the infinite-armed case, establishing the no-regret property (§ 10.1, p. 214) is not sufficient to guarantee the global optimum will ever be evaluated exactly; however, we can conclude that we evaluate points achieving objective values within any desired tolerance of the maximum.

7.8 MAXIMIZING A STATISTICAL UPPER BOUND

Effective strategies for both bandits and optimization require careful consideration of the exploration–exploitation dilemma for success. Numerous bandit algorithms have been built on the unifying principle of *optimism in the face of uncertainty*,³⁶ which has proven to be an effective heuristic for balancing these concerns. The key idea is to use any available data to both estimate the expected reward of each arm and to quantify the uncertainty in these estimates. When faced with a decision, we then always select the arm that would be optimal when allowing “the benefit of the doubt”: the arm with the highest *plausible* expected reward given the currently available information. Arms with highly uncertain reward will have a correspondingly wide range of plausible values, and this mechanism thus provides underexplored arms a so-called *exploration bonus*³⁷ commensurate with their uncertainty, encouraging exploration of plausible optimal locations.

To be more precise, assume we have gathered an arbitrary dataset \mathcal{D} , and consider an arbitrary point x . Consider the *quantile function* associated with the predictive distribution $p(\phi | x, \mathcal{D})$:³⁸

$$q(\pi; x, \mathcal{D}) = \inf \{\phi' | \Pr(\phi \leq \phi' | x, \mathcal{D}) \geq \pi\}.$$

We can interpret this function as a statistical *upper confidence bound* on ϕ : the value will exceed the bound only with tunable probability $1 - \pi$.

As a function of x , we can interpret $q(\pi; x, \mathcal{D})$ as an optimistic estimate of the entire objective function. The principle of optimism in the face of uncertainty then suggests observing where this upper confidence bound is maximized, yielding the acquisition function

$$\alpha_{\text{UCB}}(x; \mathcal{D}, \pi) = q(\pi; x, \mathcal{D}). \quad (7.18)$$

Figure 7.17 shows upper confidence bounds for our example scenario corresponding to three values of the confidence parameter π . Unlike the acquisition functions considered previously in this chapter, an upper

optimism in the face of uncertainty

³⁶ A. W. MOORE and C. G. ATKESON (1993). Memory-based Reinforcement Learning: Efficient Computation with Prioritized Sweeping. *NeurIPS 1992*.

³⁷ R. S. SUTTON (1990). Integrated Architectures for Learning, Planning, and Reacting Based on Approximating Dynamic Programming. *ICML 1990*.

³⁸ The quantile function satisfies the relation that $\phi \leq q(\pi; x, \mathcal{D})$ with probability π .

upper confidence bound

example and discussion

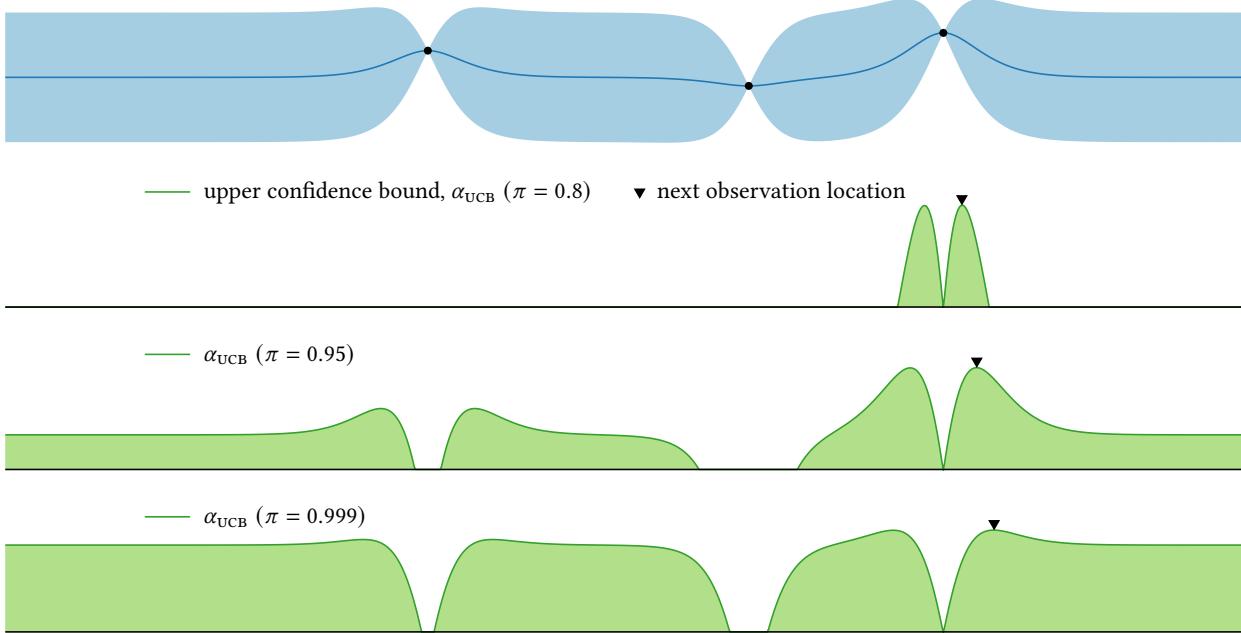


Figure 7.17: The upper confidence bound acquisition function (7.18) corresponding our running example for different values of the confidence parameter π . The vertical axis for each acquisition function is shifted to the largest observed function value. Increasing the confidence parameter leads to increasingly exploratory behavior.

correspondence between certainty parameter and exploration proclivity

adjusting the exploration parameter

confidence bound need not be nonnegative, so we shift the acquisition function in these plots so that the best-seen function value intersects the horizontal axis. We can see that relatively low confidence values ($\pi = 0.8$) give little credit to locations with high uncertainty, and exploitation is heavily favored. By increasing this parameter, our actions reflect more exploratory behavior.

In figure 7.18, we sequentially maximize the upper confidence bound on our example function to select 20 observation locations using confidence parameter $\pi = 0.999$, corresponding to the bottom and most exploratory example in figure 7.17. The global maximum was located efficiently. Notably, the observation locations chosen in the early stages of the search reflect more exploration than all the other methods we have discussed thus far.

Using an upper confidence bound policy in practice requires specifying the exploration parameter π , and it may not always be clear how best to do so. We face a similar challenge when choosing the improvement target parameter for probability of improvement, and in fact this analogy is sometimes remarkably intimate. For some models of the objective function, including Gaussian processes, a point maximizing the probability of improvement over a given threshold τ also maximizes an upper confidence bound for some confidence parameter π , and vice versa.³⁹ Therefore the sets of points obtained by maximizing these acquisition

³⁹ D. R. JONES (2001). A Taxonomy of Global Optimization Methods Based on Response Surfaces. *Journal of Global Optimization* 21(4):345–383

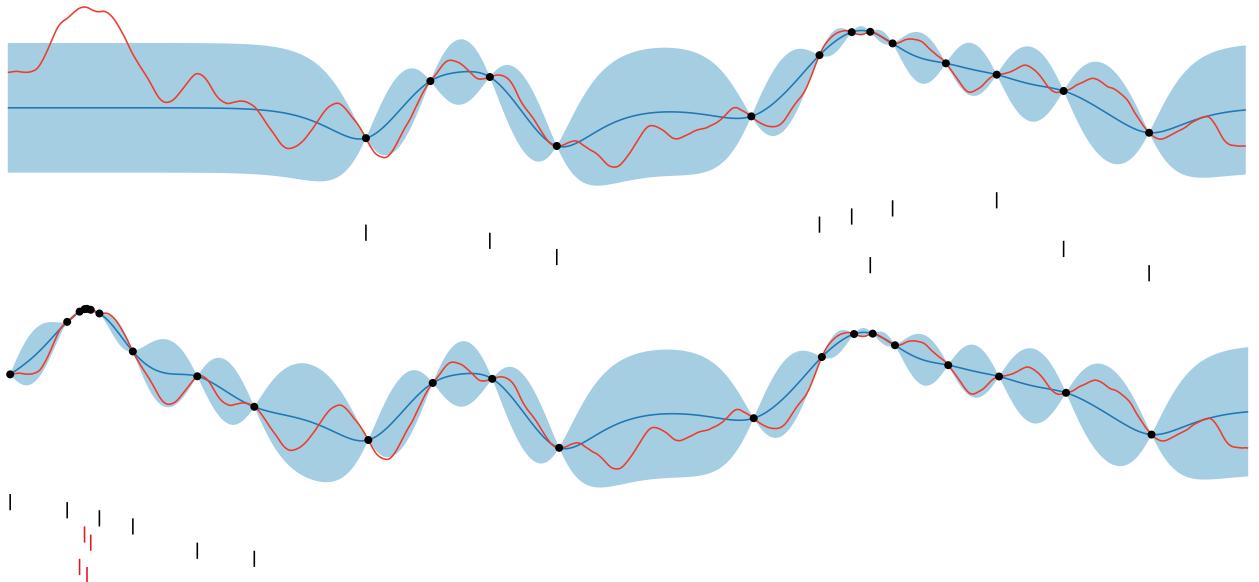


Figure 7.18: The posterior after 10 (top) and 20 (bottom) steps of the optimization policy induced by the upper confidence bound acquisition function (7.2) on our running example. The confidence parameter was set to $\pi = 0.999$. The tick marks show the points chosen by the policy, progressing from top to bottom, during iterations 1–10 (top) and 11–20 (bottom). Observations within 0.2 length scales of the optimum are marked in red; the optimum was located on iteration 15.

functions over the range of their respective parameters are identical. We will establish this relationship for Gaussian process models in the next chapter.

Little concrete advice is available for selecting the confidence parameter, as concerns such as model selection and calibration may have effects on the upper confidence bound that are hard to foresee and account for. Most authors select relatively large values in the approximate range $\pi \in (0.98, 1)$, with values of $\pi \approx 0.999$ being perhaps the most common. In line with his advice regarding probability of improvement parameter selection, JONES suggests considering a wide range of confidence values,⁴⁰ reflecting different exploration–exploitation tradeoffs. Figure 7.12 illustrates this concept by maximizing the probability of improvement for a range of improvement thresholds; by the correspondence between these policies for Gaussian process models, this is also illustrative for maximizing upper confidence bounds.

The policy realized by sequentially maximizing upper confidence bounds enjoys strong theoretical guarantees. For Gaussian process models, SRINIVAS et al. proved this policy is guaranteed to effectively maximize the objective at a nontrivial rate under reasonable assumptions.⁴¹ One of these assumptions is that the confidence parameter must increase asymptotically to 1 at a particular rate. Intuitively, the reason for this growth is that our uncertainty in the objective function will typically

equivalence of α_{PI} and α_{UCB} for GPS: § 8.4, p. 170

selecting improvement threshold: § 7.5, p. 134

⁴⁰ D. R. JONES (2001). A Taxonomy of Global Optimization Methods Based on Response Surfaces. *Journal of Global Optimization* 21(4):345–383.

theoretical analysis: chapter 10, p. 213

⁴¹ N. SRINIVAS et al. (2010). Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design. *ICML* 2010.

decrease as we continue to gather more data. As a result, we must simultaneously increase the confidence parameter to maintain a sufficient rate of exploration. This idea of slowly increasing the confidence parameter throughout optimization may be useful as a practical heuristic as well as a theoretical device.

7.9 THOMPSON SAMPLING

⁴² W. R. THOMPSON (1933). On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of Two Samples. *Biometrika* 25(3–4):285–294.

definition

mutual information with x^* : § 7.6, p. 139

alternative interpretation: optimizing a random acquisition function

⁴³ We give the sample a suggestive name!

In the early 20th century, THOMPSON proposed a simple and effective *stochastic* policy for the multi-armed bandit problem that has come to be known as *Thompson sampling*.⁴² Faced with a set of alternatives, the basic idea is to maintain a belief about which of these options is optimal in light of available information. We then design each evaluation by sampling from this distribution, yielding an adaptive stochastic policy. This procedure elegantly addresses the exploration–exploitation dilemma: sampling observations proportional to their probability of optimality automatically encourages exploitation, while the inherent randomness of the policy guarantees constant exploration. Thompson sampling can be adopted from finite-armed bandits to continuous optimization, and has enjoyed some interest in the Bayesian optimization literature.

Suppose we are at an arbitrary stage of optimization with data \mathcal{D} . The key object in Thompson sampling is the posterior distribution of the location of the global maximum x^* , $p(x^* | \mathcal{D})$, a distribution we have already encountered in our discussion on mutual information. Whereas maximizing mutual information carefully maximizes the information we expect to receive, Thompson sampling employs a considerably simpler mechanism. We choose the next observation location by sampling from this belief, yielding a nondeterministic optimization policy:

$$x \sim p(x^* | \mathcal{D}). \quad (7.19)$$

At first glance, Thompson sampling appears fundamentally different from the previous policies we have discussed, which were all defined in terms of maximizing an acquisition function. However, we may resolve this discrepancy while gaining some insight: Thompson sampling in fact designs each observation by maximizing a *random* acquisition function.

The location of the global maximum x^* is completely determined by the objective function f , and we may exploit this relationship to yield a simple two-stage implementation of Thompson sampling. We first sample a random realization of the objective function from its posterior:⁴³

$$\alpha_{\text{TS}}(x; \mathcal{D}) \sim p(f | \mathcal{D}). \quad (7.20)$$

We then optimize this function to yield the desired sample from $p(x^* | \mathcal{D})$, our next observation location:

$$x \in \arg \max_{x' \in \mathcal{X}} \alpha_{\text{TS}}(x'; \mathcal{D}).$$

From this point of view, we can interpret the sampled objective function α_{TS} as an ordinary acquisition function that is maximized as usual.

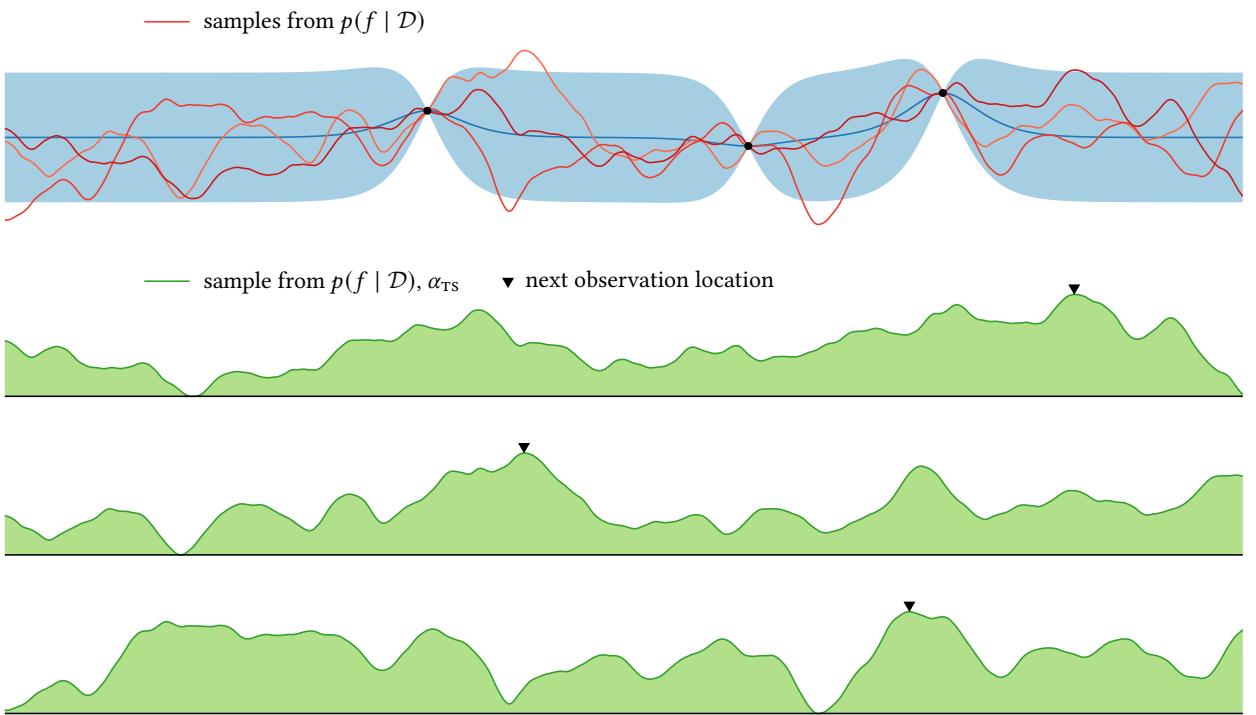


Figure 7.19: An illustration of Thompson sampling for our example optimization scenario. At the top, we show the objective function posterior $p(f | \mathcal{D})$ and three samples from this belief. Thompson sampling selects the next observation location by maximizing one of these samples. In the bottom three panels we show three possible outcomes of this process, corresponding to each of the sampled objective functions illustrated in the top panel.

Rather than representing an expected utility or a statistical upper bound, the acquisition function used in each round of Thompson sampling is a hallucinated objective function that is plausible under our belief. Whereas a Bayesian decision theoretic policy chooses the optimal action in expectation while *averaging* over the uncertain objective function, Thompson sampling chooses the optimal action for a *randomly sampled* objective function. This interpretation of Thompson sampling is illustrated for our example scenario in figure 7.19, showing three possible outcomes for Thompson sampling. In this case, two samples would exploit the region surrounding the best-seen point, and one would explore the region around the left-most observation.

Figure 7.20 shows the posterior belief of the objective after 15 rounds of Thompson sampling for our example scenario. The global maximum was located remarkably quickly. Of course, as a stochastic policy, it is not guaranteed that the behavior of Thompson sampling will resemble this outcome. In fact, this was a remarkably lucky run! The most likely locations of the optimum were ignored in the first rounds, quickly leading to the discovery of the optimum on the left. Figure 7.21 shows another

optimal decision for random sample

example and discussion

example of slow convergence

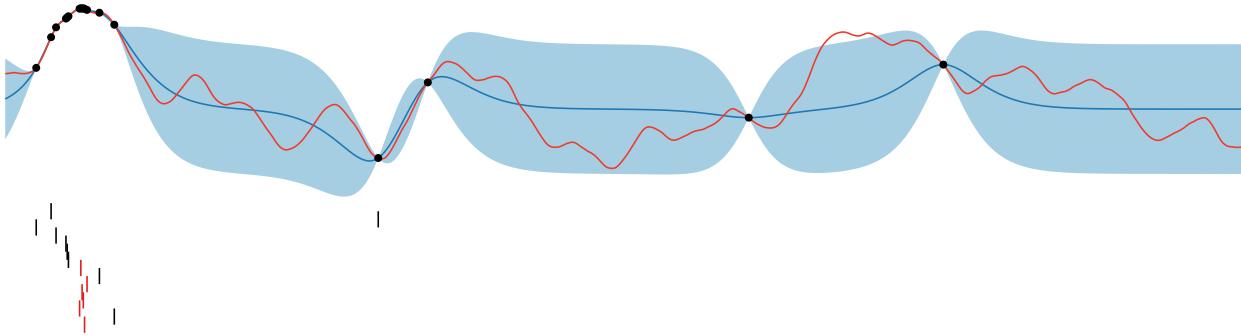


Figure 7.20: The posterior after 15 steps of Thompson sampling (7.19) on our running example. The tick marks show the points chosen by the policy, progressing from top to bottom. Observations within 0.2 length scales of the optimum are marked in red; the optimum was located on iteration 8.

run of Thompson sampling on the same scenario. The global optimum was not found nearly as quickly; however, it was eventually located after approximately 80 iterations. In 100 repetitions of the policy varying the random seed, the median iteration for discovering the optimum on this example was 38, with only 12 seeds resulting in discovery in the first 20 iterations. Despite this sometimes slow convergence, the distribution of the chosen evaluation locations nonetheless demonstrates continual management of the exploration–exploitation tradeoff.

7.10 OTHER IDEAS IN POLICY CONSTRUCTION

We have now discussed the two most pervasive approaches to policy construction in Bayesian optimization: one-step lookahead and adapting policies for multi-armed bandits. We have also introduced the most popular Bayesian optimization policies encountered in the literature, all of which stem from one of these two methods. However, there are some additional ideas worthy of discussion.

Approximate dynamic programming beyond one-step lookahead

One-step lookahead offers tremendous computational benefits, but the cost of these savings is extreme myopia in decision making. As we are oblivious to anything that might happen beyond the present observation, one-step lookahead can focus too much on exploitation. However, some less myopic alternatives have been proposed based on more complex (and more costly!) approximations to the optimal policy.

The simplest idea in this direction is to extend the lookahead horizon, and the most tractable option is of course two-step lookahead. Given an arbitrary one-step lookahead acquisition function $\alpha(x; \mathcal{D})$, the two-step analog is (5.12):

$$\alpha_2(x; \mathcal{D}) = \alpha(x; \mathcal{D}) + \mathbb{E} \left[\max_{x' \in \mathcal{X}} \alpha(x'; \mathcal{D}') \mid x, \mathcal{D} \right].$$

approximate dynamic programming: § 5.3,
p. 101

two-step lookahead

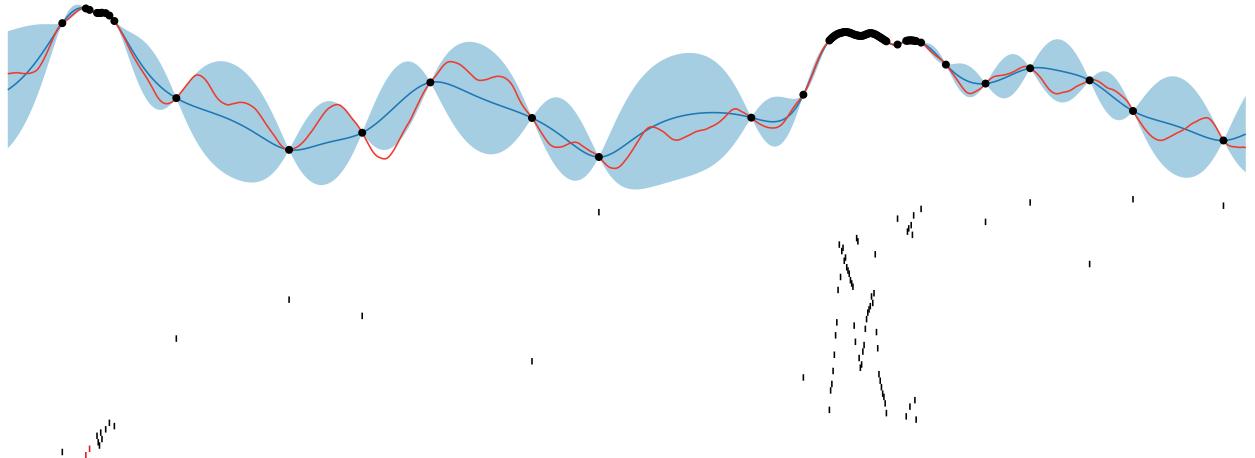


Figure 7.21: The posterior after 80 steps of Thompson sampling (7.19) on our running example, using a random seed different from figure 7.20. The global optimum was located on iteration 78.

Although this door is open for any of the decision-theoretic policies considered in this chapter, two-step expected improvement has received the most attention. OSBORNE et al. derived two-step expected improvement and demonstrated good empirical performance on some test functions compared with the one-step alternative.⁴⁴ However, it is telling that they restricted their investigation to a limited number of functions due to the inherent computational expense. GINSBOURGER and LE RICHE completed a contemporaneous exploration of two-step expected improvement and provided an explicit example showing superior behavior from the less myopic policy.⁴⁵ Recently, several authors have revisited (2+)-step lookahead and developed sophisticated implementation schemes rendering longer horizons more feasible.^{46,47}

We provided an in-depth illustration and deconstruction of two-step expected improvement for our example scenario in figures 5.2–5.3. Note that the two-step expected improvement is appreciable even for the (useless!) options of evaluating at the previously observed locations, as we can still make conscientious use of the following observation.

Figure 7.22 illustrates the progress of 20 evaluations designed by maximizing two-step expected improvement for our example scenario. Comparing with the one-step alternative in figure 7.4, the less myopic policy exhibits somewhat more exploratory behavior and discovered the optimum more efficiently – after 15 rather than 19 evaluations.

Rollout has also been considered as an approach to building nonmyopic optimization policies. Again the focus of these investigations has been on expected improvement (or the related knowledge gradient), but the underlying principles could be extended to other policies.

LAM et al. combined expected improvement with several steps of rollout, again maximizing expected improvement as the base policy.⁴⁸ The authors also proposed optionally adjusting the utility function through

⁴⁴ M. A. OSBORNE et al. (2009). Gaussian Processes for Global Optimization. *LION 3*.

⁴⁵ D. GINSBOURGER and R. LE RICHE (2010). Towards Gaussian Process-based Optimization with Finite Time Horizon. *MODA 9*.

⁴⁶ J. WU and P. I. FRAZIER (2019). Practical Two-Step Look-Ahead Bayesian Optimization. *NeurIPS 2019*

⁴⁷ S. JIANG et al. (2020b). Efficient Nonmyopic Bayesian Optimization via One-Shot Multi-Step Trees. *NeurIPS 2020*

rollout: § 5.3, p. 102

⁴⁸ R. R. LAM et al. (2016). Bayesian Optimization with a Finite Budget: An Approximate Dynamic Programming Approach. *NeurIPS 2016*.

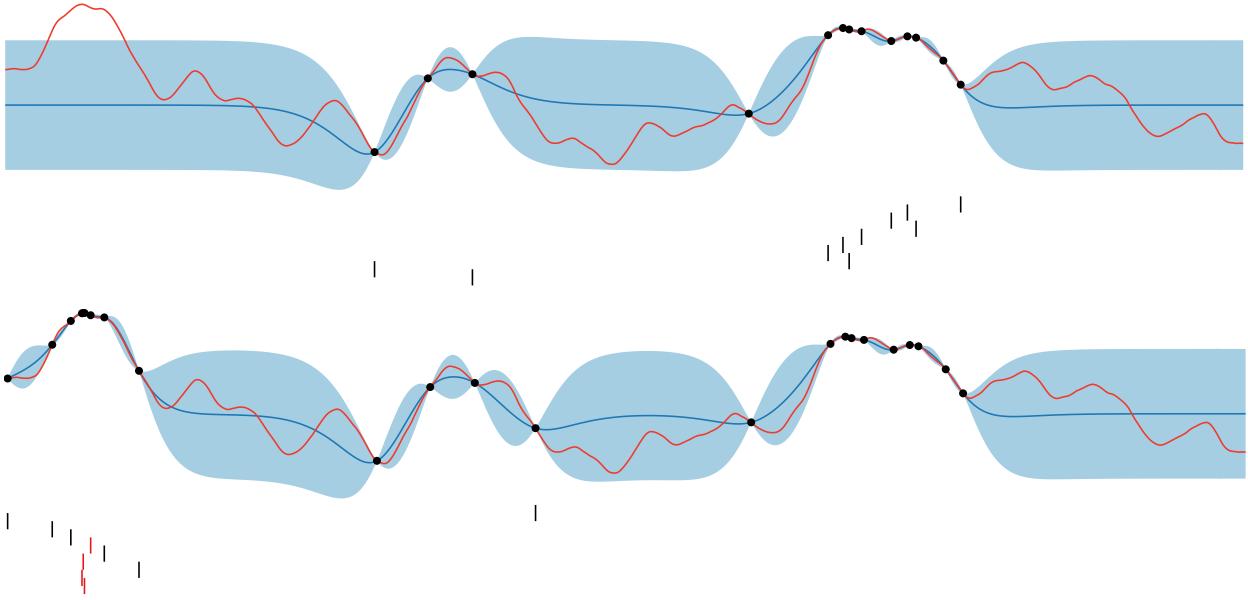


Figure 7.22: The posterior after 10 (top) and 20 (bottom) steps of the optimization policy induced by maximizing two-step expected improvement on our running example. The tick marks show the points chosen by the policy, progressing from top to bottom, during iterations 1–10 (top) and 11–20 (bottom). Observations within 0.2 length scales of the optimum are marked in red; the optimum was located on iteration 15.

⁴⁹ Such a discount factor is common in infinite-horizon decision problems:

D. P. BERTSEKAS (2017). *Dynamic Programming and Optimal Control*. Vol. 1. Athena Scientific.

⁵⁰ X. YUE and R. AL KONTAR (2020). Why Non-myopic Bayesian Optimization is Promising and How Far Should We Look-ahead? A Study via Rollout. *AISTATS 2020*.

⁵¹ See p. 125.

batch rollout: § 5.3, p. 103

⁵² J. GONZÁLEZ et al. (2016b). GLASSES: Relieving The Myopia Of Bayesian Optimization. *AISTATS 2016*.

⁵³ The policy used in GLASSES is described in

J. GONZÁLEZ et al. (2016a). Batch Bayesian Optimization via Local Penalization. *AISTATS 2016*,

as well as in § 11.3, p. 255; however, any desired alternative could also be used.

multiplicative discounting to encourage earlier rather than later progress during the rollout steps.⁴⁹ For some combinations of the rollout horizon and the discount factor, the resulting policies outperformed common one-step lookahead policies on a suite of synthetic test functions. YUE and AL KONTAR described a mechanism for dynamically choosing the rollout horizon based on the potential impact of a misspecified model,⁵⁰ exactly the issue that gave KUSHNER pause.⁵¹

The additional computational burden of rollout limited LAM et al. to a relatively short rollout horizon on the order of 4–5. Although considerably less myopic than one-step lookahead, the true decision horizon can be much greater, especially during the early stages of optimization. GONZÁLEZ et al. proposed an alternative approach based on batch rollout that can effectively look farther ahead at the expense of ignoring dependence among future decisions.⁵² The algorithm – dubbed GLASSES by the authors as it counteracts myopia – augments expected improvement with a single rollout step that designs a batch of additional observation locations of size equal to the remaining evaluation budget.⁵³ These points serve as a rough simulation of the decisions that might follow after making a proposed observation. A potential observation location is then evaluated by the expected improvement gained by simultaneously evaluating at that point as well as the constructed batch, realizing an efficient and budget-aware nonmyopic policy. GLASSES outperformed myopic and nonmyopic baselines on synthetic test functions, and the

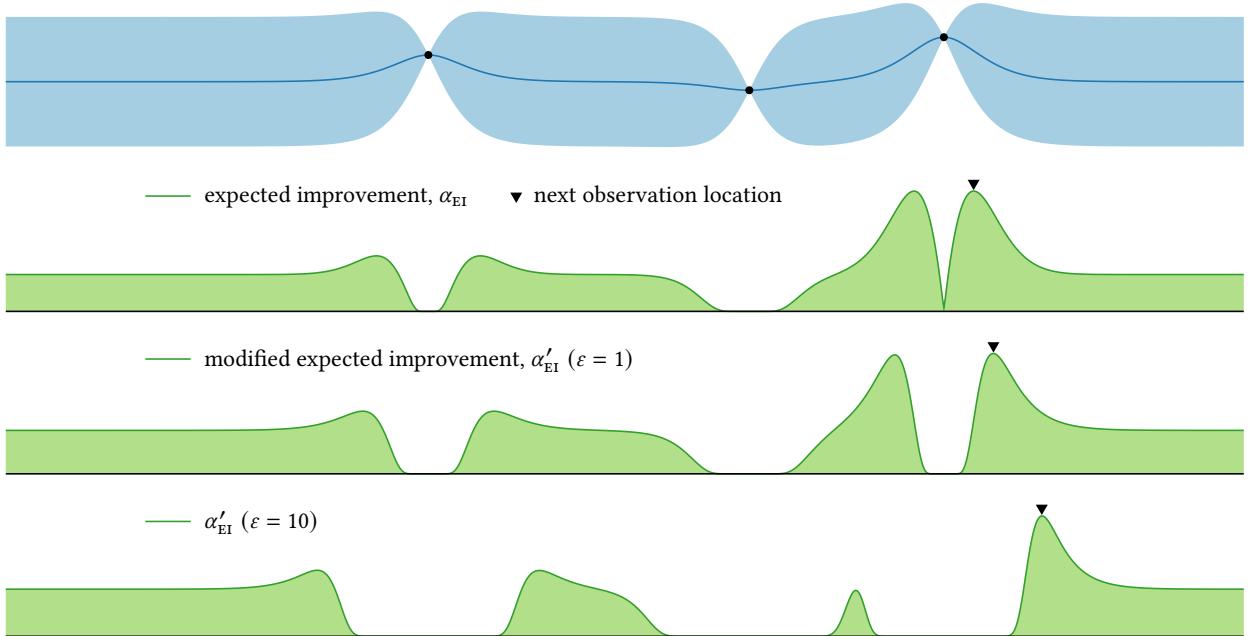


Figure 7.23: The modified expected improvement acquisition function (7.21) for our running example for different values of the target improvement ϵ . The target is expressed as a fraction of the range of the posterior mean over the space. Increasing the target improvement leads to increasingly exploratory behavior.

authors also demonstrated that dynamically setting the batch size to the remaining budget outperformed an arbitrary fixed size, suggesting that budget adaptation was important for success.

JIANG et al. continued this thread with an even more dramatic approximation dubbed BINOCULARS, potentially initiating an arms race toward increasingly nonmyopic acronyms.⁵⁴ The idea is to construct a *single* batch observation in each iteration, then select a point from this batch for evaluation. This represents an extreme computational savings over GLASSES, which must construct a batch anew for every proposed observation location. However, the method retains the same fundamental motivation: well-designed batch policies automatically induce *diversity* among batch members, encouraging exploration in the resulting sequential policy (see figure 11.3). The strong connection between the optimal batch and sequential policies (11.9–11.10) provides further motivation for this approach. JIANG et al. also conducted a study of optimization performance versus the cost of computing the policy: one-step lookahead, BINOCULARS, GLASSES, and rollout comprised the Pareto frontier, with each method increasing computational effort by an order of magnitude.

⁵⁴ S. JIANG et al. (2020a). BINOCULARS for Efficient, Nonmyopic Sequential Experimental Design. *ICML* 2020.

Pareto frontier: § 11.7, p. 267

Artificially encouraging exploration

Another more indirect approach to nonmyopic policy design is to modify a one-step lookahead acquisition function to artificially encourage more

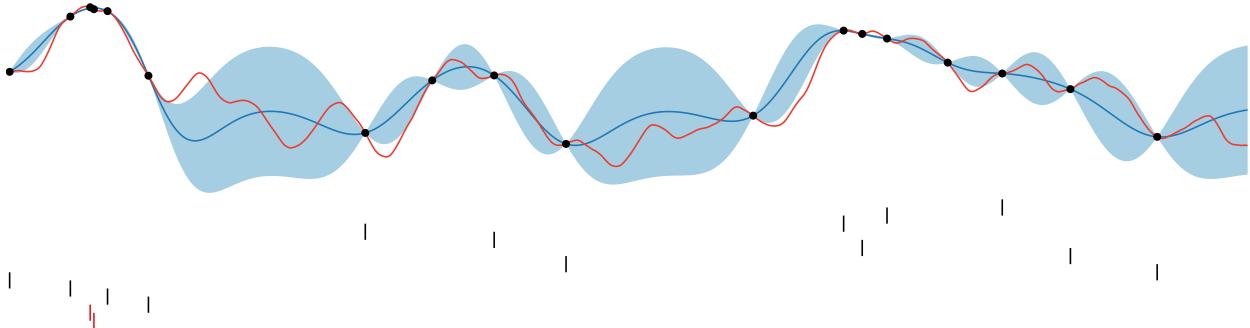


Figure 7.24: The posterior after 15 steps of the STRELTSOV and VAKILI optimization policy on our running example. The tick marks show the points chosen by the policy, progressing from top to bottom. Observations within 0.2 length scales of the optimum are marked in red; the optimum was located on iteration 14.

exploratory behavior when the remaining budget is significant. This can be motivated by the nature of the optimal policy, which maximizes a combination of immediate reward (exploitation) with expected future reward (exploration) (5.18). As the decision horizon increases, the exploration term may become increasingly dominant in this score, suggesting that optimal behavior entails early exploration of the domain followed by later exploitation and refinement.

Both the probability of improvement (7.5) and upper confidence bound acquisition functions already feature a parameter controlling the exploration-exploitation tradeoff, which can be dynamically maintained throughout optimization. This idea is quite old, reaching back to the earliest papers on Bayesian optimization. KUSHNER for example provided detailed advice on adjusting the threshold in probability of improvement to transition from early exploration to increasing levels of exploitation when appropriate.⁵⁵

In the case of exact observation, it is also possible to modify expected improvement (7.3) to incorporate a similar parameter. Rather than measuring improvement with respect to the utility of the current data $u(\mathcal{D}) = \phi^*$, we measure with respect to an inflated value $\phi^* + \varepsilon$, with no credit given for improvements less than this amount. The result is a modified expected improvement acquisition function:

$$\alpha'_{EI}(x; \mathcal{D}, \varepsilon) = \mathbb{E}[\max\{\phi - [\phi^* + \varepsilon], 0\} \mid x, \mathcal{D}]. \quad (7.21)$$

As with probability of improvement, larger improvement thresholds encourage increasing exploration, as illustrated in figure 7.23. For extreme values, the modified expected improvement drops to effectively zero except in regions with significant uncertainty. It is not obvious how this idea can be extended handle noisy observations, as the simple reward of the updated dataset may in fact decrease, raising the question of how to correctly define “sufficient improvement.”

MOCKUS proposed a scheme to set the threshold for Gaussian process models with constant mean and stationary covariance where the

⁵⁵ H. J. KUSHNER (1964). A New Method of Locating the Maximum Point of an Arbitrary Multi-peaked Curve in the Presence of Noise. *Journal of Basic Engineering* 86(1):97–106.

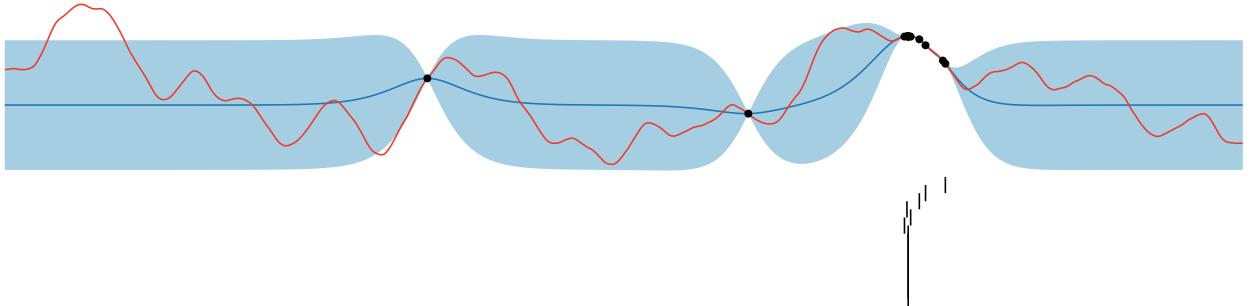


Figure 7.25: The posterior after 15 steps of two-step lookahead for cumulative reward (7.22) on our running example. The tick marks show the points chosen by the policy, progressing from top to bottom. The policy becomes stuck on iteration 7.

threshold was set dynamically based on the remaining budget, using the asymptotic behavior of the maximum of iid Gaussian random variables.⁵⁶ This can be interpreted as an approximate batch rollout policy where remaining decisions are simulated by fictitious uncorrelated observations; for some models this serves as an efficient simulation of random rollout.

The modified expected improvement (7.21) was also the basis for an unusual policy proposed by STRELTSOV and VAKILI.⁵⁷ Let $c: \mathcal{X} \rightarrow \mathbb{R}^{>0}$ quantify the cost of making an observation at any proposed location; in the simplest case we could take the cost to be constant. To evaluate the promise of making an observation at x , we solve the equation

$$\alpha'_{\text{EI}}(x; \mathcal{D}, \alpha_{\text{sv}}) = c(x)$$

for α_{sv} , which will serve as the acquisition function value at x .⁵⁸ That is, we solve for the improvement threshold that would render an observation at x cost-prohibitive in expectation, and design each observation to coincide with the last point to be ruled out when considering increasingly demanding thresholds. The resulting policy shows interesting behavior, at least on our running example; see figure 7.24. After effective initial exploration, the global optimum was located on iteration 14. The behavior is similar to the upper confidence bound approach in figure 7.18, and indeed STRELTSOV and VAKILI showed that the proposed method can be understood as a variation on this method with a location-dependent upper confidence quantile depending on observation cost and uncertainty.

Lookahead for cumulative reward?

Notably missing from our discussion on one-step lookahead policies was the cumulative reward utility function. Unfortunately, in this case one-step lookahead does not produce a particularly useful optimization policy. Suppose we adopt the cumulative reward utility function (6.7), $u(\mathcal{D}) = \sum y_i$. Then marginal gain in utility from a measurement at x is

⁵⁶ J. MOCKUS (1989). *Bayesian Approach to Global Optimization: Theory and Applications*. Kluwer Academic Publishers. [§ 2.5]

⁵⁷ S. STRELTSOV and P. VAKILI (1999). A Non-myopic Utility Function for Statistical Global Optimization Algorithms. *Journal of Global Optimization* 14(3):283–298.

⁵⁸ As α'_{EI} is monotonically decreasing with respect to α_{sv} and approaches zero as $\alpha_{\text{sv}} \rightarrow \infty$, the unique solution can be found efficiently via bisection.

posterior mean acquisition function

becoming “stuck”

two-step lookahead for cumulative reward

one-step lookahead: § 7.2, p. 124
 multi-armed bandits: § 7.7, p. 141

simply the observed value y . Therefore the expected one-step marginal gain is the posterior predictive mean:

$$\alpha(x; \mathcal{D}) = \mathbb{E}[y | x, \mathcal{D}],$$

which for zero-mean additive noise reduces to the posterior mean of f :

$$\alpha(x; \mathcal{D}) = \mu_{\mathcal{D}}(x).$$

From cursory inspection of our example scenario in figure 7.1, we can see that maximizing this acquisition function can cause the policy to become “stuck” with no chance of recovery. In our example, the posterior mean is maximized at the previously best-seen point, so the policy will select this point forevermore.

It is also interesting to consider two-step lookahead, where the acquisition function becomes (5.12):

$$\alpha_2(x; \mathcal{D}) = \mu_{\mathcal{D}}(x) + \mathbb{E}\left[\max_{x' \in \mathcal{X}} \mu_{\mathcal{D}'}(x') | x, \mathcal{D}\right].$$

After subtracting a constant (the global reward of the current data), we have

$$\alpha_2(x; \mathcal{D}) = \mu_{\mathcal{D}}(x) + \alpha_{\text{KG}}(x; \mathcal{D}), \quad (7.22)$$

the sum of the posterior mean and the knowledge gradient (7.4), reflecting both exploitation and exploration. Unfortunately even this less myopic option can become stuck due to overexploitation, as illustrated in Figure 7.25.

SUMMARY OF MAJOR IDEAS

Although we have covered a lot of ground in this chapter, there were really only two big ideas in policy construction: one-step lookahead and adopting successful policies from multi-armed bandits. However, there remains significant opportunity for novelty in this space.

8

COMPUTING POLICIES WITH GAUSSIAN PROCESSES

In the last chapter we introduced several notable Bayesian optimization policies in a model-agnostic setting, concentrating on their motivation and behavior while ignoring computational details. In this chapter we will provide further information for effectively implementing these policies. We will focus on Gaussian process models of the objective function, combined with either an exact or additive Gaussian noise observation model; this family accounts for the vast majority of models encountered in practice.

Implementing each policy in the previous chapter ultimately requires optimizing some acquisition function over the domain:¹

$$x \in \arg \max_{x' \in \mathcal{X}} \alpha(x'; \mathcal{D}).$$

We will demonstrate how to compute (or approximate) each of these acquisition functions with respect to Gaussian process models. Some will admit exact analytical expressions; when this is not possible, we will describe effective approximation schemes. In Euclidean domains, we will also show how to compute the gradient of these acquisition functions with respect to the proposed observation location, allowing efficient optimization via gradient methods. These gradient computations will sometimes be somewhat involved (but not difficult), and we will defer some details to an appendix for the sake of brevity.

The order of our presentation will differ from that in the previous chapter. Here we will begin with the acquisition functions for which exact computation is possible, then develop approximation techniques for those that remain. First, we pause to establish notation for important reoccurring quantities.

8.1 NOTATION FOR OBJECTIVE FUNCTION MODEL

As usual, let us consider an objective function $f: \mathcal{X} \rightarrow \mathbb{R}$, with a Gaussian process belief conditioned on arbitrary observations $\mathcal{D} = (\mathbf{x}, \mathbf{y})$:

$$p(f | \mathcal{D}) = \mathcal{GP}(f; \mu_{\mathcal{D}}, K_{\mathcal{D}}). \quad (8.1)$$

Our main task in this chapter will be to compute a given acquisition function at an arbitrary location $x \in \mathcal{X}$ with respect to this belief. Given a proposed location x , we will write the predictive distribution for $\phi = f(x)$ as:

$$p(\phi | x, \mathcal{D}) = \mathcal{N}(\phi; \mu, \sigma^2), \quad (8.2)$$

where the predictive mean and variance

$$\mu = \mu_{\mathcal{D}}(x); \quad \sigma^2 = K_{\mathcal{D}}(x, x) \quad (8.3)$$

depend implicitly on x . We will always treat x as given and fixed, so this convention will not lead to ambiguity.

Gaussian processes: chapter 2, p. 15
observation models: § 1.1, p. 4

¹ Even the nondeterministic Thompson sampling policy, which may be realized by optimizing a random acquisition function: § 7.9, p. 148.

gradients of common acquisition functions:
§ C.3, p. 304

Gaussian process on f , $\mathcal{GP}(f; \mu_{\mathcal{D}}, K_{\mathcal{D}})$

predictive distribution for ϕ , $\mathcal{N}(\phi; \mu, \sigma^2)$

predictive mean and variance for ϕ : μ, σ^2

We will also require the predictive distribution for the observed value y resulting from a measurement at x . In addition to the straightforward case of exact measurements, where $y = \phi$ and the predictive distribution is given above (8.2), we will also consider corruption by independent, zero-mean additive Gaussian noise:

$$p(y | \phi, \sigma_n) = \mathcal{N}(y; \phi, \sigma_n^2). \quad (8.4)$$

observation noise scale, σ_n
predictive distribution for y , $\mathcal{N}(y; \mu, s^2)$

predictive variance for y , s^2

prime notation for post-observation
quantities

general form of gradient and dependence on
parameter gradients

gradients of GP predictive distribution: § c.2,
p. 304

expected improvement: § 7.3, p. 127
simple reward: § 6.1, p. 112

Again we allow the noise scale σ_n to depend on x if desired. We will notate the resulting predictive distribution for y with:

$$p(y | x, \mathcal{D}, \sigma_n) = \mathcal{N}(y; \mu, \sigma^2 + \sigma_n^2) = \mathcal{N}(y; \mu, s^2), \quad (8.5)$$

where μ and σ^2 are the predictive moments of ϕ (8.3) and $s^2 = \sigma^2 + \sigma_n^2$ is the predictive variance for y , which again depends implicitly on x . When no distinction between exact and noisy observations is necessary, we will use the above general notation (8.5). With exact observations, the observation noise scale is identically zero, and we have $s^2 = \sigma^2$.

We will retain our convention from last chapter of indicating quantities available after acquiring a proposed observation with a prime symbol. For example $\mathbf{x}' = \mathbf{x} \cup \{x\}$ represents the updated set of observation locations after adding an observation at x , and $\mathcal{D}' = (\mathbf{x}', \mathbf{y}')$ represents the current data augmented with the observation (x, y) – a random variable.

The value of an acquisition function α at a point x will naturally depend on the distribution of the corresponding observation y , and its gradient must reflect this dependence. Applying the chain rule, the general form of the gradient will be written in terms of the gradient of the predictive parameters:

$$\frac{\partial \alpha}{\partial x} = \frac{\partial \alpha}{\partial \mu} \frac{\partial \mu}{\partial x} + \frac{\partial \alpha}{\partial s} \frac{\partial s}{\partial x}. \quad (8.6)$$

The gradients of the predictive mean and standard deviation for a Gaussian process are easily computable assuming the prior mean and covariance functions and the observation noise scale are differentiable, and general expressions are provided in an appendix.

8.2 EXPECTED IMPROVEMENT

We first consider the expected improvement acquisition function (7.2), the expected marginal gain in simple reward (6.3):

$$\alpha_{EI}(x; \mathcal{D}) = \mathbb{E}[\max \mu_{\mathcal{D}'}(\mathbf{x}') | x, \mathcal{D}] - \max \mu_{\mathcal{D}}(\mathbf{x}). \quad (8.7)$$

Remarkably, this expectation can be computed analytically for Gaussian processes with both exact and noisy observations. We will consider each case separately, as the former is considerably simpler and the latter involves minor controversy.

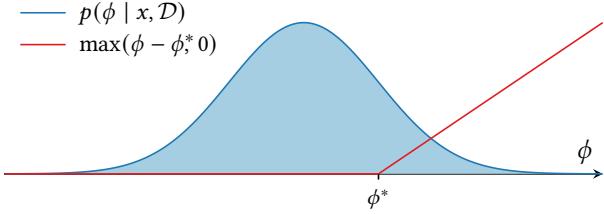


Figure 8.1: The expectation required to compute noiseless expected improvement: linear improvement for values exceeding the incumbent ϕ^* , integrated against a Gaussian distribution (8.2).

Expected improvement without noise

Expected improvement assumes a convenient form when measurements are exact (7.3):

$$\alpha_{\text{EI}}(x; \mathcal{D}) = \int \max(\phi - \phi^*, 0) \mathcal{N}(\phi; \mu, \sigma^2) d\phi. \quad (8.8)$$

Here ϕ^* is the previously best seen, *incumbent* objective function value, and $\max(\phi - \phi^*, 0)$ measures the improvement offered by observing a value of ϕ . Figure 8.1 illustrates this integral.

To proceed, we resolve the max operator to yield two integrals:

$$\alpha_{\text{EI}}(x; \mathcal{D}) = \int_{\phi^*}^{\infty} \phi \mathcal{N}(\phi; \mu, \sigma^2) d\phi - \phi^* \int_{\phi^*}^{\infty} \mathcal{N}(\phi; \mu, \sigma^2) d\phi,$$

both of which can be computed easily assuming $\sigma > 0$.² The first term is proportional to the expected value of a normal distribution truncated at ϕ^* , and the second term is the complementary normal CDF scaled by ϕ^* . The resulting acquisition function can be written conveniently in terms of the standard normal PDF and CDF:

$$\alpha_{\text{EI}}(x; \mathcal{D}) = (\mu - \phi^*) \Phi\left(\frac{\mu - \phi^*}{\sigma}\right) + \sigma \phi\left(\frac{\mu - \phi^*}{\sigma}\right). \quad (8.9)$$

Examining this expression, it is tempting to interpret its two terms as respectively encouraging exploitation (favoring points with high expected value μ) and exploration (favoring points with high uncertainty σ). Indeed, taking partial derivatives with respect to μ and σ , we have:

$$\frac{\partial \alpha_{\text{EI}}}{\partial \mu} = \Phi\left(\frac{\mu - \phi^*}{\sigma}\right) > 0; \quad \frac{\partial \alpha_{\text{EI}}}{\partial \sigma} = \phi\left(\frac{\mu - \phi^*}{\sigma}\right) > 0.$$

Expected improvement is thus monotonically increasing in both μ and σ . Increasing a point's expected value naturally makes the point more favorable for exploitation, and increasing its uncertainty makes it more favorable for exploration. Either action would increase the expected improvement. The tradeoff between these two concerns is considered automatically and is reflected in the magnitude of the derivatives above.

Maximization of expected improvement in Euclidean domains may be guided by its gradient with respect to the proposed evaluation location x . Using the results above and applying the chain rule, we have (8.6):

$$\frac{\partial \alpha_{\text{EI}}}{\partial x} = \Phi\left(\frac{\mu - \phi^*}{\sigma}\right) \frac{\partial \mu}{\partial x} + \phi\left(\frac{\mu - \phi^*}{\sigma}\right) \frac{\partial \sigma}{\partial x},$$

incumbent function value, ϕ^*

² In the degenerate case $\sigma = 0$, we simply have $\alpha_{\text{EI}}(x; \mathcal{D}) = \max(\mu - \phi^*, 0)$.

exploitation and exploration

partial derivatives with respect to predictive distribution parameters

gradient of expected improvement without noise

which expresses the gradient in terms of the implicit exploration–exploitation tradeoff parameter and the change in the predictive distribution.

Expected improvement with noise

Computing expected improvement with noisy observations is somewhat more complicated than in the noiseless case. In fact, there is not even universal agreement regarding what the definition of noisy expected improvement should be! The situation is so nuanced that JONES et al. sidestepped the issue entirely in their landmark paper:³

³ D. R. JONES et al. (1998). Efficient Global Optimization of Expensive Black-Box Functions. *Journal of Global Optimization* 13(4):455–492.

Unfortunately it is not immediately clear how to extend our *optimization algorithm* to the case of noisy functions. With noisy data, we really want to find the point where the *signal* is optimized. Similarly, our expected improvement criterion should be defined in terms of the signal component. [emphasis added by JONES et al.]

This summarizes the main challenge to optimization with noisy observations: how can we determine whether a particularly high observed value reflects a true underlying effect or is merely an artifact of noise? Several alternative definitions and heuristics have been proposed to address this question, which we will discuss further shortly.

We will first argue that seeking to maximize the simple reward utility (6.3) is *precisely* aligned with the goal outlined by JONES et al. With appropriate modeling of observation noise, the objective function posterior exactly represents our belief about the “signal component”: the latent objective f . Simple reward evaluates progress directly with respect to this belief, only ascribing merit to observations that improve the maximum of the posterior mean and thus the expected outcome of an optimal terminal recommendation. Notably, an excessively noisy observation has weak correlation with the underlying objective function value and thus yields little change in the posterior mean (2.12). Therefore even an extremely high outcome would produce only a minor improvement to the simple reward. As a result, the expected improvement of such a point would be relatively small, exactly the desired behavior.

Unfortunately, observation noise renders expected improvement somewhat more complicated than the exact case, due to the nature of the updated simple reward. After an exact observation, the simple reward can only be achieved at one of two locations: either the observed point or the incumbent. However, with noisy observations, inherent uncertainty in the objective function implies that the updated simple reward could be achieved *anywhere*, including a point that previously appeared suboptimal. We must account for this possibility in the computation, increasing its complexity. Fortunately, exact computation is still possible for Gaussian processes and additive Gaussian noise by adopting a procedure originally described by FRAZIER et al. for computing the knowledge gradient on a discrete domain;⁴ we will return to this method in our discussion on knowledge gradient in the next section.

behavior of posterior moments: § 2.2, p. 21

any point may maximize the posterior mean

⁴ P. FRAZIER et al. (2009). The Knowledge-Gradient Policy for Correlated Normal Beliefs. *INFORMS Journal on Computing* 21(4):599–613.

If we define $\mu^* = \max \mu_{\mathcal{D}}(\mathbf{x})$ to represent the simple reward of the current data, then we must compute:

$$\alpha_{EI}(x; \mathcal{D}) = \int [\max \mu_{\mathcal{D}'}(\mathbf{x}') - \mu^*] \mathcal{N}(y; \mu, s^2) dy.$$

We first reduce this computation to an expectation of the general form

$$g(\mathbf{a}, \mathbf{b}) = \int \max(\mathbf{a} + \mathbf{b}z) \phi(z) dz, \quad (8.10)$$

where $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ are arbitrary vectors and z is a standard normal random variable.⁵ Note that given the observation y , the updated posterior mean at \mathbf{x}' is a vector we may compute in closed form (2.19):

$$\mu_{\mathcal{D}'}(\mathbf{x}') = \mu_{\mathcal{D}}(\mathbf{x}') + \frac{K_{\mathcal{D}}(\mathbf{x}', x)}{s} \frac{y - \mu}{s}.$$

This update is linear in y . Applying the transformation $y = \mu + sz$ yields

$$\mu_{\mathcal{D}'}(\mathbf{x}') = \mathbf{a} + \mathbf{b}z, \quad (8.11)$$

where

$$\mathbf{a} = \mu_{\mathcal{D}}(\mathbf{x}'); \quad \mathbf{b} = \frac{K_{\mathcal{D}}(\mathbf{x}', x)}{s}, \quad (8.12)$$

and we may express expected improvement in the desired form:

$$\alpha_{EI}(x; \mathcal{D}) = g(\mathbf{a}, \mathbf{b}) - \mu^*. \quad (8.13)$$

As a function of z , $\mathbf{a} + \mathbf{b}z$ is a set of lines with intercepts and slopes given by the entries of the \mathbf{a} and \mathbf{b} vectors, respectively. In the context of expected improvement, these lines represent the updated posterior mean values for each point of interest \mathbf{x}' as a function of the z -score of the noisy observation y . See figure 8.2 for an illustration. Note that the points with the highest correlation with the proposed point have the greatest slope (8.12), as our belief at these locations will be strongly affected by the outcome.

Now $\max(\mathbf{a} + \mathbf{b}z)$ is the upper envelope of these lines, a convex piecewise linear function shown in 8.3. The interpretation of this envelope is the simple reward of the updated dataset given the z -score of the noisy observation, which will be achieved at some point in \mathbf{x}' . For this example, the updated posterior mean could be maximized at one of four locations: either at one of the points on the far right given a relatively high observation, or at a backup point farther left given a relatively low observation. Note that in the latter case the simple reward will decrease.⁶

With this geometric intuition in mind, we can deduce that g is invariant to transformations that do not alter the upper envelope. In particular, g is invariant both to reordering the lines by applying an identical permutation to \mathbf{a} and \mathbf{b} and also to the deletion of lines that never dominate. In the interest of notational simplicity, we will take advantage of these invariances and only consider evaluating g when every given line achieves

current simple reward, μ^*

reduction to evaluation of $g(\mathbf{a}, \mathbf{b})$

⁵ The dimension of \mathbf{a} and \mathbf{b} can be arbitrary as long as they are equal.

definition of \mathbf{a}, \mathbf{b}

geometric intuition of g

⁶ However, the *expected* marginal gain is always positive.

invariance to permutation

invariance to deletion of always dominated lines

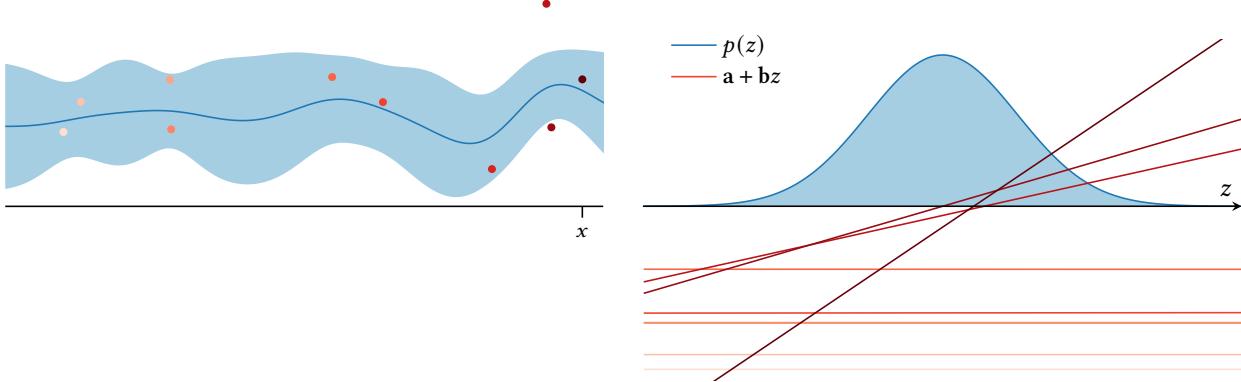


Figure 8.2: The geometric intuition of the $g(\mathbf{a}, \mathbf{b})$ function. Left: if we make a measurement at x , the z -score of the observed value completely determines the updated posterior mean at that point and all previously observed points. Right: as a function of z , the updated posterior mean at each of these points is linear; here the color of each line corresponds to the matching point on the left. The slope and intercept of each line can be determined from the posterior (8.12). Not all lines are visible.

maximal value on some interval and the lines appear in strictly increasing order of slope:

$$b_1 < b_2 < \dots < b_n.$$

⁷ Briefly, we sort the lines in ascending order of slope, then add each line in turn to a set of dominating lines, checking whether any previously added lines need to be removed and updating the intervals of dominance.

FRAZIER et al. give a simple and efficient algorithm to process a set of n lines to eliminate any always-dominated lines, reorder the remainder in increasing slope, and identify their intervals of dominance in $\mathcal{O}(n \log n)$ time.⁷ The output of this procedure is a permutation matrix \mathbf{P} , possibly with some rows deleted, such that

$$g(\mathbf{a}, \mathbf{b}) = g(\mathbf{Pa}, \mathbf{Pb}) = g(\boldsymbol{\alpha}, \boldsymbol{\beta}), \quad (8.14)$$

and the new inputs $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ satisfy the desired properties. We will assume below that the inputs have been preprocessed in such a manner.

Given a set of lines in the desired form, we may partition the real line into a collection of n intervals

$$(-\infty = c_1, c_2) \cup (c_2, c_3) \cup \dots \cup (c_n, c_{n+1} = +\infty), \quad (8.15)$$

⁸ Reordering the lines in order of increasing slope guarantees this correspondence: the line with minimal slope is always the “leftmost” in the upper envelope, etc.

such that the i th line $a_i + b_i z$ dominates on the corresponding interval (c_i, c_{i+1}) .⁸ This allows us to decompose the desired expectation (8.10) into a sum of contributions on each interval:

$$g(\mathbf{a}, \mathbf{b}) = \sum_i \int_{c_i}^{c_{i+1}} (a_i + b_i z) \phi(z) dz.$$

Finally, we may compute each integral in the sum in closed form:

$$g(\mathbf{a}, \mathbf{b}) = \sum_i a_i [\Phi(c_{i+1}) - \Phi(c_i)] + b_i [\phi(c_i) - \phi(c_{i+1})], \quad (8.16)$$

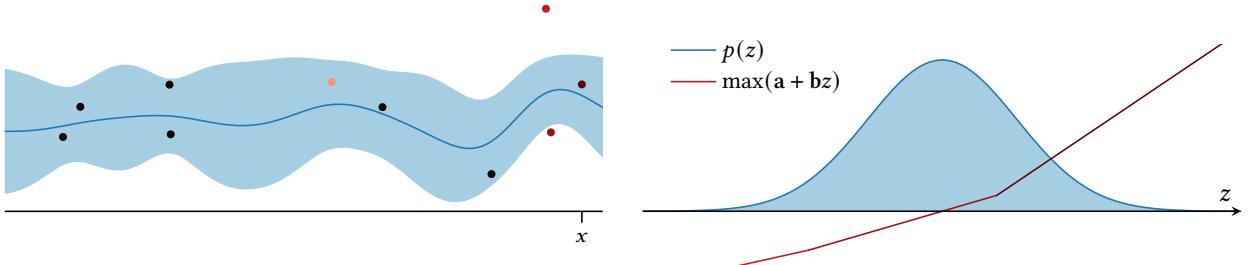


Figure 8.3: After a measurement at x , the updated simple reward can be achieved at one of four points (left), whose corresponding lines comprise the upper envelope $\max(a + bz)$ (right). The colors of the line segments on the right correspond to the possible updated maximum locations on the left. The lightest point on the left serves as a “backup option” if the observed value is low.

allowing efficient and exact computation of expected improvement in the noisy case. The main bottleneck is the modest $\mathcal{O}(n \log n)$ preprocessing step required.

This expression reverts to that for exact measurements (8.9) in the absence of noise. In that case, we have $\mathbf{b} = [\mathbf{0}, s]^\top$ and the upper envelope only contains lines corresponding to the incumbent and newly observed point, which intersect at the incumbent value $c_1 = \phi^*$. Geometrically, the situation collapses to that in figure 8.9, and it is easy to confirm (8.9) and (8.16) coincide.

Although tedious, we may compute the gradient of g and thus the gradient of expected improvement (8.13); the details are in an appendix.

compatibility with noiseless case

gradient of noisy expected improvement:
§ C.3, p. 304

Alternative formulations of noisy expected improvement

Although thematically consistent and mathematically straightforward, our approach of computing expected improvement as the expected marginal gain in simple reward is not common and may appear peculiar to some readers.

Over the years, numerous authors have grappled with the best definition of noisy expected improvement. A typical approach is to begin with the convenient formula (8.9) in the noiseless regime, then work “backwards” by identifying potential issues with its application to noisy data and suggesting a heuristic correction. This strategy of “fixing” exact expected improvement is in opposition to our ground-up approach of first defining a well-grounded utility function and only then working out the expected marginal gain. PICHENY et al. provided a survey of such approximations – representing a total of eleven different acquisition strategies – accompanied by a thorough empirical investigation.⁹ Notably, the authors were familiar with the exact computation we outlined in the previous section and recommended it for use in practice.

One popular idea in this direction is to take the formula from the exact case (8.9) and substitute a plug-in estimate for the now unknown

⁹ V. PICHENY et al. (2013b). A benchmark of kriging-based infill criteria for noisy optimization. *Structural and Multidisciplinary Optimization* 48(3):607–626.

expected improvement with plug-in estimator of ϕ^*

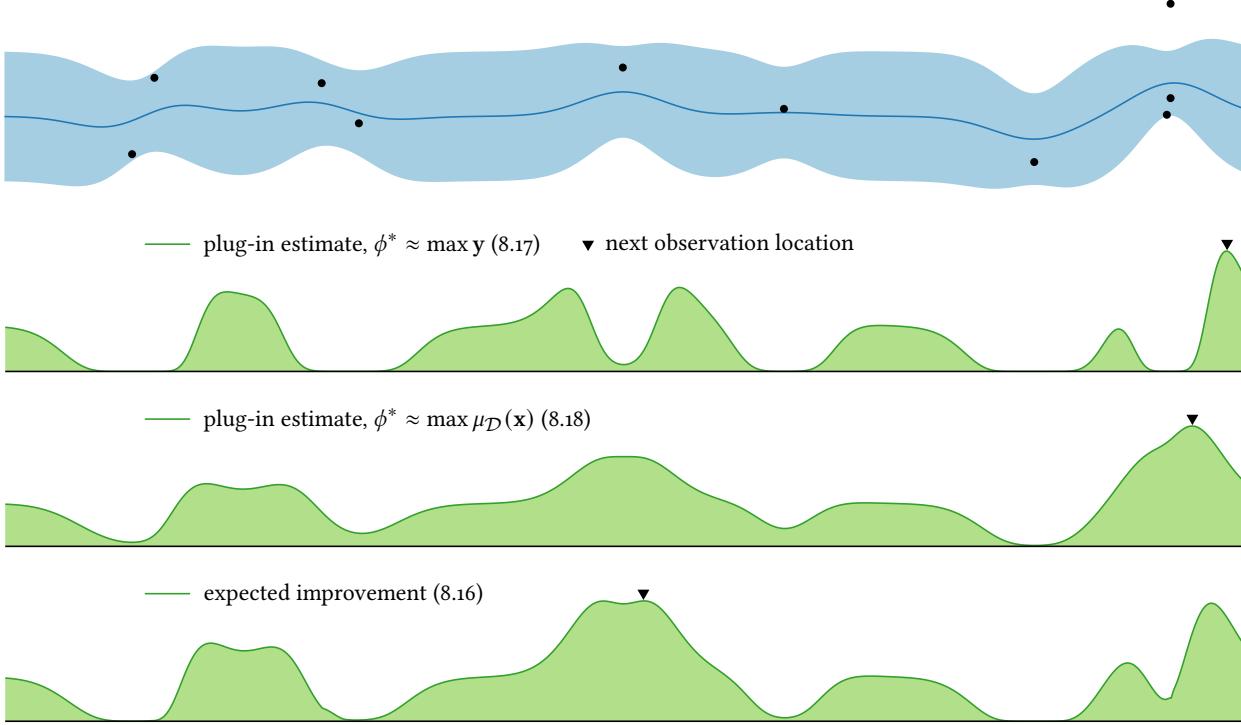


Figure 8.4: Expected improvement using different plug-in estimators (8.17–8.18) compared with the noisy expected improvement as the expected marginal gain in simple reward (8.7).

incumbent value ϕ^* . Several possibilities for this estimate have been put forward. One option is to plug in the maximum noisy observation:

$$\phi^* \approx \max y. \quad (8.17)$$

maximum noisy value “utility” function: § 6.1,
p. 113

inflating expected improvement threshold:
§ 7.10, p. 154

However, this may not always behave as expected for the same reason the maximum observed value does not serve as a sensible utility function (6.6). With very noisy data, the maximum observed value is most likely spurious rather than a meaningful goalpost. Further, as we are likely to overestimate our progress due to bias in this estimate, the resulting behavior may become excessively exploratory. The approximation will eventually devolve to expected improvement against an inflated threshold (7.21), which may overly encourage exploration; see figure 7.23. An especially spurious observation can bias the estimate in (8.17) (and our behavior) for a considerable time.

Opinions on the proposed approximation using this simple plug-in estimate (8.17) vary dramatically. PICHENY et al. discarded the idea out of hand as “naïve” and lacking robustness.¹⁰ The authors also found it empirically inferior in their investigation and described the same over-exploratory effect and explanation given above. On the other hand, NGUYEN et al. described the estimator as “standard” and concluded it was empirically preferable!¹¹

¹⁰ V. PICHENY et al. (2013b). A benchmark of kriging-based infill criteria for noisy optimization. *Structural and Multidisciplinary Optimization* 48(3):607–626.

¹¹ V. NGUYEN et al. (2017). Regret for Expected Improvement over the Best-Observed Value and Stopping Condition. *ACML 2017*.

8.2. EXPECTED IMPROVEMENT

An alternative estimator is the simple reward of the data (6.3) :^{12,13}

$$\phi^* \approx \max \mu_{\mathcal{D}}(\mathbf{x}), \quad (8.18)$$

which is less biased and may be preferable. A simple extension is to maximize other predictive quantiles,^{10,14} and HUANG et al. recommend using a relatively low quantile, specifically $\Phi(-1) \approx 0.16$, in the interest of risk aversion.

In figure 8.4 we compare noisy expected improvement with two plug-in approximations. The plug-in estimators agree that sampling on the right-hand side of the domain is the most promising course of action, but our formulation of noisy expected improvement prefers a less explored region. This decision is motivated by the interesting behavior of the posterior, which shows considerable disagreement regarding the updated posterior mean; see figure 8.6. This nuance is only revealed as our formulation reasons about the *joint* predictive distribution of \mathbf{y}' , whereas the plug-in estimators only inspect the marginals.

Another proposed approximation scheme for noisy expected improvement is reinterpolation. We fit a *noiseless* Gaussian process to imputed values of the objective function at the observed locations $\phi = f(\mathbf{x})$, then compute the exact expected improvement for this surrogate. A natural choice considered by FORRESTER et al. is to impute using the posterior mean:¹⁵

$$\phi \approx \mu_{\mathcal{D}}(\mathbf{x}),$$

resulting in the approximation (computed with respect to the surrogate):

$$\alpha_{\text{EI}}(\mathbf{x}; \mathcal{D}) \approx \alpha_{\text{EI}}(\mathbf{x}; \mathbf{x}, \phi). \quad (8.19)$$

This procedure is illustrated in figure 8.5. The resulting decision is very similar to that made by noisy expected improvement.

LETHAM et al. also promoted this basic approach, but proposed marginalizing rather than imputing the latent objective function values.¹⁶

$$\alpha_{\text{EI}}(\mathbf{x}; \mathcal{D}) \approx \int \alpha_{\text{EI}}(\mathbf{x}; \mathbf{x}, \phi) p(\phi | \mathbf{x}, \mathcal{D}) d\phi. \quad (8.20)$$

The approximate acquisition function is the expectation of the exact expected improvement if we had access to exact observations. Although this integral cannot be computed exactly, the authors described a straightforward and effective quasi-Monte Carlo approximation. LETHAM et al.'s approximation is illustrated in figure 8.6. There is good agreement with the expected improvement acquisition function from figure 8.4, *except* near the observation location chosen by that policy.

This stark disagreement is a consequence of reinterpolation: the acquisition function vanishes at previously observed locations – just as the exact expected improvement does – regardless of the observed values. Thus repeated measurements at the same location are barred, strongly encouraging exploration. HUANG et al. cite this property as undesirable,¹⁷ as repeated measurements can reinforce our understanding

approximating through reinterpolation

¹² E. VAZQUEZ et al. (2008). Global optimization based on noisy evaluations: an empirical study of two statistical approaches. *ICIP-E 2008*.

¹³ Z. WANG and N. DE FREITAS (2014). Theoretical Analysis of Bayesian Optimization with Unknown Gaussian Process Hyper-Parameters. arXiv: 1406.7758 [stat.ML].

¹⁴ D. HUANG et al. (2006b). Global Optimization of Stochastic Black-Box Systems via Sequential Kriging Meta-Models. *Journal of Global Optimization* 34(3):441–466.

¹⁵ A. I. J. FORRESTER et al. (2006). Design and Analysis of “Noisy” Computer Experiments. *AIAA Journal* 44(10):2331–2339.

¹⁶ B. LETHAM et al. (2019). Constrained Bayesian Optimization with Noisy Experiments. *Bayesian Analysis* 14(2):495–519.

reduction to zero at observed locations

¹⁷ D. HUANG et al. (2006b). Global Optimization of Stochastic Black-Box Systems via Sequential Kriging Meta-Models. *Journal of Global Optimization* 34(3):441–466.

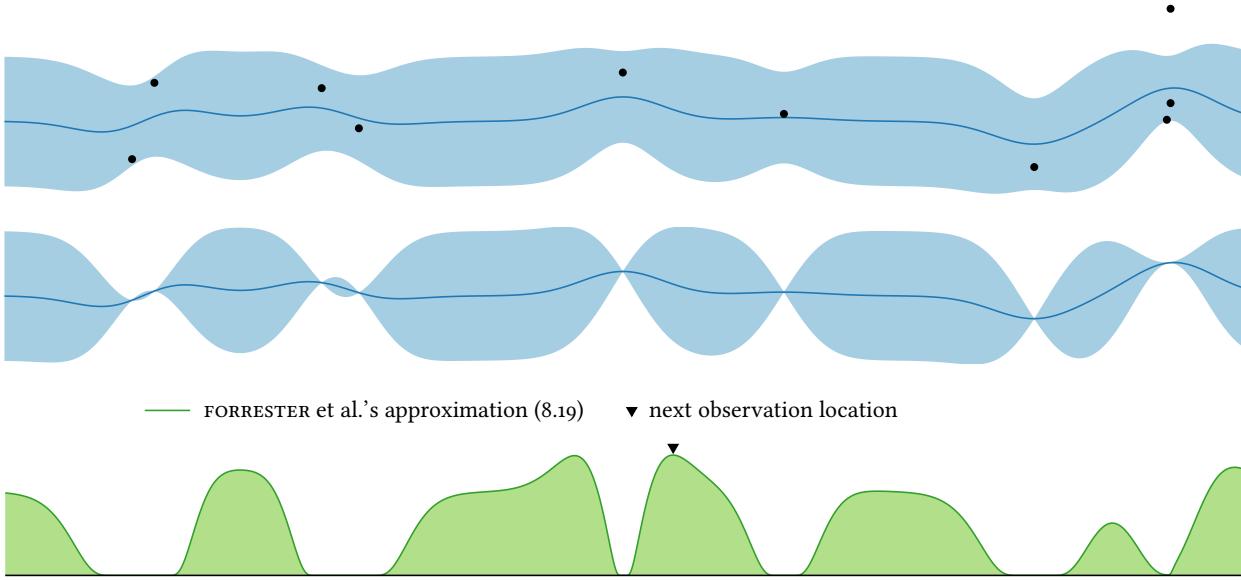


Figure 8.5: FORRESTER et al.'s approximation to noisy expected improvement (8.19). Given a Gaussian process fit to noisy data, we compute the exact expected improvement (8.9) for a noiseless Gaussian process fit to the posterior mean.

of the optimum by reducing uncertainty. LETHAM et al. on the other hand suggest the exploration boost is beneficial for optimization and point out we may reduce uncertainty through measurements in neighboring locations if desired.

A weakness shared by all these approximations is that the underlying noiseless expected improvement incorrectly assumes that our observation will reveal the *exact* objective value. In fact, observation noise is ignored entirely, as all these acquisition functions are expectations with respect to the *unobservable* quantity ϕ rather than y . This represents a disconnect between the reasoning of the optimization policy and the true nature of the observation process. HUANG et al. acknowledged this issue and proposed an *augmented expected improvement* measure accounting for observation noise¹⁸ by multiplying by the factor $1 - \sigma_n/s$,¹⁹ penalizing locations with low signal-to-noise ratios.

Ignoring the distribution of the observed value can be especially problematic with heteroskedastic noise, as shown in figure 8.7. Both the augmented¹⁸ and noisy expected improvement acquisition functions are biased towards the left-hand side of the domain, where observations reveal more information. Plug-in approximations on the other hand are oblivious to this distinction and elect to explore the noisy region on the right; see figure 8.4.

Our opinion is that all approximation schemes based on reducing to exact expected improvement should be avoided with significant noise levels, but can be a reasonable choice otherwise. With low noise, observed

¹⁸ D. HUANG et al. (2006b). Global Optimization of Stochastic Black-Box Systems via Sequential Kriging Meta-Models. *Journal of Global Optimization* 34(3):441–466.

¹⁹ This is only sensible when the signal-to-noise ratio is at least one.

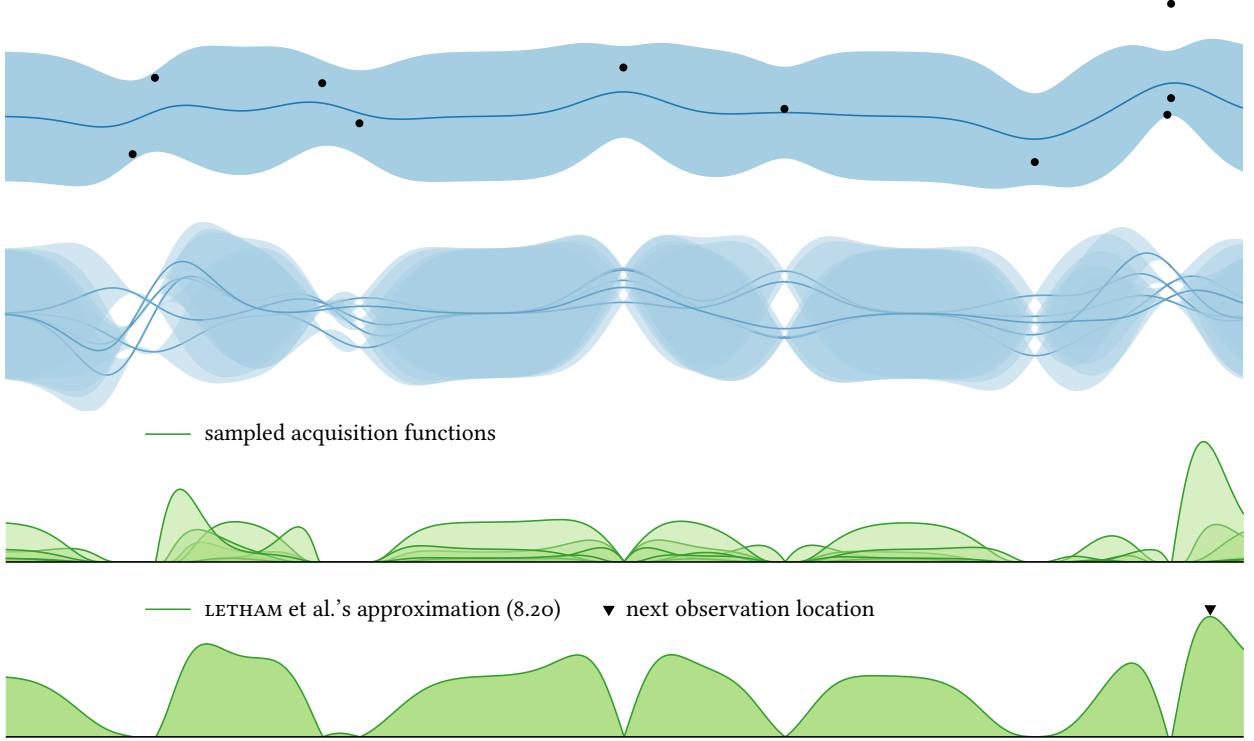


Figure 8.6: LETHAM et al.’s approximation to noisy expected improvement (8.20). We take the expectation of the exact expected improvement (8.9) for a noiseless Gaussian process fit to exact observations at the observed locations. The middle panels show realizations of the reinterpolated process and the resulting expected improvement.

values cannot stray too far from the true underlying objective function value. Thus we have $y \approx \phi$; $s \approx \sigma$, and any inaccuracy in (8.17), (8.18), or (8.20) will be minor. However, this heuristic argument breaks down in high-noise regimes.

8.3 PROBABILITY OF IMPROVEMENT

Probability of improvement (7.5) represents the probability that the simple reward of our data (6.3) will exceed a threshold τ after obtaining an observation at x :

probability of improvement: § 7.5, p. 131

$$\alpha_{\text{PI}}(x; \mathcal{D}) = \Pr(u(\mathcal{D}') > \tau \mid x, \mathcal{D}).$$

Like expected improvement, this quantity can be computed exactly for our chosen model class.

Probability of improvement without noise

The probability of improvement after an exact observation at x is simply the probability that the observed value will exceed τ (7.6). For a Gaussian

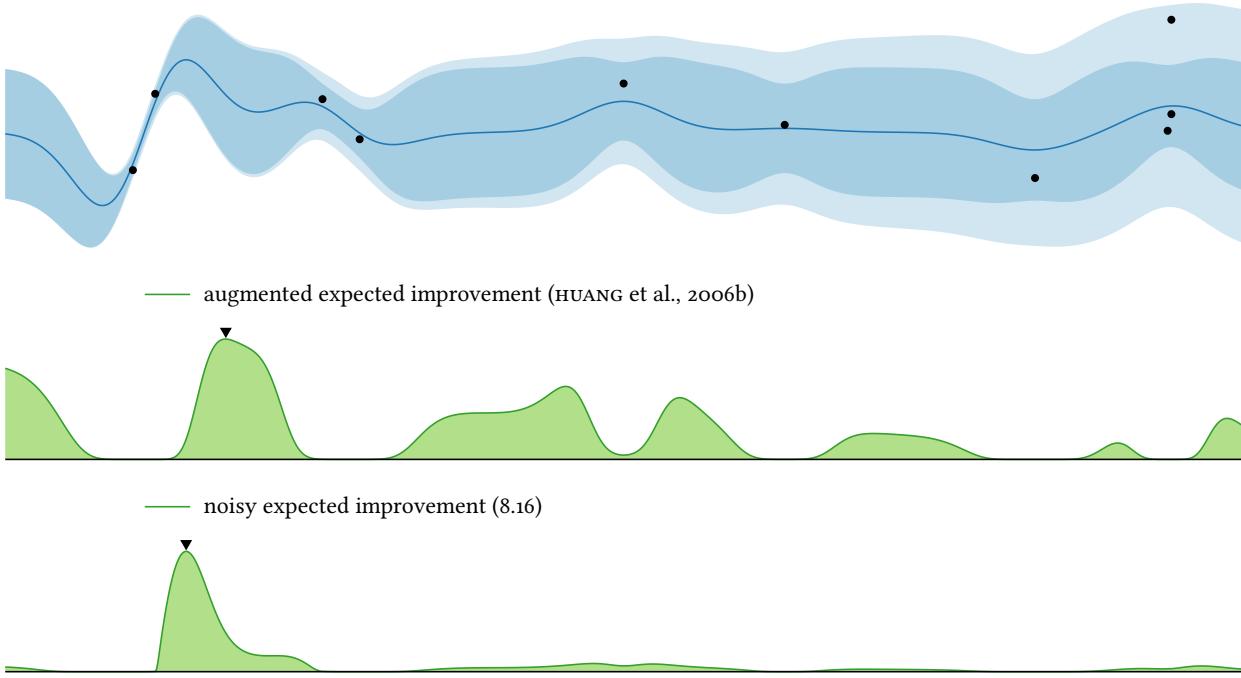


Figure 8.7: A comparison of noiseless expected improvement with a plug-in estimate for ϕ^* (8.17) and the expected one-step marginal gain in simple reward (7.2). Here the noise variance increases linearly from zero on the left-hand side of the domain to a signal-to-noise ratio of 1 on the right-hand side. The larger enveloping area shows 95% credible intervals for the noisy observation y , whereas the smaller area provides the same credible intervals for the objective function.

process, this probability is given by the complementary Gaussian CDF:

$$\alpha_{\text{PI}}(x; \mathcal{D}, \tau) = \Phi\left(\frac{\mu - \tau}{\sigma}\right). \quad (8.21)$$

equivalent acquisition function, α'_{PI}

As our policy will ultimately be determined by maximizing the probability of improvement, it is prudent to transform this expression into the simpler and better-behaved acquisition function

$$\alpha'_{\text{PI}}(x; \mathcal{D}, \tau) = \frac{\mu - \tau}{\sigma}, \quad (8.22)$$

which shares the same maxima but is slightly cheaper to compute and does not suffer from a vanishing gradient for extreme values of τ .

We may gain some insight into the behavior of probability of improvement by computing the gradient of this alternative expression (8.22) with respect to the parameters of the predictive distribution:

$$\frac{\partial \alpha'_{\text{PI}}}{\partial \mu} = \frac{1}{\sigma}; \quad \frac{\partial \alpha'_{\text{PI}}}{\partial \sigma} = \frac{\tau - \mu}{\sigma^2}.$$

We observe that probability of improvement is monotonically increasing with μ , universally encouraging exploitation. Probability of improvement

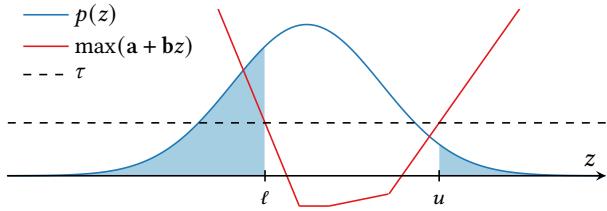


Figure 8.8: Probability of improvement with noise. The red curve represents the updated simple reward given the z -score of the observed value. The probability of improvement is the probability this value exceeds a threshold τ , the area of the shaded region.

is also increasing with σ when the target value is greater than the predictive mean, or equivalently when the probability of improvement is less than $1/2$. Therefore probability of improvement also tends to encourage exploration in this typical case. However, when the predictive mean is greater than the improvement threshold, probability of improvement is, perhaps surprisingly, *decreasing* with σ , discouraging exploration in this favorable regime. Probability of improvement favors a relatively safe option returning a certain but modest improvement over a more-risky alternative offering a potentially larger but less-certain improvement, as demonstrated in figure 7.8.

With the partial derivatives computed above, we may readily compute the gradient of (8.22) with respect to the proposed evaluation location x in the noiseless case (8.6):

$$\frac{\partial \alpha'_{\text{PI}}}{\partial x} = \frac{1}{\sigma} \left[\frac{\partial \mu}{\partial x} - \alpha'_{\text{PI}} \frac{\partial \sigma}{\partial x} \right].$$

Probability of improvement with noise

As with expected improvement, computing probability of improvement with noisy observations is slightly more complicated than with exact observations, and for the same reason. A noisy observation can affect the posterior mean at every previously observed location, and thus improvement to the simple reward can occur at any point. However, we can compute probability of improvement exactly by adapting the techniques we used to compute noisy expected improvement (8.16).

As before, the key observation is that the updated posterior mean is linear in the observed value y (8.11), allowing us to express the updated simple reward as

$$u(\mathcal{D}') = \max(a + bz),$$

where a and b are defined in (8.12) and z is the z -score of the observation. Now we can write the probability of improvement as the expectation of an indicator against a standard normal random variable; see figure 8.8:

$$\alpha_{\text{PI}}(x; \mathcal{D}) = \int [\max(a + bz) > \tau] \phi(z) dz.$$

Due to the convexity of the updated simple reward, improvement occurs on at-most two intervals:²⁰

$$(-\infty, \ell) \cup (u, \infty).$$

risk aversion of probability of improvement:
§ 7.5, p. 132

gradient of probability of improvement
without noise

²⁰ In the case that improvement is impossible or occurs at a single point, the probability of improvement is zero.

The endpoints of these intervals may be computed directly by inspecting the intersections of the lines $\mathbf{a} + \mathbf{b}z$ with the threshold:

$$\ell = \max_i \{(\tau - a_i)/b_i \mid b_i < 0\}; \quad u = \min_i \{(\tau - a_i)/b_i \mid b_i > 0\};$$

Note that one of these end points may not exist, for example if every slope in the \mathbf{b} vector were positive; see figure 8.3 for an example. In this case we may take $\ell = -\infty$ or $u = \infty$ as appropriate. Now given the endpoints (ℓ, u) , the probability of improvement may be computed in terms of the standard normal CDF:

$$\alpha_{\text{PI}}(x; \mathcal{D}) = \Phi(\ell) + \Phi(-u). \quad (8.23)$$

gradient of noisy probability of improvement:
§ C.3, p. 305

upper confidence bound: § 7.8, p. 145

exploration parameter, β

²¹ A word of warning: the exploration parameter is sometimes denoted $\sqrt{\beta}$, so that β is a weight on the variance σ^2 instead.

gradient of upper confidence bound

We may compute the gradient of this noisy formulation of probability of improvement; details are given in the appendix.

8.4 UPPER CONFIDENCE BOUND

Computing an upper confidence bound (7.18) for a Gaussian process is trivial. Given a confidence parameter $\pi \in (0, 1)$, we must compute a pointwise upper bound for the objective function with that confidence:

$$\alpha_{\text{UCB}}(x; \mathcal{D}, \pi) = q(\pi; x, \mathcal{D}), \quad (8.24)$$

where q is the quantile function of the predictive distribution. For a Gaussian process, this quantile takes the simple form

$$\alpha_{\text{UCB}}(x; \mathcal{D}, \pi) = \mu + \beta\sigma, \quad (8.25)$$

where $\beta = \Phi^{-1}(\pi)$ depends on the confidence level and can be computed from the inverse Gaussian CDF. This acquisition function is normally parameterized directly in terms of β rather than the confidence level π . In this case β can be interpreted as an “exploration parameter,” as higher values clearly reward uncertainty more than smaller values.²¹ No special care is required in computing (8.25), and its gradient with respect to the proposed observation location x can also be computed easily:

$$\frac{\partial \alpha_{\text{UCB}}}{\partial x} = \frac{\partial \mu}{\partial x} + \beta \frac{\partial \sigma}{\partial x}.$$

Correspondence with probability of improvement

²² D. R. JONES (2001). A Taxonomy of Global Optimization Methods Based on Response Surfaces. *Journal of Global Optimization* 21(4):345–383.

²³ Z. WANG et al. (2016a). Optimization as Estimation with Gaussian Processes in Bandit Settings. *AISTATS 2016*.

For Gaussian process models, we can derive an intimate correspondence between the probability of improvement and upper confidence bound acquisition functions. Namely, a point maximizing probability of improvement for a given target also maximizes some statistical upper bound of the objective and vice versa, a fact that has been noted by several authors, including JONES²² and, independently, WANG et al.²³

To establish this result, let an arbitrary exploration parameter β be given. Consider a point optimizing the upper confidence bound:

$$x \in \arg \max_{x' \in \mathcal{X}} \mu + \beta\sigma,$$

and define

$$\tau(\beta) = \max_{x' \in \mathcal{X}} \mu + \beta\sigma$$

equivalent improvement target, $\tau(\beta)$

to be its optimal value. Then the following equalities are satisfied at x :

$$\tau(\beta) = \mu + \beta\sigma; \quad \beta = \frac{\tau(\beta) - \mu}{\sigma}.$$

Now it is easy to show that x also minimizes the score

$$x \in \arg \min_{x' \in \mathcal{X}} \frac{\tau(\beta) - \mu}{\sigma}, \quad (8.26)$$

because if there were some other point x' with

$$\frac{\tau(\beta) - \mu'}{\sigma'} < \beta,$$

then we would have

$$\mu' + \beta\sigma' > \tau(\beta),$$

contradicting the optimality of x . Therefore x also maximizes probability of improvement with target $\tau(\beta)$ (8.22).

It is important to note that the value of this target $\tau(\beta)$ is data- and model-dependent and may change from iteration to iteration. Therefore sequentially maximizing an upper confidence bound with a fixed exploration parameter is not equivalent to maximizing probability of improvement with a fixed improvement target.

data-dependence of relationship

8.5 APPROXIMATE COMPUTATION FOR ONE-STEP LOOKAHEAD

Unfortunately we have exhausted the acquisition functions for which exact computation is possible with Gaussian process models. However, we can still proceed effectively with appropriate approximations. We will begin by discussing the implementation of arbitrary one-step lookahead policies when exact computation is not possible.

Recall that one-step lookahead entails maximizing the expected marginal gain to a utility function $u(\mathcal{D})$ after making an observation at a proposed location x :

$$\alpha(x; \mathcal{D}) = \int [u(\mathcal{D}') - u(\mathcal{D})] \mathcal{N}(y; \mu, s^2) dy.$$

one-step lookahead: § 5.3, p. 101

If this integral is intractable, we must resort to analytic approximation or numerical integration to evaluate and optimize the acquisition function. Fortunately in the case of sequential optimization, this is a one-dimensional integral that can be approximated using standard tools.

It will be convenient below to introduce simple notation for the marginal gain in utility resulting from a putative observation (x, y) . We will write

$$\Delta(x, y) = u(\mathcal{D}') - u(\mathcal{D})$$

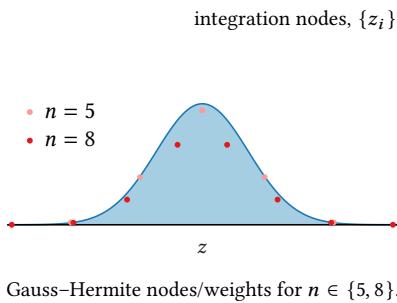
marginal gain from observation (x, y) ,
 $\Delta(x, y)$

for this quantity, leaving the dependence on \mathcal{D} implicit. Now we seek to approximate

$$\alpha(x; \mathcal{D}) = \int \Delta(x, y) \mathcal{N}(y; \mu, s^2) dy. \quad (8.27)$$

Gauss–Hermite quadrature

24 P. J. DAVIS and P. RABINOWITZ (1984). *Methods of Numerical Integration*. Academic Press.

Gauss–Hermite nodes/weights for $n \in \{5, 8\}$.

25 T. S. SHAO et al. (1964). Tables of Zeros and Gaussian Weights of Certain Associated Laguerre Polynomials and the Related Generalized Hermite Polynomials. *Mathematics of Computation* 18(88):598–616.

26 We take

$$y = \mu + \sqrt{2}sz \Leftrightarrow z = \frac{y - \mu}{\sqrt{2}s};$$

Note we must account for the normalization factor $(\sqrt{2}\pi\sigma)^{-1}$ of the Gaussian distribution, hence the constant that appears.

approximating gradient via Gauss–Hermite quadrature: § C.3, p. 305

knowledge gradient: § 7.4, p. 129

We recommend using off-the-shelf quadrature methods to estimate the expected marginal gain (8.27).²⁴ The most natural approach is *Gauss–Hermite quadrature*, a classical approach for approximating integrals of the form

$$I = \int h(z) \exp(-z^2) dz. \quad (8.28)$$

Like all numerical integration methods, Gauss–Hermite quadrature entails measuring the integrand h at a set of n points, called *nodes*, $\{z_i\}$, then approximating the integral by a weighted sum of the measured values:

$$I \approx \sum_{i=1}^n w_i h(z_i).$$

Remarkably, this estimator is *exact* when the integrand is a polynomial of degree less than $2n - 1$. This guarantee provides some guidance for selecting the order of the quadrature rule depending on how well the marginal gain Δ may be approximated by a polynomial over the range of plausible observations. Tables of integration nodes and quadrature weights for n up to 64 are readily available, although such a high order should rarely be needed.²⁵

Through a simple transformation,²⁶ we can rewrite the arbitrary Gaussian expectation (8.27) in the Gauss–Hermite form (8.28):

$$\int \Delta(x, y) \mathcal{N}(y; \mu, s^2) dy = \frac{1}{\sqrt{\pi}} \int \Delta(x, \mu + \sqrt{2}sz) \exp(-z^2) dz.$$

If we define appropriately renormalized weights

$$\bar{w}_i = w_i / \sqrt{\pi},$$

then we arrive at the following approximation to the acquisition function:

$$\alpha(x; \mathcal{D}) \approx \sum_{i=1}^n \bar{w}_i \Delta(x, y_i); \quad y_i = \mu + \sqrt{2}z_i s. \quad (8.29)$$

We may also extend this scheme to approximate the gradient of the acquisition function as well; details are provided in the accompanying appendix.

8.6 KNOWLEDGE GRADIENT

The knowledge gradient is the expected one-step gain in the global reward (6.5). If we define

$$\mu^* = \max_{x \in \mathcal{X}} \mu_{\mathcal{D}}(x)$$

to be the utility of the current dataset, then we must compute:

$$\alpha_{\text{KG}}(x; \mathcal{D}) = \int \left[\max_{x' \in \mathcal{X}} \mu_{\mathcal{D}'}(x') - \mu^* \right] p(y | x, \mathcal{D}) dy. \quad (8.30)$$

The global optimization in the expectation renders the knowledge gradient nontrivial to compute in most cases, so we must resort to quadrature or other approximation.

Exact computation in discrete domains

FRAZIER et al. first proposed the knowledge gradient for optimization on a *discrete* domain $\mathcal{X} = \{1, 2, \dots, n\}$.²⁷ In this case, the objective is simply a vector $\mathbf{f} \in \mathbb{R}^n$ and a Gaussian process belief about the objective is simply a multivariate normal distribution:

$$p(\mathbf{f} | \mathcal{D}) = \mathcal{N}(\mathbf{f}; \boldsymbol{\mu}, \Sigma).$$

The knowledge gradient now reduces to the expected marginal gain in the maximum of the posterior mean vector:

$$\alpha_{\text{KG}}(x; \mathcal{D}) = \int \left[\max \boldsymbol{\mu}' - \mu^* \right] p(y | x, \mathcal{D}) dy.$$

We may compute this expectation in closed form following our analysis of noisy expected improvement, which was merely a slight adaptation of FRAZIER et al.'s approach.

The updated posterior mean vector $\boldsymbol{\mu}'$ after observing an observation (x, y) is linear in the observed value y :

$$\boldsymbol{\mu}' = \boldsymbol{\mu} + \frac{\Sigma_x}{s} \frac{y - \mu_x}{s},$$

where μ_x is the entry of $\boldsymbol{\mu}$ corresponding to the index x , and Σ_x is similarly the corresponding column of Σ . If we define

$$\mathbf{a} = \boldsymbol{\mu}; \quad \mathbf{b} = \frac{\Sigma_x}{s},$$

then we may rewrite the knowledge gradient in terms of the g function introduced in the context of noisy expected improvement (8.10):

$$\alpha_{\text{KG}}(x; \mathcal{D}) = g(\mathbf{a}, \mathbf{b}) - \mu^*.$$

We may evaluate this expression exactly in $\mathcal{O}(n \log n)$ time following our previous discussion. Thus we may compute the knowledge gradient policy with a discrete domain in $\mathcal{O}(n^2 \log n)$ time per iteration by exhaustive computation of the acquisition function.

²⁷ P. FRAZIER et al. (2009). The Knowledge-Gradient Policy for Correlated Normal Beliefs. *INFORMS Journal on Computing* 21(4):599–613.

computation of noisy expected improvement:
§ 8.2, p. 160

computation in terms of $g(\mathbf{a}, \mathbf{b})$

computation of g : § 8.2, p. 161

Approximation via numerical quadrature

Although the knowledge gradient acquisition function can be computed exactly over discrete domains, the situation becomes significantly more

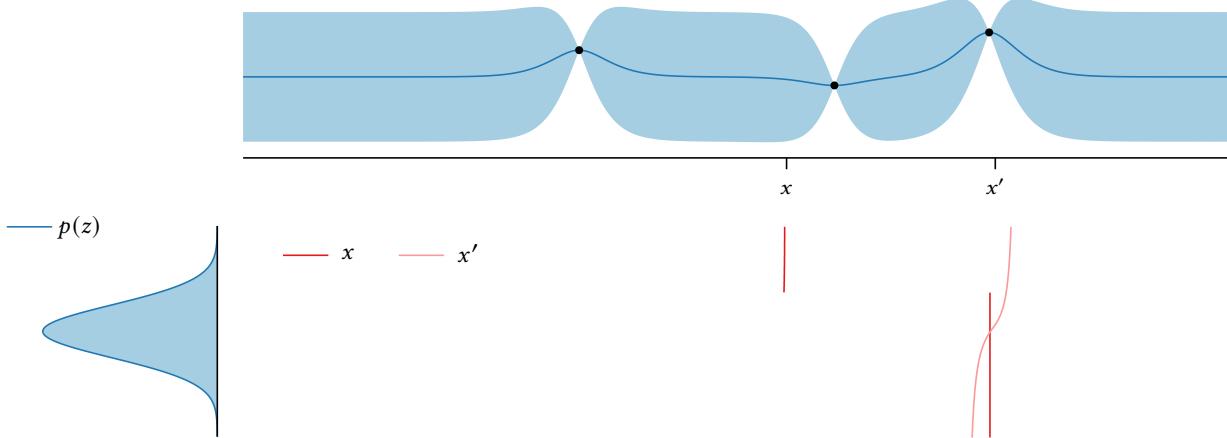


Figure 8.9: The complex behavior of the updated global reward. Above: the location of the posterior mean maximum given the z -score of an observation at two points, x and x' . An observation at x' always results in shoring up the existing maximum, whereas an observation at x reveals a new maximum given a sufficiently high observation. Below: the marginal gain in the global reward as a function of the z -score of an observation at these locations.

complicated in continuous domains. The culprit is the nonconvex global optimization in (8.30), which makes the knowledge gradient intractable except in a few special cases.

Some Gaussian processes give rise to convenient structure in the posterior distribution facilitating computation of the knowledge gradient. For example, in one dimension the Wiener or Ohrstein–Uhlenbeck processes satisfy a Markov property guaranteeing the posterior mean is always maximized at an observed location. As a result the simple reward and global reward are always equal, and the knowledge gradient reduces to expected improvement. This structure was often exploited in the early literature on Bayesian optimization.^{28,29}

Figure 8.9 give some insight into the complexity of the knowledge gradient integral (8.30). We illustrate the possible results from adding an observation at two locations as a function of the z -score of the observed value. For the point on the left, the behavior is similar to expected improvement: a sufficiently high value moves the maximum of the posterior mean to that point; otherwise, we retain the incumbent. The marginal gain in utility is a piecewise linear function corresponding to these two outcomes, as in expected improvement. The point on the right displays entirely different behavior – the marginal gain in utility is smooth, and nearly any outcome would be beneficial.

exact computation for special cases

²⁸ H. J. KUSHNER (1964). A New Method of Locating the Maximum Point of an Arbitrary Multi-peak Curve in the Presence of Noise. *Journal of Basic Engineering* 86(1):97–106.

²⁹ J. MOCKUS (1972). Bayesian Methods of Search for an Extremum. *Avtomatika i Vychislitel'naya Tekhnika (Automatic Control and Computer Sciences)* 6(3):53–62.

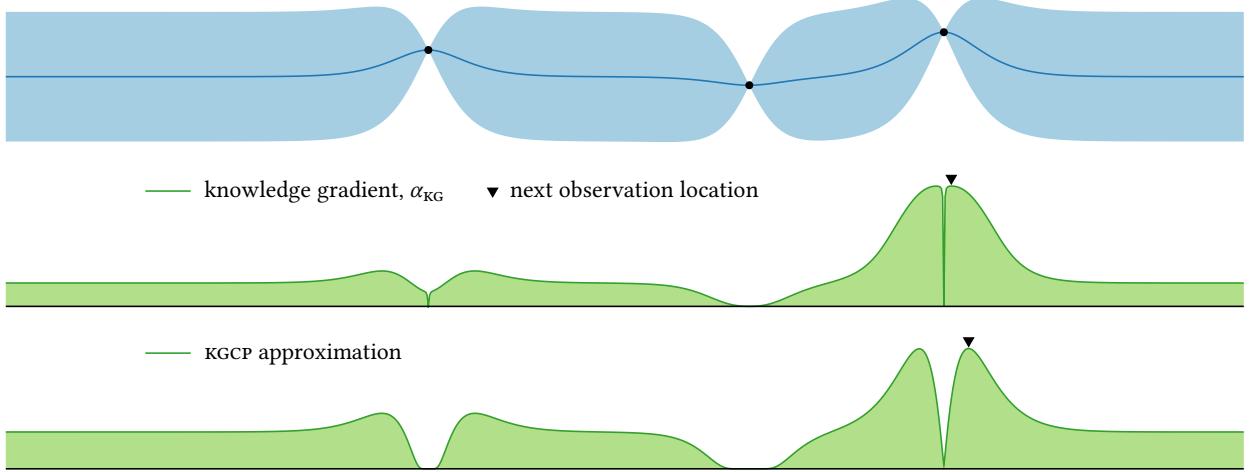


Figure 8.10: A comparison of the knowledge gradient acquisition function and the KGCP approximation for an example scenario. In this case the KGCP approximation reverts to expected improvement.

When exact computation is not possible, we must result to numerical integration or an analytic approximation when adopting the knowledge gradient policy. The former path was explored in depth by Wu et al.³⁰ To compute the acquisition function at a given point, we can compute a high-accuracy approximation following the numerical techniques outlined in the previous section. Order- n Gauss–Hermite quadrature for example, would approximate with (8.29):

$$\alpha_{\text{KG}}(x; \mathcal{D}) \approx \frac{1}{\sqrt{\pi}} \sum_{i=1}^n w_i \Delta(x, y_i); \quad \Delta(x, y_i) = \max_{x' \in \mathcal{X}} \mu_{\mathcal{D}'}(x') - \mu^*.$$

With some care, we can also approximate the gradient of the knowledge gradient in this scheme; details are given in an appendix.

³⁰ J. Wu et al. (2017). Bayesian Optimization with Gradients. *NeurIPS 2017*.

approximating gradient of knowledge gradient: § C.3, p. 306

Knowledge gradient for continuous parameters approximation

SCOTT et al. suggested an alternative and lightweight approximation scheme for the knowledge gradient the authors called the *knowledge gradient for continuous parameters* (KGCP).³¹ The KGCP approximation entails replacing the domain of the maximization in the definition of the global reward utility, normally the entire domain \mathcal{X} , with a conveniently chosen discrete set: the already observed locations x and the proposed new observation location x . Let x' represent this set. We approximate the current and future utility with

$$u(\mathcal{D}) \approx \max \mu_{\mathcal{D}}(x'); \quad u(\mathcal{D}') \approx \max \mu_{\mathcal{D}'}(x'),$$

yielding the approximation

$$\alpha_{\text{KG}}(x; \mathcal{D}) \approx \mathbb{E}[\max \mu_{\mathcal{D}'}(x') | x, \mathcal{D}] - \max \mu_{\mathcal{D}}(x'). \quad (8.31)$$

³¹ W. Scott et al. (2011). The Correlated Knowledge Gradient for Simulation Optimization of Continuous Parameters Using Gaussian Process Regression. *SIAM Journal on Optimization* 21(3):996–1026.

comparison with expected improvement

This expression is almost identical to expected improvement (8.7)! In fact, if we define

$$\mu^* = \max_{\mathbf{x}} \mu_{\mathcal{D}}(\mathbf{x}),$$

then a simple manipulation gives:

$$\alpha_{\text{KG}}(x; \mathcal{D}) \approx \alpha_{\text{EI}}(x; \mathcal{D}) - \max(\mu - \mu^*, 0).$$

Effectively, the KGCP approximation is a simple adjustment to expected improvement where we punish points for already having large expected value. From the point of view of the global reward, these points already represent success and their large expected values are already reflected in the current utility. Therefore, we should not necessarily waste precious evaluations confirming what we already believe.

gradient of KGCP approximation

The gradient of the KGCP approximation may also be computed in terms of the gradient of expected improvement and the posterior mean:

$$\frac{\partial \alpha_{\text{KG}}}{\partial x} \approx \frac{\partial \alpha_{\text{EI}}}{\partial x} - [\mu > \mu^*] \frac{\partial \mu}{\partial x}.$$

example and discussion

In the case of our example scenario, the posterior mean never exceeds the highest observed point. Therefore, the KGCP approximation to the knowledge gradient globally reduces to the expected improvement; compare with the expected improvement in figure 7.3. Comparing with true knowledge gradient, we can see that this approximation cannot necessarily be trusted to be unconditionally faithful. However, the KGCP approximation has the advantage of efficient computation and may be used as a drop-in replacement for expected improvement that may offer a slight boost in performance when the global reward utility is preferred to simple reward.

8.7 THOMPSON SAMPLING

Thompson sampling: § 7.9, p. 148

Thompson sampling designs each observation by sampling a point proportional to its probability of maximizing the objective (7.19):

$$x \sim p(x^* | \mathcal{D}).$$

³² One notable example is the Wiener process, where x^* famously has an arcsine distribution:

P. LÉVY (1948). *Processus stochastiques et mouvement brownien*. Gauthier-Villars.

A major barrier to Thompson sampling with Gaussian processes is the complex nature of this distribution. Except in a small number of special cases,³² this distribution cannot be computed analytically. Figure 8.11 illustrates the complicated nature of this distribution for our running example, which was only revealed via brute-force sampling. However, a straightforward implementation strategy is to maximize a draw from the objective function posterior, which assumes the role of an acquisition function:

$$\alpha_{\text{TS}}(x; \mathcal{D}) \sim p(f | \mathcal{D}).$$

The global optimum of α_{TS} is then a sample from the desired distribution.

This procedure in fact yields a sample from the *joint* distribution of the location and value of the optimum, $p(x^*, f^* | \mathcal{D})$, as the value of

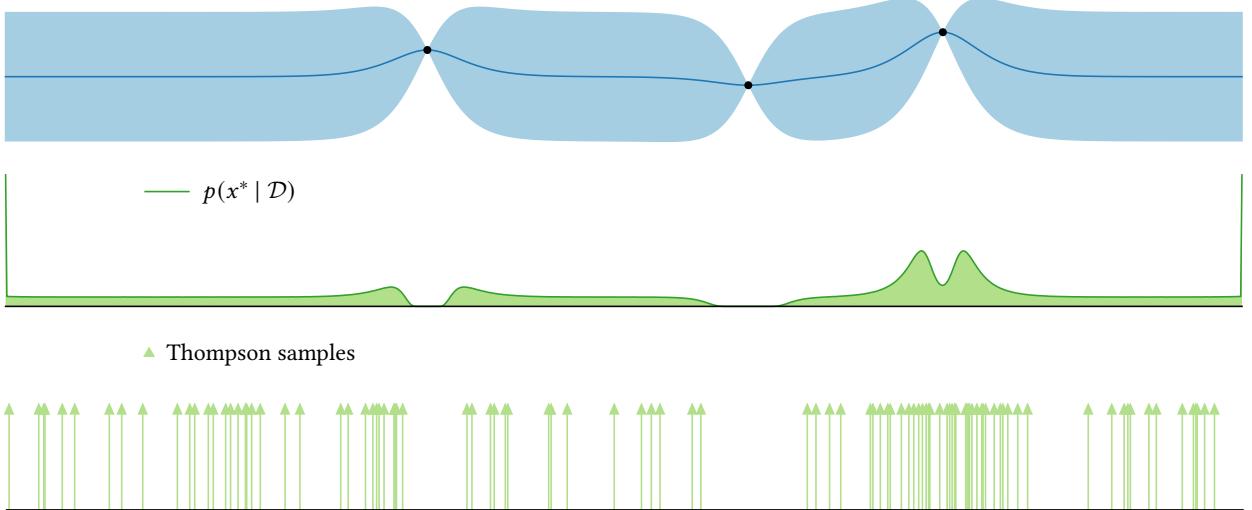


Figure 8.11: The distribution of the location of the global maximum, $p(x^* | \mathcal{D})$, for an example scenario, and 100 samples drawn from this distribution.

the sampled objective function α_{TS} at its maximum provides a sample from $p(f^* | x^*, \mathcal{D})$; see the margin. We discuss Thompson sampling now because the ability to sample from these distributions will be critical for computing mutual information, our focus in the following sections.

Exhaustive sampling

In “small” domains, we can realize Thompson sampling via brute force. If the domain can be exhaustively covered by a sufficiently small set of points ξ (for example, with a dense grid or a low-discrepancy sequence) then we can simply sample the associated objective function values $\phi = f(\xi)$ and maximize:³³

$$x = \arg \max \phi; \quad \phi \sim p(\phi | \xi, \mathcal{D}).$$

The distribution of ϕ is multivariate normal, making sampling easy:

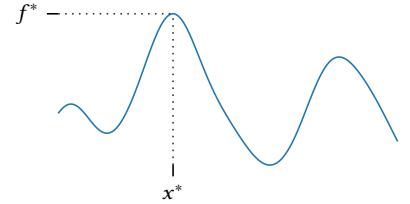
$$p(\phi | \xi, \mathcal{D}) = \mathcal{N}(\phi; \mu, \Sigma); \quad \mu = \mu_{\mathcal{D}}(\xi); \quad \Sigma = K_{\mathcal{D}}(\xi, \xi).$$

The running time of this procedure grows quickly with the size of ξ , although sophisticated numerical methods enable scaling to roughly 50 000 points.³⁴

Figure 8.11 shows 100 Thompson samples for our example scenario generated via exhaustive sampling, taking ξ to be a grid of 1000 points.

On-demand sampling

An alternative to exhaustive sampling is to use off-the-shelf optimization routines to maximize a draw from the objective function posterior we



Maximizing a draw from a Gaussian process naturally samples from $p(x^*, f^* | \mathcal{D})$

³³ We are taking a slight liberty with notation here as we have previously used ϕ for $f(x)$, the latent objective function values at the observed locations. However, for the remainder of the we will be assuming the general, potentially noisy case where our data will be written $\mathcal{D} = (x, y)$ and we will have no need to refer to $f(x)$.

³⁴ G. PLEISS et al. (2020). Fast Matrix Square Roots with Applications to Gaussian Processes and Bayesian Optimization. *NeurIPS 2020*.

Algorithm 8.1: On-demand sampling.

$\mathcal{D}_{\text{TS}} \leftarrow \mathcal{D}$	► initialize fictitious dataset with current data
repeat	
given request for observation at x :	
$\phi \leftarrow p(\phi x, \mathcal{D}_{\text{TS}})$	► sample value at x
$\mathcal{D}_{\text{TS}} \leftarrow \mathcal{D}_{\text{TS}} \cup (x, \phi)$	► update fictitious dataset
yield ϕ	
until external optimizer terminates	

build progressively on demand. Namely, when an optimization routine requests an evaluation of α_{TS} at a point x , we sample an objective function value at that location:

$$\phi \sim p(\phi | x, \mathcal{D}),$$

then augment our dataset with the simulated observation (x, ϕ) . We proceed in this manner until the optimizer terminates. This procedure avoids simulating the entire objective function while guaranteeing the joint distribution of the provided evaluations is correct via the chain rule of probability. Pseudocode is provided in Algorithm 8.1 in the form of a generator function; the “yield” statement returns a value while maintaining state. Using rank-one updates to update the posterior, the computational cost of generating each evaluation scales quadratically with the size of the fictitious dataset \mathcal{D}_{TS} .

If desired, we may use gradient methods to optimize the generated sample by sampling from the joint posterior of the function value and its gradient at each requested location:

$$p(\phi, \nabla\phi | x, \mathcal{D}_{\text{TS}})$$

However, in high dimensions, the additional cost required to condition on these gradient observations may become excessive. In d dimensions, returning gradients effectively reduces the number of function evaluations we can allow for the optimizer by a factor of $(d + 1)$ if we wish to maintain the same total computational effort.

Sparse spectrum approximation for stationary covariance functions

If the prior covariance function K of the Gaussian process is stationary, we may use a *sparse spectrum approximation*³⁵ to the posterior Gaussian process to dramatically accelerate Thompson sampling.³⁶

Consider a stationary covariance function K on \mathbb{R}^d with spectral density κ . The key idea in sparse spectrum approximation is to interpret the characterization in (3.10) as an expectation with respect to the spectral density:³⁷

$$K(\mathbf{x} - \mathbf{x}') = K(\mathbf{0}) \mathbb{E}_\xi [\exp(2\pi i (\mathbf{x} - \mathbf{x}')^\top \xi)], \quad (8.32)$$

and approximate via Monte Carlo integration. We first sample a set of m frequencies, called *spectral points*, from the spectral density: $\{\xi_i\} \sim \kappa(\xi)$.

low-rank updates: § 9.1, p. 202
running time

using gradient methods

- 35 M. LÁZARO-GREDILLA et al. (2010). Sparse Spectrum Gaussian Process Regression. *Journal of Machine Learning Research* 11(Jun):1865–1881.
- 36 J. M. HERNÁNDEZ-LOBATO et al. (2014). Predictive Entropy Search for Efficient Global Optimization of Black-box Functions. *NeurIPS 2014*.

stationarity: § 3.2, p. 50

37 Recall the convention of writing a stationary covariance with respect to a single input, $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x} - \mathbf{x}')$.

spectral points, $\{\xi_i\}$

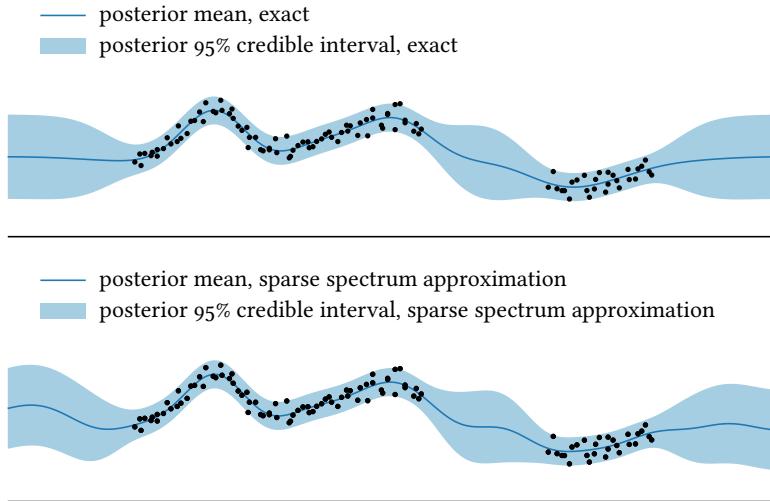


Figure 8.12: Sparse spectrum approximation. Top: the exact posterior belief (about noisy observations rather than the latent function) for a Gaussian process conditioned on 200 observations. Bottom: a sparse spectrum approximation using 100 spectral points sampled from the spectral density.

To enforce the symmetry around the origin inherent to the spectral density (theorem 3.2), we augment each sample ξ_i with its negation, $-\xi_i$.

Using these samples for a Monte Carlo approximation to (8.32) has the effect of approximating a Gaussian process $\mathcal{GP}(f; \mu, K)$ with a *finite-dimensional* Gaussian process:

$$f(\mathbf{x}) \approx \mu(\mathbf{x}) + \boldsymbol{\beta}^\top \boldsymbol{\psi}(\mathbf{x}). \quad (8.33)$$

Here $\boldsymbol{\beta}$ is a vector of normally distributed weights and $\boldsymbol{\psi}: \mathbb{R}^d \rightarrow \mathbb{R}^{2m}$ is a feature representation determined by the spectral points:

$$\psi_{2i-1}(\mathbf{x}) = \cos(2\pi \xi_i^\top \mathbf{x}); \quad \psi_{2i}(\mathbf{x}) = \sin(2\pi \xi_i^\top \mathbf{x}).$$

For the additive Gaussian noise model, if \mathbf{N} is the covariance of the noise contributions to the observed values \mathbf{y} , the posterior moments of the weight vector in this approximation are:

$$\begin{aligned} \mathbb{E}[\boldsymbol{\beta} | \mathcal{D}] &= \boldsymbol{\Psi}^\top (\boldsymbol{\Psi} \boldsymbol{\Psi}^\top + a^{-2} \mathbf{N})^{-1} (\mathbf{y} - \boldsymbol{\mu}); \\ \text{cov}[\boldsymbol{\beta} | \mathcal{D}] &= a^2 [\mathbf{I} - \boldsymbol{\Psi}^\top (\boldsymbol{\Psi} \boldsymbol{\Psi}^\top + a^{-2} \mathbf{N})^{-1} \boldsymbol{\Psi}], \end{aligned}$$

where $a^2 = K(\mathbf{0})/m$, $\boldsymbol{\mu}$ is the prior mean of \mathbf{y} , and $\boldsymbol{\Psi}$ is a matrix whose rows comprise the feature representations of the observed locations.

The sparse spectrum approximation allows us to generate a posterior sample of the objective function in time $\mathcal{O}(nm^2)$ by drawing a sample from the weight posterior and appealing to the representation in (8.33). The resulting sample is nontrivial to maximize due to the nonlinear nature of the representation, but with the sampled weights in hand, we can generate requested function and gradient evaluations in constant time, enabling exhaustive optimization via gradient methods.

A sparse spectrum approximation is illustrated in figure 8.12, using a total of 100 spectral points. The approximation is quite reasonable near data and acceptable during extrapolation as well.

weight vector, $\boldsymbol{\beta}$
feature representation, $\boldsymbol{\psi}$

noise covariance matrix, \mathbf{N}

training features, $\boldsymbol{\Psi}$

sampling from a multivariate normal distribution: § A.2, p. 295

8.8 MUTUAL INFORMATION WITH x^*

mutual information: § 7.6, p. 135

Mutual information measures the expected information gain (6.3) provided by an observation at x about some random variable of interest ω , which can be expressed in two equivalent forms (7.14–7.15):

$$\alpha_{\text{MI}}(x; \mathcal{D}) = H[\omega | \mathcal{D}] - \mathbb{E}_y[H[\omega | \mathcal{D}'] | x, \mathcal{D}] \quad (8.34)$$

$$= H[y | x, \mathcal{D}] - \mathbb{E}_{\omega}[H[y | \omega, x, \mathcal{D}] | x, \mathcal{D}]. \quad (8.35)$$

In the context of Bayesian optimization, the most natural choices for ω are the location x^* and value f^* of the global optimum, both of which were discussed in the previous chapter. Unfortunately, a Gaussian process belief on the objective function in general induces complex distributions for these quantities, which makes even *approximating* mutual information somewhat involved. However, several effective approximation schemes are available for both options. We will consider each in turn, beginning with x^* .

Direct form of mutual information

Of the two equivalent formulations of mutual information above, the former (8.34) is perhaps the most natural, as it reasons directly about changes in the entropy of the variable of interest. Initial work on mutual information with x^* considered approximations to this direct expression³⁸ – including the work coining the now-common moniker *entropy search* for information-theoretic optimization policies³⁹ – but this approach has fallen out of favor in preference for the latter formulation (8.35), discussed below.

The main computational difficulty in approximating (8.34) is the unwieldy (and potentially high-dimensional) distribution $p(x^* | \mathcal{D})$; see even the simple example in figure 8.11. There is no general method for computing the entropy of this distribution in closed form, and overcoming this barrier requires complex approximation schemes. The usual approach is to make a discrete approximation to this distribution via a set of carefully maintained so-called *representer points*, then reason about changes in the entropy of this surrogate distribution via further layers of approximation.

Predictive form of mutual information

³⁸ J. VILLEMONTEIX et al. (2009). An informational approach to the global optimization of expensive-to-evaluate functions. *Journal of Global Optimization* 44(4):509–534.

³⁹ P. HENNIG and C.J. SCHULER (2012). Entropy Search for Information-Efficient Global Optimization. *Journal of Machine Learning Research* 13(Jun):1809–1837.

⁴⁰ J. M. HERNÁNDEZ-LOBATO et al. (2014). Predictive Entropy Search for Efficient Global Optimization of Black-box Functions. *NeurIPS 2014*.

$$\alpha_{x^*}(x; \mathcal{D}) = H[y | x, \mathcal{D}] - \mathbb{E}[H[y | x, x^*, \mathcal{D}] | x, \mathcal{D}] \quad (8.36)$$

Compared to the direct form (8.34), this formulation is attractive as all entropy computations are restricted to the one-dimensional predictive distribution $p(y | x, \mathcal{D})$ rather than the inconvenient distribution $p(x^* | \mathcal{D})$. The authors named their approach *predictive entropy search* to highlight this feature.

Let us consider the evaluation of (8.36) for a Gaussian process model with additive Gaussian noise. To begin, the first term is simply the differential entropy of a one-dimensional Gaussian distribution (8.5) and may be computed in closed form (A.17):

$$H[y | x, \mathcal{D}] = \frac{1}{2} \log(2\pi e s^2).$$

Unfortunately, the second term is significantly more complicated to work with, and we can identify two primary challenges.

*Approximating expectation with respect to x^**

First we must compute an expectation with respect to the location of the global optimum x^* . The complexity of its distribution limits our options, but Monte Carlo integration remains a viable option. We use Thompson sampling to generate samples of the optimal location $\{x_i^*\}_{i=1}^n \sim p(x^* | \mathcal{D})$ and estimate the expected updated entropy with:

$$\mathbb{E}[H[y | x, x^*, \mathcal{D}] | x, \mathcal{D}] \approx \frac{1}{n} \sum_{i=1}^n H[y | x, x_i^*, \mathcal{D}].$$

When the prior covariance function is stationary, HERNÁNDEZ-LOBATO et al. further propose to exploit the efficient approximate Thompson sampling scheme via sparse spectrum approximation described in the last section. When feasible, this reduces the cost of drawing the samples, but it is not necessary for correctness.

Gaussian approximation to conditional predictive distribution

Now we must address the predictive distribution conditioned on the location of the global optimum, $p(y | x, x^*, \mathcal{D})$, which we may express as the result of marginalizing the latent objective value ϕ :

$$p(y | x, x^*, \mathcal{D}) = \int p(y | x, \phi) p(\phi | x, x^*, \mathcal{D}) d\phi. \quad (8.37)$$

It is unclear how we can condition our belief on the objective function given knowledge of the optimum, and – as the resulting posterior on ϕ will be non-Gaussian – how we can resolve the resulting integral (8.37).

A key insight is that if the predictive distribution $p(\phi | x, x^*, \mathcal{D})$ were Gaussian, we could compute (8.37) in closed form,⁴¹ suggesting a promising path forward. In particular, consider an arbitrary Gaussian approximation:

$$p(\phi | x, x^*, \mathcal{D}) \approx \mathcal{N}(\phi; \mu_*, \sigma_*^2), \quad (8.38)$$

whose parameters depend x^* as their subscripts indicate. Plugging into (8.37), we may estimate the predictive variance of y given x^* with (A.15):

$$\text{var}[y | x, x^*, \mathcal{D}] \approx \sigma_*^2 + \sigma_n^2 = s_*^2, \quad (8.39)$$

computing first term

Thompson sampling for GPS: § 8.7, p. 176

sparse spectrum approximation for
Thompson sampling: § 8.7, p. 178

⁴¹ In this case, the integral in (8.37) becomes the convolution of two Gaussians, which may be interpreted as the distribution of the sum of independent Gaussian random variables – see § A.2, p. 296.

approximate variance of y given x^* , s_*^2

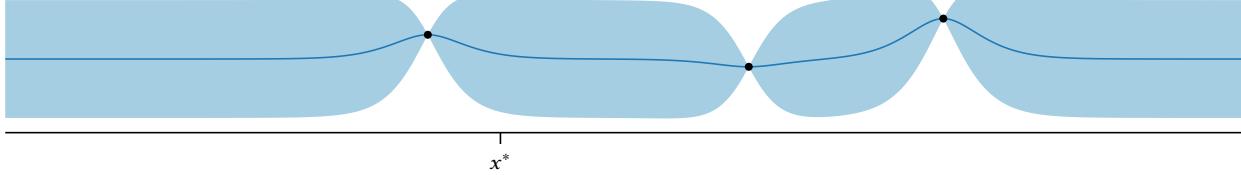


Figure 8.13: The example scenario we will consider for illustrating the predictive entropy search approximation to $p(f | x, x^*, \mathcal{D})$, using the marked location for x^* .

and thus its differential entropy with (A.17):

$$H[y | x, x^*, \mathcal{D}] \approx \frac{1}{2} \log(2\pi e s_*^2). \quad (8.40)$$

After simplification, the resulting approximation to (8.36) becomes:

$$\alpha_{x^*}(x; \mathcal{D}) \approx \alpha_{\text{PES}}(x; \mathcal{D}) = \log s - \frac{1}{n} \sum_{i=1}^n \log s_{*i}. \quad (8.41)$$

Approximation via Gaussian expectation propagation

To realize a complete algorithm, we need some way of finding a suitable Gaussian approximation to $p(\phi | x, x^*, \mathcal{D})$, and HERNÁNDEZ-LOBATO et al. describe one effective approach. The high-level idea is to begin with the objective function posterior $p(f | \mathcal{D})$, impose a series of constraints implied by knowledge of x^* , then approximate the desired posterior via approximate inference. We will describe the procedure for an arbitrary putative optimum x^* illustrating each step of the approximation for the example scenario and assumed optimum location shown in figure 8.13. For the moment, we will proceed with complete disregard for computational efficiency, and return to the question of implementation shortly.

Ensuring x^ is a local optimum*

We first condition our belief on x^* being a *local* optimum by insisting that the point satisfy the second partial derivative test. Let ∇^* and \mathbf{H}^* respectively represent the gradient and Hessian of the objective function at x^* . Local optimality implies that the gradient is zero and the Hessian is negative definite:⁴²

$$\nabla^* = \mathbf{0}; \quad (8.42)$$

$$\mathbf{H}^* < \mathbf{0}. \quad (8.43)$$

Enforcing the gradient constraint (8.42) is straightforward, as we can directly condition on a gradient observation. We show the result for our example scenario in the top panel of figure 8.14.

The Hessian constraint (8.43) however is nonlinear and we must resort to approximate inference. HERNÁNDEZ-LOBATO et al. approximate

gradient, Hessian at x^* : ∇^*, \mathbf{H}^*

⁴² We discussed a simpler, one-dimensional analog of this task in § 2.8, p. 39.

gradient constraint: $\nabla^* = \mathbf{0}$

first Hessian constraint: diagonal entries are negative

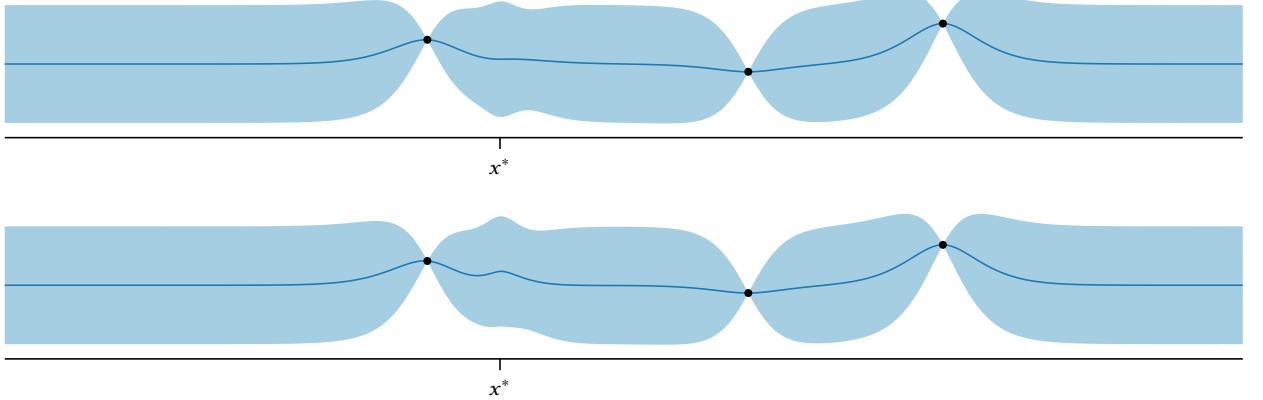


Figure 8.14: Top: the posterior for our example after conditioning on the derivative being zero at x^* (8.42). Bottom: the approximate posterior after conditioning on the second derivative being negative at x^* (8.44).

this condition by breaking it into two complimentary components. First we compel every diagonal entry of the Hessian to be negative. Letting $\mathbf{h}^* = \text{diag } \mathbf{H}^*$; we assume:

$$\forall i: h_i^* < 0. \quad (8.44)$$

We then fix the off-diagonal entries of the Hessian to values of our choosing. One simple option is to set all off-diagonal entries to zero:

$$\text{upper } \mathbf{H}^* = \mathbf{0}, \quad (8.45)$$

which combined with the diagonal constraint guarantees negative definiteness.⁴³ The combination of these conditions (8.44–8.45) is stricter than mere negative definiteness, as we eliminate all degrees of freedom for the off-diagonal entries. However, an advantage of this approach is that we can enforce the off-diagonal constraint via exact conditioning.

To proceed we dispense with the constraints we can condition on exactly. Let \mathcal{D}' represent our dataset augmented with the gradient (8.42) and off-diagonal Hessian (8.45) observations. The joint distribution of the latent objective value ϕ , the purportedly optimal value $\phi^* = f(x^*)$, and the diagonal of the Hessian \mathbf{h}^* given this additional information is multivariate normal:

$$p(\phi, \phi^*, \mathbf{h}^* | x^*, \mathcal{D}') = \mathcal{N}(\phi, \phi^*, \mathbf{h}^*; \boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*). \quad (8.46)$$

We will now subject this initial belief to a series of factors corresponding to desired nonlinear constraints. These factors will be compatible with Gaussian expectation propagation, which we will use to finally derive the desired Gaussian approximation to the posterior (8.38). To begin, the Hessian diagonal constraint (8.44) contributes one factor for each entry:

$$\prod_i [h_i^* < 0]. \quad (8.47)$$

diagonal of Hessian at x^* , \mathbf{h}^*

second Hessian constraint: off-diagonal entries are fixed

⁴³ HERNÁNDEZ-LOBATO et al. point out this may not be faithful to the model and suggest the alternative of matching the off-diagonal entries of the Hessian of the objective function sample that generated x^* . However, this does not guarantee negative definiteness without tweaking (8.44).

Gaussian EP: § B.2, p. 298
truncating a variable with EP: § B.2, p. 301

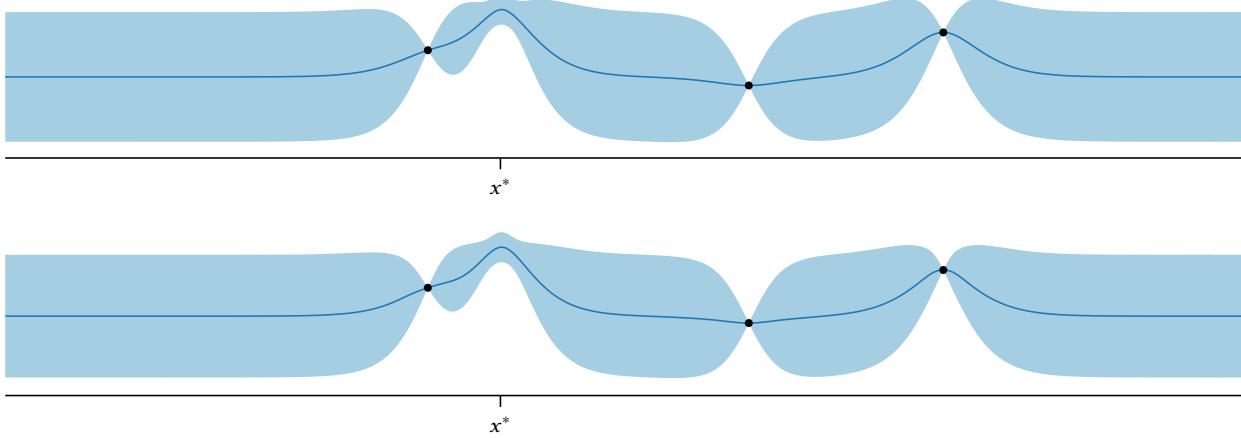


Figure 8.15: Top: the approximate posterior after conditioning on ϕ^* exceeding the function values at previously measured locations (8.49). Bottom: the approximate posterior after conditioning on ϕ^* dominating elsewhere (8.50).

⁴⁴ For this and the following demonstrations, we show the expectation propagation approximation to the entire objective function posterior; this is not required to approximate the marginal predictive distribution and is only for illustration.

Our approximate posterior after incorporating (8.47) and performing expectation propagation is shown in the bottom panel of 8.14.⁴⁴

Ensuring x^* is a global optimum

Our belief now reflects our desire that x^* be a *local* maximum; however, we wish for x^* to be the *global* maximum. Global optimality is not easy to enforce, as it entails infinitely many constraints bounding the objective at every point in the domain. HERNÁNDEZ-LOBATO et al. instead approximate this condition with optimality at the most relevant locations: the already-observed points \mathbf{x} and the proposed point x .

ensuring $\phi^* > \max \phi$

To enforce that ϕ^* exceed the objective function values at the observed points, we could theoretically add $\phi = f(\mathbf{x})$ to our prior (8.46), then add one factor for each observation: $\prod_j [\phi_j < \phi^*]$. However, this approach requires an increasing number of factors as we gather more data, rendering expectation propagation (and thus the acquisition function) increasingly expensive. Further, factors corresponding to obviously suboptimal observations are uninformative and simply represent extra work for no benefit.

Instead, we enforce this constraint through a single factor truncating with respect to the maximal value of ϕ : $[\phi^* < \max \phi]$. In general, this threshold will be a random variable unless our observations are noiseless. Fortunately, expectation propagation enables tractable approximate truncation at an *unknown*, Gaussian-distributed threshold. Define

$$\mu_{\max} = \mathbb{E}[\max \phi \mid \mathcal{D}]; \quad \sigma_{\max}^2 = \text{var}[\max \phi \mid \mathcal{D}]; \quad (8.48)$$

⁴⁵ C. E. CLARK (1961). The Greatest of a Finite Set of Random Variables. *Operations Research* 9(2): 145–162.

these moments can be approximated via either sampling or an assumed density filtering approach described by CLARK.⁴⁵ Taking a moment-matched Gaussian approximation to $\max \phi$ and integrating yields the

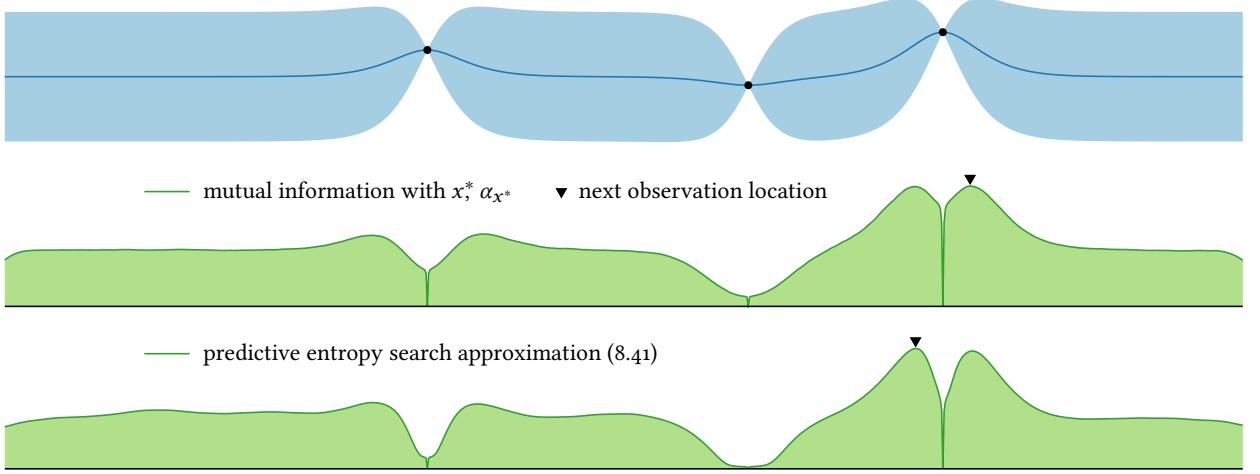


Figure 8.16: The predictive entropy search approximation (8.41) to the mutual information with x^* acquisition function (8.36) using the 100 Thompson samples from figure 8.11.

factor (B.13):⁴⁶

$$\Phi\left(\frac{\phi^* - \mu_{\max}}{\sigma_{\max}}\right). \quad (8.49)$$

⁴⁶ With high noise, this approximation could be improved slightly by computing moments of $\max \phi$ given \mathcal{D}' , at additional expense.

We show the approximate posterior for our running example after incorporating this factor in the top panel of figure 8.15. The probability mass at x^* has shifted up dramatically.

Finally, to constrain ϕ^* to dominate the objective function value at a point of interest x , we add one additional factor:

$$[\phi < \phi^*]. \quad (8.50)$$

ensuring $\phi^* > \phi$
enforcing order with EP: § B.2, p. 302

We obtain the final approximation to $p(\phi | x, x^*, \mathcal{D})$ by combining the prior in (8.46) with the factors (8.47, 8.49–8.50):

$$\mathcal{N}(\phi, \phi^*, \mathbf{h}^*; \boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*) [\phi < \phi^*] \Phi\left(\frac{\phi^* - \mu_{\max}}{\sigma_{\max}}\right) \prod_i [h_i^* < 0],$$

approximating with Gaussian expectation propagation, and deriving the marginal belief about ϕ . The resulting final approximate posterior for our example scenario is shown in the bottom panel of figure 8.15. Our predictive uncertainty has been moderated in response to the final constraint (8.50).

With the ability to approximate the latent predictive posterior, we can now compute the predictive entropy search acquisition function (8.41) via Thompson sampling and following the above procedure for each sample. Figure 8.16 shows the approximation computed with 1000 Thompson samples. The approximation is excellent and induces a near-optimal decision. Any deviation from the truth mostly reflects bias in the Thompson sample distribution.

Efficient implementation

A practical realization of predictive entropy search can benefit from careful precomputation and reuse of partial results. We will outline an efficient implementation strategy, beginning with three steps of one-time initial work.

1. Estimate the moments of $\max \phi$ (8.48).
2. Generate a set of Thompson samples $\{x_i^*\}$.
3. For each sample x^* , derive the joint distribution of the function value, gradient, and Hessian at x^* . Let $\mathbf{z}^* = (\phi^*, \mathbf{h}^*, \text{upper } \mathbf{H}^*, \nabla^*)$ represent a vector comprising these random variables, with those that will be subjected to expectation propagation first. We compute:

$$p(\mathbf{z}^* | x^*, \mathcal{D}) = \mathcal{N}(\mathbf{z}^*; \boldsymbol{\mu}^*, \Sigma^*).$$

⁴⁷ In the interest of numerical stability, the inverse of \mathbf{V}_* should not be stored directly; for relevant practical advice see:

C. E. RASMUSSEN and C. K. I. WILLIAMS (2006). *Gaussian Processes for Machine Learning*. MIT Press.

J. P. CUNNINGHAM et al. (2011). Gaussian Probabilities and Expectation Propagation. arXiv: 1111.6832 [stat.ML].

Find the marginal belief over (ϕ^*, \mathbf{h}^*) and use Gaussian expectation propagation to approximate the posterior after incorporating the factors (8.47, 8.49). Let vectors $\tilde{\boldsymbol{\mu}}$ and $\tilde{\sigma}^2$ denote the site parameters at termination. Finally, precompute⁴⁷

$$\mathbf{V}_*^{-1} = \left[\Sigma^* + \begin{bmatrix} \tilde{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right]^{-1}; \quad \boldsymbol{\alpha}^* = \mathbf{V}_*^{-1} \left[\begin{bmatrix} \tilde{\boldsymbol{\mu}} \\ \mathbf{0} \end{bmatrix} - \boldsymbol{\mu}^* \right],$$

where $\tilde{\Sigma} = \text{diag } \tilde{\sigma}^2$. These quantities do not depend on x and will be repeatedly reused during prediction.

After completing the preparations above, suppose a proposed observation location x is given. For each sample x^* , we compute the joint distribution of ϕ and ϕ^* :

$$p(\phi, \phi^* | x, x^*, \mathcal{D}) = \mathcal{N}(\phi, \phi^*; \boldsymbol{\mu}, \Sigma),$$

and derive the approximate posterior given the exact gradient (8.42) and off-diagonal Hessian (8.45) observations and the factors (8.47, 8.49). Defining

$$\mathbf{K} = \begin{bmatrix} \mathbf{k}^\top \\ \mathbf{k}_*^\top \end{bmatrix} = \text{cov}([\phi, \phi^*]^\top, \mathbf{z}^*),$$

the desired distribution is $\mathcal{N}(\phi, \phi^*; \mathbf{m}, \mathbf{S})$, where:⁴⁸

$$\mathbf{m} = \begin{bmatrix} m \\ m^* \end{bmatrix} = \boldsymbol{\mu} + \mathbf{K}\boldsymbol{\alpha}^* \quad \mathbf{S} = \begin{bmatrix} \varsigma^2 & \rho \\ \rho & \varsigma_*^2 \end{bmatrix} = \Sigma - \mathbf{K}\mathbf{V}_*^{-1}\mathbf{K}^\top. \quad (8.51)$$

We now apply the prediction constraint (8.50) with one final step of expectation propagation. Define (B.7, B.11–B.12):

$$\int p(\phi, \phi^* | \mathbf{z}^*) p(\mathbf{z}^*) d\mathbf{z}^*$$

⁴⁸ This can be derived by marginalizing \mathbf{z}^* according to its approximate posterior from step 3 above:

$$\bar{\mu} = m - m^*; \quad \tilde{\sigma}^2 = \varsigma^2 - 2\rho + \varsigma_*^2;$$

$$z = -\frac{\bar{\mu}}{\tilde{\sigma}}; \quad \alpha = -\frac{\phi(z)}{\Phi(z)\tilde{\sigma}}; \quad \gamma = -\frac{\tilde{\sigma}}{\alpha} \left(\frac{\phi(z)}{\Phi(z)} + z \right)^{-1}.$$

8.9. MUTUAL INFORMATION WITH f^*

The final approximation to the predictive variance of ϕ given x^* is

$$\sigma_*^2 = \zeta^2 - (\zeta^2 - \rho)^2 / \gamma, \quad (8.52)$$

from which we can compute the contribution to the acquisition function from this sample with (8.39–8.40).

Although it may seem unimaginable at this point, in Euclidean domains we may compute the gradient of the predictive entropy search acquisition function; see the accompanying appendix for details.

gradient of predictive entropy search
acquisition function: § c.3, p. 307

8.9 MUTUAL INFORMATION WITH f^*

Finally, we consider the computation of the mutual information between the observed value y and the value of the global maximum f^* (8.53). Several authors have considered this acquisition function in its predictive form (8.35), which is the most convenient choice for Gaussian process models:

$$\alpha_{f^*}(x; \mathcal{D}) = H[y | x, \mathcal{D}] - \mathbb{E}[H[y | x, f^*, \mathcal{D}] | x, \mathcal{D}]. \quad (8.53)$$

Unfortunately, this expression cannot be computed exactly due to the complexity of the distribution $p(f^* | \mathcal{D})$; see figure 8.17 for an example. However, several effective approximations have been proposed, including *max-value entropy search* (MES)⁴⁹ and *output-space predictive entropy search* (OPES).⁵⁰

The issues we face in estimating (8.53), and the strategies we use to overcome them, largely mirror those in predictive entropy search. To begin, the first term is the differential entropy of a Gaussian and may be computed exactly:

$$H[y | x, \mathcal{D}] = \frac{1}{2} \log(2\pi e s^2). \quad (8.54)$$

The second term, however, presents some challenges, and the available approximations to (8.53) diverge in their estimation approach. We will discuss the MES and OPES approximations in parallel, as they share the same basic strategy and only differ in some details along the way.

*Approximating expectation with respect to f^**

The first complication in evaluating the second term of (8.53) is that we must compute an expectation with respect to f^* . Although Thompson sampling and simple Monte Carlo approximation is one way forward, we can exploit the fact that f^* is one dimensional to pursue more sophisticated approximations. One convenient and rapidly converging strategy is to design n samples $\{f_i^*\}_{i=1}^n$ to be equally spaced quantiles of f^* , then use the familiar estimator

$$\mathbb{E}[H[y | x, f^*, \mathcal{D}] | x, \mathcal{D}] \approx \frac{1}{n} \sum_{i=1}^n H[y | x, f_i^*, \mathcal{D}]. \quad (8.55)$$

mutual information with f^* : § 7.6, p. 140

49 Z. WANG and S. JEGELKA (2017). Max-value Entropy Search for Efficient Bayesian Optimization. *ICML 2017*.

50 M. W. HOFFMAN and Z. GHAHRAMANI (2015). Output-Space Predictive Entropy Search for Flexible Global Optimization. *Bayesian Optimization Workshop, NeurIPS 2015*.

51 R. E. CAFLISCH (1998). Monte Carlo and quasi-Monte Carlo methods. *Acta Numerica* 7:1–49.

52 The chosen samples are the first n points from a base- n van de Corput sequence:

J. G. VAN DE CORPUT (1935). Verteilungsfunktionen: Erste Mitteilung. *Proceedings of the Koninklijke Nederlandse Akademie van Wetenschappen* 38:813–821

mapped to the desired distribution via the inverse CDF. These points have the minimum possible discrepancy for a set of size n .

This represents a *quasi*-Monte Carlo⁵¹ approximation that converges considerably faster than simple Monte Carlo integration.⁵² Further, this scheme only requires the ability to estimate quantiles of f^* .

*Approximating quantiles of f^**

WANG and JEGETKA proposed estimating the quantiles of f^* via a convenient analytic approximation. To build this approximation, they proposed selecting a set of *representer points* ξ and approximating the distribution of f^* with that of the maximum function value restricted to these locations. Let $\phi = f(\xi)$ and define the random variable $\phi^* = \max \phi$. We estimate:

$$p(f^* | \mathcal{D}) \approx p(\phi^* | \xi, \mathcal{D}).$$

Unfortunately, the distribution of the maximum of dependent Gaussian random variables is intractable in general, even if the dimension is finite,^{53,54} so we must resort to further approximation.

We could proceed in several ways, but WANG and JEGETKA propose using an exhaustive, dense set of representer points (for example, covering the domain with a low-discrepancy sequence) and then making the simplifying assumption that the associated function values are *independent*. If the marginal belief at each representer point is

$$p(\phi_i | \xi_i, \mathcal{D}) = \mathcal{N}(\phi_i; \mu_i, \sigma_i^2),$$

then we may approximate the cumulative distribution function of ϕ^* with a product of normal CDFs:

$$\Pr(\phi^* < z | \xi, \mathcal{D}) \approx \prod_i \Phi\left(\frac{z - \mu_i}{\sigma_i}\right). \quad (8.56)$$

We can now estimate the quantiles of f^* via numerical inversion of (8.56); Brent's⁵⁵ method applied to the log CDF would offer rapid convergence. The resulting approximation for our example scenario is shown in figure 8.17 using 100 equally spaced representer points covering the domain. In general, the independence assumption tends to overestimate the maximal value, a finding WANG and JEGETKA explained heuristically by appealing to SLEPIAN's lemma.^{56,57}

A diametrically opposing alternative to using dense representer points with a crude approximation (8.56) would be using a few, carefully chosen representer points with a better approximation to the distribution of ϕ^* . For example, we could generate a set of Thompson samples, $\xi_i \sim p(x^* | \mathcal{D})$, then approximate the quantiles of ϕ^* by repeatedly sampling their corresponding function values ϕ . Figure 8.17 shows such an approximation for our example scenario using 100 Thompson samples. The agreement with the true distribution is excellent, but the computational cost was significant compared to the independent approximation. However, in high-dimensional spaces where dense coverage of the domain is not possible, such a direct approach may become appealing.

53 C. E. CLARK (1961). The Greatest of a Finite Set of Random Variables. *Operations Research* 9(2): 145–162.

54 A. M. ROSS (2010). Computing Bounds on the Expected Maximum of Correlated Normal Variables. *Methodology and Computing in Applied Probability* 12(1):111–138.

55 R. P. BRENT (1973). *Algorithms for Minimization Without Derivatives*. Prentice-Hall. [chapter 4]

56 D. SLEPIAN (1962). The One-Sided Barrier Problem for Gaussian Noise. *The Bell System Technical Journal* 41(2):463–501.

57 This result requires the posterior covariance function to be positive everywhere, which is not guaranteed even if true for the prior covariance function. However, it is “usually true” for typical models in high-dimensional spaces.

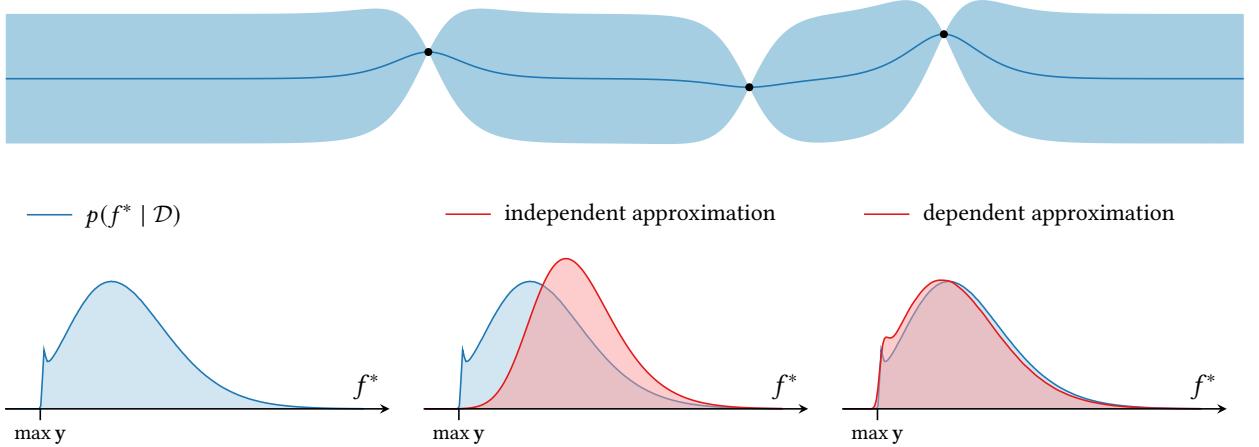


Figure 8.17: Left: the true distribution $p(f^* | \mathcal{D})$ for our running scenario, estimated using exhaustive sampling. Middle: WANG and JEGELKA's approximation to $p(f^* | \mathcal{D})$ (8.56) using a grid of 100 equally spaced representer points. Right: an approximation to $p(f^* | \mathcal{D})$ from sampling the function values at the 100 Thompson samples from figure 8.11.

Predictive entropy with exact observations

Regardless of the exact estimator we use to approximate the second term of the acquisition function, we will need to compute the entropy of the predictive distribution for y given the optimal value f^* :

$$p(y | x, f^*, \mathcal{D}) = \int p(y | x, \phi) p(\phi | x, f^*, \mathcal{D}) d\phi. \quad (8.57)$$

Remarkably, unlike in the case of predictive entropy search, the latent predictive distribution of ϕ given f^* can be derived easily as a normal distribution with upper tail truncated at f^* :

$$p(\phi | x, f^*, \mathcal{D}) = \mathcal{T}\mathcal{N}(\phi; \mu, \sigma^2 (-\infty, f^*)). \quad (8.58)$$

Even more remarkably, the entropy of this distribution can be computed in closed form:

$$H[\phi | x, f^*, \mathcal{D}] = \frac{1}{2} \left[\log(2\pi e \sigma^2 \Phi(z)^2) - z \frac{\phi(z)}{\Phi(z)} \right]; \quad z = \frac{f^* - \mu}{\sigma}. \quad (8.59)$$

In the absence of observation noise, this result is sufficient to realize a complete algorithm by combining (8.59) with (8.53–8.55). After simplification, this yields the final MES approximation, which ignores the effect of observation noise in the predictive distribution:

$$\alpha_{\text{MES}}(x; \mathcal{D}) = \frac{1}{2n} \sum_{i=1}^n \left[z_i \frac{\phi(z_i)}{\Phi(z_i)} - \log \Phi(z_i)^2 \right]; \quad z_i = \frac{f_i^* - \mu}{\sigma}. \quad (8.60)$$

Figure 8.18 illustrates this approximation for our running example using the independent approximation to $p(f^* | \mathcal{D})$ (8.56). The approximation is faithful and induces a near-optimal decision; any (slight) inaccuracy is entirely due to the independence assumption.

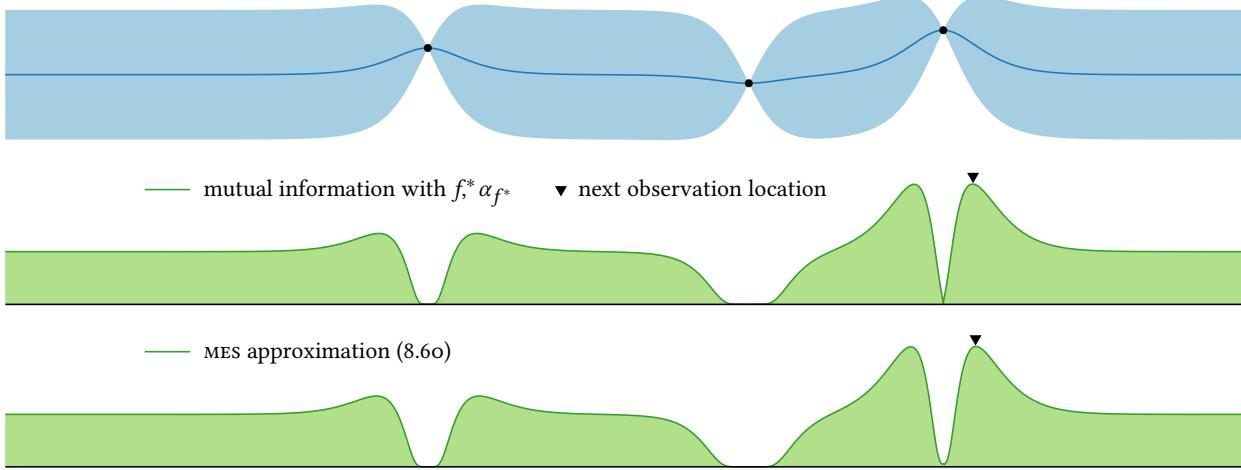


Figure 8.18: An approximation to the mutual information between the observed value y and the value of the global optimum $\alpha_{f^*} = I(y; f^* | x, \mathcal{D})$ for our running example using the independent approximation to $p(f^* | \mathcal{D})$ (8.56) and numerical integration.

Approximating predictive entropy with noisy observations

direct computation of predictive entropy

⁵⁸ S. TURBAN (2010). *Convolution of a truncated normal and a centered normal variable*. Technical report. Columbia University.

In the case of additive Gaussian noise, however, the predictive distribution of y (8.57) becomes the convolution of a centered normal distribution and a truncated normal distribution; see figure 8.19 for an illustration of a particularly noisy scenario. TURBAN provides a closed form for the resulting probability density function:⁵⁸

$$p(y | x, f^*, \mathcal{D}) = \gamma \Phi\left(\frac{\alpha(y) - y + f^*}{\beta}\right) \exp\left(-\frac{(y - \mu)^2}{2s^2}\right),$$

where

$$\alpha(y) = \frac{\sigma_n^2(y - \mu)}{s^2}; \quad \beta = \frac{\sigma_n \sigma}{s}; \quad \gamma = \frac{1}{\sqrt{2\pi s} \Phi(z)},$$

and z is defined in (8.59). Unfortunately this distribution does not admit a closed-form expression for its entropy, although we can approximate the entropy effectively via numerical quadrature.

Alternatively, we can appeal to analytic approximation to estimate the predictive entropy, and both WANG and JEGELKA and HOFFMAN and GHAHRAMANI take this approach.⁵⁹ The MES approximation simply ignores the observation noise in the acquisition function (8.60) and instead computes the mutual information between the latent function value ϕ and f^* . This approach overestimates the true mutual information by assuming exact observations, but serves as a reasonable approximation when the signal-to-noise ratio is high. This approximation is illustrated in figure 8.20 for a scenario with relatively high noise; although the approximation is not perfect, the chosen point is close to optimal.

analytic approximation to predictive entropy

⁵⁹ However, the closed form for the predictive distribution above was not discussed by either and perhaps deserves further consideration.

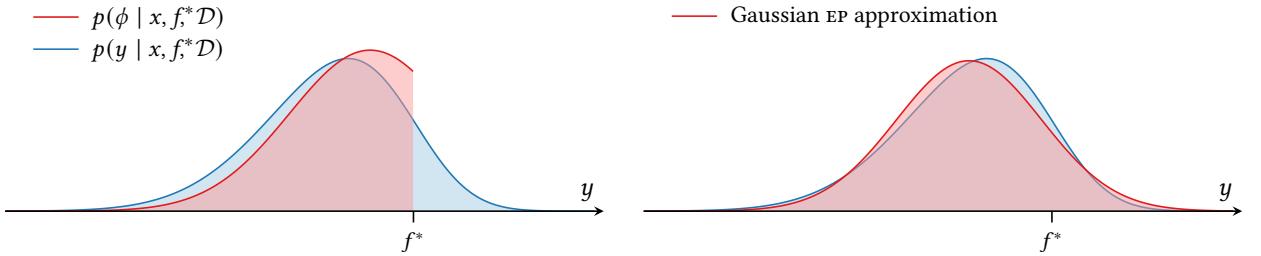


Figure 8.19: Left: an example of the latent predictive distribution $p(\phi | x, f^* \mathcal{D})$, which takes the form of a truncated normal distribution (8.58), and the resulting predictive distribution $p(y | x, f^* \mathcal{D})$, which is a convolution with a centered normal distribution accounting for Gaussian observation noise. Right: a Gaussian expectation propagation approximation to the predictive distribution (8.62).

Hoffman and Ghahramani instead approximate the latent predictive distribution (8.58) using Gaussian expectation propagation:

truncating a variable with EP: § B.2, p. 301

$$p(\phi | x, f^* \mathcal{D}) \approx \mathcal{N}(\phi; \mu_*, \sigma_*^2),$$

where

$$\sigma_*^2 = \sigma^2 \left[1 - z \frac{\phi(z)}{\Phi(z)} - \frac{\phi(z)^2}{\Phi(z)^2} \right] \quad (8.61)$$

is the variance of the truncated normal latent predictive distribution (8.58) and z is defined in (8.59). We may now approximate the differential entropy of y with

$$\begin{aligned} \text{var}[y | x, f^* \mathcal{D}] &\approx \sigma_*^2 + \sigma_n^2 = s_*^2; \\ H[y | x, f^* \mathcal{D}] &\approx \frac{1}{2} \log(2\pi e s_*^2). \end{aligned} \quad (8.62)$$

This is similar to the strategy used in predictive entropy search, although computing the approximate latent posterior is considerably easier in this case, only requiring the single factor $[\phi < f^*]$. Figure 8.19 shows the resulting Gaussian approximation to the predictive distribution for an example point; the approximation is excellent. Estimating the expectation over f^* with (8.55) and simplifying, the final OPES approximation to (8.53) is

$$\alpha_{\text{OPES}}(x; \mathcal{D}) = \log s - \frac{1}{n} \sum_{i=1}^n \log s_{*i}, \quad (8.63)$$

where s_{*i} is the approximate predictive standard deviation corresponding to the sample f_i^* . This takes the same form as the predictive entropy search approximation to the mutual information with x^* (8.41), although the approximations to the predictive distribution are of course different. The OPES approximation (8.63) is shown for a high-noise scenario in figure 8.20; the approximation is almost perfect and in this case yields the optimal decision.

Both the MES and OPES approximations to the mutual information can be differentiated with respect to the proposed observation location.

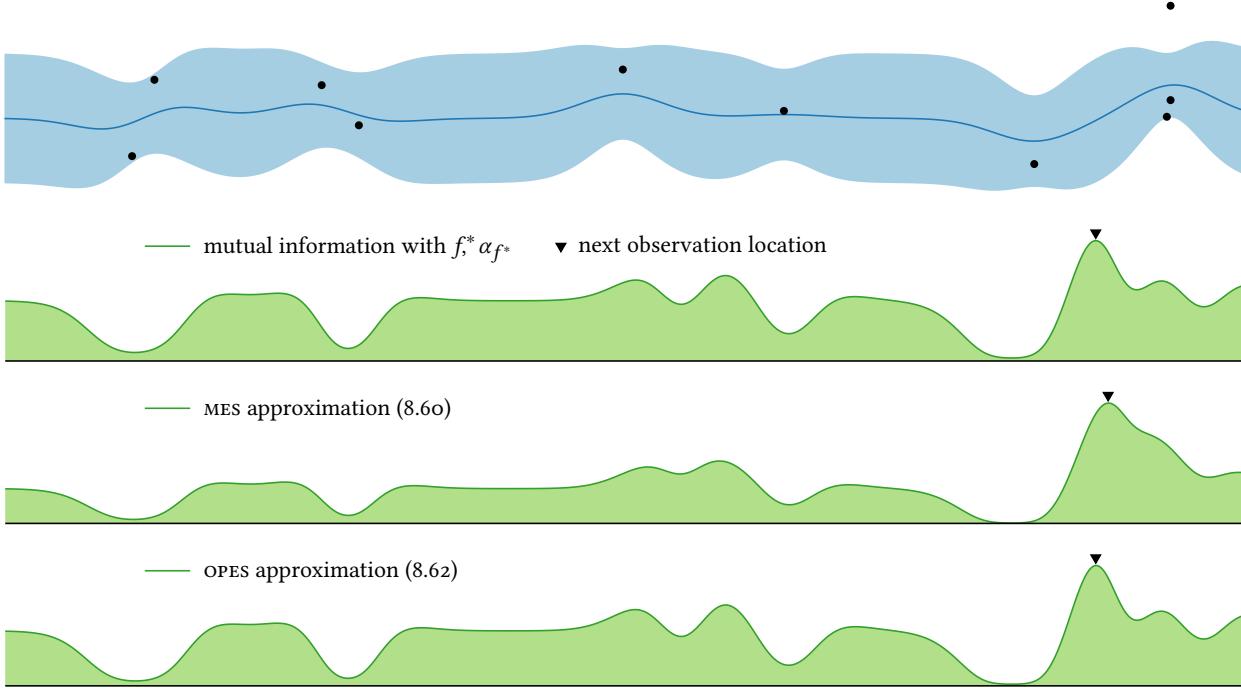


Figure 8.20: The MES and OPES approximations to the mutual information with f^* for an example high-noise scenario with unit signal-to-noise ratio. Both use the independent representer point approximation to $p(f^* | \mathcal{D})$ (8.56).

For the former, we simply differentiate (8.60):

$$\frac{\partial \alpha_{\text{MES}}}{\partial x} = \frac{1}{2n\sigma} \sum_{i=1}^n \frac{\phi(z_i)}{\Phi(z_i)} \left[\frac{\partial \mu}{\partial x} + z_i \frac{\partial \sigma}{\partial x} \right] \left[1 + z_i \frac{\phi(z_i)}{\Phi(z_i)} + z_i^2 \right];$$

note the $\{f_i^*\}$ samples will in general not depend on x , so we do not need to worry about any dependence of their distribution on the observation location. The details for the OPES approximation are provided in the appendix.

gradient of OPES acquisition function: § c.3,
p.307

model averaging: § 4.4, p. 74

8.10 AVERAGING OVER A SPACE OF GAUSSIAN PROCESSES

Throughout this chapter we have assumed our belief regarding the objective function is given by a single Gaussian process. However, especially when our dataset is relatively small, we may seek robustness to model misspecification by averaging over multiple plausible models. Here we will provide some guidance for policy computation when performing Bayesian model averaging over a space of Gaussian processes. Although this may seem straightforward, there is some nuance involved.

Below we will consider a space of Gaussian processes indexed by θ , which determines both the moments the objective function prior and any relevant parameters of the observation noise process. Bayesian

model averaging over this space yields marginal posterior and predictive distributions:

$$p(f \mid \mathcal{D}) = \int p(f \mid \mathcal{D}, \theta) p(\theta \mid \mathcal{D}) d\theta; \quad (8.64)$$

$$p(y \mid x, \mathcal{D}) = \int p(y \mid x, \mathcal{D}, \theta) p(\theta \mid \mathcal{D}) d\theta, \quad (8.65)$$

which are integrated against the model posterior $p(\theta \mid \mathcal{D})$ (4.7). Both of these distributions are in general intractable, but we developed several viable approximations in chapter 3, all of which approximate the objective function posterior (8.64) with a mixture of Gaussian processes and the posterior predictive distribution (8.65) with a mixture of Gaussians.

model posterior, $p(\theta \mid \mathcal{D})$

Noiseless expected improvement and probability of improvement

When observations are exact, the marginal gain in utility underlying both expected improvement and probability of improvement depends *only* on the objective function value ϕ and the value of the incumbent ϕ^* :

$$\Delta_{EI}(x, \phi) = \max(\phi - \phi^*, 0); \quad \Delta_{PI}(x, \phi) = [\phi > \phi^*].$$

As a result, the formulas we derived for these acquisition functions given a GP belief on the objective (8.9, 8.21) depend only on the moments of the (Gaussian) predictive distribution $p(\phi \mid x, \mathcal{D})$.⁶⁰ With a Gaussian *mixture* approximation to the predictive distribution, the expected marginal gain,

$$\alpha(x; \mathcal{D}) = \int \Delta(x, \phi) p(\phi \mid x, \mathcal{D}) d\phi,$$

is simply a weighted combination of these results by linearity of expectation.

⁶⁰ In fact, a Gaussian *process* belief is not required at all, only Gaussian predictive distributions. This will be important in the next section.

Model-dependent utility functions

For the remaining decision-theoretic acquisition functions, noisy expected improvement and probability of improvement, knowledge gradient, and mutual information with any relevant random variable ω , the situation is somewhat more complicated, because the underlying utility functions *depend on the model of the objective function*. The first three depend on the posterior mean function, and the last depends on the posterior belief $p(\omega \mid \mathcal{D})$, both of which are induced by our belief regarding the objective function. To make this dependence explicit, for a model θ in our space of interest, let us respectively notate the utility, marginal gain in utility, and expected marginal gain in utility with:

$$u(\mathcal{D}; \theta); \quad \Delta(x, y; \theta); \quad \alpha(x; \mathcal{D}, \theta).$$

There are two natural ways we might address this dependence when averaging over models. One is to seek to maximize the *expected* utility

dependence on objective function model

average model-conditional utility, $\mathbb{E}u$

marginal gain in $\mathbb{E}u$, $\mathbb{E}\Delta$
expected marginal gain in $\mathbb{E}u$, $\mathbb{E}\alpha$

marginal posterior mean

marginal belief about ω

utility of marginal model, $u\mathbb{E}$
(expected) marginal gain in utility of
marginal model: $\Delta\mathbb{E}$, $\alpha\mathbb{E}$

of the data, averaged over the choice of model:

$$\mathbb{E}u(\mathcal{D}) = \int u(\mathcal{D}; \theta) p(\theta | \mathcal{D}) d\theta. \quad (8.66)$$

Writing the marginal gain in expected utility as $\mathbb{E}\Delta$, we may derive an acquisition function via one-step lookahead:

$$\begin{aligned} \mathbb{E}\alpha(x; \mathcal{D}) &= \int \mathbb{E}\Delta(x, y) p(y | x, \mathcal{D}) dy \\ &= \int \left[\int \Delta(x, y; \theta) p(y | x, \mathcal{D}, \theta) dy \right] p(\theta | \mathcal{D}) d\theta \\ &= \int \alpha(x; \mathcal{D}, \theta) p(\theta | \mathcal{D}) d\theta. \end{aligned} \quad (8.67)$$

As hinted by its notation, this is simply the expectation of the conditional acquisition functions, which we can approximate via standard methods.

Although this approach is certainly convenient, it may overestimate optimization progress as utility is *only measured under the assumption of a perfectly identified model* – the utility function is “blind” to model uncertainty. An arguably more appealing alternative is to evaluate utility with respect to the marginal objective function model (8.64) from the start, defining simple and global reward with respect to the marginal posterior mean:

$$\int \mu_{\mathcal{D}}(x; \theta) p(\theta | \mathcal{D}) d\theta, \quad (8.68)$$

and information gain about ω with respect to its marginal belief:

$$\int p(\omega | \mathcal{D}, \theta) p(\theta | \mathcal{D}) d\theta. \quad (8.69)$$

Let us notate a utility function defined in this manner with $u\mathbb{E}(\mathcal{D})$, contrasting with its post hoc averaging equivalent, $\mathbb{E}u(\mathcal{D})$ (8.66). Similarly, let us notate its marginal gain with $\Delta\mathbb{E}$ and its expected marginal gain with:

$$\alpha\mathbb{E}(x; \mathcal{D}) = \int \Delta\mathbb{E}(x, y) p(y | x, \mathcal{D}) dy. \quad (8.70)$$

Example and discussion

We may shed some light on the differences between these approaches with a barebones example. Let us work with the knowledge gradient, and consider a pair of simple models for an objective function on the interval $\mathcal{X} = [-1, 1]$: either $f = x$ or $f = -x$, with equal probability.⁶¹

The model-conditional global reward in either case is 1, and thus the expected utility (8.66) *does not depend on the data*: $\mathbb{E}u(\mathcal{D}) \equiv 1$. What a strange set of affairs – although we know the maximal *value* of the objective a priori, we need data to tell us where it is! For this particular model space, an optimal recommendation for either model is maximally *suboptimal* if that model is incorrect.

⁶¹ These models may be interpreted as degenerate Gaussian processes with covariance $K \equiv 0$.

In contrast, the global reward of the marginal model *does* depend on the data via the model posterior. The global reward is monotonic in the probability of the model most favored by the data, π :

$$u\mathbb{E}(\mathcal{D}) = 2\pi - 1.$$

A priori, the utility is $u\mathbb{E}(\emptyset) = 0$ – the marginal mean function is $\mu \equiv 0$, and no compelling recommendation is possible until we determine the correct model.

As the expected utility $\mathbb{E}u$ (8.66) is independent of the data, it does not lead to an especially insightful policy. In contrast, the marginal gain in $u\mathbb{E}$ (8.70) can differentiate potential observation locations via their expected impact on the model posterior. We illustrate this acquisition function in the margin given an empty dataset and assuming moderate additive Gaussian noise. As one might hope, it prefers evaluating on the boundary of the domain, where observations are expected to provide more information regarding the model and thus greater improvement in model-marginal utility.

Computing expected gain in marginal utility

Unfortunately, the expected gain in utility for the marginal model (8.70) does not simplify as before (8.66), as we must account for the effect of the observed data on the model posterior in the utility function. However, we may sketch a Monte Carlo approximation using samples from the model posterior, $\{\theta_i\} \sim p(\theta | \mathcal{D})$.

For utility functions based on the marginal posterior mean (8.68), the resulting approximation to the expected marginal gain (8.70) is a weighted sum of Gaussian expectations of $\Delta\mathbb{E}$, each of which we may approximate via Gauss–Hermite quadrature. To approximate $\Delta\mathbb{E}$ for a putative observation (x, y) , a simple sequential Monte Carlo approximation to the updated model posterior would reweight each sample θ_i by $p(y | x, \mathcal{D}, \theta_i)$, then approximate the updated marginal posterior mean by a weighted sum.

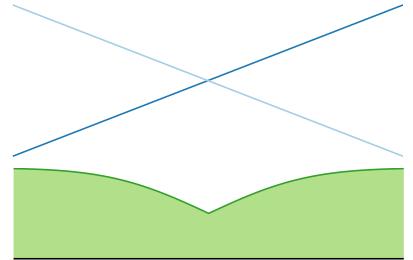
To approximate the predictive form of mutual information with x^* or f^* (8.36, 8.53):

$$H[y | x, \mathcal{D}] - \mathbb{E}_\omega[H[y | x, \omega, \mathcal{D}] | x, \mathcal{D}],$$

we first note that the first term is the entropy of a Gaussian mixture, which we can approximate via quadrature. We may approximate the expectation in the second term via Thompson sampling (see below); for each of these samples, we may approximate the updated predictive distribution as before, reweighting the mixture by $p(\omega | \mathcal{D}, \theta)$.⁶²

Upper confidence bound and Thompson sampling

We may compute any desired upper confidence bound of a Gaussian process mixture at x by bisecting the cumulative distribution function of ϕ , which is a weighted sum of Gaussian CDFs.



Above: the two possible objectives for our comparison of the $\mathbb{E}u$ and $u\mathbb{E}$. Below: the expected marginal gain in $u\mathbb{E}$, which prefers sampling on the boundary to reveal more information regarding the model. The expected marginal gain in $\mathbb{E}u$ is constant.

Gauss–Hermite quadrature: § 8.5, p. 171

⁶² This requires one final layer of approximation, which is produced as a byproduct of expectation propagation in the case of x^* (B.6), and may be dealt with without much trouble in the case of the univariate random variable f^* .

Thompson sampling for GPS: § 8.7, p. 176

Finally, to perform Thompson sampling from the marginal posterior, we first sample from the model posterior, $\theta \sim p(\theta | \mathcal{D})$; the conditional posterior $p(f | \mathcal{D}, \theta)$ is then a GP, and we may proceed via the previous discussion.

8.11 ALTERNATIVES TO GAUSSIAN PROCESSES

Although Gaussian processes are without question the most prominent objective function model used in Bayesian optimization, we are of course free to use any other model when prudent. Below we briefly outline some notable alternative model classes that have received some attention in the context of Bayesian optimization and comment on any issues arising in computing common policies with these surrogates.

Random forests

63 L. BREIMAN (2001). Random Forests. *Machine Learning* 45(1):5–32.

64 M. FERNÁNDEZ-DELGADO et al. (2014). Do we Need Hundreds of Classifiers to Solve Real World Classification Problems? *Journal of Machine Learning Research* 15(90):3133–3181.

65 F. HUTTER et al. (2014). Algorithm runtime prediction: Methods & evaluation. *Artificial Intelligence* 206:79–111.

66 F. HUTTER et al. (2011). Sequential Model-Based Optimization for General Algorithm Configuration. *LION* 5.

67 HUTTER et al. then fit a single Gaussian distribution to this mixture via moment matching, although this is not strictly necessary.

68 A similar approach can also be used to estimate arbitrary predictive quantiles:

N. MEINSHAUSEN (2006). Quantile Regression Forests. *Journal of Machine Learning Research* 7(35):983–999.

*Random forests*⁶³ are a popular model class renown for their excellent off-the-shelf performance,⁶⁴ offering good generalization, strong resistance to overfitting, and efficient training and prediction. Of particular relevance for optimization, random forests are adept at handling high-dimensional data and categorical and conditional features, and may be a better choice than Gaussian processes for objectives featuring any of these characteristics.

Algorithm configuration is one setting where these capabilities are critical: complex algorithms such as compilers or SAT solvers often have complex configuration schemata with many mutually dependent parameters, and it can be difficult to build nontrivial covariance functions for such inputs. Random forests require no special treatment in this setting and have delivered impressive performance in predicting algorithmic performance measures such as runtime.⁶⁵ They are thus a natural choice for Bayesian optimization of these same measures.⁶⁶

Classical random forests are not particularly adept at quantifying uncertainty in predictions off-the-shelf. Seeking more nuanced uncertainty quantification, HUTTER et al. proposed a modification of the vanilla model wherein leaves store both the mean (as usual) and the standard deviation of the training data terminating there.⁶⁵ We then estimate the predictive distribution with a mixture of Gaussians with moments corresponding to the predictions of the member trees.^{67,68} Figure 8.21 compares the predictions of a Gaussian process and a random forest model on a toy dataset. Although they differ in their extrapolatory behavior, the models make very similar predictions otherwise.

To realize an optimization policy with a random forest, HUTTER et al. suggested approximating acquisition functions depending only on marginal predictions – such as (noiseless) expected improvement or probability of improvement – by simply plugging this Gaussian approximation into the expressions derived in this chapter (8.9, 8.22). Either can be computed easily from a Gaussian mixture predictive distribution as well due to linearity of expectation.

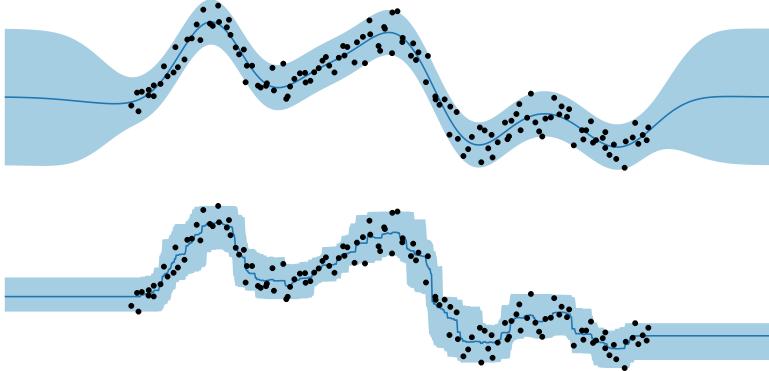


Figure 8.21: The predictions of a Gaussian process model (above) and a random forest model comprising 100 regression trees (below) for an example dataset; the credible intervals of the latter are not symmetric as the predictive distribution is estimated with a Gaussian mixture.

Conditional density estimation

BERGSTRA et al. described a lightweight Bayesian optimization algorithm – scaling only linearly with the number of observations – that operates by maximizing expected improvement (7.3) for a simple objective function model.⁶⁹ In fact, the model underlying this algorithm is so simple that it is *incomplete* in the sense that we cannot compute the predictive distribution $p(y | x, \mathcal{D})$ at all without further modeling! Nonetheless, we can still (indirectly) maximize the expected improvement to yield an effective algorithm.

We begin each iteration of this algorithm by choosing some reference value y^* that we wish to exceed with our next observation; we will then select the next observation location to maximize the expected improvement over this threshold:⁷⁰

$$\mathbb{E}[\max(y - y^*, 0) | x, \mathcal{D}]. \quad (8.71)$$

We will discuss the selection of y^* further shortly. We proceed by forming two conditional probability density estimates referencing this threshold:

$$g(x) \approx p(x | y > y^*, \mathcal{D}); \quad \ell(x) \approx p(x | y \leq y^*, \mathcal{D}); \quad (8.72)$$

that is, g estimates the probability density of observation locations exceeding the threshold, and ℓ of locations failing to. Regardless of the exact nature of these estimates, BERGSTRA et al. showed that maximizing the expected improvement over y^* (8.71) is equivalent to maximizing their ratio:⁷¹

$$\alpha(x; \mathcal{D}) = g(x)/\ell(x). \quad (8.73)$$

We are free to design the density estimates (8.72) however we please, but BERGSTRA et al. proposed one simple and efficient scheme. We construct separate kernel (Parzen) density estimates⁷² $\{g_j, \ell_j\}$ along each dimension of the domain,⁷³ setting the kernel bandwidth to ensure that unexplored regions beyond the range of the data have nontrivial density. For this approach to be feasible, we must have enough observations to estimate both densities; BERGSTRA et al. addressed this issue by taking y^*

⁶⁹ J. BERGSTRA et al. (2011). Algorithms for Hyperparameter Optimization. *NeurIPS 2011*.

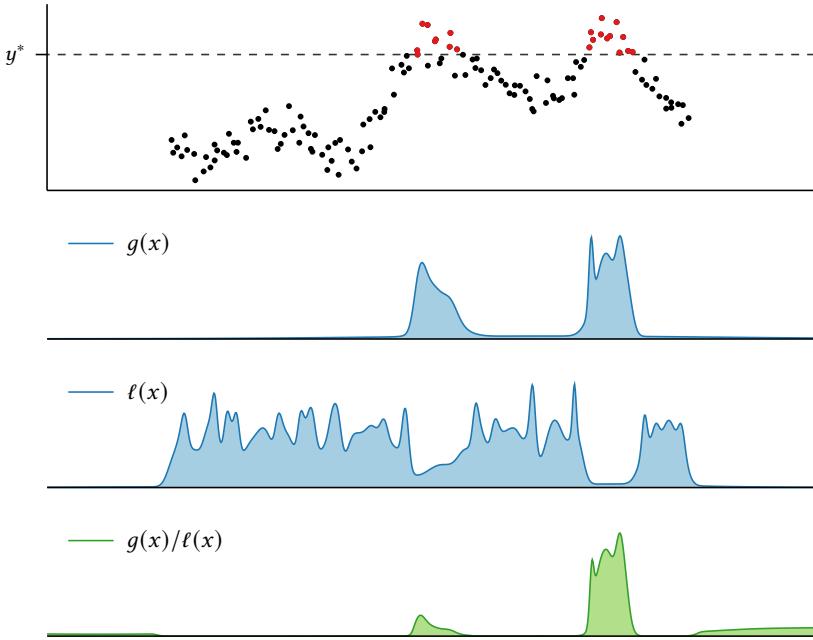
⁷⁰ Although we have spent considerable effort arguing *against* maximizing the “improvement” over a noisy observation – at least in the presence of high noise – we will see that this choice provides considerable computational benefit.

⁷¹ This is *not* the expected improvement over y^* , which would require further modeling of $p(y)$ to compute. However, the expected improvement is monotonic in this quantity, so maximizing yields the same policy.

⁷² B. W. SILVERMAN (1986). *Density Estimation for Statistics and Data Analysis*. Chapman & Hall.

⁷³ BERGSTRA et al.’s focus was hyperparameter tuning, where conditional dependence among variables can lead to complex structure in the domain. BERGSTRA et al. considered domains with *tree* structure, where each node represents a potential variable that can be assigned provided the assignments of its parents. Due to this structure, they called their modeling approach *tree Parzen estimation*. This does not rule out its use on simpler domains.

Figure 8.22: BERGSTRA et al.’s Parzen estimation optimization policy. The top panel shows an example dataset along with a threshold y^* set to be the 85th percentile of the observed values. The next two panels illustrate the central kernel density estimates (8.72) for the density of observation locations above and below this threshold. The “wiggleness” of these estimates stems from the kernel bandwidth scheme proposed by BERGSTRA et al. The bottom panel shows the ratio of these densities, which is monotonic with the expected improvement over y^* (8.73).



to be a relatively (but not exceedingly) high quantile of the observed data, using the 85th percentile in their experiments. Given these estimates, BERGSTRA et al. realized a policy by first generating a set of proposals $\{x_i\}$ by sequentially sampling each component x_{ij} from the relevant “enthusiastic” density g_j (8.72).⁷⁴ Finally, we evaluate the proposal maximizing the expected improvement (8.73).

Figure 8.22 illustrates the key components of this algorithm for a toy dataset in one dimension. The next observation location cannot be determined from the acquisition function due to the stochastic nature of the policy, but we can identify several likely regions offering different tradeoffs between exploration and exploitation.

Bayesian neural networks

Given the predictive power and modeling flexibility offered by neural networks, it is no surprise that (Bayesian) neural networks have seen serious consideration in Bayesian optimization.

SNOEK et al. took an early step in this direction and reported success with a relatively simple construction.⁷⁵ The authors first trained a typical (non-Bayesian) deep neural network using traditional empirical loss minimization, then replaced the final layer with Bayesian linear regression. The weights in all but the final layer were then fixed for the remainder of the procedure; the resulting model can thus be interpreted as Bayesian linear regression using highly complex neural basis functions. As this is in fact a Gaussian process,⁷⁶ we may appeal to the computational details outlined in the remainder of this chapter to compute any policy

⁷⁴ For a tree-structured domain, sampling from the root down.

⁷⁵ J. SNOEK et al. (2015). Scalable Bayesian Optimization Using Deep Neural Networks. *ICML 2015*.

⁷⁶ See the result in (3.6), setting $K \equiv 0$ as there is no nonlinear component in the model.

we desire. One beneficial side effect of this model is that the cost of inference for Bayesian linear regression is only linear with the number of observations.

A somewhat more involved (and, arguably, “more Bayesian”) approach was proposed by SPRINGENBERG et al.,⁷⁷ who combined a parametric objective function model $f(x; \mathbf{w})$ – the output of a neural network with input x and weights \mathbf{w} – with an additive Gaussian noise observation model:

$$p(y | x, \mathbf{w}, \sigma) = \mathcal{N}(y; f(x; \mathbf{w}), \sigma^2).$$

Bayesian inference for this model proceeds by selecting a prior $p(\mathbf{w}, \sigma)$ and computing the posterior from the observed data. The highly non-linear nature of neural networks renders this posterior intractable, but SPRINGENBERG et al. described a stochastic Hamiltonian Monte Carlo scheme for drawing samples $\{\mathbf{w}_i, \sigma_i\}_{i=1}^s$ from the posterior. We may then use these samples to form a Gaussian mixture approximation to the predictive distribution:⁷⁸

$$p(y | x, \mathcal{D}) \approx \frac{1}{s} \sum_{i=1}^s \mathcal{N}(y; f(x; \mathbf{w}_i), \sigma_i^2).$$

As in our discussion on random forests, this is sufficient to compute policies that only require access to Gaussian marginal predictions such as (noiseless) expected improvement and probability of improvement (8.9, 8.22).

SUMMARY OF MAJOR IDEAS

In this chapter we considered the computation of the popular optimization policies described in the last chapter for Gaussian process models of an objective function with an exact or additive Gaussian noise observation model. The acquisition functions for most of these policies represent the one-step expected marginal gain to some underlying utility function:

$$\alpha(x; \mathcal{D}) = \int \Delta(x, y) \mathcal{N}(y; \mu, s^2) dy,$$

where $\Delta(x, y)$ is the gain in utility resulting from the observation (x, y) (8.27). When Δ is a (not necessarily continuous) piecewise linear function of y , this integral can be resolved analytically in terms of the standard normal CDF. This is the case for the expected improvement and probability of improvement acquisition functions, both with and without observation noise.

However, when Δ is a more complicated function of the putative observation, we must in general rely on approximate computation to resolve this integral. When the predictive distribution is normal – as in the model class considered in this chapter – *Gauss–Hermite quadrature* provides a useful and sample-efficient approximation via a weighted average of carefully chosen integration nodes. This allows us to address some more complex acquisition functions such as the knowledge gradient.

⁷⁷ J. T. SPRINGENBERG et al. (2016). Bayesian Optimization with Robust Bayesian Neural Networks. *NeurIPS 2016*.

⁷⁸ Following HUTTER et al.’s approach to optimization with random forests, SPRINGENBERG et al. fit a single Gaussian distribution to this mixture via moment matching, but this is optional.

computation of expected improvement and probability of improvement: §§ 8.2–8.3, p. 167

approximate computation for one-step lookahead: § 8.5, p. 171

computation of knowledge gradient: § 8.6, p. 172

Thompson sampling: § 8.7, p. 176
 approximate computation of mutual information: §§ 8.8–8.9, p. 180

averaging over a space of GPS: § 8.10, p. 192

alternatives to GPS: § 8.11, p. 196

The computation of mutual information with x^* or f^* entails an expectation with respect to these random variables, which cannot be approximated using simple quadrature schemes. Instead, we must rely on schemes such as Thompson sampling – a notable policy in its own right – to generate samples and proceed via (simple or quasi-) Monte Carlo integration, and in some cases, further approximations to the conditional predictive distributions resulting from these samples.

Finally, Bayesian optimization is of course not limited to a single Gaussian process belief on the objective function, which may be objectionable even when Gaussian processes are the preferred model class due to uncertainty in hyperparameters or model structure. Averaging over a space of Gaussian processes is possible – with some care – by adopting a Gaussian process mixture approximation to the marginal objective function posterior and relying on results from the single GP case. If desired, we may also abandon the model class entirely and compute policies with respect to an alternative such as random forests or Bayesian neural networks.

9

IMPLEMENTATION

There is a rich and mature software ecosystem available for Gaussian process modeling and Bayesian optimization, and it is relatively easy to build sophisticated optimization routines using off-the-shelf libraries. However, successful implementation of the underlying algorithms requires attending to some nitty-gritty details to ensure optimal performance, and what may appear to be simple equations on the page can be challenging to realize in a limited-precision environment. In this chapter we will provide a brief overview of the computational details that practitioners should be aware of when designing Bayesian optimization algorithms, even when availing themselves of existing software libraries.

9.1 GAUSSIAN PROCESS INFERENCE

As is typical with nonparametric models, the computational cost of Gaussian process inference grows (considerably!) with the number of observations, and it is important to understand the nature of this growth and be aware of methods for scaling to large-scale data when necessary.

The primary computational bottleneck in Gaussian process inference is solving systems of linear equations that scale with the number of observed values. Consider the general case of exact inference where we condition a Gaussian process $\mathcal{GP}(f; \mu, K)$ on the observation of a length- n vector of values \mathbf{y} with marginal distribution and cross-covariance function (2.6):

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}; \mathbf{m}, \mathbf{C}); \quad \kappa(x) = \text{cov}[\mathbf{y}, \phi | x]. \quad (9.1)$$

The posterior is a Gaussian process with moments (2.10):

$$\begin{aligned} \mu_{\mathcal{D}}(x) &= \mu(x) + \kappa(x)^T \mathbf{C}^{-1} (\mathbf{y} - \mathbf{m}); \\ K_{\mathcal{D}}(x, x') &= K(x, x') - \kappa(x)^T \mathbf{C}^{-1} \kappa(x'). \end{aligned} \quad (9.2)$$

Evaluating either the posterior mean or the posterior covariance requires solving a linear system with respect to the observation covariance \mathbf{C} . Inference with non-Gaussian observations using Monte Carlo sampling or Gaussian approximate inference also entails solving linear systems with respect to this matrix (2.34, 2.38–2.40).

Direct computation via Cholesky decomposition

NEAL outlined a straightforward implementation based on direct numerical methods that suffices for the small-to-moderate datasets typical in Bayesian optimization.¹ We take advantage of the fact that \mathbf{C} is symmetric and positive definite and precompute and store its Cholesky factorization:²

$$\mathbf{L} = \text{chol}(\mathbf{C}); \quad \mathbf{LL}^T = \mathbf{C}.$$

This computation requires one-time $\mathcal{O}(n^3)$ work and represents the bulk of effort spent in this implementation. Note that despite having the same

solving linear systems with respect to the observation covariance \mathbf{C}

number of observed values, $n = |\mathbf{y}|$

¹ R. M. NEAL (1998). Regression and Classification Using Gaussian Process Priors. In: *Bayesian Statistics 6*.

² G. H. GOLUB and C. F. VAN LOAN (2013). *Matrix Computations*. Johns Hopkins University Press. [§ 4.2]

asymptotic running time as the general-purpose LU decomposition, the Cholesky factorization exploits symmetry to run twice as fast.

With the Cholesky factor in hand, we may solve an arbitrary linear system $\mathbf{C}\mathbf{x} = \mathbf{b}$ in time $\mathcal{O}(n^2)$ via forward–backward substitution by rewriting as $\mathbf{L}\mathbf{L}^\top\mathbf{x} = \mathbf{b}$ and twice applying a triangular solver. We exploit this fact to additionally precompute the vector

$$\boldsymbol{\alpha} = \mathbf{C}^{-1}(\mathbf{y} - \mathbf{m})$$

appearing in the posterior mean. After this initial precomputation of \mathbf{L} and $\boldsymbol{\alpha}$, we may compute the posterior mean

$$\mu_{\mathcal{D}}(\mathbf{x}) = \mu(\mathbf{x}) + \kappa(\mathbf{x})^\top \boldsymbol{\alpha}$$

on demand in linear time and the posterior covariance

$$K_{\mathcal{D}}(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}, \mathbf{x}') - [\mathbf{L}^{-1}\kappa(\mathbf{x})]^\top [\mathbf{L}^{-1}\kappa(\mathbf{x}')]^\top$$

in quadratic time. We may also efficiently compute the log marginal likelihood (4.8) of the data in linear time:³

$$\log p(\mathbf{y} \mid \mathbf{x}) = -\frac{1}{2}[(\mathbf{y} - \mathbf{m})^\top \boldsymbol{\alpha} + 2 \sum_i \log L_{ii} + n \log 2\pi].$$

Low-rank updates to the Cholesky factorization for sequential inference

Optimization is an inherently sequential procedure, where our dataset grows incrementally as we gather new observations. In this setting, we can accelerate sequential inference with a *fixed* Gaussian process⁴ by replacing direct computation of the Cholesky decomposition in favor of a fast incremental updates to previously computed Cholesky factors.

For the sake of argument, suppose we have computed the Cholesky factor $\mathbf{L} = \text{chol } \mathbf{C}$ for a set of n observations \mathbf{y} with covariance \mathbf{C} . Suppose we then receive k additional observations \mathbf{v} , resulting in the augmented observation vector

$$\mathbf{y}' = \begin{bmatrix} \mathbf{y} \\ \mathbf{v} \end{bmatrix}.$$

The covariance matrix of \mathbf{y}' is formed by appending the previous covariance \mathbf{C} with new rows/columns:

$$\text{cov}[\mathbf{y}'] = \mathbf{C}' = \begin{bmatrix} \mathbf{C} & \mathbf{X}_1^\top \\ \mathbf{X}_1 & \mathbf{X}_2 \end{bmatrix}.$$

The Cholesky factor of the updated covariance matrix has the form

$$\text{chol } \mathbf{C}' = \begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{\Lambda}_1 & \mathbf{\Lambda}_2 \end{bmatrix}.$$

Note that the upper-left block is simply the previously computed Cholesky factor, which we can reuse. The new blocks may be computed as

$$\mathbf{\Lambda}_1 = \mathbf{X}_1 \mathbf{L}^{-\top}; \quad \mathbf{\Lambda}_2 = \text{chol}[\mathbf{X}_2 - \mathbf{\Lambda}_1 \mathbf{\Lambda}_1^\top].$$

The first of these blocks can be computed efficiently using a triangular solver with the previous Cholesky factor, and the second block only requires factoring a $(k \times k)$ matrix.

This low-rank update requires $\mathcal{O}(kn^2)$ work to compute, compared to the $\mathcal{O}(n^3 + kn^2)$ work required to compute the Cholesky decomposition of C' from scratch; asymptotically, the low-rank update is $\mathcal{O}(n/k)$ times faster. In particular, if we begin with an empty dataset and sequentially apply this update for a total of n observations (in any order and with any number of observations at a time) the *total* cost of inference would be $\mathcal{O}(n^3)$. This is equivalent to the cost of one-time inference with the full dataset, so the update scheme is as efficient as one could hope for.

amortized analysis

Ill-conditioned covariance matrices

When an optimization policy elects to make an observation in the pursuit of “exploitation,” the value observed is (by design!) highly correlated with at least one existing observation. Although these decisions may be well-grounded, the resulting highly correlated observations can wreak havoc on numerical linear algebra routines.

To sketch the problems that may arise, consider the following scenario. Suppose the covariance function is stationary, that observations are corrupted by additive Gaussian noise with variance σ_n^2 , and that the signal-to-noise ratio is large. Now, if two locations in the dataset correspond to highly correlated values, the corresponding rows/columns in the observation covariance matrix C will be nearly equal, and C will thus be nearly singular: one or more eigenvalue will be near zero, and in extreme cases, some may even become (numerically) negative. This poor conditioning⁵ can cause loss of precision when solving a linear system via the Cholesky decomposition, and a negative eigenvalue will cause the Cholesky routine to fail altogether.

When necessary, we may sidestep these issues with a number of “tricks.” One simple solution is to add a small (in terms of σ_n^2) multiple of the identity to the observation covariance, replacing C by $C + \varepsilon I$. Numerically, this shifts the singular values of C by ε , improving conditioning; practically, this caps the signal-to-noise ratio by increasing the noise floor. When high correlation is a result of small spatial separation, OSBORNE et al. suggested replacing problematic observations of the objective function with (noisy) observations of directional derivatives,⁶ which are only weakly correlated with the corresponding function values.⁷ Another option is to appeal to iterative numerical methods, discussed below, which actually *benefit* from having numerous eigenvalues clustered around zero.

⁵ Numerical conditioning of C is usually evaluated by its *condition number*, defined in terms of its singular values $\sigma = \{\sigma_i\}$:

$$\kappa = \frac{\max \sigma}{\min \sigma}.$$

Given infinite precision, the singular values are equal to the (nonnegative) eigenvalues of C , but very poor conditioning can cause negative (numerical) eigenvalues.

⁶ That is, we replace observations of $f(x)$ and $f(x')$ with one of $f(x)$ and one of the directional derivative in the direction of $x - x'$, evaluated at the midpoint $(x + x')/2$.

⁷ M. A. OSBORNE et al. (2009). Gaussian Processes for Global Optimization. *LION 3*.

Iterative numerical methods

Direct methods can handle datasets of perhaps a few tens of thousands of observations before the cubic scaling becomes too much to bear. In most settings where Bayesian optimization would be considered, the

cost of observation will preclude obtaining a dataset anywhere near this size, in which case we need not consider the issue further. However, when necessary, we may appeal to more complex approximate inference schemes to scale to larger datasets.

One line of work in this direction is to solve the linear systems arising in the posterior with iterative rather than direct numerical methods. The method of *conjugate gradients* is especially well-suited as it is designed for symmetric positive-definite systems such as appear in the GP posterior.⁸ The main idea behind the conjugate gradient method is to reinterpret the solution of the linear system $Cx = b$ as the solution of a related and particularly well-behaved convex optimization problem.⁹ With this insight, we may derive a simple procedure to construct a sequence of vectors $\{x_i\}$ guaranteed to converge in finite time to the desired solution. For a system of size n , each iteration of this procedure requires only $\mathcal{O}(n^2)$ work; the most expensive operation is a single matrix–vector multiplication with C .

The method of conjugate gradients is guaranteed to converge (up to round-off error) after n iterations, but this does not offer any speedup over direct methods. However, when C has a well-behaved spectrum – that is, it is well-conditioned and/or has clustered eigenvalues – the sequence converges rapidly. In many cases we may terminate after only $k \ll n$ iterations with an accurate estimate of the solution, for an effectively quadratic running time of $\mathcal{O}(kn^2)$. Although many covariance matrices arising in practice are not necessarily well-conditioned, we may use a technique known as *preconditioning* to transform a poorly conditioned matrix to speed up convergence, with only minor overhead.^{10,11}

⁸ M. R. HESTENES and E. STIEFEL (1952). Methods of Conjugate Gradients for Solving Linear Systems. *Journal of Research of the National Bureau of Standards* 49(6):409–436.

⁹ If $Cx = b$, then $|Cx - b| = 0$. As this norm is nonnegative and only vanishes at x , squaring gives

$$x = \arg \min_y y^\top Cy - 2b^\top y.$$

Thus x is also the solution of an unconstrained quadratic optimization problem, which is convex as C is symmetric positive definite.

¹⁰ G. H. GOLUB and C. F. VAN LOAN (2013). *Matrix Computations*. Johns Hopkins University Press. [§ 11.5]

¹¹ K. CUTAJAR et al. (2016). Preconditioning Kernel Matrices. *JCML* 2016.

¹² M. N. GIBBS (1997). Bayesian Gaussian Processes for Regression and Classification. PhD thesis. University of Cambridge.

¹³ J. R. GARDNER et al. (2018). GPyTorch: Black-box Matrix–Matrix Gaussian Process Inference with GPU Acceleration. *NeurIPS* 2018.

The use of conjugate gradients for GP inference can be traced back to the doctoral work of GIBBS.¹² Numerous authors have provided enhancements in the intervening years, and there is a now a substantial body of related work. A good starting point is the work of GARDNER et al., who provide a review of the literature and the key ideas from numerical linear algebra required for large-scale GP inference.¹³ The authors refine these tools to exploit modern massively parallel hardware and build an accompanying software package scaling inference to hundreds of thousands of observations.

Sparse approximations

An alternative approach for scaling to large datasets is *sparse approximation*. Here rather than approximating the linear algebra arising in the exact posterior, we approximate the posterior distribution itself with a Gaussian process admitting tractable computation with direct numerical methods. A large family of sparse approximations have been proposed, which differ in their details but share the same general approach.

As we have seen, specifying an *arbitrary* Gaussian distribution for a set of values jointly Gaussian distributed with a function of interest induces a GP posterior consistent with that belief (2.38). This is a powerful tool, used in approximate inference to optimize the fit of the induced

Gaussian process to a true, intractable posterior. Sparse approximation methods make use of this property as well, but to achieve computational rather than mathematical tractability. The idea is to craft a Gaussian belief for a sufficiently small set of values such that the induced posterior is a faithful, but tractable approximation to the true posterior.

For this discussion it is important that we explicitly account for any observation noise that may be present in the values we wish to condition on, as only *independent* noise – that is, with diagonal error covariance – is suitable for sparse approximation. Consider conditioning a Gaussian process on a vector of n values \mathbf{z} , for large n . We will assume the observation model $\mathbf{z} = \mathbf{y} + \boldsymbol{\epsilon}$, where \mathbf{y} is a vector of jointly Gaussian distributed values as in (9.1), and $\boldsymbol{\epsilon}$ is a vector of independent, zero-mean Gaussian measurement noise with diagonal covariance matrix \mathbf{N} .

The first step of sparse approximation is to identify a set of $m \ll n$ values \mathbf{v} , called *inducing values*, whose distribution can in some sense capture most of the information in the full dataset. We will discuss the selection of inducing values shortly; for the moment we assume an arbitrary set has been chosen. The joint prior distribution of the observed and inducing values is Gaussian:

$$p(\mathbf{v}, \mathbf{z}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{v} \\ \mathbf{z} \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{m} \end{bmatrix}, \begin{bmatrix} \Sigma & \mathbf{K}^\top \\ \mathbf{K} & \mathbf{C} + \mathbf{N} \end{bmatrix}\right), \quad (9.3)$$

and we will write the cross-covariance function for \mathbf{v} as

$$k(x) = \text{cov}[\mathbf{v}, \phi | x].$$

Conditioning (9.3) on \mathbf{z} would yield the true Gaussian posterior on the inducing values, but computing this posterior would be intractable. In a sparse approximation, we instead *prescribe* a computationally tractable posterior for \mathbf{v} informed by the available observations:

$$p(\mathbf{v} | \mathbf{z}) \approx q(\mathbf{v} | \mathbf{z}) = \mathcal{N}(\mathbf{v}; \tilde{\boldsymbol{\mu}}, \tilde{\Sigma}).$$

This assumed distribution then induces (hence the moniker inducing values) a Gaussian process posterior on f :

$$p(f | \mathbf{z}) \approx \int p(f | \mathbf{v}) q(\mathbf{v} | \mathbf{z}) d\mathbf{v} = \mathcal{GP}(f; \boldsymbol{\mu}_{\mathcal{D}}, K_{\mathcal{D}}),$$

which represents the sparse approximation. After transliteration of notation, the posterior moments take the same form as (2.37), and the cost of computation now scales according to the number of inducing values.

To complete this approximation scheme, we must specify a procedure for identifying a set of inducing values \mathbf{v} as well as the approximate posterior $q(\mathbf{v} | \mathbf{z})$, and it is in these details that the various available methods differ. The inducing values are usually taken to be function values at a set of locations ξ called *pseudo-* or *inducing points*, taking $\mathbf{v} = f(\xi)$.¹⁴ These inducing points are fictitious and do not need to coincide with any actual observations. Once a suitable parameterization

observed values, \mathbf{z} ; $|\mathbf{z}| = n$

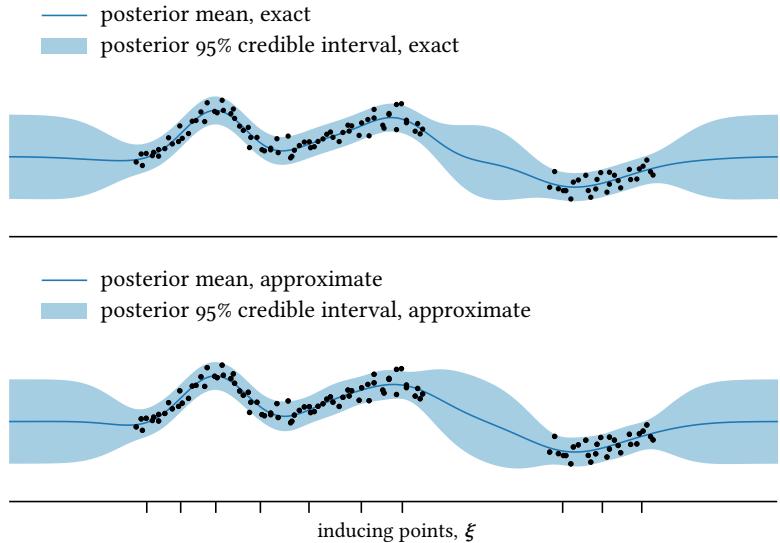
diagonal noise covariance, \mathbf{N}

inducing values, \mathbf{v} ; $|\mathbf{v}| = m \ll n$

¹⁴ This is not strictly necessary. We could consider using other values such as derivatives as inducing values for added flexibility.

pseudopoints, inducing points

Figure 9.1: Sparse approximation. Top: the exact posterior belief (about noisy observations rather than the latent function) for a Gaussian process conditioned on 200 observations. Bottom: a sparse approximation (9.5–9.7) using ten inducing values corresponding to the indicated inducing points, designed to minimize the KL divergence between the induced and true posterior distributions (9.4).



of the inducing values is chosen, we usually design them – as well as their inducing distribution – by optimizing a measure of fit between the true posterior distribution and the resulting approximation.

TITSIAS introduced a variational approach that has gained prominence.¹⁵ The idea is to minimize the Kullback–Leibler divergence between the true and induced posteriors on \mathbf{y} and the inducing values \mathbf{v} :

$$D_{\text{KL}}[q(\mathbf{y}, \mathbf{v} | \mathbf{z}) \| p(\mathbf{y}, \mathbf{v} | \mathbf{z})], \quad (9.4)$$

where

$$q(\mathbf{y}, \mathbf{v} | \mathbf{z}) = p(\mathbf{y} | \mathbf{v}) q(\mathbf{v} | \mathbf{z})$$

is the implied joint distribution in the approximate posterior. Once optimal inducing values are determined, the optimal inducing distribution is as well, and examining the resulting approximate posterior gives insight into the typical behavior of sparse approximations. The approximate posterior mean is

$$\mu_{\mathcal{D}}(x) = \mu(x) + [\mathbf{K}\Sigma^{-1}k(x)]^\top (\mathbf{K}\Sigma^{-1}\mathbf{K}^\top + \mathbf{N})^{-1}(\mathbf{z} - \mathbf{m}). \quad (9.5)$$

Although this expression nominally entails solving a linear system of size n , the low-rank-plus-diagonal structure of the matrix $\mathbf{K}\Sigma^{-1}\mathbf{K}^\top + \mathbf{N}$ allows the system to be solved in time $\mathcal{O}(nm^2)$, merely linear in the size of the dataset.¹⁶ With this favorable cost, sparse approximation can scale Gaussian process inference to millions of observations without issue.

Comparing this expression (9.5) with the true posterior mean:

$$\mu_{\mathcal{D}}(x) = \mu(x) + \kappa(x)^\top (\mathbf{C} + \mathbf{N})^{-1}(\mathbf{z} - \mathbf{m}),$$

we can identify a key approximation to the covariance structure of the GP prior. Namely, we approximate the covariance function with:

$$K(x, x') \approx k(x)^\top \Sigma^{-1} k(x') = \text{cov}[x, \mathbf{v}] \text{ cov}[\mathbf{v}, \mathbf{v}]^{-1} \text{ cov}[\mathbf{v}, x']. \quad (9.6)$$

¹⁵ M. TITSIAS (2009). Variational Learning of Inducing Variables in Sparse Gaussian Processes. *AISTATS 2009*.

¹⁶ For example, we may appeal to the Woodbury identity and recognize that dealing with the diagonal matrix \mathbf{N} is trivial. This is the reason why we restricted the present discussion to diagonal error covariance.

that is, we assume that all covariance between function values is moderated through the inducing values. This is a popular approximation scheme known as the *Nyström method*.¹⁷ Importantly, however, we note that the posterior mean does still reflect the information contained in the entire dataset through the true residuals ($\mathbf{z} - \mathbf{m}$) and noise \mathbf{N} . The approximate posterior covariance also reflects this approximation:

$$\begin{aligned} K(x, x') - [\mathbf{K}\Sigma^{-1}k(x)]^\top (\mathbf{K}\Sigma^{-1}\mathbf{K}^\top + \mathbf{N})^{-1} [\mathbf{K}\Sigma^{-1}k(x')] \\ \approx K(x, x') - \kappa(x)^\top (\mathbf{C} + \mathbf{N})^{-1} \kappa(x'). \end{aligned} \quad (9.7)$$

Figure 9.1 illustrates a sparse approximation for a toy example following this approach. Here the inducing values were taken to be the function values at ten inducing points, and both the inducing points and their distribution were designed to minimize the KL divergence between the true and induced posterior distributions (9.4). The approximation is faithful: the posterior mean is nearly identical to the true mean and the posterior credible intervals only display some minor differences from the truth. Increasing the number of inducing values would naturally improve the approximation.

Sparse approximation for Gaussian processes has a long history, beginning in earnest with investigation into the Nyström approximation (9.6).^{18,19} In addition to the variational approach mentioned above, an approximation known as the *fully independent training conditional* (FITC) approximation has also received significant attention²⁰ and gives rise to a similar approximate posterior. HENSMAN et al. provided a variational sparse approximation for non-Gaussian observation models, allowing for scaling general GP latent models to large datasets.²¹

9.2 OPTIMIZING ACQUISITION FUNCTIONS

In our discussion on computing optimization policies with Gaussian processes, we considered the *pointwise* evaluation of common acquisition functions and their gradient with respect to the proposed observation location. However, we realize an optimization policy via the global *optimization* of an acquisition function:

$$x \in \arg \max_{x' \in \mathcal{X}} \alpha(x'; \mathcal{D}). \quad (9.8)$$

Every common Bayesian acquisition function is nonconvex in general, so we must resort to some generic global optimization routine for this inner optimization. Some care is required to guarantee success in this optimization, as the behavior of a typical acquisition function can make it a somewhat unusual objective function. In particular, consider a prototypical Gaussian process model combining:

- a constant mean function (3.1),
- a stationary covariance function decaying to zero as $|x - x'| \rightarrow \infty$, and
- independent, homoskedastic observation noise (2.16).

¹⁷ C. K. I. WILLIAMS and M. SEEGER (2000). Using the Nyström Method to Speed Up Kernel Machines. *NeurIPS 2000*.

¹⁸ C. K. I. WILLIAMS and M. SEEGER (2000). Using the Nyström Method to Speed Up Kernel Machines. *NeurIPS 2000*.

¹⁹ A. J. SMOLA and B. SCHÖLKOPF (2000). Sparse Greedy Matrix Approximation for Machine Learning. *ICML 2000*.

²⁰ E. Snelson and Z. Ghahramani (2005). Sparse Gaussian Processes using Pseudo-inputs. *NeurIPS 2005*.

²¹ J. HENSMAN et al. (2015). MCMC for Variationally Sparse Gaussian Processes. *NeurIPS 2015*.

computing common policies with Gaussian processes: chapter 8, p. 157

stationarity: § 3.2, p. 50

For such a model, the prior predictive distribution $p(y | x)$ is identical regardless of location. However, the *posterior* predictive distribution $p(y | x, \mathcal{D})$ also degenerates to the prior for locations sufficiently far from observed locations, due to the decay of the covariance function. In these regions, the gradients of the posterior predictive parameters effectively vanish:

$$\frac{\partial \mu}{\partial x} \approx 0; \quad \frac{\partial s}{\partial x} \approx 0.$$

As a result, the gradient of the acquisition function (8.6) vanishes as well! This vanishing gradient is especially problematic in high-dimensional spaces, where the acquisition function will be flat on an overwhelming fraction of the domain unless the prior encodes absurdly long-scale correlations, a consequence of the unshakable *curse of dimensionality*. Thus, the acquisition function will only exhibit interesting behavior in the neighborhood of previous observations – where the posterior predictive distribution is nontrivial – and it is here we should spend most of our effort during optimization.

curse of dimensionality: § 3.5, p. 61

Optimization approaches

There are two common lines of attack for optimizing acquisition functions in Bayesian optimization. One approach is to use an off-the-shelf derivative-free global optimization method such as the “dividing rectangles” (DIRECT) algorithm of JONES et al.²² or a member of the covariance matrix adaptation evolution strategy (CMA-ES) family of algorithms.²³ Although popular, we argue that neither is a particularly good choice in situations where the acquisition function may devolve into effective flatness as described above. However, an algorithm of this class may be reasonable in modest dimension.

An alternative is multistart *local* optimization, making use of the gradients computed in the previous chapter for rapid convergence. To ensure success, we must carefully select starting points to ensure that the relevant regions of the domain are searched. JONES (the same JONES of the DIRECT algorithm) recognized the problem of vanishing gradients described above and suggested a simple heuristic for selecting local optimization starting points by enumerating and pruning the midpoints between all pairs of observed points.²⁴ A more brute-force approach is to measure the acquisition function on a exhaustive covering of the domain – generated for example by a low-discrepancy sequence – then begin local searches from the highest values seen. This can be effective if the initial set of points is dense enough to probe the neighborhoods of previous observations; otherwise, it would be prudent to augment with a locally motivated approach such as JONES’s.

Multistart local optimization has the advantage of being embarrassingly parallel. Further, both global and multistart local optimization of the acquisition function can be treated as *anytime* algorithms that constantly improve their proposed observations until we are ready to act.²⁵

²² D. R. JONES et al. (1993). Lipschitzian Optimization Without the Lipschitz Constant. *Journal of Optimization Theory and Application* 79(1): 157–181.

²³ N. HANSEN (2016). The cma Evolution Strategy: A Tutorial. arXiv: 1604.00772 [cs.LG].

²⁴ D. R. JONES (2001). A Taxonomy of Global Optimization Methods Based on Response Surfaces. *Journal of Global Optimization* 21(4):345–383.

²⁵ E. BROCHU et al. (2010). A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning. arXiv: 1012.2599 [cs.LG].

Optimization in latent spaces

A common approach for modeling on high-dimensional domains is to apply some mapping from the domain to some lower-dimensional representation space, then construct a Gaussian process on that space. With such a model, it is tempting to optimize an acquisition function on the latent space rather than on the original domain so as to at least partially sidestep the curse of dimensionality. This can be an effective approach when we can faithfully “decode” from the latent space back into the original domain for evaluation, a process that can require careful thought even when the latent embedding is *linear* (3.30) due to the nonbijective nature of the map.²⁶

Fortunately, modern neural embedding techniques such as (variational) autoencoders provide a decoding mechanism as a natural side effect of their construction. When we are fortunate enough to have sufficient unlabeled data to learn a useful unsupervised representation prior to optimization, we may simply optimize in the latent space and feed each chosen observation location through the decoder. GÓMEZ-BOMBARELLI et al. for example applied Bayesian optimization to *de novo* molecular design.²⁷ Their model combined a pretrained variational autoencoder for molecular structures with a Gaussian process on the latent embedding space; the autoencoder was trained on a large precompiled database of known molecules.²⁸ A welcome side effect of this construction was that, by optimizing the acquisition function over the continuous embedding space, the decoding process had the freedom to generate novel structures not seen in the autoencoder’s training data.

We may also pursue this approach even when we begin optimization without any data at all. For example, MORICONI et al. demonstrated success jointly learning a nonlinear low-dimensional representation as well as a corresponding decoding stage throughout optimization on the fly.²⁹

Optimization on combinatorial domains

Combinatorial domains present a challenge for Bayesian optimization as the optimization policy (9.8) requires combinatorial optimization of the acquisition function. It is difficult to provide concrete advice for such a situation, as the details of the domain may in some cases suggest a natural path forward. That said, we can provide some ideas. One potential solution is outlined above: when the domain is a space of discrete structured objects such as graphs (say, molecules), we might find a useful *continuous* embedding of the domain and simply work there instead.²⁷

When this is not possible, we note that our previous sketch of how a “typical” acquisition function behaves with a “typical” model – with its most interesting behavior near the observed data – carries over to combinatorial domains and may lead to useful heuristics. For example, rather than enumerating the domain (presumably impossible) or sam-

modeling functions on high-dimensional domains: § 3.5, p. 61

²⁶ Detailed advice for this setting is provided by:

Z. WANG et al. (2016b). Bayesian Optimization in a Billion Dimensions via Random Embeddings. *Journal of Artificial Intelligence Research* 55:361–387.

linear embeddings: § 3.5, p. 62

neural embeddings: § 3.5, p. 61

de novo molecular design: p. 310

²⁷ R. GÓMEZ-BOMBARELLI et al. (2018). Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules. *ACS Central Science* 4(2):268–276.

²⁸ T. STERLING and J. J. IRWIN (2015). ZINC 15 – Ligand Discovery for Everyone. *Journal of Chemical Information and Modeling* 55(11):2324–2337.

²⁹ R. MORICONI et al. (2020). High-dimensional Bayesian optimization using low-dimensional feature spaces. *Machine Learning* 109(9–10):1925–1943.

IMPLEMENTATION

- ³⁰ As evaluated by the covariance function.
- ³¹ R. GARNETT et al. (2010). Bayesian Optimization for Sensor Set Selection. *IPSN 2010*.
- ³² R. BAPTISTA and M. POLOCZEK (2018). Bayesian Optimization of Combinatorial Structures. *ICML 2018*.
- ³³ C. OH et al. (2019). Combinatorial Bayesian Optimization using the Graph Cartesian Product. *NeurIPS 2019*.
- ³⁴ J. KIM et al. (2021). Bayesian optimization with approximate set kernels. *Machine Learning* 110(5):857–879.

pling from the domain (presumably yielding poor coverage), we might instead curate a small list of candidate points offering options for both exploitation and exploration. We could perhaps generate a list of points similar to³⁰ the thus-far best-seen points, encouraging exploitation, and augment with a small set of points constructed to cover the domain, encouraging exploration. This approach was used for example by GARNETT et al. for Bayesian set function optimization.³¹

Other approaches are of course possible. For example, several authors have constructed GP models for particular families of combinatorial spaces whose structure simplifies the (possibly approximate but near-optimal) optimization of particular acquisition functions induced from the model.^{32,33,34}

The benefit of imperfect optimization?

We conclude this discussion with one paradoxical remark. In some cases it may actually be advantageous to *avoid* perfectly optimizing an acquisition function. As many common policies are based on extremely myopic (one-step) reasoning, imperfect optimization may yield an exploration benefit not encoded in the acquisition function itself. It might be challenging to study this effect formally and offer practical advice, but this phenomenon – and the effect it may have on empirical results – deserves thoughtful consideration.

9.3 STARTING AND STOPPING OPTIMIZATION

Finally, we briefly consider the part of optimization that happens *outside* the application of an optimization policy: initialization and termination.

Initialization

Theoretically, one can begin a Bayesian optimization routine with a completely empty dataset $\mathcal{D} = \emptyset$ and then use an optimization policy to design every observation, and indeed this has been our working model of Bayesian optimization since sketching the basic idea in algorithm 1.1. However, Bayesian optimization policies are informed by an underlying belief about the objective function, which can be significantly misinformed when too little data is available, especially when relying on point estimation for model selection rather than accounting for (significant!) uncertainty in model hyperparameters and/or model structures.

Due to the sequential nature of optimization and the dependence of each decision on the data observed in previous iterations, it can be wise to use a model-*independent* procedure to design a small number of initial observations before beginning optimization in earnest. This procedure can be as simple as random sampling³⁵ or a space-filling design such as a low-discrepancy sequence or Latin hypercube design.³⁶ When repeatedly solving related optimization problems, we may even be able to *learn* how to initialize Bayesian optimization routines from experience. Some au-

model averaging: § 4.4, p. 74

- ³⁵ These approaches are discussed and evaluated in:

M. W. HOFFMAN and B. SHAHRIARI (2014). Modular mechanisms for Bayesian optimization. *Bayesian Optimization Workshop, NeurIPS 2014*.

- ³⁶ D. R. JONES et al. (1998). Efficient Global Optimization of Expensive Black-Box Functions. *Journal of Global Optimization* 13(4):455–492.

thors have proposed sophisticated “warm start” initialization procedures for hyperparameter tuning using so-called *metafeatures* characterizing the datasets under consideration.³⁷

Termination

In many applications of Bayesian optimization, we assume a preallocated budget on the number of observations we will make and simply terminate optimization when that budget is expended. However, we may also treat termination as a *decision* and adaptively determine when to stop based on collected data. Of course, in practice we are free to design a stopping rule however we see fit, but we can outline some possible options.

Especially when using a policy grounded in decision theory, it is natural to terminate optimization when the maximum of our chosen acquisition function drops below some threshold c , which may depend on x :

$$\max_{x \in \mathcal{X}} [\alpha(x; \mathcal{D}) - c(x)] < 0. \quad (9.9)$$

For acquisition functions derived from decision theory, such a stopping rule may be justified theoretically: we stop when the expected gain from the optimal observation is no longer worth the cost of acquisition.³⁸ A majority of the stopping rules described in the literature assume this form, with the threshold c often being determined dynamically based on the scale of observed data.³⁹ DAI et al. combined a stopping rule of this form with an otherwise non-decision-theoretic policy (GP-UCB) and showed that its asymptotic performance in terms of expected regret was not adversely affected despite the mismatch in motivation between the policy and stopping rule.⁴⁰

It may also be prudent to consider purely data-dependent stopping rules in order to avoid undue expense arising from miscalibrated models fruitlessly continuing optimization based on incorrect beliefs. For example, ACERBI and MA proposed augmenting a bound on the total number of observations with an early stopping option if no optimization progress is made over a given number of observations.⁴¹

SUMMARY OF MAJOR IDEAS

- Gaussian process inference requires solving a system of linear equations whose size grows with the number of observed values (9.2).
- A direct implementation via the Cholesky decomposition is a straightforward option, but scales cubically with the number of observed values.
- Scaling to larger datasets is possible by appealing to iterative numerical methods such as the method of *conjugate gradients*, or to *sparse approximation*, where we approximate an intractable posterior Gaussian process with a Gaussian process conditioned on carefully designed, fictitious observations at locations called *inducing points*. These methods can scale inference to hundreds of thousands of observations or more.

³⁷ M. FEURER et al. (2015). Initializing Bayesian Hyperparameter Optimization via Meta-Learning. *AAAI 2015*.

optimal stopping rules: § 5.4, p. 103

³⁸ See in particular our discussion of the one-step optimal stopping rule, p. 104.

³⁹ Some early (but surely not the earliest!) examples:

D. D. COX and S. JOHN (1992). A Statistical Method for Global Optimization. *SMC 1992*.

D. R. JONES et al. (1998). Efficient Global Optimization of Expensive Black-Box Functions. *Journal of Global Optimization* 13(4):455–492.

⁴⁰ Z. DAI et al. (2019). Bayesian Optimization Meets Bayesian Optimal Stopping. *ICML 2019*.

⁴¹ L. ACERBI and W. J. MA (2017). Practical Bayesian Optimization for Model Fitting with Bayesian Adaptive Direct Search. *NeurIPS 2017*.

direct computation via Cholesky decomposition: § 9.1, p. 201

iterative numerical methods: § 9.1, p. 203
sparse approximation: § 9.1, p. 204

IMPLEMENTATION

the potential for vanishing gradients: § 9.2,
p. 207

initialization: § 9.3, p. 210

termination: § 9.3, p. 211

- When optimizing acquisition functions, it is important to be aware of the potential for vanishing gradients in extrapolatory regions of the domain and plan accordingly.
- It is usually a good idea to begin optimization with a small set of observations designed in a model-agnostic fashion in order to begin with a somewhat informed model of the objective function.
- When dynamic termination is desired, simple schemes based on thresholding the acquisition function can be effective.

10

annotated bibliography of applications: p. 309

THEORETICAL ANALYSIS

The Bayesian optimization procedures we have developed throughout the course of this book have demonstrated excellent empirical performance in a huge array of practical settings. However, good empirical performance may not be enough to satisfy those who value rigor over results. Fortunately, many Bayesian optimization procedures are also backed by strong theoretical guarantees on their performance. The literature on this topic is now quite vast, and convergence has been studied by different authors in different ways, sometimes involving slight nuances interpretation. We will take an in-depth look at this topic in this chapter, covering the most common lines of attack and outlining the state-of-the-art in results.

A running theme throughout this chapter will be understanding how various measures of optimization error decrease asymptotically as an optimization policy is repeatedly executed. To facilitate this discussion, we will universally use τ in this chapter to indicate dataset size, or equivalently to indicate the number of steps an optimization procedure is assumed to have run. As we will primarily be interested in asymptotic results as $\tau \rightarrow \infty$, this convention does not rule out the possibility of starting optimization with some arbitrary dataset whose size does not depend on τ . We will use subscripts to indicate dataset size when necessary, notating the dataset comprising the first τ observations with:

$$\mathcal{D}_\tau = (\mathbf{x}_\tau, \mathbf{y}_\tau) = \{(x_i, y_i)\}_{i=1}^\tau.$$

When studying the convergence of a global optimization algorithm, we must be careful to define exactly what we *mean* by “convergence.” In general, deriving a convergence result entails:

- choosing some measure of optimization error,
- choosing some space of possible objective functions, and
- establishing some guarantee for the chosen error on the chosen function space, such as an asymptotic bound on the worst- or average-case error in the large-sample limit $\tau \rightarrow \infty$.

There is a great deal of freedom in the last of these steps, and we will discuss several important results and proof strategies later in this chapter. However, there are well-established conventions for the first two of these steps, which we will introduce in the following two sections. We begin with the notion of *regret*, which provides a natural measure of optimization error.

size of dataset, τ

dataset after τ observations, \mathcal{D}_τ

what does it mean to converge?

regret: below

useful spaces of objective functions: § 10.2,
p. 215

10.1 REGRET

Regret is a core concept in the analysis of optimization algorithms, Bayesian or otherwise. The role of regret is to quantify optimization progress in a manner suitable for establishing convergence to the global optimum and studying the rate of this convergence. There are several definitions

THEORETICAL ANALYSIS

- 1 One might propose an alternative definition of error by measuring how closely the observed *locations* approach a global maximum x^* rather than by how closely the observed *values* approach the value of the global optimum f^* . However, it turns out this is both less convenient for analysis and harder to motivate: in practice it is the value of the objective that we care about the most, and discovering a near-optimal value is a success regardless of its distance to the global optimum.
- 2 As outlined in chapter 5 (p. 87), optimal actions maximize *expected* utility in the face of uncertainty. Many observations may not result in progress, even though designed with the best of intentions.
- 3 For example, it is easy to develop a space-filling design that will eventually locate the global optimum of any continuous function through sheer dumb luck – but it won't do so very quickly!

of regret used in different contexts, all based on the same idea: comparing the objective function values visited during optimization to the globally optimal value, f^* . The larger this gap, the more “regret” we incur in retrospect for having invested in observations at suboptimal locations.¹

Regret is an unavoidable consequence of decision making under uncertainty – without foreknowledge of the global optimum, we must of course spend some time searching for it. As a result, even what may be *optimal* actions in the face of uncertainty may seem disappointing in retrospect.² However, such actions are necessary in order to learn about the environment and inform future decisions. This reasoning gives rise to the classic tension between exploration and exploitation in policy design: although exploration may not yield immediate progress, it enables future success, and if we are careful, reduces future regret. However, exploration alone is not sufficient to realize a compelling optimization strategy,³ as we must also exploit what we have learned and adapt our behavior accordingly. An ideal algorithm thus explores *efficiently* enough that its regret can at least be limited, and establishing bounds on this regret is a fundamental goal in theoretical analysis.

Most analysis is performed in terms of either *simple* or *cumulative* regret, defined below. These notions are closely related, and which is best in a given situation usually comes down to what is most convenient.

Simple regret

Let \mathcal{D}_τ represent some set of (potentially noisy) observations gathered during optimization,⁴ and let $\phi_\tau = f(\mathbf{x}_\tau)$ represent the objective function values at the observed locations. The simple regret associated with this data is the difference between the global maximum of the objective and the maximum restricted to the observed locations:⁵

$$r_\tau = f^* - \max_{x \in \mathcal{X}} \mu_{\mathcal{D}_\tau}(x). \quad (10.1)$$

- 4 The notions of regret introduced here do not depend on the observed values y but only on the underlying objective values ϕ . We are making a tacit assumption that y is sufficiently informative about ϕ for this to be sensible.
- 5 Occasionally a slightly different definition of simple regret is used, analogous to the global reward (6.5). Here we replace $\max \phi$ with the maximum of the posterior mean given the data:

$$r_\tau = f^* - \max_{x \in \mathcal{X}} \mu_{\mathcal{D}_\tau}(x).$$

This distinction is rarely relevant.

convergence goal: show $r_\tau = o(1)$

instantaneous regret, ρ

cumulative regret of \mathcal{D}_τ , R_τ

It is immediate from its definition that simple regret is nonnegative and vanishes only if the data contain a global optimum. With this in mind, a common goal is to show that the simple regret of data obtained by some policy approaches zero, implying the policy will eventually (and perhaps efficiently) identify the global optimum, up to vanishing error.

Cumulative regret

To define cumulative regret, we first introduce the *instantaneous regret* ρ corresponding to an observation at some point x , which is the difference between the global maximum of the objective and the function value ϕ :

$$\rho = f^* - \phi. \quad (10.2)$$

The *cumulative regret* for a dataset \mathcal{D}_τ is then the total instantaneous regret incurred:

$$R_\tau = \sum_i \rho_i = \tau f^* - \sum_i \phi_i.$$

Relationship between simple and cumulative regret

Simple and cumulative regret are analogous to the simple (6.3) and cumulative (6.7) reward utility functions, and any intuition regarding these utilities transfers to their regret counterparts.⁶

However, these two definitions of regret are not directly comparable, even on the same data – for starters, simple regret is nonincreasing as more data is collected, whereas cumulative regret is nondecreasing. However, suitable normalization allows some useful comparison.⁷ Namely, consider the *average*, rather than cumulative, regret:

$$\frac{R_\tau}{\tau} = f^* - \frac{1}{\tau} \sum_i \phi_i.$$

As the mean of a vector is a lower bound on its maximum, we may derive an upper bound on simple regret in terms of cumulative regret: (10.1):

$$r_\tau \leq \frac{R_\tau}{\tau}. \quad (10.3)$$

In this light, a common goal is to show that an optimization algorithm has the so-called *no-regret property*, which means that its average regret vanishes with increasing data:

$$\lim_{\tau \rightarrow \infty} \frac{R_\tau}{\tau} = 0. \quad (10.4)$$

Equivalently, a policy achieves no regret if its cumulative regret grows sublinearly with the dataset size τ . This is sufficient to prove convergence in terms of simple regret as well by appealing to the squeeze theorem:

$$0 \leq r_\tau \leq \frac{R_\tau}{\tau} \rightarrow 0.$$

Although no regret guarantees convergence for simple regret, the reverse is not necessarily the case. It is easy to find counterexamples – consider for example modifying a no-regret policy to select a *fixed* suboptimal point every other observation. The simple regret would still vanish, but constant instantaneous regret on alternating iterations would prevent sublinear cumulative regret. Thus the no-regret property is somewhat stronger than convergence in terms of the simple regret alone. From another perspective, simple regret is more tolerant of exploration: we need only visit the global optimum *once* to converge in terms of simple regret, whereas effectively *all* observations must eventually become effectively optimal to achieve the no-regret property.

⁶ Note that neither notion of regret represents a valid utility function itself, as in general both f^* and ϕ would be random variables in that context.

⁷ For a deeper discussion on the connections between simple and cumulative regret, see:

S. BUBECK et al. (2009). Pure Exploration in Multi-armed Bandits Problems. *ALT 2009*.

convergence goal, no-regret property:
 $R(\mathcal{D}_\tau) = o(\tau)$

no-regret property implies convergence for simple regret

convergence in simple regret does not imply the no-regret property

10.2 USEFUL FUNCTION SPACES FOR STUDYING CONVERGENCE

Identifying the “right” function space to consider when studying convergence is a subtle decision, as we must strike a balance between generality and practicality. One obvious choice would be the space of *all* continuous

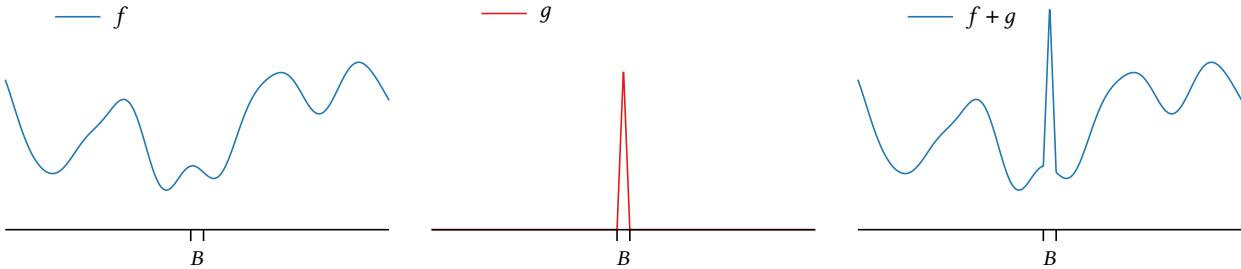


Figure 10.1: Modifying a continuous function f (left) to feature a “needle” on some ball B . We construct a continuous function g vanishing on the compliment of B (middle), then add this “correction” to f (right).

functions. However, it turns out this space is too large to be of much theoretical interest, as it contains adversarial functions that are arbitrarily hard to optimize. Nonetheless, it is easy to characterize convergence (in terms of the simple regret) on this space, and results of this type have been established for select Bayesian optimization algorithms.

We may gain some traction by considering a family of more plausible, “nice” functions whose complexity can be controlled enough to guarantee rapid convergence. In particular, choosing a Gaussian process prior for the objective function implies strong correlation structure, and this insight leads to natural function spaces to study. This has become the standard approach in modern analysis, and we will consider it shortly.

Convergence on all continuous functions

The largest reasonable space we might want to consider is the space of all continuous functions. However, continuous functions can be poorly behaved from the point of view of optimization, and we cannot hope for strong convergence guarantees as a result. There is simply too much freedom for functions to “hide” their optima in inconvenient places.

To begin, we establish a simple characterization of convergence on all continuous functions in terms of eventual density of observation.

Theorem. *Let \mathcal{X} be a compact metric space. An optimization policy converges in terms of simple regret on all continuous functions $f: \mathcal{X} \rightarrow \mathbb{R}$ if and only if the set of eventually observed points $\mathbf{x} = \bigcup_{i=1}^{\infty} \{x_i\}$ is always dense in \mathcal{X} .*

The proof is instructive as it shows what can go wrong with general continuous functions. First, if the set of eventually observed points is dense in \mathcal{X} , we may construct a sequence of observations converging to a global maximum x^* .⁸ The associated function values $\{\phi_i\}$ then converge to f^* by continuity, and thus $r_\tau \rightarrow 0$.

If density fails⁹ for some continuous function f – and thus there is some ball $B \subset \mathcal{X}$ that will never contain an observation – then we may foil the policy with a “needle in a haystack.” We construct a continuous

⁸ For example by taking x_i within a ball of radius $1/i$ around x^* .

⁹ Some care would be required to make the following argument rigorous for stochastic observations and/or policies, but the spirit would remain intact.

function g vanishing on the complement of B and achieving arbitrarily high values on B .¹⁰ Adding the needle to f creates a continuous function $f + g$ with arbitrary – and, once the needle is sufficiently tall, never observed – maximum. As f and $f + g$ agree outside B , the policy cannot distinguish between these two functions, so we can find continuous functions for which the policy has arbitrarily high simple regret.

We can use this strategy of building adversarial “needles in haystacks” to further show that, even if we can guarantee convergence on all continuous functions, we cannot hope to demonstrate *rapid* convergence in simple regret, at least not in the worst case. Unless the domain is finite, running any policy for any finite number of iterations will leave unobserved “holes” that we can fill in with arbitrarily tall “needles,” and thus the worst-case regret will be unbounded at every stage of the algorithm.

Convergence results for all continuous functions on the unit interval

Establishing the density criterion above has proven difficult for Bayesian optimization algorithms in general spaces with arbitrary models. However, restricting the domain to the unit interval $\mathcal{X} = [0, 1]$ and limiting the objective function and observation models to certain well-behaved combinations has yielded universal convergence guarantees for some Bayesian procedures.

For example, KUSHNER sketched a proof of convergence in simple regret when maximizing probability of improvement on the unit interval, when the objective function model is the Wiener process and observations are exact or corrupted by additive Gaussian noise.¹¹ ŽILINSKAS later provided a proof in the noiseless case.¹² This particular model exhibits a Markov property that enables relatively straightforward analysis by characterizing its behavior on each of the subintervals subdivided by the data. This structure enables a simple proof strategy by assuming some subinterval is never subdivided and arriving at a contradiction.

Convergence guarantees for this policy (under the name the “P-algorithm”) on the unit interval have also been established for smoother models of the objective function – that is, with differentiable rather than merely continuous sample paths¹³ – including the once-integrated Wiener process.¹⁴ Convergence rates for these algorithms for general continuous functions have also been derived.^{13, 14, 15} In light of the above discussion, these convergence rates are not in terms of simple regret but rather in terms of shrinkage in the size of the subinterval containing the global optimum, and thus the distance from the closest observed location to the global optimum. The proofs of these results all relied on exact observation of the objective function.

Convergence on all continuous functions on the unit interval has also been established for maximizing expected improvement, in the special case of exact observation and a Wiener process model on the objective.¹⁶ The proof strategy again relied on the special structure of the posterior to show that the observed locations will always subdivide the domain such that no “hole” is left behind.

¹⁰ For example, we may scale the distance to the complement of B .

worst-case simple regret unbounded

¹¹ H. J. KUSHNER (1964). A New Method of Locating the Maximum Point of an Arbitrary Multi-peak Curve in the Presence of Noise. *Journal of Basic Engineering* 86(1):97–106.

¹² A. G. ŽILINSKAS (1975). Single-Step Bayesian Search Method for an Extremum of Functions of a Single Variable. *Kibernetika (Cybernetics)* 11(1):160–166.

¹³ J. CALVIN and A. ŽILINSKAS (1999). On the Convergence of the P-Algorithm for One-Dimensional Global Optimization of Smooth Functions. *Journal of Optimization Theory and Applications* 102(3):479–495.

¹⁴ J. M. CALVIN and A. ŽILINSKAS (2001). On Convergence of a P-Algorithm Based on a Statistical Model of Continuously Differentiable Functions. *Journal of Global Optimization* 19(3):229–245.

¹⁵ J. M. CALVIN (2000). Convergence Rate of the P-Algorithm for Optimization of Continuous Functions. In: *Approximation and Complexity in Numerical Optimization: Continuous and Discrete Problems*.

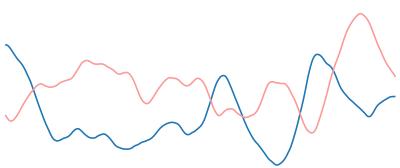
¹⁶ M. LOCATELLI (1997). Bayesian Algorithms for One-Dimensional Global Optimization. *Journal of Global Optimization* 10(1):57–76.

Convergence on “nice” continuous functions

In the pursuit of stronger convergence results, we may abandon the space of all continuous functions and focus on some pace of suitably well-behaved objectives. By limiting the complexity of the functions we consider, we can avoid complications arising from adversarial examples like “needles in haystacks.” This can be motivated from the basic assumption that optimization is to be feasible at all, which is not the case when facing an adversary creating arbitrarily difficult problems. Instead, we may seek strong performance guarantees when optimizing “plausible” objective functions. Conveniently, a Gaussian process model on the objective $\mathcal{GP}(f; \mu, K)$ gives rise to several paths forward.

Sample paths of a Gaussian process

sample path continuity: § 2.5, p. 28



Sample paths of a stationary Gaussian process with Matérn covariance, $\nu = 5/2$ (3.14) show more regularity than arbitrary continuous functions.

Bayesian (Bayes) regret

In a Bayesian analysis, it is natural to assume that the objective function is a sample path from the Gaussian process used to model it. Assuming sample path continuity, sample paths of a Gaussian process are much better behaved than general continuous functions. The covariance function provides regularization on the behavior of sample paths via the induced correlations among function values; see the figure in the margin. As a result, we can ensure that functions with exceptionally bad behavior (such as hidden “needles”) are also exceptionally rare.

The sample path assumption provides a known distribution for the function values at observed locations (2.2–2.3), allowing us to derive expected-case results. An important concept here is the *Bayesian* (or simply *Bayes*) *regret* of a policy, which is the expected value of the (simple or cumulative) regret incurred when following the policy, say for τ steps:

$$\mathbb{E}[r_\tau]; \quad \mathbb{E}[R_\tau]. \quad (10.5)$$

This expectation is taken with respect to uncertainty in the objective function f , the observation locations \mathbf{x} (in the case of a stochastic policy), and the observed values \mathbf{y} .

The worst-case alternative

The GP sample path assumption is not always desirable, for example in the context of a frequentist (that is, worst-case) analysis of a Bayesian optimization algorithm. This is not as contradictory as it may seem, as Bayesian analyses can lack robustness to model misspecification – a certainty in practice. An alternative is to assume that the objective function lies in some explicit space of “nice” functions \mathcal{H} , then find worst-case convergence guarantees for the Bayesian algorithm on inputs satisfying this regularity assumption. For example, we might seek to bound the worst-case expected (simple or cumulative) regret for a function in this space after τ decisions:

$$\bar{r}_\tau[\mathcal{H}] = \sup_{f \in \mathcal{H}} \mathbb{E}[r_\tau]; \quad \bar{R}_\tau[\mathcal{H}] = \sup_{f \in \mathcal{H}} \mathbb{E}[R_\tau]. \quad (10.6)$$

worst-case regret on \mathcal{H} after τ steps: $\bar{r}(\tau, \mathcal{H})$,
 $\bar{R}(\tau, \mathcal{H})$

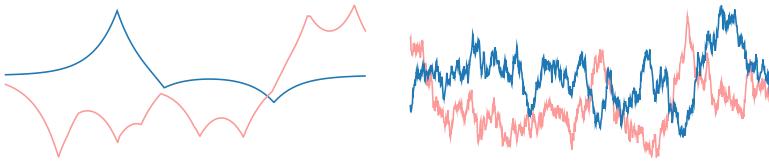


Figure 10.2: Left: functions in the RKHS corresponding to the Matérn covariance function with $v = 1/2$ (3.11). Right: sample paths from a GP with the same covariance.

Here the expectation is over any uncertainty in the observed locations \mathbf{x} and the observed values \mathbf{y} ,¹⁷ but uncertainty in the objective function – presumably the most troubling factor in a Bayesian analysis – is replaced by a pessimistic bound on the functions in \mathcal{H} . In such a result, the policy may still refer to a model $p(f)$ in making decisions, but we are not rewarded in terms of this belief.

Reproducing kernel Hilbert spaces

Corresponding to every covariance function K is a natural companion function space, its *reproducing kernel Hilbert space* (RKHS) \mathcal{H}_K . There is a strong connection between this space and a centered Gaussian process with the same covariance function, $\mathcal{GP}(f; \mu \equiv 0, K)$. Namely, consider the set of functions of the form

$$x \mapsto \sum_{i=1}^n \alpha_i K(x_i, x), \quad (10.7)$$

where $\{x_i\} \subset \mathcal{X}$ is an arbitrary finite set of input locations with corresponding real-valued weights $\{\alpha_i\}$.¹⁸ Note this is *precisely* the set of all possible posterior mean functions for the Gaussian process arising from exact inference.¹⁹ The RKHS \mathcal{H}_K is then the completion of this space endowed with the inner product

$$\left\langle \sum_{i=1}^n \alpha_i K(x_i, x), \sum_{j=1}^m \beta_j K(x'_j, x) \right\rangle = \sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j K(x_i, x'_j). \quad (10.8)$$

That is, the RKHS is roughly the set of functions “as smooth as” a posterior mean function of the corresponding GP, according to a notion of “explainability” by the covariance function.

It turns out that belonging to the RKHS \mathcal{H}_K is a stronger regularity assumption than being a sample path of the corresponding GP. In fact, unless the RKHS is finite-dimensional (which is not normally the case), sample paths from a Gaussian process almost surely do *not* lie in the corresponding RKHS:^{20,21} $\Pr(f \in \mathcal{H}_K) = 0$. However, the posterior mean function of the same process *always* lies in the RKHS by the above construction. Figure 10.2 illustrates a striking example of this phenomenon: sample paths from a stationary GP with Matérn covariance function with $v = 1/2$ (3.11) are *nowhere* differentiable, whereas members of the corresponding RKHS are *almost everywhere* differentiable. Effectively, the process of averaging over sample paths “smooths out” their erratic behavior in the posterior mean, and elements of the RKHS exhibit similar smoothness.

¹⁷ In the special case of exact observation and a deterministic policy, such a bound would entail no probabilistic elements.

reproducing kernel Hilbert space corresponding to K, \mathcal{H}_K

¹⁸ Equivalently, this is the span of the set of covariance functions with one input held fixed: $\{x \mapsto K(x, x') \mid x' \in \mathcal{X}\}$.

¹⁹ Inspection of the general posterior mean functions in (2.14, 2.19) reveals they can always be (and can *only* be) written in this form.

²⁰ M. N. LUKIĆ and J. H. BEDER (2001). Stochastic Processes with Sample Paths in Reproducing Kernel Hilbert Spaces. *Transactions of the American Mathematical Society* 353(10):3945–3969.

²¹ However, sample paths *do* lie in a “slightly larger” RKHS we can determine from the covariance function:

M. KANAGAWA et al. (2018). Gaussian Processes and Kernel Methods: A Review on Connections and Equivalences. arXiv: 1807 . 02582 [stat.ML] [theorem 4.12]

Reproducing kernel Hilbert space norm

RKHS norm, $\|f\|_{\mathcal{H}_K}$

Associated with a RKHS \mathcal{H}_K is a norm $\|f\|_{\mathcal{H}_K}$ that can be interpreted as a measure of function complexity with respect to the covariance function K . The RKHS norm derives from the pre-completion inner product (10.8), and we can build intuition for the norm by drawing a connection between that inner product and familiar concepts from Gaussian process regression. The key is the characterization of the pre-completion function space (10.7) as the space of all possible posterior mean functions for the corresponding centered Gaussian process $\mathcal{GP}(f; \mu \equiv 0, K)$.

RKHS norm of posterior mean function

To be more explicit, consider the posterior mean after observing a dataset of *exact* observations $\mathcal{D} = (\mathbf{x}, \mathbf{y})$, inducing the posterior mean (2.14):

$$\mu_{\mathcal{D}}(x) = K(x, \mathbf{x})\Sigma^{-1}\mathbf{y},$$

where $\Sigma = K(\mathbf{x}, \mathbf{x})$. Then the (squared) RKHS norm of the posterior mean is (10.8):²²

$$\|\mu_{\mathcal{D}}\|_{\mathcal{H}_K}^2 = \langle \mu_{\mathcal{D}}, \mu_{\mathcal{D}} \rangle = \mathbf{y}^\top \Sigma^{-1} \mathbf{y}. \quad (10.9)$$

That is, the RKHS norm of the posterior mean is the Mahalanobis norm of the observed data \mathbf{y} under their Gaussian prior distribution (2.2–2.3).

We have actually seen this score before: it appears in the log marginal likelihood of the data under the Gaussian process (4.8), where we interpreted it as a score of data fit.²³ This reveals a deep connection between the RKHS norm and the associated Gaussian process: posterior mean functions arising from “more unusual” observations from the point of view of the GP have higher complexity from the point of view of the RKHS. Whereas the Gaussian process judges the *data* \mathbf{y} to be complex via the marginal likelihood, the RKHS judges the resulting *posterior mean* $\mu_{\mathcal{D}}$ to be complex via the RKHS norm – but these are simply two ways of interpreting the same scenario.²⁴

The role of the RKHS norm in quantifying function complexity suggests a natural space of objective functions to work with when seeking worst-case results (10.6). We take the *RKHS ball* of radius B , the space of functions with complexity bounded by B in the RKHS:

$$\mathcal{H}_K[B] = \{f \mid \|f\|_{\mathcal{H}_K} \leq B\}. \quad (10.10)$$

The radius B is left as a parameter that is absorbed into derived bounds.

10.3 RELEVANT PROPERTIES OF COVARIANCE FUNCTIONS

A sizable majority of the results discussed in the remainder of this chapter concern optimization performance on one of the function spaces discussed in the previous section: sample paths of a centered Gaussian process with covariance function K or functions in the corresponding RKHS \mathcal{H}_K . Given the fundamental role the covariance function plays in determining sample path behavior, it should not be surprising that the nature of the covariance function also has profound influence on optimization performance.

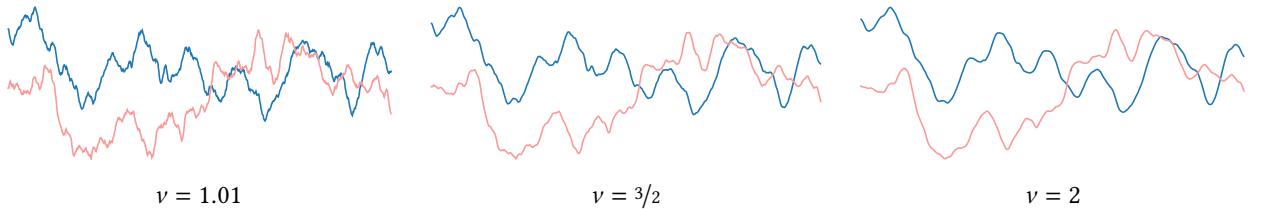


Figure 10.3: Sample paths from a centered GP with Matérn covariance with smoothness parameter v ranging from just over 1 (left) to 2 (right). The random seed is shared so that paths of the same color are comparable. All samples are once differentiable, but some are smoother than others.

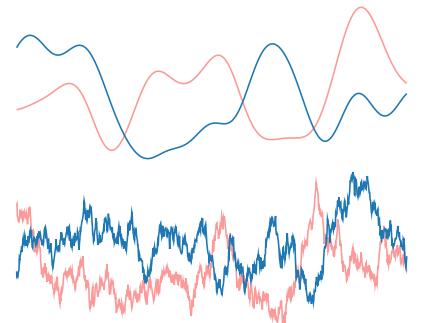
One might intuitively expect that optimizing smoother functions should be easier than optimizing rougher functions, as a rougher function gives the optimum more places to “hide;” see the figure in the margin. This intuition turns out to be correct. The key insight is that rougher functions require more *information* to describe than smoother ones. As each observation we make is limited in how much information it can reveal regarding the objective function, rougher objectives require more observations to learn with a similar level of confidence, and to optimize with a similar level of success. We can make this intuition precise through the concept of *information capacity*, which bounds the rate at which we can learn about the objective through noisy observations and serves as a fundamental measure of function complexity appearing in numerous analyses.

Smoothness of sample paths

The connection between sample path smoothness and the inherent difficulty of learning is best understood for the Matérn covariance family and the limiting case of the squared exponential covariance. Combined, these covariance functions allow us to model functions with a *continuum* of smoothness. To this end, there is a general form of the Matérn covariance function modeling functions of any finite smoothness, controlled by a parameter $v > 0$:²⁵

$$K_M(d; v) = \frac{2^{1-v}}{\Gamma(v)} (\sqrt{2v}d)^v K_v(\sqrt{2v}d),$$

where $d = |x - x'|$ and K_v is the modified Bessel function of the second kind. Sample paths from a centered Gaussian process with this covariance are $\lceil v \rceil - 1$ times continuously differentiable, but the smoothness of sample paths is not as granular as a simple count of derivatives. Rather, the parameter v allows us to fine-tune sample path smoothness as desired.²⁶ Figure 10.3 illustrates sample paths generated from a Matérn covariance with a range of smoothness from $v = 1.01$ to $v = 2$. All of these samples are exactly once differentiable, but we might say that the $v = 1.01$ samples are “just barely” so, and that the $v = 2$ samples are “very nearly” twice differentiable.



Sample paths from a smooth GP (above) and a rough one (below). The rough samples have more degrees of freedom – and far more local maxima – and we might conclude they are harder to optimize as a result.

Matérn and squared exponential covariance functions: § 3.3, p. 51

25 The given expression has unit length and output scale; if desired, we can introduce parameters for these following § 3.4: smoothness parameter, v

26 This can be made precise through the coefficient of Hölder continuity in the “final” derivative of K_M , which controls the smoothness of the final derivative of sample paths. Except when v is an integer, the Matérn covariance with parameter v belongs to the Hölder space

$$\mathcal{C}^{\alpha, \beta}; \quad \alpha = \lfloor 2v \rfloor; \quad \beta = 2v - \lfloor 2v \rfloor,$$

and we can expect the final derivative of sample paths to be $(v - \lfloor v \rfloor)$ -Hölder continuous:

A. D. BULL (2011). Convergence Rates of Efficient Global Optimization Algorithms. *Journal of Machine Learning Research* 12(88):2879–2904.

Taking the limit $v \rightarrow \infty$ recovers the squared exponential covariance function K_{SE} . This serves as the extreme end of the continuum, modeling functions with infinitely many continuous derivatives. Together, the Matérn and square exponential covariance allow us to model functions with any smoothness $v \in (0, \infty]$.

Information capacity

We now require some way to relate the complexity of sample path behavior to our ability to learn about an unknown function. Information theory provides an answer through the concept of *information capacity*, the maximum rate of information transfer through a noisy observation mechanism.

In the analysis of Bayesian optimization algorithms, the central concern is how efficiently we can learn about a GP-distributed objective function $\mathcal{GP}(f; \mu, K)$ through a set of τ noisy observations \mathcal{D}_τ . The information capacity of this observation process, as a function of the number of observations τ , provides a fundamental bound on our ability to learn about f . For this discussion, let us adopt the common observation model of independent and homoskedastic additive Gaussian noise with scale $\sigma_n > 0$ (2.16). In this case, information capacity is a function of:

- the covariance K , which determines the information content of f , and
- the noise scale σ_n , which limits the amount of information obtainable through a single observation.

The information regarding f contained in an arbitrary set of observations $\mathcal{D} = (\mathbf{x}, \mathbf{y})$ can be quantified by the *mutual information* (A.16):

$$I(\mathbf{y}; f) = H[\mathbf{y}] - H[\mathbf{y} \mid \boldsymbol{\phi}] = \frac{1}{2} \log |\mathbf{I} + \sigma_n^{-2} \Sigma|, \quad (10.11)$$

²⁷ For more on information gain, see § 6.3, p. 115. Thus far we have primarily concerned ourselves with information gain regarding x^* or f^* ; here we are reasoning about the function f itself.

information capacity, γ_τ

where $\Sigma = K(\mathbf{x}, \mathbf{x})$. Note that the entropy of \mathbf{y} given $\boldsymbol{\phi}$ does not depend on the actual value of $\boldsymbol{\phi}$, and thus the mutual information $I(\mathbf{y}; f)$ is also the *information gain* about f provided by the data.²⁷ The *information capacity* (also known as the *maximum information gain*) of this observation process is now the maximum amount of information about f obtainable through *any* set of τ observations:

$$\gamma_\tau = \sup_{|\mathbf{x}|=\tau} I(\mathbf{y}; f). \quad (10.12)$$

Known bounds on information capacity

The information capacity of a GP observation process (10.12) is commonly invoked in theoretical analyses of algorithms making use of this model class. Unfortunately, working with information capacity is somewhat unwieldy for two reasons. First, as mentioned above, information capacity is a function of the covariance function K , which can be verified through the explicit formula in (10.11). Thus performance guarantees in terms of information capacity require further analysis to derive explicit

results for a particular choice of model. Second, information capacity of any given model is in general NP-hard to compute due to the difficulty of the set function maximization in (10.12).²⁸

For these reasons, the typical strategy is to derive *agnostic* convergence results in terms of the information capacity of an arbitrary observation process, then seek to derive bounds on the information capacity for notable covariance functions such as for the Matérn family and the squared exponential covariance. A common proof strategy for bounding the information gain is to relate the information capacity to the spectrum of the covariance function, with faster spectral decay yielding stronger bounds on information capacity. The first explicit bounds on information capacity for the Matérn and squared exponential covariances were provided by SRINIVAS et al.,²⁹ and the bounds for the Matérn class have since been sharpened using similar techniques.^{30,31}

For a compact domain $\mathcal{X} \subset \mathbb{R}^d$ and fixed noise scale σ_n , we have the following asymptotic bounds on the information capacity. For the Matérn covariance function with smoothness parameter v , we have:³¹

$$\gamma_\tau = \mathcal{O}(\tau^\alpha (\log \tau)^{1-\alpha}), \quad \alpha = \frac{d}{2v+d}; \quad (10.13)$$

and for the squared exponential covariance, we have:²⁹

$$\gamma_\tau = \mathcal{O}((\log \tau)^{d+1}). \quad (10.14)$$

These results embody our stated goal of characterizing smoother sample paths (as measured by v) as being inherently less complex than rougher sample paths. The information capacity decreases steadily as $v \rightarrow \infty$, eventually dropping to only logarithmic growth in τ for the squared exponential covariance. The correct interpretation of this result is that the smoother sample paths require *less* information to describe, and thus the maximum amount of information one *could* learn is limited compared to rougher sample paths.

Bounding the sum of predictive variances

A key result linking information capacity to an optimization policy is the following.²⁹ Suppose some optimization policy selected an arbitrary sequence of τ observation locations $\{x_i\}$, and let $\{\sigma_i\}$ be the corresponding predictive standard deviations at the time of selection. By repeatedly applying block determinant identities, we can rewrite the information gain (10.11) from observing y in terms of the marginal predictive variances:

$$I(y; f) = \frac{1}{2} \sum_{i=1}^{\tau} \log\left(1 + \frac{\sigma_i^2}{\sigma_n^2}\right). \quad (10.15)$$

Now assume that the prior covariance function is bounded: $K(x, x) \leq M$. Noting that $z^2 / \log(1 + z^2)$ is increasing for $z > 0$ and that $\sigma_i^2 \leq M$, the following inequality holds for every observation:

$$\sigma_i^2 \leq \frac{M}{\log(1 + \sigma_n^{-2} M)} \log\left(1 + \frac{\sigma_i^2}{\sigma_n^2}\right).$$

²⁸ C.-W. KO et al. (1995). An Exact Algorithm for Maximum Entropy Sampling. *Operations Research* 43(4):684–691.

²⁹ N. SRINIVAS et al. (2010). Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design. *ICML 2010*.

³⁰ D. JANZ et al. (2020). Bandit optimisation of functions in the Matérn kernel RKHS. *AISTATS 2020*.

³¹ S. VAKILI et al. (2021). On Information Gain and Regret Bounds in Gaussian Process Bandits. *AISTATS 2021*.

predictive standard deviation of ϕ_i , σ_i

bound on prior variance $K(x, x)$, M
posterior variance is nonincreasing; § 2.2,
p. 22.

- ³² If an explicit leading constant is desired, we have $\sum_i \sigma_i^2 \leq c\gamma_\tau$ with

$$c = \frac{2M}{\log(1 + \sigma_n^{-2}M)}.$$

We can now bound the sum of the predictive variances in terms of the information capacity:³²

$$\sum_{i=1}^{\tau} \sigma_i^2 = \mathcal{O}(\gamma_\tau). \quad (10.16)$$

This bound will repeatedly prove useful below.

10.4 BAYESIAN REGRET WITH OBSERVATION NOISE

We have now covered the background required to understand – or at least to meaningfully interpret – most of the theoretical convergence results appearing in the literature. In the remainder of the chapter we will summarize some notable results built upon these ideas. The literature is vast, and navigation can be challenging. Broadly, we can categorize these results according to certain dichotomies in their approach and focus, listed below.

analysis: frequentist or Bayesian?

observations: noisy or exact?

asymptotic behavior with logarithmic factors suppressed, \mathcal{O}^*

- The first is whether the result is regarding the worst-case (frequentist) regret or expected-cased (Bayesian) regret. In the former case, we assume the objective function lies in some RKHS \mathcal{H}_K with bounded norm (10.10) and seek to bound the worst-case regret on this space (10.6). In the latter case, we assume the objective function is a sample path from a Gaussian process $\mathcal{GP}(f; \mu, K)$ and seek to bound the expected regret (10.5).
- The second is the assumption and treatment of observation noise. Stronger guarantees can often be derived in the noiseless setting, as we can learn much faster about the objective function. Observation noise may also be modeled somewhat differently in the frequentist and Bayesian settings.

In this and the following sections, we will provide an overview of convergence results for all combinations of these choices: frequentist and Bayesian guarantees, with and without noise. For each of these cases, we will discuss both upper bounds on regret, which provide guarantees for the performance of specific Bayesian optimization algorithms, and also lower bounds on regret, which provide algorithm-agnostic bounds on the best possible performance. Here we will begin with results for Bayesian regret in the noisy setting.

To facilitate the discussion, we will adopt the notation \mathcal{O}^* (sometimes written $\tilde{\mathcal{O}}$ in other texts) to describe asymptotic bounds in which logarithmic factors of τ are suppressed:

$$f(\tau) = \mathcal{O}(g(\tau)(\log \tau)^k) \implies f(\tau) = \mathcal{O}^*(g(\tau)).$$

Common assumptions

In this section, we will assume that the objective function $f: \mathcal{X} \rightarrow \mathbb{R}$ is a sample path from a centered Gaussian process $\mathcal{GP}(f; \mu \equiv 0, K)$, and that observation noise is independent, homoskedastic Gaussian noise with scale $\sigma_n > 0$ (2.16). The domain \mathcal{X} will at various times be either a

finite set (as a stepping stone toward the continuous case) or a compact and convex subset of a d -dimensional cube: $\mathcal{X} \subset [0, m]^d$. We will also be assuming that the covariance function is continuous and bounded on \mathcal{X} : $K(x, x) \leq 1$. Since the covariance function is guaranteed to be bounded anyway (as \mathcal{X} is compact) this simply fixes the scale without loss of generality.

assumption: $K(x, x) \leq 1$ is bounded

Upper confidence bound

In a landmark paper, SRINIVAS et al. derived sublinear cumulative regret bounds for the Gaussian process upper confidence bound (GP-UCB) policy in the Bayesian setting with noise.³³ The authors considered policies of the form (8.25):

$$x_i = \arg \max_{x \in \mathcal{X}} \mu + \beta_i \sigma, \quad (10.17)$$

where x_i is the point chosen in the i th iteration of the policy, μ and σ are shorthand for the posterior mean and standard deviation of $\phi = f(x)$ given the data available at time i , \mathcal{D}_{i-1} , and β_i is a time-dependent exploration parameter. The authors were able to demonstrate that if this exploration parameter is carefully tuned over the course of optimization, then the cumulative regret of the policy can be asymptotically bounded in terms of the information capacity.

We will discuss this result and its derivation in some detail below, as it demonstrates important proof strategies that will be repeated throughout this section and the next.

upper confidence bound: § 7.8, p. 145, § 8.4, p. 170

³³ N. SRINIVAS et al. (2010). Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design. *ICML 2010*.

information capacity: § 10.3, p. 222

Regret bound on finite domains

To proceed, we first assume the domain \mathcal{X} is finite; we will lift this to the continuous case shortly via a secondary argument.³⁴

The construction of the UCB policy suggests the following *confidence ellipsoid* condition is likely to hold for any point $x \in \mathcal{X}$ at any time i :

$$\phi \in [\mu - \beta_i \sigma, \mu + \beta_i \sigma] = \mathcal{B}_i(x). \quad (10.18)$$

In fact, we can show that for appropriately chosen confidence parameters $\{\beta_i\}$, every such confidence ellipsoid is *always* valid, with high probability.

This may seem like a strong claim, but it is simply a consequence of the exponentially decreasing tails of the Gaussian distribution. At time i , we can use tail bounds on the Gaussian CDF to bound the probability of a given confidence ellipsoid failing in terms of β_i ,³⁵ then use the union bound to bound the probability of (10.18) failing anywhere at time i .³⁶ Finally, we show that by increasing the confidence parameter β_i over time – so that the probability of failure decreases suitably quickly – the probability of failure anywhere and at any time is small. SRINIVAS et al. showed in particular that for any $\delta \in (0, 1)$, taking

$$\beta_i^2 = 2 \log\left(\frac{i^2 \pi^2 |\mathcal{X}|}{6\delta}\right) \quad (10.19)$$

³⁴ The general strategy for this case was established in the linear bandit setting in:

V. DANI et al. (2008). Stochastic Linear Optimization Under Bandit Feedback. *COLT 2008*.

step 1: show confidence ellipsoids (10.18) are universally valid with high probability

³⁵ SRINIVAS et al. use

$$\Pr(\phi \notin \mathcal{B}_i(x)) \leq \exp(-\beta_i^2/2).$$

³⁶ Using the above bound, the probability of failure anywhere at time i is at most

$$|\mathcal{X}| \exp(-\beta_i^2/2).$$

THEORETICAL ANALYSIS

step 2: bound instantaneous regret by width of chosen confidence ellipsoid

the lower bound for every x holds for x^*

the upper bound for x_i holds for every x

³⁷ D. RUSSO and B. VAN ROY (2014). Learning to Optimize via Posterior Sampling. *Mathematics of Operations Research* 39(4):1221–1243.

step 3: bound cumulative regret via bound on instantaneous regret

³⁸ We have $M = 1$ as a bound on $K(x, x)$ according to our common assumptions (p. 224).

sample path differentiability: § 2.6, p. 30

³⁹ Hölder continuity of the derivative process covariance is sufficient; this holds for the Matérn covariance with $\nu > 1$.

guarantees that the confidence ellipsoids (10.18) are universally valid with probability at least $1 - \delta$.

With this claim established, we are actually almost finished. The next key insight is that if the confidence ellipsoids in (10.18) are universally valid, then we can bound the instantaneous regret of the policy in every iteration by noting that the confidence ellipsoid of the chosen point always contains the global optimum:

$$f^* \in \mathcal{B}_i(x_i). \quad (10.20)$$

To show this, we first note that the lower bound of every confidence ellipsoid applies to f^* , as all ellipsoids are valid and f^* is the global maximum. Further, the upper bound of the chosen confidence ellipsoid is valid for every function value, as it is the loosest possible upper bound by definition (10.17).

The result in (10.20) allows us to bound, with high probability, the instantaneous regret (10.2) in every iteration by the width of the confidence ellipsoid of the chosen point:

$$\rho_i \leq 2\beta_i \sigma_i, \quad (10.21)$$

RUSSO and VAN ROY interpret this bound on the instantaneous regret as guaranteeing that regret can only be high when we also *learn* a great deal about the objective function to compensate (10.15).³⁷

Finally, we bound the cumulative regret. Assuming the confidence ellipsoids (10.18) are universally valid, we may bound, with high probability, the sum of the squared instantaneous regret up to time τ by

$$\sum_{i=1}^{\tau} \rho_i^2 \leq 4 \sum_{i=1}^{\tau} \beta_i^2 \sigma_i^2 \leq 4\beta_{\tau}^2 \sum_{i=1}^{\tau} \sigma_i^2 = \mathcal{O}(\beta_{\tau}^2 \gamma_{\tau}). \quad (10.22)$$

From left-to-right, we plug in (10.21), note that $\{\beta_i\}$ is nondecreasing (10.19), and appeal to the information capacity bound on the sum of predictive variances (10.16).³⁸ Plugging in β_{τ} and appealing to the Cauchy-Schwartz inequality gives

$$R_{\tau} = \mathcal{O}^*(\sqrt{\tau \gamma_{\tau} \log |\mathcal{X}|}) \quad (10.23)$$

with probability at least $1 - \delta$, using the $\{\beta_i\}$ sequence in (10.19).

Extending to continuous domains

The above analysis can be extended to continuous domains via a discretization argument. The proof is technical, but the technique is general and potentially useful in other settings, so we provide a sketch.

Let us assume that the domain $\mathcal{X} \subset [0, m]^d$ is convex and compact. If the covariance function of our Gaussian process is smooth enough, its sample paths will be continuously differentiable;³⁹ as \mathcal{X} is compact, the sample paths will in fact be *Lipschitz* continuous with (random)

Lipschitz constant L . SRINIVAS et al. assume that the covariance function is sufficiently smooth to ensure the value of this Lipschitz constant has an exponential tail bound of the following form for $\lambda > 0$:

$$\Pr(L > \lambda) \leq da \exp(-\lambda^2/b^2), \quad (10.24)$$

where $a, b > 0$; this ensures we can control the value of L with high probability. We will say more about this assumption shortly.

To proceed, we analyze an algorithm that, in each iteration, discretizes the domain with a grid becoming increasingly fine over time. Rather than maximizing the upper confidence bound on all of \mathcal{X} (10.17), we maximize on the chosen grid instead. As the domain is always finite, we can hook into the analysis of the finite case to reason about this algorithm. Of course, this is not the true GP-UCB algorithm due to the discretization, but we can rely on Lipschitz continuity of sample paths to bound any extra regret caused by the discretization.

SRINIVAS et al. demonstrate in particular that a slight correction to the confidence parameter sequence from the finite case⁴⁰ (10.19) to include dependence on the parameters relevant to the discretization argument (those controlling the measure of \mathcal{X} , $\{m, d\}$, and the Lipschitz tail bound parameters, $\{a, b\}$) (10.24) is sufficient to ensure

$$R_\tau = \mathcal{O}^*(\sqrt{\tau \gamma_\tau d}) \quad (10.25)$$

with high probability. Thus, assuming the Gaussian process sample paths are smooth enough, the cumulative regret of GP-UCB on continuous domains grows at rate comparable to the discrete case.

Plugging in the information capacity bounds in (10.13–10.14) and dropping the dependence on dimension, we have the following high-probability regret bounds for specific covariance functions. For the Matérn covariance, we have

$$R_\tau = \mathcal{O}^*(\tau^\alpha), \quad \alpha = \frac{v+d}{2v+d}, \quad (10.26)$$

and for the squared exponential we have

$$R_\tau = \mathcal{O}^*(\sqrt{\tau}). \quad (10.27)$$

The growth is sublinear for all values of the smoothness parameter v , so the algorithm achieves no regret with high probability.

The Lipschitz tail bound condition

Finally, we address the exponential tail bound assumption on the Lipschitz constant (10.24). Despite its seeming strength in controlling the behavior of sample paths, it is actually a fairly weak assumption on the Gaussian process. In fact, all that is needed is that sample paths be continuous differentiability. In this case, each coordinate of the gradient,

⁴⁰ Specifically, SRINIVAS et al. took:

$$\begin{aligned} \beta_i^2 = 2 \log\left(\frac{i^2 \pi^2 |\mathcal{X}|}{3\delta}\right) \\ + 2d \log(i^2 dbm \sqrt{\log(4da/\delta)}). \end{aligned}$$

THEORETICAL ANALYSIS

- ⁴¹ A. W. VAN DER VAART and J. A. WELLNER (1996). *Weak Convergence and Empirical Processes with Applications to Statistics*. Springer-Verlag. [proposition A.2.1]

$f_i = \partial f / \partial x_i$, is a sample path continuous Gaussian process. By compactness, each f_i is then almost surely bounded on \mathcal{X} . Now the Borell-TIS inequality ensures that for any $\lambda > 0$ we have:⁴¹

$$\Pr(\max |f_i| > \lambda) \leq 4 \exp(-\lambda^2/b_i^2), \quad b_i = 2\sqrt{2} \mathbb{E}[\max |f_i|]. \quad (10.28)$$

Taking a union bound then establishes a bound of the desired form (10.24). In particular, this argument applies to the Matérn covariance with $v > 1$, which has continuously differentiable sample paths. Thus GP-UCB can achieve sublinear cumulative regret for all but the roughest of sample paths.

However, the bound in (10.28) does not immediately lead to an *effective* algorithm due to the difficulty of controlling the expected maximum $\mathbb{E}[\max |f_i|]$. Things simplify considerably if sample paths are *twice* differentiable, in which case GHOSAL and ROY showed it is possible to derive bounds of the above form (10.28) for each of the coordinates of the gradient process with *explicitly* computable constants,⁴² yielding an effective algorithm for the Matérn covariance with $v > 2$.

Intricate arguments regarding the objective function's Lipschitz constant are only necessary due to its randomness in the Bayesian setting. In the frequentist setting, where f is in some RKHS ball $\mathcal{H}_K[B]$, we can immediately derive a *hard* upper bound L in terms of K and B .⁴³

Thompson sampling

- ⁴² S. GHOSAL and A. ROY (2006). Posterior Consistency of Gaussian Process Prior for Nonparametric Binary Regression. *The Annals of Statistics* 34(5):2413–2429. [lemma 5]

- ⁴³ N. DE FREITAS et al. (2012a). Regret Bounds for Deterministic Gaussian Process Bandits. arXiv: 1203.2177 [cs.LG] [lemma 1]

Thompson sampling: § 7.9, p. 148, § 8.7, p. 176

- ⁴⁴ D. RUSSO and B. VAN ROY (2014). Learning to Optimize via Posterior Sampling. *Mathematics of Operations Research* 39(4):1221–1243.

RUSSO and VAN ROY developed a general approach for transforming Bayesian regret bounds for a wide class of UCB-style algorithms into regret bounds for analogous Thompson sampling algorithms.⁴⁴

Namely, consider a UCB policy selecting a sequence of points $\{x_i\}$ for observation by maximizing a sequence of “upper confidence bounds,” which here can be *any* deterministic functions of the observed data, regardless of their statistical validity. Let $\{u_i\}$ be the sequence of upper confidence bounds associated with the selected points at the time of their selection, and let $\{u_i^*\}$ be the sequence of upper confidence bounds associated with a given global optimum, x^* .

Now consider a corresponding Thompson sampling policy selecting $x_i \sim p(x^* \mid \mathcal{D}_{i-1})$. Note that, given the observed data, x_i and x^* are identically distributed by design; because the upper confidence bounds are deterministic functions of the observed data, u_i and u_i^* are identically distributed as well. This observation allows us to express the expected cumulative regret of the Thompson sampling policy entirely in terms of the upper confidence bounds:

$$\mathbb{E}[R_\tau] = \sum_{i=1}^\tau \mathbb{E}[u_i - \phi_i] + \sum_{i=1}^\tau \mathbb{E}[f^* - u_i^*]. \quad (10.29)$$

This is RUSSO and VAN ROY's main result.

With a bit of extra work, we may use this result to bound the expected regret of Gaussian process Thompson sampling (GP-TS) in terms of the

upper confidence bounds used in GP-UCB (10.17). The argument closely follows SRINIVAS et al.'s proof strategy for GP-UCB: we bound the regret on finite domains, then rely on the discretization argument above to lift the result to continuous domains.

Let us assume \mathcal{X} is finite and otherwise assume the common assumptions for this section. RUSSO and VAN ROY showed how to bound each of the terms in (10.29) in this case. First, they show that the confidence parameter sequence (compare with the analogous sequence for GP-UCB (10.19))

$$\beta_i^2 = 2 \log\left(\frac{(i^2 + 1)|\mathcal{X}|}{\sqrt{2\pi}}\right)$$

is sufficient to bound the second term by a constant: $\sum_i \mathbb{E}[f^* - u_i^*] \leq 1$.⁴⁵ The first term can then be bounded in terms of the information capacity following our previous discussion (10.16, 10.22):

$$\sum_{i=1}^{\tau} \mathbb{E}[u_i - \phi_i] = \sum_{i=1}^{\tau} \beta_i \sigma_i \leq \beta_{\tau} \sum_{i=1}^{\tau} \sigma_i = \mathcal{O}^*(\sqrt{\tau \gamma_{\tau} \log |\mathcal{X}|}). \quad (10.30)$$

This resulting bound on the GP-TS regret matches that for the GP-UCB algorithm (10.23). However, note that whereas the GP-UCB bound is a high-probability bound on the *actual* cumulative regret (with the desired probability determining the width of the confidence bounds (10.19)), the result for Thompson sampling bounds the *expected* cumulative regret unconditionally.

common assumptions: p. 224

⁴⁵ This is again a consequence of the rapidly decaying tails of the Gaussian distribution.

Upper bounds on simple regret

We may use the general bound on simple regret in terms of the average regret in (10.3) to derive asymptotic bounds on the simple regret of GP-UCB and GP-TS policies. Dropping dependence on the domain size, we have

$$r_{\tau} = \mathcal{O}^*(\sqrt{\gamma_{\tau}/\tau})$$

$$r_{\tau} \leq \frac{R_{\tau}}{\tau}$$

for both algorithms, where this is to be understood as either a high-probability (GP-UCB) or expected-case (GP-TS) result.

Plugging in the information capacity bounds (10.13–10.14), we have the following bounds for the simple regret for specific covariance functions. For the Matérn covariance, we have

$$r_{\tau} = \mathcal{O}^*(\tau^{\alpha}), \quad \alpha = -\frac{v}{2v+d},$$

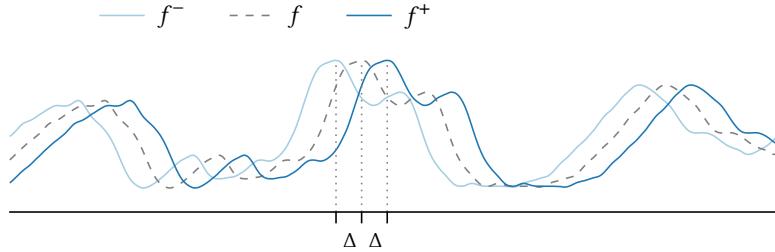
and for the squared exponential we have

$$r_{\tau} = \mathcal{O}^*(1/\sqrt{\tau}).$$

Lower bounds and tightness of existing algorithms

We have now derived upper bounds on the regret of particular Bayesian optimization algorithms in Bayesian setting with noise. A natural question is whether we can derive corresponding algorithm-agnostic *lower*

Figure 10.4: A sketch of SCARLETT’s proof strategy. Given access to a reference function f , which of its translations by Δ , f^- or f^+ is the objective function?



regret bounds establishing the fundamental difficulty of optimization in this setting, which might hint at how much room for improvement there may be.

Lower bounds on Bayesian regret are not easy to come by. As we will see, the frequentist setting offers considerable flexibility in deriving lower bounds, as we can construct explicit *objectives* in a given RKHS, then prove that these adversarial examples are difficult to optimize. In the Bayesian setting, the objective function is random, so we must instead seek explicit *distributions* over objective functions with enough structure that we can bound the expected regret, a much more challenging task.

That said, nontrivial lower bounds have been derived in this setting, most notably on the unit interval $\mathcal{X} = [0, 1]$. Under the assumption of a stationary Gaussian process with twice differentiable sample paths, SCARLETT demonstrated the following high-probability lower bound on the expected cumulative regret of any optimization algorithm:⁴⁶

$$\mathbb{E}[R_\tau] = \Omega(\sqrt{\tau}). \quad (10.31)$$

The key idea behind this result is to identify multiple plausible objective functions with identical distributions such that an optimization algorithm is prone to “prefer the wrong function,” and in doing so, incur high regret.

To illustrate the construction, we imagine an objective function on the unit interval is generated *indirectly* by the following process. We first realize an initial sample path f on the larger domain $[-\Delta, 1 + \Delta]$, for some $\Delta > 0$.⁴⁷ We then take the objective function to be one of the following translations of f with equal probability:

$$f^+: x \mapsto f(x + \Delta); \quad f^-: x \mapsto f(x - \Delta). \quad (10.32)$$

As the prior process is assumed to be stationary, this procedure does not change the distribution of the objective function.

To proceed, we consider optimization of the generated objective *given access to the initial reference function f* .⁴⁸ We now frame optimization in terms of using a sequence of noisy observations to determine which of the two possible translations (10.32) represents the objective function; see figure 10.4. *Fano’s inequality*,⁴⁹ an information-theoretic lower bound on error probability in adaptive hypothesis testing, then allows us to show that there is a significant probability that the data collected by *any* algorithm have better cumulative regret for the *wrong* translation.

⁴⁶ J. SCARLETT (2018). Tight Regret Bounds for Bayesian Optimization in One Dimension. *ICML 2018*.

⁴⁷ This introduces the additional minor requirement that the covariance function be defined on this larger domain.

⁴⁸ This is a so-called *genie argument*: as the extra information can be simply ignored if desired, it cannot possibly impede optimization.

⁴⁹ J. SCARLETT and V. CEVHAR (2021). An Introductory Guide to Fano’s Inequality with Applications in Statistical Estimation. In: *Information-Theoretic Methods in Data Science*.

Finally, we may use assumed properties of the objective function to show that when this happens, the cumulative regret is significant.⁵⁰

On the rougher side of the spectrum, WANG et al. used the same proof strategy to bound the expected regret – both simple and cumulative – of any algorithm maximizing a (nondifferentiable) sample path of Brownian motion (from the Wiener process) on the unit interval:⁵¹

$$\mathbb{E}[r_\tau] = \Omega(1/\sqrt{\tau \log \tau}); \quad \mathbb{E}[R_\tau] = \Omega(\sqrt{\tau \log \tau}). \quad (10.33)$$

Both WANG et al. and SCARLETT also described straightforward, but not necessarily practical, optimization algorithms to provide corresponding upper bounds on the unit interval. The algorithms are based on a simple branch-and-bound scheme whereby the domain is adaptively partitioned based on confidence bounds computed from available data. For relatively smooth sample paths,⁵² SCARLETT’s algorithm achieves cumulative regret

$$\mathbb{E}[R_\tau] = \mathcal{O}(\sqrt{\tau \log \tau})$$

with high probability. This is within a factor of $\sqrt{\log \tau}$ of the corresponding lower bound (10.31), so there is not too much room for improvement on the unit interval. Note that both GP-UCB and GP-TS with the squared exponential covariance match this rate (10.14, 10.25), but SCARLETT’s algorithm is better for the Matérn family with $\nu > 2$ (10.13). For Brownian motion sample paths, WANG et al. established the following upper bounds:

$$\mathbb{E}[r_\tau] = \mathcal{O}(\log \tau / \sqrt{\tau}); \quad \mathbb{E}[R_\tau] = \mathcal{O}(\sqrt{\tau \log \tau}), \quad (10.34)$$

which are within a factor of $(\log \tau)^{\frac{3}{2}}$ of the corresponding lower bounds (10.33), again fairly tight.

10.5 WORST-CASE REGRET WITH OBSERVATION NOISE

We now turn our attention to worst-case (frequentist) results analogous to the expected-case (Bayesian) results discussed in the previous section. Here we no longer view the objective function as a random sample path from a Gaussian process, but rather as a fixed, but unknown function lying in some reproducing kernel Hilbert space \mathcal{H}_K with norm bounded by some constant B (10.10). We also relax the assumption of independent Gaussian observation errors with more flexible and agnostic conditions. We then seek to bound the worst-case (simple or cumulative) regret (10.6) of a given algorithm when presented such an input, which we will notate with $\bar{r}_\tau[B]$ and $\bar{R}_\tau[B]$, respectively.

Worst-case regret is in some ways easier to deal with than Bayesian regret, as we no longer need to address uncertainty in the objective function. However, the analysis is also complicated by model misspecification, as the model assumptions underlying the Bayesian optimization algorithms are no longer valid.⁵³ We must therefore seek other methods of analysis.

⁵⁰ Specifically, we need that the global maximum of the reference function is unlikely to be “pushed off the edge” of the domain during translation, and that sample paths are likely to have locally quadratic behavior in a neighborhood of the global optimum. The lower bound on regret then holds with probability depending on these events. For more discussion, see:

N. DE FREITAS et al. (2012b). Exponential Regret Bounds for Gaussian Process Bandits with Deterministic Observations. *ICML 2012*.

⁵¹ Z. WANG et al. (2020b). Tight Regret Bounds for Noisy Optimization of a Brownian Motion. arXiv: 2001.09327 [cs.LG].

⁵² Again $\nu > 2$ is sufficient for the Matérn family.

worst-case regret on $\mathcal{H}_K[B]$: $\bar{r}_\tau[B], \bar{R}_\tau[B]$

⁵³ Recall sample paths of a centered Gaussian process with covariance K do *not* lie inside \mathcal{H}_K ; see p. 219.

assumption: $f \in \mathcal{H}_K[B]$

assumption: $K(x, x) \leq 1$ is bounded

assumption: errors have zero mean
conditioned on their history

assumption: scale of errors is limited
sub-Gaussian distribution

finite-domain Bayesian analysis of GP-UCB:
p. 225

Common assumptions

In this section, we will assume that the objective function $f: \mathcal{X} \rightarrow \mathbb{R}$, where \mathcal{X} is compact, lies in a reproducing kernel Hilbert space corresponding to the covariance function K and has norm bounded by some constant $B: f \in \mathcal{H}_K[B]$ (10.10). As in the previous section, we will also assume that the covariance function is continuous and bounded by unity on $\mathcal{X}: K(x, x) \leq 1$.

We will also assume that a sequence of observations at $\{x_i\}$ will take the form $y_i = \phi_i + \varepsilon_i$, and that the distribution of the errors respects mild regularity conditions. First, we will assume each ε_i has mean zero conditioned on its history:

$$\mathbb{E}[\varepsilon_i | \boldsymbol{\varepsilon}_{<i}] = 0,$$

where $\boldsymbol{\varepsilon}_{<i}$ is the vector of errors occurring before time i . We will also make assumptions regarding the scale of the errors. The most typical assumption is that the distribution of ε_i is σ_n -sub-Gaussian conditioned on its history, that is, that the tail of the conditional distribution shrinks at least as quickly as a Gaussian distribution with variance σ_n^2 :

$$\forall c > 0: \Pr(|\varepsilon_i| > c | \boldsymbol{\varepsilon}_{<i}) \leq 2 \exp(-\frac{1}{2}c^2/\sigma_n^2).$$

This condition is satisfied, for example, by a distribution bounded on the interval $[-\sigma_n, \sigma_n]$ and by any Gaussian distribution with standard deviation of at most σ_n .

Complementary with the above assumptions, the Bayesian optimization algorithms that we will analyze model the function with the centered Gaussian process $\mathcal{GP}(f; \mu \equiv 0, K)$ and assume independent Gaussian observation noise with scale σ_n .

Upper confidence bound and Thompson sampling

Both the Gaussian process upper confidence bound (GP-UCB) and Thompson sampling (GP-TS) algorithms have worst-case regret bounds with rates comparable to the Bayesian setting.

The primary strategy for deriving regret bounds for the GP-UCB algorithm is to prove concentration inequalities regarding the deviation of the posterior mean of the Gaussian process from the objective function in the assumed RKHS. A prototypical result is to show that for some sequence of confidence parameters $\{\beta_i\}$, the following confidence ellipsoid assumption (10.18) is, with high probability, universally valid:

$$\phi \in [\mu - \beta_i \sigma, \mu + \beta_i \sigma]. \quad (10.35)$$

If this holds, we can then bound the cumulative regret in terms of the information capacity following our previous analysis. An advantage of the frequentist setting is that we can avoid complicated arguments lifting this result from discrete to continuous domains: the assumption of the objective function having bounded RKHS norm offers enough regularity that we can show (10.35) is universally valid even in continuous domains.

SRINIVAS et al. proved a concentration inequality of the above form for a parameter sequence $\{\beta_i\}$ depending on the RKHS norm B and the information capacity γ_i , assuming that the errors $\{\varepsilon_i\}$ are almost surely bounded, a stronger assumption than sub-Gaussianity.⁵⁴ This result was improved by ABBASI-YADKORI in the context of regression and later applied in the context of optimization by CHOWDHURY and GOPALAN.^{55,56} Namely, under our common assumptions and the assumption of σ_n -sub-Gaussian errors, the confidence parameter sequence

$$\beta_i = B + \sigma_n \sqrt{2(\gamma_{i-1} + \log(1/\delta))} \quad (10.36)$$

provides universally valid confidence ellipsoids with probability at least $1 - \delta$. Note that these confidence parameters are much larger than needed in the Bayesian setting (see (10.19) and footnote 40), due in part to the model mismatch between the Gaussian process model and the smoother objective function lying in the RKHS. Following the same argument leading up to (10.22), universally valid confidence ellipsoids allow us to show that, with high probability, the worst-case cumulative regret is⁵⁷

$$\bar{R}_\tau[B] = \mathcal{O}^*(\sqrt{\tau}\gamma_\tau) \quad (10.37)$$

CHOWDHURY and GOPALAN also extended this argument to provide worst-case regret bounds for Gaussian process Thompson sampling. Although the instantaneous regret in each iteration of Thompson sampling can no longer be bounded by the width of the confidence ellipsoids due to its stochastic nature, the authors showed that the algorithm nonetheless samples “good” points with low regret sufficiently often that the worst-case cumulative regret is, with high probability, limited to

$$\bar{R}_\tau[B] = \mathcal{O}^*(\sqrt{d}\tau\gamma_\tau). \quad (10.38)$$

Plugging in the information capacity bounds in (10.13–10.14) and dropping the \sqrt{d} factor in (10.38), we have the following high-probability regret bounds for specific covariance functions. For the Matérn covariance, we have the following bound for GP-UCB and GP-TS:

$$\bar{R}_\tau[B] = \mathcal{O}^*(\tau^\alpha); \quad \alpha = \frac{2\nu + 3d}{4\nu + 2d}, \quad (10.39)$$

and for the squared exponential we have:

$$\bar{R}_\tau[B] = \mathcal{O}^*(\sqrt{\tau}). \quad (10.40)$$

This is sublinear only when $d < 2\nu$, so the best-known results in the frequentist case only guarantee convergence for sufficiently smooth objective functions.

There may be room to improve these results. In particular, there is a gap of $\sqrt{\gamma_\tau}$ between the above bounds and the corresponding results in the Bayesian case. There is good reason to believe these bounds are not tight. No such gap on the worst-case exists in the finite-dimensional case under otherwise similar conditions,⁵⁸ and the lower bounds described below also match the lower (and upper) bounds in the Bayesian setting.

⁵⁴ N. SRINIVAS et al. (2010). Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design. *ICML 2010*. [theorem 6]

⁵⁵ Y. ABBASI-YADKORI (2012). Online Learning for Linearly Parameterized Control Problems. PhD thesis. University of Alberta. [theorem 3.11, remark 3.13]

⁵⁶ S. R. CHOWDHURY and A. GOPALAN (2017). On Kernelized Multi-armed Bandits. *ICML 2017*. [theorem 2]

⁵⁷ If we wish to retain explicit dependence on the RKHS norm B , we have

$$\bar{R}_\tau[B] = \mathcal{O}^*(\sqrt{\tau\gamma_\tau}(B + \sqrt{\gamma_\tau})).$$

⁵⁸ M. VALKO et al. (2013). Finite-Time Analysis of Kernelised Contextual Bandits. *UAI 2013*.

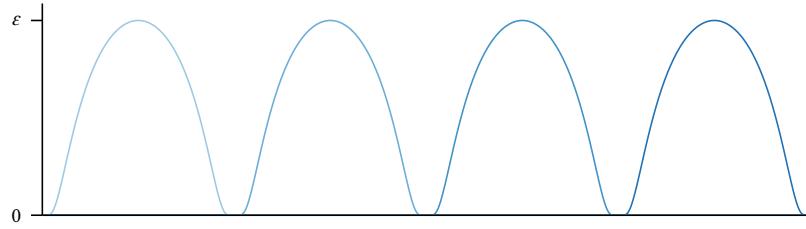


Figure 10.5: Four smooth objective functions with disjoint support.

needles in haystacks: see p. 216

⁵⁹ J. SCARLETT et al. (2017). Lower Bounds on Regret for Noisy Gaussian Process Bandit Optimization. *COLT 2017*.



The bump function used in SCARLETT et al.’s lower bound analysis, the Fourier transform of a smooth function with compact support as in figure 10.5.

⁶⁰ This particular bump function is the prototypical smooth function with compact support:

$$x \mapsto \begin{cases} \exp\left(-\frac{1}{1-|x|^2}\right) & |x| < 1; \\ 0 & \text{otherwise.} \end{cases}$$

Lower bounds and tightness of existing algorithms

Lower bounds on regret are much easier to come by in the frequentist setting than in the Bayesian setting, as we have considerable freedom to construct *explicit* adversarial objective functions that are provably difficult to optimize.

A common strategy to this end is to take inspiration from the “needle in a haystack” trick discussed earlier. We construct a large set of suitably well behaved “needles” that have disjoint (or near disjoint) support, then argue that there will always be some needle “missed” by an algorithm with insufficient budget to distinguish all the functions. Figure 10.5 shows a motivating example with four translations of a smooth bump function with height ϵ and mutually disjoint support. Given any set of three observations – regardless of how they were chosen – the cumulative regret for at least one of these functions would be 3ϵ .

This construction embodies the spirit of most of the lower bound arguments appearing in the literature. To yield a full proof, we must show how to construct a large number of suitable needles with bounded RKHS norm. For a *stationary* process, this can usually be accomplished by scaling and translating a suitable bump-shaped function to cover the domain. We also need to bound the regret of an algorithm given an input chosen from this set; here, we can keep the set of potential objectives larger than the optimization budget and appeal to pigeonhole arguments.

In the frequentist setting with noise, the strongest known lower bounds are due to SCARLETT et al.⁵⁹ The function class considered in the analysis was scaled and translated versions of a function similar to (in fact, precisely the Fourier transform of) the bump function in figure 10.5;⁶⁰ this function has the advantage of having “nearly compact” support while having finite RKHS norm in the entire Matérn family. For optimization on the unit cube $\mathcal{X} = [0, 1]^d$ with the Matérn covariance function, the authors were able to establish a lower bound on the cumulative regret of any algorithm of

$$\bar{R}_\tau[B] = \Omega(\tau^\alpha); \quad \alpha = \frac{\nu + d}{2\nu + d},$$

and for the squared exponential covariance

$$\bar{R}_\tau[B] = \Omega(\sqrt{\tau}).$$

These bounds are within log factors of the best-known upper bounds in the *Bayesian* setting (10.26–10.27), but there is a sizable gap (on the order

of $\sqrt{\gamma_\tau}$) between this lower bound and the best-known upper bounds (10.39–10.40) in the worst-case setting.

SCARLETT et al. also provide lower bounds on the worst-case *simple* regret $\bar{r}_\tau[B]$, in terms of the expected time required to reach a given level of regret. Inverting these bounds in terms of the simple regret at a given time yields rates that are as expected in light of the relation in (10.3).

$$r_\tau \leq \frac{R_\tau}{\tau}$$

10.6 THE EXACT OBSERVATION CASE

The arguments outlined in the previous sections all ultimately depend on the information capacity of noisy observations of a Gaussian process through the pivotal result (10.16), which allows us to control the width of confidence ellipsoids in terms of the information capacity. After noting that the instantaneous regret of an upper confidence bound algorithm is in turn bounded by the width of these confidence ellipsoids (10.21), we may piece together a bound on its cumulative regret (10.22).

Unfortunately, this line of attack breaks down with exact observations, as the information capacity of the now *deterministic* observation process is no longer well defined.⁶¹ However, all is not lost – we can appeal to different techniques to find much *stronger* bounds on the width of confidence ellipsoids induced by exact observations, and thereby establish much faster rates of convergence than in the noisy case.

For convenience, as in the rest of this chapter, throughout this section we will assume the domain $\mathcal{X} \subset \mathbb{R}^d$ is compact and that the covariance function is bounded by unity: $K(x, x) \leq 1$.

Bounding the posterior standard deviation

As before, the key idea in deriving regret bounds with exact observations is to bound the width of confidence ellipsoids (10.35) derived from the posterior process. In this light, let us consider the behavior of a Gaussian process $\mathcal{GP}(f; \mu, K)$ on an objective $f: \mathcal{X} \rightarrow \mathbb{R}$ as it is conditioned on a set of exact observations. Given some set of observation locations $\mathbf{x} \subset \mathcal{X}$, we seek to bound the *maximum* posterior standard deviation of the process:

$$\bar{\sigma}_{\mathbf{x}} = \max_{x \in \mathcal{X}} \sqrt{K_D(x, x)},$$

in terms of properties of \mathbf{x} . Given such a bound, we may then bound the cumulative regret of a UCB-style algorithm via previous arguments.

Thankfully, bounds of this type have enjoyed a great deal of attention, and strong results are available.⁶² Most of these results bound the posterior standard deviation of a Gaussian process in terms of the so-called *fill distance*, a measure of how densely a given set of observations fills the domain. For a set of observation locations $\mathbf{x} \subset \mathcal{X}$, its fill distance is defined to be

$$\delta_{\mathbf{x}} = \max_{x \in \mathcal{X}} \min_{x_i \in \mathbf{x}} |x - x_i|, \quad (10.41)$$

the largest distance from a point in the domain to the closest observation.

⁶¹ This can be seen by taking $\sigma_n \rightarrow 0$ in (10.15).

maximum posterior standard deviation given exact observations at \mathbf{x} , $\bar{\sigma}_{\mathbf{x}}$

⁶² The literature on this topic is substantial, but the following references provide good entry points:

M. KANAGAWA et al. (2018). Gaussian Processes and Kernel Methods: A Review on Connections and Equivalences. arXiv: 1807.02582 [stat.ML]. [§ 5.2]

H. WENDLAND (2004). *Scattered Data Approximation*. Cambridge University Press. [chapter 11]

fill distance of observation locations \mathbf{x} , $\delta_{\mathbf{x}}$

Intuitively, we should expect the maximum posterior standard deviation induced by a set of observation locations \mathbf{x} to shrink with its fill distance. This is indeed the case, and particularly nice results for the rate of this shrinkage are available for the Matérn family. Namely, for finite v , we have

$$\bar{\sigma}_{\mathbf{x}} \leq c \delta_{\mathbf{x}}^v, \quad (10.42)$$

where the constant c does not depend on \mathbf{x} . This bound is obviously only useful once the fill distance becomes less than the (unit) length scale of the covariance function: $\delta_{\mathbf{x}} < 1$. Once this is the case, the posterior standard deviation shrinks rapidly with the fill distance, even more so as sample paths become smoother in the limit $v \rightarrow \infty$.⁶³ This result will be instrumental in several results discussed below.

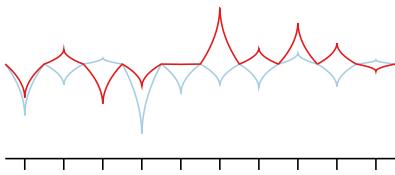
Bayesian regret with deterministic observations

⁶³ Carefully taking this limit yields the following bound for the squared exponential covariance:

$$\bar{\sigma}_{\mathbf{x}} \leq \exp(c \log(\delta_{\mathbf{x}})/\sqrt{\delta_{\mathbf{x}}}). \quad (10.43)$$

⁶⁴ S. GRÜNEWÄLDER et al. (2010). Regret Bounds for Gaussian Process Bandit Problems. *AISTATS 2010*.

⁶⁵ The upper bound is within a factor of $\sqrt{\log \tau}$ of WANG et al.’s upper bound for the *noisy* optimization of Brownian motion sample paths, where $\alpha = d = 1$ (10.34). The lower bound is identical to that case (10.33).



Samples from an example Gaussian process used in deriving GRÜNEWÄLDER et al.’s lower bound. Here 10 $(1/2)$ -Hölder continuous “needles” are scaled and translated to cover the unit interval, with one centered on each tick mark. Each needle is then scaled by an independent normal random variable and summed, yielding a Gaussian process with the desired properties.

⁶⁶ P. MASSART (2007). *Concentration Inequalities and Model Selection: Ecole d’Été de Probabilités de Saint-Flour XXXIII – 2003*. Vol. 1896. Springer–Verlag.

GRÜNEWÄLDER et al. provided bounds on Bayesian regret in the exact observation case, although their results are only relevant for particularly rough objectives.⁶⁴ The authors considered a function on the unit cube $f: [0, 1]^d \rightarrow \mathbb{R}$ with distribution $\mathcal{GP}(f; \mu, K)$, where the mean function is Lipschitz continuous and the covariance function is α -Hölder continuous. The latter is an exceptionally weak smoothness assumption on the prior process: the Matérn covariance with parameter v is Hölder continuous with $\alpha = \min(2v, 1)$, and any distinction in smoothness is lost beyond $v > 1/2$, where sample paths are not yet even differentiable. Under these assumptions, GRÜNEWÄLDER et al. proved the following bounds on the Bayesian simple regret:⁶⁵

$$\begin{aligned} \mathbb{E}[r_\tau] &= \Omega(\tau^\alpha / \sqrt{\log \tau}); \quad c = -\frac{\alpha}{2d}; \\ \mathbb{E}[r_\tau] &= \mathcal{O}(\tau^\alpha \sqrt{\log \tau}). \end{aligned}$$

The lower bound derives from constructing an explicit Gaussian process with Hölder continuous covariance function whose “needle in a haystack” nature makes it difficult to optimize by any strategy. Specifically, for any budget τ and Hölder continuity parameter α , the authors show how to cover the unit cube with 2τ disjoint and α -Hölder continuous “needles” of compact support. We may then create a Gaussian process by scaling each of these needles by an independent normal random variable and summing. This construction is perhaps most clearly demonstrated visually, and the marginal figure shows an example on the unit interval with $\tau = 5$ and $\alpha = 1/2$. As the needles are disjoint with independently random heights, no policy with a budget of τ can determine more than half of the weights, and we must therefore pay a penalty in terms of expected simple regret.

The upper bound derives from analyzing the performance of a naïve grid strategy via classical concentration inequalities.⁶⁶ It is remarkable that a nonadaptive strategy would yield such a small gap in regret compared to the corresponding lower bound, but this result is probably best

understood as illustrating the inherent difficulty of optimizing rough functions rather than any inherent aptitude of grid search.

To underscore this remark, we may turn to a result of DE FREITAS et al., who showed that we may optimize sample paths of sufficiently smooth Gaussian processes with *exponentially* decreasing expected simple regret in the deterministic setting.⁶⁷ This result relies on a few technical properties of the Gaussian process in question, in particular that it have a unique global optimum and that it exhibit “nice” behavior in the neighborhood of the global optimum.⁶⁸ A centered Gaussian process with covariance function in the Matérn family exhibits the required smoothness assumptions if $\nu > 2$.

The algorithm analyzed by the authors was a simple branch-and-bound policy based on GP-UCB, wherein the domain is recursively subdivided into finer and finer divisions. After each division, we identify the regions that could still contain the global optimum based on the current confidence ellipsoids, then evaluate on these regions such that the fill distance (10.41) is sufficiently small before the next round of subdivision. Here DE FREITAS et al. rely on the bound in (10.43) (with $\nu = 2$) to ensure that the confidence ellipsoids induced by observed data shrink rapidly throughout this procedure.

At some point in this procedure, we will (with high probability) find ourselves having rejected all but the local neighborhood of the global optimum, at which point the assumed “nice” behavior of the sample path guarantees rapid convergence thereafter. Specifically, the authors demonstrate that at some point the *instantaneous* regret of this algorithm will converge exponentially:

$$\mathbb{E}[\rho_\tau] = \mathcal{O}\left(\exp\left(-\frac{c\tau}{(\log \tau)^{d/4}}\right)\right),$$

for some constant $c > 0$ depending on the process but not on τ . This condition implies rapid convergence in terms of simple regret as well (as we obviously have $r_\tau \leq \rho_\tau$) and *bounded* cumulative regret after we have entered the converged regime.⁶⁹

Evidently optimization is *much* easier with deterministic observations than noisy ones, at least in terms of Bayesian regret. Intuitively, the reason for this discrepancy is that noisy observations may compel us to make repeated measurements in the same region in order to shore up our understanding of the objective function, whereas this is never necessary with exact observations.

Worst-case regret with deterministic observations

Worst-case bounds have also been derived for regret arising from deterministic observations. BULL for example derived tight upper and lower bounds on the worst-case simple regret in this setting for covariance functions in the Matérn family with finite smoothness $\nu < \infty$:

$$\bar{r}_\tau[B] = \Theta(\tau^{-\nu/d}). \quad (10.44)$$

⁶⁷ N. DE FREITAS et al. (2012b). Exponential Regret Bounds for Gaussian Process Bandits with Deterministic Observations. *ICML 2012*

uniqueness of global maxima: § 2.7, p. 34

⁶⁸ Here “nice” means that if the global optimum lies on the interior of the domain, the objective function must have locally quadratic behavior, and if the global optimum lies on the boundary of the domain, the objective function must not intersect the boundary with gradient zero at the maximum.

⁶⁹ This is not inconsistent with GRÜNEWÄLDER et al.’s lower bound: here we assume smoothness consistent with $\nu > 2$ in the Matérn class, whereas the adversarial Gaussian process constructed in the lower bound is only as smooth as $\nu = 1/2$.

THEORETICAL ANALYSIS

lower bound: 2τ needles with bounded RKHS norm and disjoint support

upper bound: shrinking fill distance limits error

⁷⁰ H. WENDLAND (2004). *Scattered Data Approximation*. Cambridge University Press [chapter 11]

convergence rates for expected improvement

As with previous results, the lower bound derives from an adversarial construction where we construct 2τ “needles” with bounded RKHS norm and disjoint support and argue that any policy with a budget of τ will necessarily incur significant regret. The smooth bump function illustrated in figure 10.5 suffices for this construction, after suitable scaling and translation.

The upper bound is achieved by a naïve grid strategy: as the domain $\mathcal{X} \subset \mathbb{R}^d$ is compact, it is also bounded, and thus we may construct grids with fill distance (10.41) shrinking as $\mathcal{O}(\tau^{-1/d})$. Now we may appeal to kernel interpolation results effectively equivalent to the result in (10.42) to show the simple regret of grid search in this setting decreases as $\mathcal{O}(\tau^{-v/d})$.⁷⁰

This result is perhaps not as exciting as it could be, as it demonstrates that no adaptive algorithm can perform (asymptotically) better than grid search in this setting. However, we may reasonably seek similar guarantees for algorithms that are also effective in practice. BULL was able to show that maximizing expected improvement yields worst-case simple regret

$$\bar{r}_\tau[B] = \mathcal{O}^*(\tau^{-\min(v,1)/d}),$$

which is near optimal for $v \leq 1$. BULL also showed that augmenting expected improvement with occasional random exploration akin to an ϵ -greedy policy improves its performance to near optimal (10.44) for any finite smoothness:

$$\bar{r}_\tau[B] = \mathcal{O}^*(\tau^{-v/d}).$$

The added randomness effectively guarantees the fill distance of observations shrinks quickly enough that we can still rely on posterior contraction arguments, and this strategy could be useful for analyzing other policies.

These results were refined by VAKILI et al., who proved the following upper bounds on the worst-case *cumulative* regret for the GP-UCB algorithm using a Matérn kernel with finite smoothness $v < \infty$ on the d -dimensional unit ball:⁷¹

$$\bar{R}_\tau[B] = \begin{cases} \mathcal{O}(\tau^\alpha); & \alpha = 1 - \frac{v}{d}; \quad v < d; \\ \mathcal{O}(\log \tau); & v = d; \\ \mathcal{O}(1); & v > d. \end{cases} \quad (10.45)$$

These results are compatible with BULL’s bounds on the simple regret (10.44), but offer more detail in the very smooth regime $v \geq d$.

The proof of these bounds relied on two key observations. First, we note that with exact observations we can guarantee that confidence ellipsoids (10.35) are universally valid for a function in $\mathcal{H}_K[B]$ by simply taking the *constant* confidence parameter $\beta = B$.⁷² Next we note that we can bound the posterior standard deviation of the point chosen at time i in terms of its distance to the nearest previous observation:⁷³

$$\sigma_i \leq c \min_{j < i} |x_i - x_j|^v$$

⁷¹ S. VAKILI et al. (2020). Regret Bounds for Noise-Free Bayesian Optimization. arXiv: 2002.05096 [stat.ML].

⁷² Compare to the noisy case where we had to inflate the confidence parameters over time (10.36), and for this result in a related context see also:

N. DE FREITAS et al. (2012a). Regret Bounds for Deterministic Gaussian Process Bandits. arXiv: 1203.2177 [cs.LG]. [lemma 2]

⁷³ VAKILI et al. provide a clever proof by appealing to the result in (10.42): we consider the maximum standard deviation on the ball centered on the closest observation, which has fill distance at most $\min|x_i - x_j|$. The standard deviation using observations from the entire space can only be smaller.

for some constant c . These two results immediately allow us to bound the cumulative regret of GP-UCB as:⁷⁴

$$\bar{R}_\tau = \sum_i \rho_i \leq 2\beta \sum_i \sigma_i \leq 2\beta + 2\beta c \sum_{i>1} \min_{j< i} |x_i - x_j|^v$$

The final step of the proof was to provide bounds for the term

$$\sum_{i>1} \min_{j< i} |x_i - x_j|^v$$

and VAKILI et al. provide a geometric argument giving rise to (10.45).

VAKILI et al. also provided upper bounds on the cumulative regret of Gaussian process Thompson sampling, which match (10.45) after scaling by an additional factor of $\sqrt{d \log \tau}$.

⁷⁴ Recall the instantaneous regret is bounded by the width of the corresponding confidence ellipsoid (10.21): $\rho_i \leq 2\beta\sigma_i$. The 2β term on the right hand side comes from special handling of the first observation and the bound $K(x, x) \leq 1$.

10.7 THE EFFECT OF UNKNOWN HYPERPARAMETERS

We have now outlined a plethora of convergence results: upper bounds on the regret of specific algorithms and algorithm-agnostic lower bounds, in the expected and worst case, with and without observation noise. However, all of these results assumed *intimate* knowledge about the objective function being optimized: in Bayesian analysis, we assumed the objective is sampled from the prior used to model it, and for worst-case analysis, we assumed the objective lay in the corresponding RKHS. Of course, neither of these assumptions is likely to hold in any practical setting. Further, in practice the model used to reason about the objective is typically inferred from observed data and constantly updated throughout optimization, but the analysis discussed thus far has assumed the model is not only perfectly informed, but also fixed.

It turns out that violations to these (rather implausible!) assumptions can be disastrous for convergence. For example, consider this prototypical Bayesian optimization algorithm: we model the objective function with an automatic relevance determination (ARD) version of a Matérn covariance function, learn its output (3.20) and length scales (3.25) via maximum likelihood estimation, and design observation locations by maximizing expected improvement. Although this setup has proven remarkably effective in practice,⁷⁵ BULL proved that it can actually fail miserably in the frequentist setting: we may construct functions in the RKHS of any such covariance function (that is, with any desired parameters) “fooling” this algorithm into having high regret with high probability.⁷⁶

Estimation of hyperparameters can also cause problems with convergence in the Bayesian setting. For example, LOCATELLI considered optimization on the unit interval $\mathcal{X} = [0, 1]$ from exact observations using a Wiener process prior and maximizing expected improvement. In this relatively straightforward model, the only hyperparameter under consideration was an output scale (3.20). LOCATELLI showed that for a fixed output scale, this procedure converges for *all* continuous functions.⁷⁷ However, when the output scale is learned from data – even in a

modeling with Gaussian processes: § 3, p. 45

automatic relevance determination: § 3.4, p. 56

ML estimation of hyperparameters: § 4.3, p. 73

expected improvement: § 7.3, p. 127

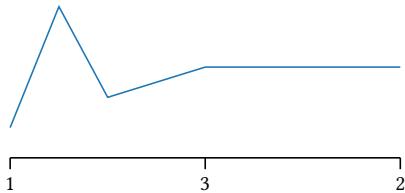
⁷⁵ J. SNOEK et al. (2012). Practical Bayesian Optimization of Machine Learning Algorithms. *NeurIPS 2012*.

⁷⁶ A. D. BULL (2011). Convergence Rates of Efficient Global Optimization Algorithms. *Journal of Machine Learning Research* 12(88):2879–2904. [theorem 3]

⁷⁷ M. LOCATELLI (1997). Bayesian Algorithms for One-Dimensional Global Optimization. *Journal of Global Optimization* 10(1):57–76

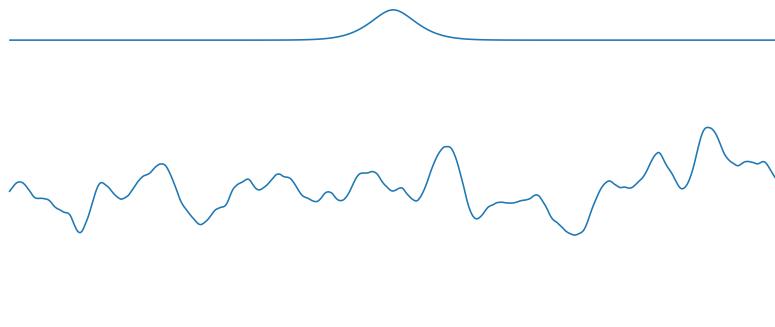
THEORETICAL ANALYSIS

Figure 10.6: The problem with estimating hyperparameters. Above: a function in the RKHS for the $\nu = 5/2$ Matérn covariance with unit norm. Below: a sample path from a Gaussian process with the same covariance. Unless we're lucky, we may never find the “hump” in the former function and thus may never build a reasonable belief regarding the objective.



A function given by LOCATELLI for which expected improvement combined with a Wiener process prior will not converge if its output scale is marginalized. The first three observations are at the indicated locations; after iteration 3, the left hand side of the domain is forever ignored.

⁷⁸ Recall that sample paths almost surely do not lie in \mathcal{H}_K , assuming it is infinite dimensional, and see figure 10.2 and surrounding text.



fully Bayesian manner where it is endowed with a prior distribution and marginalized in the predictive distribution – there are extremely simple (piecewise linear!) functions for which this algorithm does *not* recover the optimum. An example is shown in the margin.

In both of these examples, the roadblock to convergence when estimating hyperparameters is a mismatch between the objective function model used by the Bayesian optimization algorithm and the true objective. In the frequentist setting, there is an inherent tension as functions lying in the RKHS \mathcal{H}_K are *not* representative of sample paths from the Gaussian process $\mathcal{GP}(f; \mu, K)$.⁷⁸

Functions in the RKHS are much smoother than sample paths and may feature relatively “flat” regions of little variation, which may in turn lead to poorly fit models. See figure 10.6 for a striking example. LOCATELLI’s piecewise linear counterexample also reflects extreme model misspecification – the function is far too smooth to resemble Brownian motion in any meaningful statistical sense, and we should not be surprised that maximum likelihood estimation leads to an inconsistent algorithm.

A line of continuing work has been to derive *robust* algorithms that do not require perfect prior knowledge regarding the objective function to provide strong convergence guarantees. There are several ideas in this direction, all of which employ some mechanism to ensure that the space of objective functions considered by the algorithm does not ever become too “small.” For example, both BULL and LOCATELLI addressed their counterexamples with schemes wherein the output scale used during optimization is never allowed to shrink so rapidly that true objective function is left behind.⁷⁹

In the frequentist setting, one common strategy is to replace the assumption that the objective function lay in some particular RKHS with the assumption that it lay in some parametric family of RKHSes indexed by a set of hyperparameters. We then slowly expand the space of functions considered in the algorithm over the course of the algorithm such that the objective function is guaranteed to eventually – and forever thereafter – be well explained. In particular, consider augmenting an isotropic covariance function from the Matérn family K with an output scale (3.20) λ and a vector of length scales ℓ (3.25). We have the remarkable property that if a function f is in the RKHS ball of radius B (10.10)

⁷⁹ BULL simply inflated the output scale estimated by maximum likelihood by a factor of $\sqrt{7}$, and LOCATELLI chose a prior on the output scale that had no support below some minimum threshold.

for some setting of these hyperparameters:

$$f \in \mathcal{H}_K[B; \lambda, \ell],$$

then the same function is *also* in the RKHS ball for any larger output scale and for any vector of shorter (as in the lexicographic order) length scales:

$$f \in \mathcal{H}_K[B; \lambda', \ell']; \quad \lambda' \geq \lambda; \quad \ell' \leq \ell.$$

With this in mind, a natural idea for deriving theoretical convergence results (but not necessarily for realizing a *practical* algorithm!) is to ignore any data-dependent scheme for setting hyperparameters and instead simply slowly increase the output scale and slowly decrease the length scales over the course of optimization, so that at some point any function lying in any such RKHS will eventually be captured. This scheme has been used to provide convergence guarantees for both expected improvement⁸⁰ and GP-UCB⁸¹ with unknown hyperparameters.

SUMMARY OF MAJOR IDEAS

- Convergence analysis for Bayesian optimization algorithms entails selecting a measure of optimization performance, a space of objective functions to consider, and deriving bounds on the asymptotic growth of the chosen performance measure on the chosen function space.
- Optimization performance is almost always assessed via one of two related notations of *regret*, both of which depend on the difference between the function value at measured locations and the global optimum.
- Although tempting, the space of all continuous functions is too large to be of much interest in analysis: we may construct “needles in haystacks” to foil any algorithm by any amount. Instead, we may study spaces of objective functions with more plausible behavior. A Gaussian process model suggests two natural possibilities: sample paths from the process and the reproducing kernel Hilbert space (RKHS) associated with the process, the closure of the space of all possible posterior mean functions.
- Putting these pieces together, a Bayesian analysis of regret assumes the objective is a sample path from a Gaussian process and seeks asymptotic bounds on the expected regret (10.5). A frequentist analysis, on the other hand, assumes the objective function is a fixed, but unknown function in a given RKHS and seeks asymptotic bounds on the worst-case regret (10.5).
- Most upper bounds on regret derive from a proof strategy where we identify some suitable set of predictive credible intervals that are universally valid with high probability. We may then argue that the cumulative regret is bounded in terms of the total width of these intervals. This strategy lends itself most naturally to analyzing the Gaussian process upper confidence bound (GP-UCB) policy, but also yields bounds for Gaussian process Thompson sampling (GP-TS) due to strong theoretical connections between these algorithms (10.29).

⁸⁰ Z. WANG and N. DE FREITAS (2014). Theoretical Analysis of Bayesian Optimization with Unknown Gaussian Process Hyper-Parameters. arXiv: 1406.7758 [stat.ML].

⁸¹ F. BERKENKAMP et al. (2019). No-Regret Bayesian Optimization with Unknown Hyperparameters. *Journal of Machine Learning Research* 20(50):1–24.

simple and cumulative regret: § 10.1, p. 213

useful function spaces for studying convergence: § 10.2, p. 215

upper bounds:

Bayesian regret with noise: § 10.4, p. 224
worst-case regret with noise: § 10.5, p. 231
Bayesian regret without noise: § 10.6, p. 236
worst-case regret without noise: § 10.6, p. 237

this argument is carefully laid out for Bayesian regret with noise in § 10.4, p. 225

THEORETICAL ANALYSIS

information capacity: § 10.3, p. 222

bounding the posterior standard deviation
with exact observations: § 10.6, p. 235

lower bounds:

Bayesian regret with noise: § 10.4, p. 229

worst-case regret with noise, § 10.5, p. 234

Bayesian regret without noise: § 10.6, p. 236

worst-case regret without noise: § 10.6, p. 237

- In the presence of noise, a key quantity is the *information capacity*, the maximum information about the objective that can be revealed by a set of noisy observations. Bounds on this quantity yield bounds on the sum of predictive variances (10.16) and thus cumulative regret. With exact observations, we may derive bounds on credible intervals by relating the *fill distance* (10.41) of the observations to the maximum standard deviation of the process, as in (10.42).
- To derive lower bounds, we may construct explicit problem examples and prove they are difficult to optimize. For Bayesian regret, we seek distributions over objective functions that are provably difficult to optimize well. For worst-case regret, we seek objective functions in a given RKHS and that are difficult to optimize; here the “needles in haystacks” idea again proves useful.

11

EXTENSIONS AND RELATED SETTINGS

Thus far we have focused exclusively on a simple model problem (algorithm 1.1): sequential optimization of a single objective with either a fixed evaluation budget or known observation costs. These assumptions are convenient for study and often reasonable in practice, but not all optimization scenarios fit neatly into this mold. Numerous extensions of this setup have received serious attention in the literature, and we will provide an overview of the most important of these in this chapter.

A running theme throughout this discussion will be adapting the decision-theoretic framework developed in chapter 5 to derive policies for each of these new settings. In that chapter, we derived optimal optimization policies for our model problem, but we want to stress that the overarching approach to decision making can be extended to effectively any scenario.

Namely, we may derive an optimal policy for *any* sequential experimental design problem by following a simple recipe, an abstraction of our previous presentation:

1. Identify the action space \mathcal{A} of each decision.
2. Define preferences over outcomes with a utility function $u(\mathcal{D})$.
3. Identify the uncertain elements ψ relevant for each decision and determine how to compute the posterior belief given data, $p(\psi | \mathcal{D})$.
4. Compute the one-step marginal gain α_1 (5.8) for every action.¹
5. Derive the optimal policy by induction on the horizon (5.15–5.16).

decision theory for optimization: chapter 5,
p. 87

general procedure for optimal policies

¹ Conveniently, this step also yields the one-step lookahead approximate policy as a side effect.

We will sketch how this scheme can be realized for notable optimization settings not yet addressed. Each will involve additional complexity in building effective policies due to additional complexity in the problem formulation, but we will demonstrate how we can adapt our existing tools by addressing the appropriate steps of the above general procedure. We will also provide a survey of proposed approaches to each of the optimization extensions we will consider, regardless of whether they are grounded in decision theory.

11.1 UNKNOWN OBSERVATION COSTS

In our discussion on cost-aware optimization, we assumed that observation costs were prescribed by a *known* cost function c , which enabled a simple mechanism for dynamic termination through careful accounting. However, in some scenarios the cost of an observation may not be known *a priori*, but rather only determined through the course of acquisition. These unknown costs reflect additional uncertainty that must be accounted for in policy design, as the utility of the returned data depends critically on their values. The natural solution is to perform inference about observation costs throughout optimization and use the resulting beliefs to guide our decisions, just as we do with an unknown

cost-aware optimization: § 5.4, p. 103
cost function, c

objective function. To adapt the cost-aware policies from chapter 5 to this setting, we must revisit steps 3–4 of the above procedure.

Inference for unknown cost function

At a high level, reasoning about an unknown cost function given observations is no different from reasoning about an unknown objective function, and we can apply any suitable regression model for this task. The details of this inference will obviously depend on the situation, but we can outline one rather general approach.

notation for observation costs

cost function, c

observed cost at x, z
value of cost function at $x, \kappa = c(x)$

cost observation model, $p(z | x, \kappa)$

cost function prior, $p(c)$
cost function posterior, $p(c | \mathcal{D})$

predictive distribution for cost, $p(z | x, \mathcal{D})$

First let us define some notation for cost observations mirroring our notation for the objective function. Suppose that evaluation costs are determined, perhaps stochastically, by an underlying cost function $c: \mathcal{X} \rightarrow \mathbb{R}$ we wish to infer. Suppose further that an evaluation at a point x now returns both a measured value y , whose distribution depends on the corresponding objective function value $\phi = f(x)$, and an observation cost z , whose distribution depends on the corresponding cost function value $\kappa = c(x)$. We will accumulate these values throughout optimization in an augmented dataset of observations and their costs, $\mathcal{D} = (\mathbf{x}, \mathbf{y}, \mathbf{z})$.

If we can safely model cost observations as conditionally independent of objective observations given the chosen location, then we may follow our modeling strategy for the objective function and assume that each observed cost is generated by an observation model $p(z | x, \kappa)$. This allows for modeling a wide range of different scenarios, including nondeterministic costs.² Now we can proceed with inference about the cost function as usual. We choose a suitable (possibly parametric) prior process $p(c)$, which we condition on observed costs to form the posterior $p(c | \mathcal{D})$. Finally, we can form the predictive distribution for the cost of making an observation at an arbitrary point x by marginalizing the latent cost function:

$$p(z | x, \mathcal{D}) = \int p(z | x, \kappa) p(\kappa | x, \mathcal{D}) d\kappa.$$

In some applications, observation costs may be nontrivially correlated with the objective function. As an extreme example, consider a common problem in *algorithm configuration*,³ where the goal is to design the parameters of a complex algorithm so as to minimize its expected running time. Here the cost of evaluating a proposed configuration might be reasonably defined to be proportional to its running time. Up to scaling, the observation cost is precisely equal to the objective! To model such correlations, we could define a joint prior $p(f, c)$ over the cost and objective functions, as well as a joint observation model $p(y, z | x, \phi, \kappa)$. We could then continue as normal, computing expected utilities with respect to the joint predictive distribution

$$p(y, z | x, \mathcal{D}) = \iint p(y, z | x, \phi, \kappa) p(\phi, \kappa | x, \mathcal{D}) d\phi d\kappa,$$

a setup offering considerable flexibility in modeling.

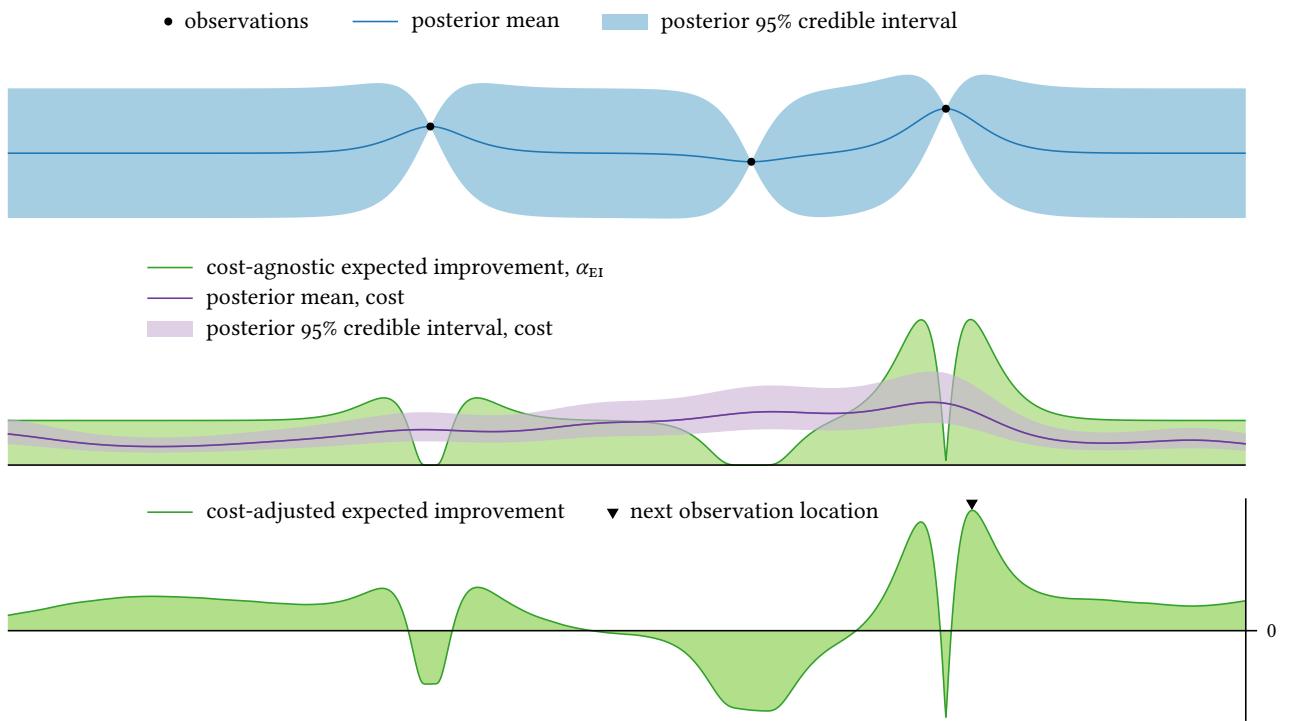


Figure 11.1: Decision making with uncertain costs. The middle panel shows the cost-agnostic expected improvement acquisition function along with a belief about an uncertain cost function, here assumed to be independent of the objective. The bottom panel shows the cost-adjusted expected improvement, marginalizing uncertainty in the objective and cost function (11.1).

Decision making with unknown costs

The approach outlined above suffices to maintain a belief about the potential cost of observations proposed throughout optimization, but we still must account for this uncertainty in the optimization policy. Thankfully, this is relatively straightforward in our decision-theoretic framework: we simply compute expected utility accounting for all relevant uncertainty as usual, here to include cost. Given an arbitrary dataset \mathcal{D} , now augmented with observation costs, consider the one-step expected marginal gain in utility:

$$\alpha_1(x; \mathcal{D}) = \mathbb{E}[u(\mathcal{D}_1) | x, \mathcal{D}] - u(\mathcal{D}).$$

Computing this expectation now requires integrating over both the unknown measurement y and the unknown observation cost z :

$$\mathbb{E}[u(\mathcal{D}_1) | x, \mathcal{D}] = \iint u(\mathcal{D} \cup \{x, y, z\}) p(y, z | x, \mathcal{D}) dy dz. \quad (11.1)$$

Although there is some slight added complexity in this computation, the optimal policy otherwise remains exactly as derived in (5.15–5.17).⁴

⁴ Of course, all nested expectations must also be taken with respect to unknown observation costs!

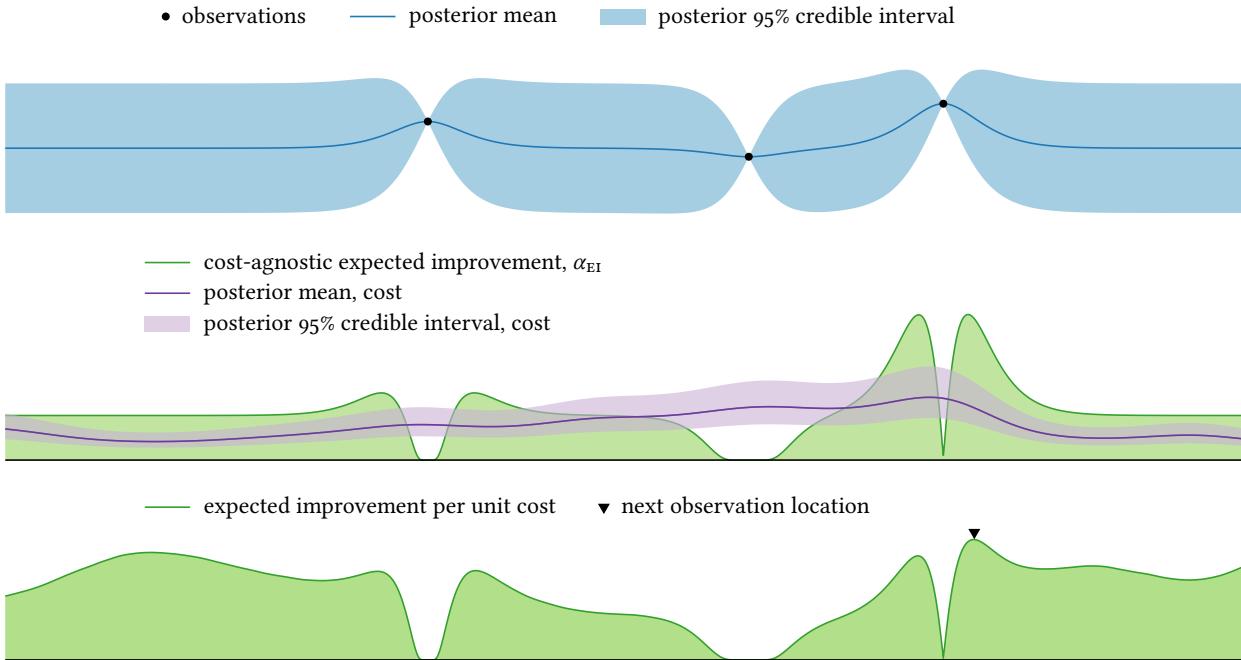


Figure 11.2: Expected gain per unit cost. The middle panel shows the cost-agnostic expected improvement acquisition function along with a belief about an uncertain cost function, here assumed to be independent of the objective. The bottom panel shows the expected improvement per unit cost (11.1).

Figure 11.1 illustrates this policy, combining expected improvement with an independent uncertain cost function.

Expected gain per unit cost

Other approaches to dealing with unknown costs are also of course possible. SNOEK et al. proposed one notable option that has gained some popularity based on a heuristic common in *anytime algorithms*⁵ – that is, algorithms that seek to maximize the instantaneous *rate* of improvement under the premise that the procedure may be terminated at any time.⁶ The idea is simple: we can approximate the expected immediate marginal gain in utility *per unit cost* from an observation at x by

$$\mathbb{E}\left[\frac{u(\mathcal{D}_1) - u(\mathcal{D})}{\kappa} \mid x, \mathcal{D}\right] \approx \frac{\alpha_1(x; \mathcal{D})}{\mathbb{E}[\kappa \mid x, \mathcal{D}]}, \quad (11.2)$$

where α_1 is the expected gain in data utility. This is a first-order approximation to the expected gain-per-cost (which is not the ratio of their respective expectations, even in the independent case) that could be further refined if desired,⁷ but works well in practice. The motivating example for SNOEK et al. was hyperparameter tuning of machine-learning algorithms with unknown training costs, and the simple heuristic of maximizing “expected improvement per (expected) second” delivered

⁵ S. ZILBERSTEIN (1996). Using Anytime Algorithms in Intelligent Systems. *AI Magazine* 17(3):73–83.

⁶ J. SNOEK et al. (2012). Practical Bayesian Optimization of Machine Learning Algorithms. *NeurIPS 2012*.

⁷ Assuming independence between the objective and cost, a second-order expansion is just as easy to compute:

$$\frac{\alpha_1(x; \mathcal{D})}{\mathbb{E}[\kappa \mid x, \mathcal{D}]} + \frac{\alpha_1(x; \mathcal{D}) \text{ var}[\kappa \mid x, \mathcal{D}]}{\mathbb{E}[\kappa \mid x, \mathcal{D}]^3}.$$

promising results in their experiments. This heuristic has since appeared in other contexts.⁸

Figure 11.2 illustrates this policy with the sample example in figure 11.2. The chosen decision closely matches the decision reached in figure 11.1. It is interesting to compare the behavior of the two acquisition functions on both sides of the domain: whereas these regions are not especially exciting in the additive cost approach, they are appealing from the anytime view – although they are expected to give only modest improvement, they are also inexpensive.

⁸ G. MALKOMES et al. (2016). Bayesian optimization for automated model selection. *NeurIPS 2016*.

11.2 CONSTRAINED OPTIMIZATION AND UNKNOWN CONSTRAINTS

Many optimization problems feature *constraints* restricting allowable solutions in a potentially complex, even uncertain manner. To this end we may extend our running optimization problem (1.1) to incorporate arbitrary constraints; a common formulation is:

$$x^* \in \arg \max_{x \in \mathcal{X}} f(x) \quad \text{subject to} \quad \forall i: g_i(x) \leq 0,$$

where the functions $\{g_i\}: \mathcal{X} \rightarrow \mathbb{R}$ comprise a set of inequality constraints. The subset of the domain where all constraints are satisfied is known as the *feasible region*:

$$\mathcal{F} = \{x \in \mathcal{X} \mid \forall i: g_i(x) \leq 0\}.$$

inequality constraints, $\{g_i\}$

feasible region, \mathcal{F}

uncertain constraints

In some situations, the value of some or all of the constraint functions may in fact be *unknown* a priori and only revealed through experimentation, complicating policy design considerably. As an example, consider a business optimizing the parameters of a service to maximize revenue, subject to constraints on customer response. If customer response is measured experimentally – for example via a focus group – we cannot know the feasibility of a proposed solution until after the objective has been measured, and even then only with limited confidence.

Further, even if the constraint functions can be computed exactly on demand, constrained optimization of an uncertain objective function is not entirely straightforward. In particular, an observation of the objective at an *infeasible* point may yield useful information regarding behavior on the feasible region, and could represent an optimal decision if that information were compelling enough. Thus simply restricting the action space to the feasible region may not be the best approach to policy design.⁹ Instead, to derive effective policies for constrained optimization, we must reconsider steps 2–4 of our general approach.

⁹ R. B. GRAMACY and H. K. H. LEE (2011). Optimization Under Unknown Constraints. In: *Bayesian Statistics 9*.

Modeling constraint functions

To allow for uncertain constraint functions, we begin by modeling the joint observation process of the objective and constraint functions. As with an uncertain cost function, we will assume that each observation of

the objective is accompanied by some information regarding constraint satisfaction at the chosen location. Modeling this process is trivial when the constraint functions are known *a priori*, but otherwise may require some care.

It is difficult to provide concrete advice as information about constraint satisfaction may assume different forms, ranging from exact observation of the constraint functions to mere *binary* indicators of constraint satisfaction. In some situations, we may even face stochastic constraint processes where the feasibility of a given location is only achieved with some unknown probability. Fortunately, our discussion in the previous section regarding modeling an uncertain cost function is general enough to handle any of these situations by choosing an appropriate joint prior processes and observation model for the objective and constraint functions.

Defining a utility function

Next we must define a utility function appropriate for constrained optimization. Most of the utility functions in chapter 6 can be suitably modified to this end.

simple reward: § 6.1, p. 112

We can realize a constrained version of simple reward by considering the expected utility of a risk-neutral terminal recommendation following optimization. As before, the resulting utility will depend on the action space used for the terminal decision. Perhaps the simplest option would be to limit the recommendation to the *feasible* observed locations (if known!):

$$u(\mathcal{D}) = \max_{x \in \mathbf{x} \cap \mathcal{F}} \mu_{\mathcal{D}}(x), \quad (11.3)$$

resulting in a natural adaptation of the simple reward (6.3). If the entire feasible region is known, we might instead allow recommending any feasible point, giving rise to an adaptation of the global simple reward (6.5):

$$u(\mathcal{D}) = \max_{x \in \mathcal{F}} \mu_{\mathcal{D}}(x). \quad (11.4)$$

Finally, in the case of uncertain constraints, we might limit our recommendation to those points believed to be feasible with sufficient confidence:¹⁰

$$u(\mathcal{D}) = \max_{x \in \mathcal{F}(\delta)} \mu_{\mathcal{D}}(x); \quad \mathcal{F}(\delta) = \{x \mid \Pr(x \in \mathcal{F} \mid \mathcal{D}) \geq 1 - \delta\}. \quad (11.5)$$

With some care, we could also modify other utility functions to be aware of a (possibly uncertain) feasible region, although the variations on simple reward above have received the most attention.

Deriving a policy

After selecting a model for the constraint functions and a utility function for our observations, we can derive a policy for constrained optimization following the standard procedure of induction on the decision horizon.

¹⁰ M. A. GELBART et al. (2014). Bayesian Optimization with Unknown Constraints. *UAI 2014*.

A policy that has received particular attention is the result of one-step lookahead with (11.3).^{10,11,12} If we assume that the constraint functions are conditionally independent of the objective function given the observation location, then the one-step expected marginal gain in utility becomes

$$\alpha(x; \mathcal{D}) = \alpha'_{\text{EI}}(x; \mathcal{D}, \mu^*) \Pr(x \in \mathcal{F} | x, \mathcal{D}). \quad (11.6)$$

This is simply the expected improvement, measured with respect to the *feasible incumbent* value $\mu^* = u(\mathcal{D})$ (7.21, 11.3), weighted by the probability of feasibility, a natural policy we might arrive at via purely heuristic arguments. GELBART et al. point out this acquisition function has a slight pathology (also present with unconstrained expected improvement): the utility degenerates when no feasible observations are available, and (11.6) becomes ill-defined. In this case, the authors propose simply maximizing the probability of feasibility:

$$\alpha(x; \mathcal{D}) = \Pr(x \in \mathcal{F} | x, \mathcal{D}).$$

The expected feasible improvement (11.6) encodes a strong preference for evaluating on the feasible region only, and in the case where the constraint functions are all known, the resulting policy will *never* evaluate outside the feasible region.¹³ This is a natural consequence of the one-step nature of the acquisition function: an infeasible observation cannot yield any immediate improvement to the utility (11.3) and thus cannot be one-step optimal. However, this behavior might be seen as undesirable given our previous comment that infeasible observations may yield valuable information about the objective on the feasible region.

If we wish to realize a policy more open to observing outside the feasible region, there are several paths forward. A less-myopic policy built on the same utility (11.3) is one option; even two-step lookahead could elect to obtain an infeasible measurement. Another possibility is one-step lookahead with a more broadly defined utility such as (11.4–11.5), which can see the merit of infeasible observations through more global evaluation of success.

To encourage infeasible observations when prudent, GRAMACY and LEE proposed a score they called the *integrated expected conditional improvement*:¹⁴

$$\iint [\alpha_{\text{EI}}(x'; \mathcal{D}) - \alpha_{\text{EI}}(x'; \mathcal{D}')] \Pr(x' \in \mathcal{F} | x', \mathcal{D}) p(y | x, \mathcal{D}) dx' dy,$$

where \mathcal{D}' is the putative updated dataset and the location x' is integrated over the domain. This is a measure of the expected impact of a measurement on the entire acquisition surface over the feasible region, which can effectively capture the potential impact of an infeasible observation when it is useful. A similar approach was taken by PICHENY, who integrated the change in probability of improvement against the feasibility probability.¹⁵ Although these approaches can be heuristically motivated, the required integrals over the acquisition surfaces are intractable, and no obvious approximations are available beyond standard methods.

¹¹ M. SCHONLAU et al. (1998). Global versus Local Search in Constrained Optimization of Computer Models. In: *New Developments and Applications in Experimental Design*.

¹² J. R. GARDNER et al. (2014). Bayesian Optimization with Inequality Constraints. *ICML 2014*.

observations in the infeasible region

¹³ In this case the policy would be equivalent to redefining the domain to encompass only the feasible region \mathcal{F} and maximizing the unmodified expected improvement (7.2).

¹⁴ R. B. GRAMACY and H. K. H. LEE (2011). Optimization Under Unknown Constraints. In: *Bayesian Statistics 9*.

¹⁵ V. PICHENY (2014). A Stepwise uncertainty reduction approach to constrained global optimization. *AISTATS 2014*.

decoupled observations

¹⁶ M. A. GELBART et al. (2014). Bayesian Optimization with Unknown Constraints. *UAI 2014*.

¹⁷ J. M. HERNÁNDEZ-LOBATO et al. (2016b). A General Framework for Constrained Bayesian Optimization using Information-based Search. *Journal of Machine Learning Research* 17:1–53.

predictive entropy search: § 8.8, p. 180

synchronous vs. asynchronous batch construction

asynchronous batch observations: § 11.4,
p. 260batch of observation locations, \mathbf{x}
corresponding observed values, \mathbf{y} action space for batch observations, $\mathcal{A} = \mathcal{X}^b$ expected one-step marginal gain from batch
observation, β_1

GELBART et al. considered a variation on the constrained optimization problem discussed above wherein observations of the objective and constraints can be “decoupled” – that is, when we can elect to measure any of these functions independent of the others, expanding the action space of the decision problem¹⁶ The authors noted the expected feasible improvement (11.6) displayed undesirable behavior in this scenario and proposed an alternative policy based on mutual information, which was later refined and expanded into a fully fledged information theoretic policy for constrained optimization (with or without decoupled observations) based on predictive entropy search.¹⁷

11.3 SYNCHRONOUS BATCH OBSERVATIONS

Many optimization settings allow for the possibility of making multiple observations in parallel. In fact, some settings such as high-throughput screening for scientific discovery practically *demand* parallel experiments due to the growing capacity of sophisticated automated instruments. Numerous batch policies have been proposed for Bayesian optimization to harness this capability, including variants of virtually every popular sequential policy.

Here we can distinguish two settings: *synchronous* and *asynchronous* batch construction. In both cases, multiple experiments must be designed to run in parallel. The distinguishing factor is that in the synchronous case, the results from each entire batch of experiments are obtained before designing the next, whereas in the asynchronous case, each time an experiment completes, we may immediately design a new one in light of those still pending. Bayesian decision theory naturally offers one possible approach to both of these scenarios. Here we will focus on the synchronous case here and discuss the asynchronous case in the next section.

Decision-theoretic batch construction

Consider an optimization scenario where in each iteration we may design a batch of b points $\mathbf{x} = \{x_1, x_2, \dots, x_b\}$ for simultaneous evaluation, resulting in a corresponding vector of measured values \mathbf{y} obtained before our next action. The design of each batch represents a decision with action space $\mathcal{A} = \mathcal{X}^b$, a modification to step 3 of our general procedure.

We proceed by computing the one-step expected gain in utility from a proposed batch measurement \mathbf{x} . We will call the corresponding batch acquisition function β_1 to distinguish it from its sequential analog α_1 :

$$\beta_1(\mathbf{x}; \mathcal{D}) = \mathbb{E}[u(\mathcal{D}_1) | \mathbf{x}, \mathcal{D}] - u(\mathcal{D}).$$

Here \mathcal{D}_1 represents the data available after the batch observation is resolved: $\mathcal{D}_1 = \mathcal{D} \cup \{(\mathbf{x}, \mathbf{y})\}$. Computing this expected marginal gain is an expectation with respect to the unknown values \mathbf{y} :

$$\beta_1(\mathbf{x}; \mathcal{D}) + u(\mathcal{D}) = \int u(\mathcal{D}_1) p(\mathbf{y} | \mathbf{x}, \mathcal{D}) d\mathbf{y}, \quad (11.7)$$

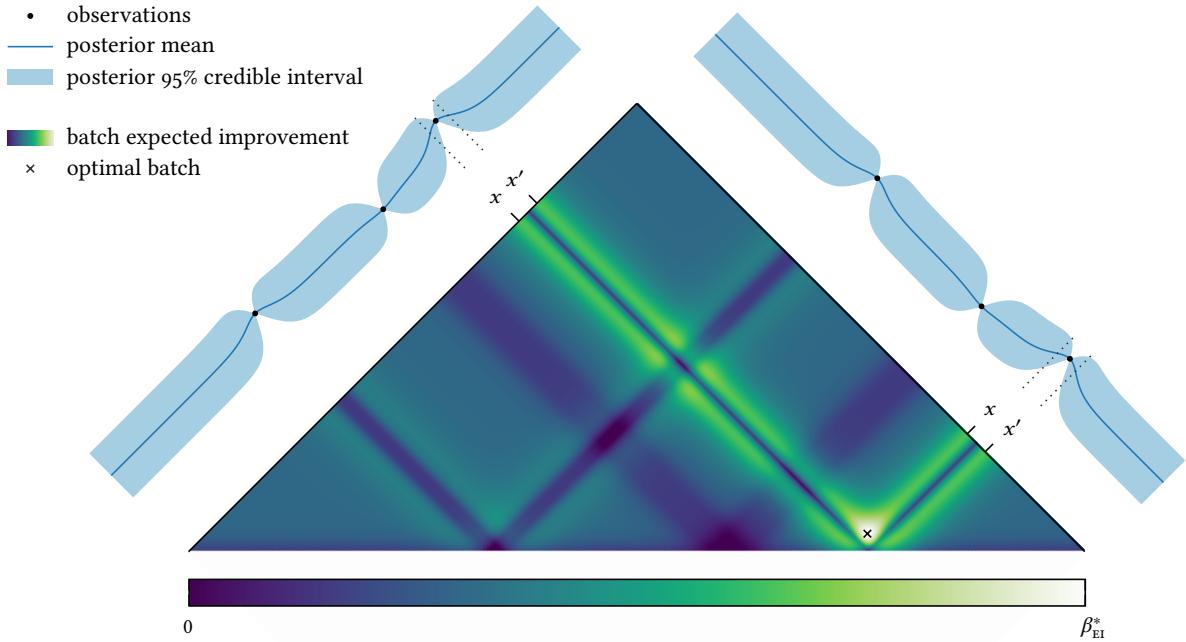


Figure 11.3: Optimal batch selection. The heatmap shows the expected one-step marginal gain in simple reward (6.1) from adding a batch of two points – corresponding in location to the belief about the objective plotted along the margins – to the current dataset. Note that the expected marginal gain is symmetric. The optimal batch will observe on both sides of the previously best-seen point. In this example, incorporating either one of the selected points alone would also yield relatively high expected marginal gain.

and an optimal batch with decision horizon $\tau = 1$ maximizes this score:

$$\mathbf{x} \in \arg \max_{\mathbf{x}' \in \mathcal{X}^b} \beta_1(\mathbf{x}'; \mathcal{D}). \quad (11.8)$$

Finally, we can derive the optimal batch policy for a fixed evaluation budget by induction on the horizon, accounting for the expanded action space for each future decision:

$$\begin{aligned} \mathbf{x} &\in \arg \max_{\mathbf{x}' \in \mathcal{X}^b} \underbrace{\beta_1(\mathbf{x}'; \mathcal{D}) + \mathbb{E}[\beta_{\tau-1}^*(\mathcal{D}_1) | \mathbf{x}', \mathcal{D}]}_{= \beta_\tau(\mathbf{x}'; \mathcal{D})}; \\ \beta_\tau^*(\mathcal{D}) &= \max_{\mathbf{x}' \in \mathcal{X}^b} \beta_\tau(\mathbf{x}'; \mathcal{D}). \end{aligned}$$

If desired, we could also allow for variable-cost observations and the option of dynamic termination by accounting for costs and including a termination option in the action space. Another compelling possibility would be to consider *dynamic* batch sizes by expanding the action space further and assigning an appropriate size-dependent cost function for proposed batch observations.

Optimal batch selection is illustrated in figure 11.3 for designing a batch of two points with horizon $\tau = 1$. We compute the expected one-step gain in utility (11.7) – here the simple reward (6.1), analogous to

optimal batch policy with fixed evaluation budget

variable costs and termination option

dynamic batch sizes

example of optimal batch policy

expected improvement α_{EI} (7.2) – for every possible batch and observe where the score is maximized. The optimal batch evaluates on either side of the previously best-seen point, achieving distributed exploitation. The expected marginal gain surface has notably complex structure, for example expressing a strong preference for batches containing at least one of the chosen locations over any purely exploratory alternative, as well as severely punishing batches containing an observation too close to an existing one.

connection to b -step lookahead

We may gain some insight into the optimal batch policy by decomposing the expected batch marginal gain in terms of corresponding quantities from the optimal sequential policy. Let us first consider the expected marginal gain from selecting a batch of two points, $\mathbf{x} = \{x, x'\}$, resulting in observed values $\mathbf{y} = \{y, y'\}$. Let \mathcal{D}' represent the current data augmented with the *single* observation (x, y) . We may rewrite the marginal gain from the batch observation (\mathbf{x}, \mathbf{y}) as a telescoping sum with terms corresponding to the impact of each individual observation:

$$u(\mathcal{D}_1) - u(\mathcal{D}) = [u(\mathcal{D}_1) - u(\mathcal{D}')] + [u(\mathcal{D}') - u(\mathcal{D})],$$

which allows us to rewrite the one-step expected batch marginal gain as:

$$\beta_1(\mathbf{x}; \mathcal{D}) = \alpha_1(x; \mathcal{D}) + \mathbb{E}[\alpha_1(x'; \mathcal{D}') | \mathbf{x}, \mathcal{D}].$$

This expression is remarkably similar to the optimal two-step expected sequential marginal gain (5.12):

$$\alpha_2(x; \mathcal{D}) = \alpha_1(x; \mathcal{D}) + \mathbb{E}\left[\max_{x' \in \mathcal{X}} \alpha_1(x'; \mathcal{D}') | x, \mathcal{D}\right].$$

The main difference is that in the batch setting, we must commit to both observation locations *a priori*, whereas in the sequential setting, we can design our second observation optimally given the outcome of the first.

We can extend this relationship to the general case. Temporarily adopting compact notation, a horizon- b optimal sequential decision satisfies:

$$x \in \arg \max \left\{ \alpha_1 + \mathbb{E} \left[\max \left\{ \alpha_1 + \mathbb{E} \left[\max \left\{ \alpha_1 + \cdots \right\} \right] \right\} \right] \right\}, \quad (11.9)$$

and the optimal one-step batch of size b satisfies:

$$\mathbf{x} \in \arg \max \left\{ \alpha_1 + \mathbb{E} \left[\alpha_1 + \mathbb{E} \left[\alpha_1 + \cdots \right] \right] \right\}. \quad (11.10)$$

Clearly the expected utility gained from making b optimal sequential decisions surpasses the expected utility from a single optimal batch the same size: the sequential policy benefits from designing each successive observation location adaptively, whereas the batch policy must make all decisions simultaneously and cannot benefit from replanning. This unavoidable difference in performance is called the *adaptivity gap* in the analysis of algorithms.¹⁸

Unfortunately, working with the larger action space inherent to batch optimization requires significant computational effort. First, computing

¹⁸ J. VONDRAK (2005). Probabilistic Methods in Combinatorial and Stochastic Optimization. Ph.D. thesis. Massachusetts Institute of Technology.

adaptivity gap

```

input: dataset  $\mathcal{D}$ , batch size  $b$ , acquisition function  $\alpha$ 
 $\mathcal{D}' \leftarrow \mathcal{D}$                                 ▶ initialize fictitious dataset
for  $i = 1 \dots b$  do
     $x_i \leftarrow \arg \max_{x \in \mathcal{X}} \alpha(x; \mathcal{D}')$     ▶ select the next batch member
     $\hat{y}_i \leftarrow \text{SIMULATE-OBSERVATION}(x_i, \mathcal{D}')$ 
     $\mathcal{D}' \leftarrow \mathcal{D}' \cup \{(x_i, \hat{y}_i)\}$            ▶ update fictitious dataset
end for
return  $\mathbf{x}$ 

```

Algorithm 11.1: Sequential simulation.

the expected marginal gain (11.7) is more expensive than in the sequential setting (5.8), as we must now integrate with respect to the *joint* distribution over outcomes $p(\mathbf{y} \mid \mathbf{x}, \mathcal{D})$. Thus even evaluating the acquisition function at a single point entails more work. Additionally, finding the optimal decision (11.8) requires optimizing this score over a significantly larger domain than in the sequential analog, a nontrivial task due to its potentially complex and multimodal nature – see figure 11.3.

Despite these computational difficulties, synchronous batch Bayesian optimization has enjoyed significant attention from the research community. We can identify two recurring research thrusts: deriving general strategies for extending arbitrary sequential policies to batch policies and deriving batch extensions of specific sequential policies. We provide a brief survey below.

Batch construction via sequential simulation

Sequential simulation is an efficient strategy for creating batch policies by simulating multiple steps of an existing sequential policy. Pseudocode for this procedure is listed in Algorithm 11.1. Given a sequential acquisition function α , we choose the first batch member by maximization:

$$x_1 \in \arg \max_{x \in \mathcal{X}} \alpha(x; \mathcal{D}),$$

and commit to this choice. We now augment our dataset with the chosen point and a *fictitious* observed value, forming $\mathcal{D}_1 = \mathcal{D} \cup \{(x_1, \hat{y}_1)\}$, then maximize the acquisition function again to choose the second point:

$$x_2 \in \arg \max_{x \in \mathcal{X}} \alpha(x; \mathcal{D}_1).$$

We proceed in this manner until the desired batch size has been reached. Sequential simulation entails b optimization problems on \mathcal{X} rather than a single problem on \mathcal{X}^b , which can be a significant computational savings.

When the sequential policy represents one-step lookahead, sequential simulation can be regarded as a natural greedy approximation to the one-step batch policy via the decomposition in (11.10): whereas the optimal batch policy must maximize this score *jointly*, sequential simulation maximizes the score *pointwise*, fixing each point once chosen.

This procedure requires some mechanism for generating fictitious observations as we build the batch. GINSBOURGER et al. described two

first batch member, x_1

fictitious observation at x_1 , \hat{y}_1
second batch member, x_2

greedy approximation of one-step lookahead

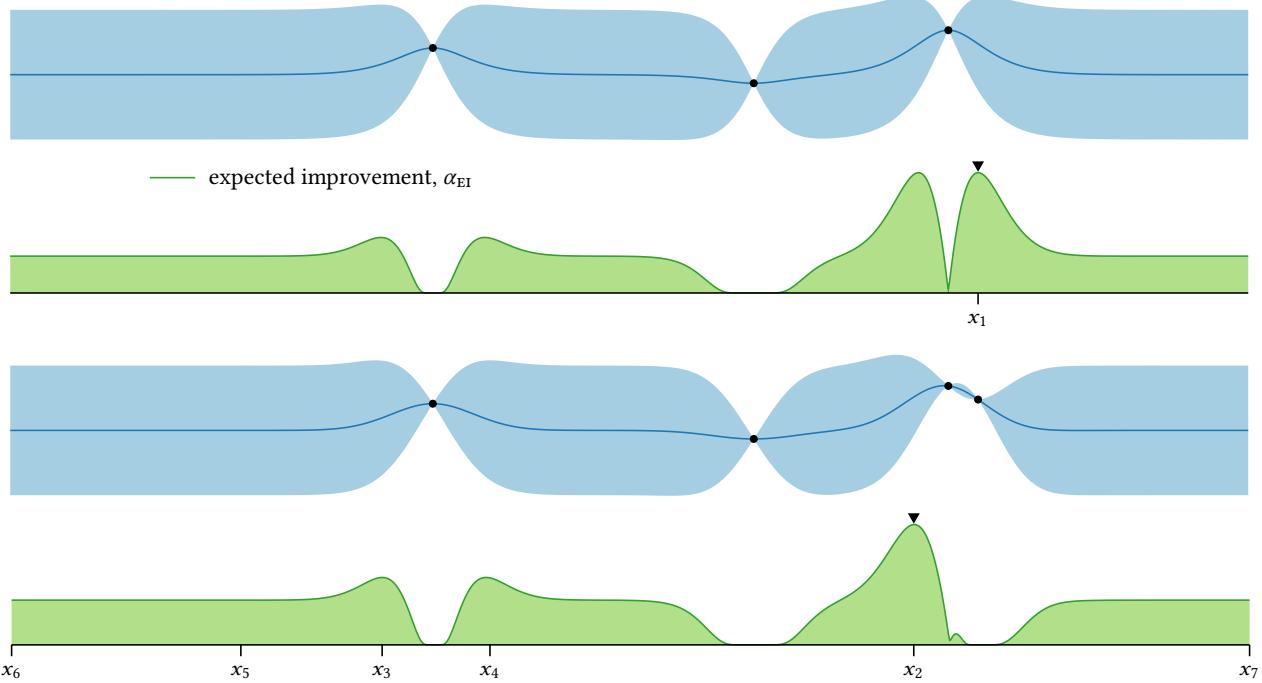


Figure 11.4: Sequential simulation using the expected improvement (7.2) policy and the kriging believer (11.11) imputation strategy. The first point is selected by maximizing expected improvement, and we condition the model on an observation equal to the posterior mean (top panel). We then maximize the updated expected improvement, condition, and repeat as desired. The bottom panel indicates the locations of several further points selected in this manner.

¹⁹ D. GINSBOURGER et al. (2010). Kriging is well-suited to parallelize optimization. In: *Computational Intelligence in Expensive Optimization Problems*.

simple heuristics that have been widely adopted.¹⁹ Perhaps the most natural option is to impute the expected value of each observation, a heuristic GINSBOURGER et al. dubbed the *kriging believer* strategy:

$$\hat{y} = \mathbb{E}[y | x, \mathcal{D}]. \quad (11.11)$$

kriging believer heuristic

constant liar heuristic

example and discussion

This has the effect of fixing the posterior mean of the objective function throughout simulation. An even simpler option is to impute a constant value independent of the chosen point, which the authors called the *constant liar* strategy:

$$\hat{y} = c. \quad (11.12)$$

Although this might seem silly, this has the advantage of being model independent and has demonstrated surprisingly good performance in practice. Three natural options for the constant, ranging from the most optimistic to most pessimistic, are to impute the maximum, mean, or minimum of the known observed values y .

Seven steps of sequential simulation with the expected improvement acquisition function (7.2) and the kriging believer strategy (11.11) are demonstrated in figure 11.4. The selected points appear reasonable: the

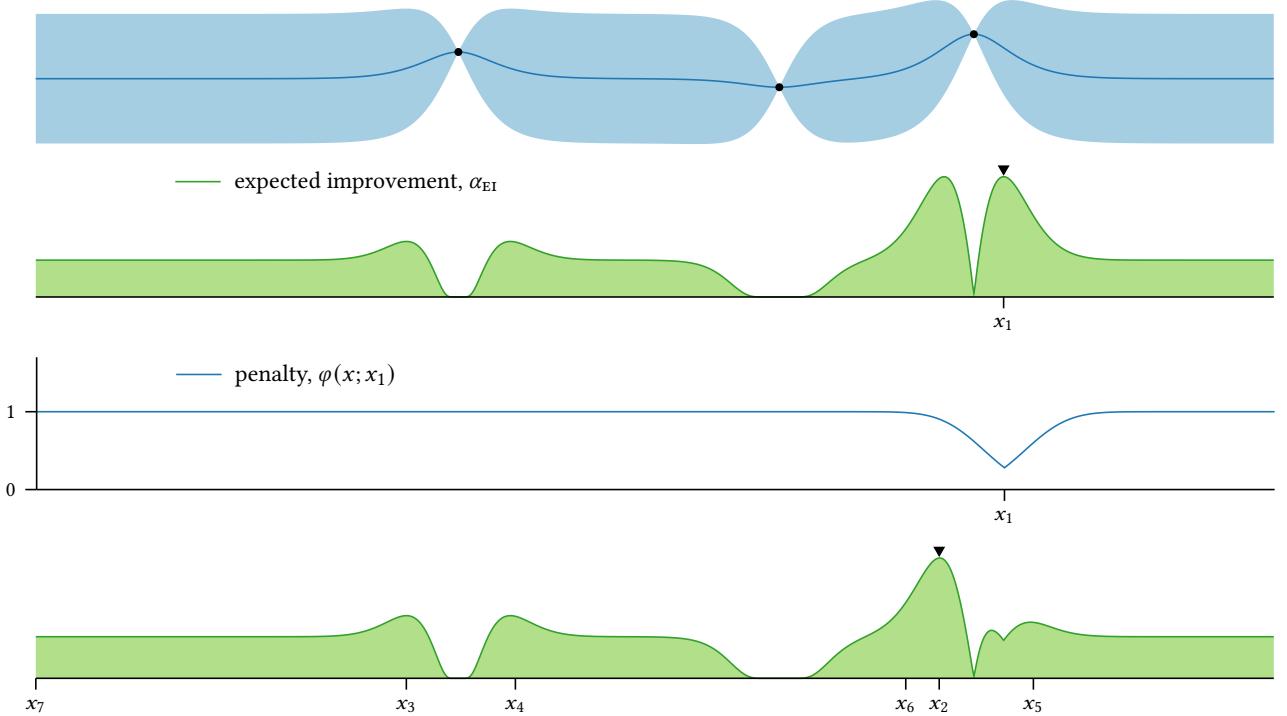


Figure 11.5: Batch construction via local penalization of the expected improvement acquisition function (7.2). We select the first point by maximizing the expected improvement, after which the acquisition function is multiplied by a penalty factor discouraging future batch members in that area. We then maximize the updated acquisition function, penalize, and repeat as desired. The bottom panel indicates the locations of several further points selected in this manner.

first two exploit the best-seen observation (and are near optimal for a batch size of two; see figure 11.3), the next two exploit another local optimum, and the remainder explore the domain.

Batch construction via local penalization

GONZÁLEZ et al. proposed another general mechanism for extending a given sequential policy (defined by the acquisition function α) to the batch setting.²⁰ Like sequential simulation, we select the first point by maximizing the acquisition function:

$$x_1 \in \arg \max_{x \in \mathcal{X}} \alpha(x; \mathcal{D}).$$

We then incorporate a multiplicative penalty $\varphi(x; x_1)$ into the acquisition function discouraging future batch members from being in a neighborhood of the initial point. This penalty is designed to avoid redundancy between batch members without disrupting the differentiability of the original acquisition function. GONZÁLEZ et al. describe one simple and effective penalty function from an estimate of the global maximum and

²⁰ J. GONZÁLEZ et al. (2016a). Batch Bayesian Optimization via Local Penalization. *AISTATS 2016*.

penalty from selecting x_1 , $\varphi(x; x_1)$

²¹ Any continually differentiable function on a compact set is Lipschitz continuous.

Lipschitz constant of the objective function.²¹ We now select the second batch member by maximizing the penalized acquisition function

$$x_2 \in \arg \max_{x \in \mathcal{X}} \alpha(x; \mathcal{D}) \varphi(x; x_1),$$

after which we apply another penalty and continue in this manner as desired. This process is illustrated in figure 11.5. Comparing with sequential simulation, the first batch members are very similar; however, there is some divergence in the final stages, with local penalization preferring to revisit the local optimum on the right.

Like sequential simulation, local penalization also entails b optimization problems on \mathcal{X} and is in fact even faster than sequential simulation, as the objective function model does not need to be updated along the way.

Approximation via Monte Carlo integration

approximate computation for sequential one-step lookahead: § 8.5, p. 171

If we wish to proceed via joint optimization of (11.7) rather than one of the above heuristics, we will often face the complication that the expectation with respect to the observed values y is intractable. We can offer some advice for approximating this quantity and its gradient for a Gaussian process model coupled with an exact or additive Gaussian noise observation model. This abbreviated discussion will closely follow the approach for the sequential case.

First we write the one-step marginal gain (11.7) as an expectation of the marginal gain in utility $\Delta(\mathbf{x}, \mathbf{y}) = u(\mathcal{D}') - u(\mathcal{D})$ with respect to a multivariate normal belief on the observations (2.20):

$$\beta(\mathbf{x}; \mathcal{D}) = \int \Delta(\mathbf{x}, \mathbf{y}) \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}, \mathbf{S}) d\mathbf{y}. \quad (11.13)$$

²² Gauss–Hermite quadrature, as we recommended in the one-dimensional case, does not scale well with dimension.

We may approximate this expectation via Monte Carlo integration²² by sampling from this belief. It is convenient to do so by sampling n vectors $\{\mathbf{z}_i\}$ from a *standard* normal distribution, then transforming these via

$$\mathbf{z}_i \mapsto \boldsymbol{\mu} + \boldsymbol{\Lambda} \mathbf{z}_i = \mathbf{y}_i, \quad (11.14)$$

where $\boldsymbol{\Lambda}$ is the Cholesky factor of \mathbf{S} : $\boldsymbol{\Lambda} \boldsymbol{\Lambda}^\top = \mathbf{S}$. Monte Carlo estimation now gives

$$\beta(\mathbf{x}; \mathcal{D}) \approx \frac{1}{n} \sum_{i=1}^n \Delta(\mathbf{x}, \mathbf{y}_i). \quad (11.15)$$

As in the sequential case, we may reuse these samples to approximate the gradient of the acquisition function with respect to the proposed observation locations, under mild assumptions (c.3–c.4):

$$\frac{\partial \beta}{\partial \mathbf{x}} \approx \frac{1}{n} \sum_{i=1}^n \frac{\partial \Delta}{\partial \mathbf{x}}(\mathbf{x}, \mathbf{y}_i); \quad \frac{\partial \Delta}{\partial x_j}(\mathbf{x}, \mathbf{y}_i) = \left[\frac{\partial \Delta}{\partial x_j} \right]_{\mathbf{y}} + \frac{\partial \Delta}{\partial \mathbf{y}} \left[\frac{\partial \boldsymbol{\mu}}{\partial x_j} + \frac{\partial \boldsymbol{\Lambda}}{\partial x_j} \mathbf{z}_i \right], \quad (11.16)$$

where, in the second expression, we account for the dependence of the \mathbf{y} samples on \mathbf{x} through the transformation (11.14). The gradient of the Cholesky factor $\mathbf{\Lambda}$ can be computed efficiently via automatic differentiation.²³

We will spend the remainder of this section discussing the details of explicit batch extensions of popular sequential policies.

Expected improvement and knowledge gradient

Both the expected improvement and knowledge gradient policies have been extended to the batch case, closely following the decision-theoretic approach outlined above.

Unlike its sequential counterpart, computation of batch expected improvement for Gaussian process models is rather involved, even in the case of exact observation. The primary challenge is evaluating the expectation of a multivariate truncated normal distribution, a computation whose difficulty increases rapidly with the batch size. There are exact formulas to compute batch expected improvement²⁴ and its gradient²⁵ based on the moment-generating function for the truncated multivariate normal derived by TALLIS;²⁶ however, these formulas require b evaluations of the b -dimensional and b^2 evaluations of the $(b - 1)$ -dimensional multivariate normal CDF, itself a notoriously difficult computation that can only effectively be approximated via Monte Carlo methods in dimension greater than four. This limits the utility of the direct approach to relatively small batch sizes, perhaps $b \leq 10$.

For larger batch sizes, GINSBOURGER et al. proposed sequential simulation, and found the constant liar strategy (11.12) using the “optimistic” estimate $\hat{y} = \max \mathbf{y}$ to deliver good empirical performance in simulation.²⁷ WANG et al. proposed an efficient alternative: joint optimization of batch expected improvement via multistart stochastic gradient ascent using the Monte Carlo estimators in (11.15–11.16).²⁸ The authors demonstrated this procedure scaling up to batch sizes of $b = 128$ with impressive performance and runtime compared with exact computation and sequential simulation.

A similar approach for approximating the batch knowledge gradient policy was described by WU and FRAZIER.²⁹ The overall approach is effectively the same: multistart stochastic gradient ascent relying on (11.15–11.16). An additional complication in estimating the batch knowledge gradient is the global optimization inherent to the global reward (6.5). WU and FRAZIER suggest a discretization approach where the global reward is estimated on a dynamically managed discrete set of points drawn via Thompson sampling from the objective function posterior.

The batch expected improvement acquisition function is illustrated for an example scenario in figure 11.3, and batch knowledge gradient for the same scenario in figure 11.6. In both cases, the optimal batch measures on either side of the previously best-seen point; however, knowledge gradient is more tolerant of batches containing only one observation in this region.

²³ I. MURRAY (2016). Differentiation of the Cholesky decomposition. arXiv: 1602.07527 [stat.CO].

expected improvement: § 7.3 p. 127
knowledge gradient: § 7.4 p. 129

²⁴ C. CHEVALIER and D. GINSBOURGER (2013). Fast Computation of the Multi-points Expected Improvement with Applications in Batch Selection. *LION* 7.

²⁵ S. MARMIN et al. (2015). Differentiating the Multipoint Expected Improvement for Optimal Batch Design. *MOD* 2015.

²⁶ G. M. TALLIS (1961). The Moment Generating Function of the Truncated Multi-normal Distribution. *Journal of the Royal Statistical Society Series B (Methodological)* 23(1):223–229.

²⁷ D. GINSBOURGER et al. (2010). Kriging is well-suited to parallelize optimization. In: *Computational Intelligence in Expensive Optimization Problems*.

²⁸ J. WANG et al. (2020a). Parallel Bayesian Global Optimization of Expensive Functions. *Operations Research* 68(6):1850–1865.

²⁹ J. WU and P. I. FRAZIER (2016). The Parallel Knowledge Gradient Method for Batch Bayesian Optimization. *NeurIPS* 2016.

Thompson sampling: 7.9, p. 148

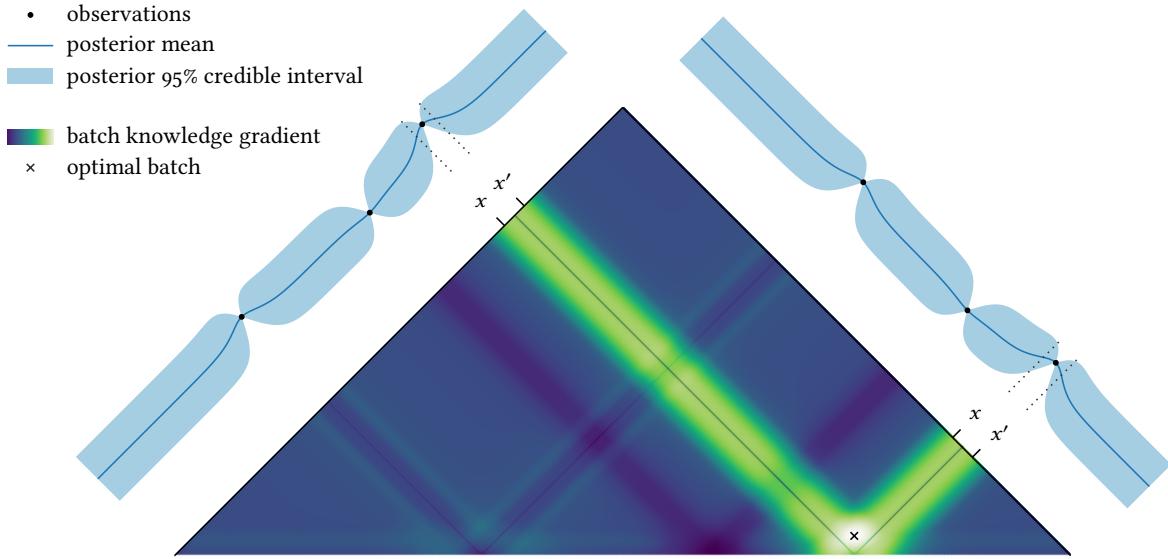


Figure 11.6: The batch knowledge gradient acquisition function for an example scenario. The optimal batch exploits the local optimum, but any batch containing at least one point in that neighborhood is near-optimal.

*Mutual information with x^**

mutual information with x^* : § 7.6, p. 139

predictive entropy search: § 8.8, p. 180

³⁰ A. SHAH and Z. GHAHRAMANI (2015). Parallel Predictive Entropy Search for Batch Global Optimization of Expensive Objective Functions. *NeurIPS 2015*.

Compared with the difficulty faced in computing (or even approximating) batch analogs of the expected improvement and knowledge gradient acquisition functions, extending the predictive entropy search acquisition function to the batch setting is relatively straightforward.³⁰ The mutual information between the observed values y and the location of the global maximum x^* is (compare with (8.36)):

$$\beta_{x^*}(\mathbf{x}; \mathcal{D}) = H[\mathbf{y} | \mathbf{x}, \mathcal{D}] - \mathbb{E}[H[\mathbf{y} | \mathbf{x}, x^*, \mathcal{D}] | \mathbf{x}, \mathcal{D}]. \quad (11.17)$$

The first term is the differential entropy of a multivariate normal and may be computed in closed form (A.16). The second term is somewhat difficult to approximate, but no innovation is required in the batch setting beyond the machinery already developed for the sequential case. We may approximate the expectation with respect to x^* via Thompson sampling and may approximate $p(y | \mathbf{x}, x^*, \mathcal{D})$ as a multivariate normal following the expectation propagation approach described previously.

Probability of improvement

probability of improvement: § 7.5, p. 131

³¹ See § 7.5, p. 134.

³² D. R. JONES (2001). A Taxonomy of Global Optimization Methods Based on Response Surfaces. *Journal of Global Optimization* 21(4):345–383.

Batch probability of improvement has received relatively little attention, but JONES proposed one simple option in the context of threshold selection.^{31,32} The idea is to find the optimal sequential decisions using a range of improvement thresholds, representing a spectrum exploration-exploitation tradeoffs. JONES then recommends a simple clustering procedure to remove redundant points, resulting in a batch (of variable size)

reflecting diversity in location and behavior. A compelling aspect of this procedure is that it is naturally *nonmyopic*, as each batch is explicitly constructed to address both immediate and long-term gain.

This approach is illustrated in figure 7.12; depending on the aggressiveness the pruning procedure, the constructed batch would contain 2–4 points chosen from the visible clusters.

Upper confidence bound

Due to the equivalence between the probability of improvement and upper confidence bound policies for Gaussian processes (8.22, 8.26), the procedure proposed by JONES described above may also be used to realize a simple batch upper confidence bound policy for that model class. In this case, we would design each batch by maximizing the upper confidence bound for a range of confidence parameters, clustering, and pruning.

Several more direct batch upper confidence bound policies have been developed, all variations on a theme. DESAUTELS et al. proposed a strategy – dubbed simply *batch upper confidence bound* (BUCB) – based on sequential simulation with the kriging believer strategy (11.11).³³ Batch diversity is automatically encouraged: each point added to the batch globally reduces the upper confidence bound, most dramatically at the locations with the most strongly correlated function values.

The BUCB algorithm was later refined by several authors to encourage more exploration, which can improve both empirical and theoretical performance. Like BUCB, we seed each batch with the maximum of the upper confidence bound (8.25). We now identify the so-called “relevant region” of the domain, defined to be the set of locations whose upper confidence bound exceeds the global maximum of the *lower* confidence bound.³⁴ The intuition behind this region is that the objective value at any point in its complement is – with high probability – lower than at the point maximizing the lower confidence bound and can thus be discarded with some confidence; see the illustration in the margin.

With the relevant region in hand, we design the remaining batch members to promote maximal information gain about the objective on this region. For a Gaussian process, this is intimately related to a diversity-encouraging distribution known as a *k-determinantal point process* (*k*-DPP) built from the posterior covariance function.³⁵ We may proceed by either a simple greedy procedure³⁶ or via more nuanced maximization or sampling using methods developed for *k*-DPPs.³⁷

These schemes are all backed by strong theoretical analysis, including sublinear cumulative regret bounds under suitable conditions.

Thompson sampling

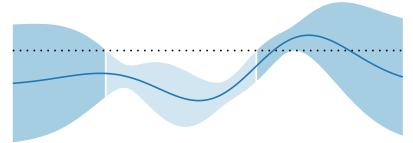
The stochastic nature of Thompson sampling enables trivial batch construction by drawing b independent samples of the location of the global maximum (7.19):

$$\{x_i\} \sim p(x^* | \mathcal{D}).$$

maximizing an upper confidence bound: § 7.8,
p. 145

33 T. DESAUTELS et al. (2014). Parallelizing Exploration–Exploitation Tradeoffs in Gaussian Process Bandit Optimization. *Journal of Machine Learning Research* 15(119):4053–4103.

34 N. DE FREITAS et al. (2012b). Exponential Regret Bounds for Gaussian Process Bandits with Deterministic Observations. *ICML* 2012.



The (disconnected) relevant region (darker blue) for an example Gaussian process. Points outside the region (lighter blue) are unlikely to maximize the objective function.

35 A. KULESZA and B. TASKAR (2012). Determinantal Point Processes for Machine Learning. *Foundations and Trends in Machine Learning* 5(2–3):123–286.

36 E. CONTAL et al. (2013). Parallel Gaussian Process Optimization with Upper Confidence Bound and Pure Exploration. *ECCML PKDD* 2013.

37 T. KATHURIA et al. (2016). Batched Gaussian Process Bandit Optimization via Determinantal Point Processes. *NeurIPS* 2016.

- 38 J. M. HERNÁNDEZ-LOBATO et al. (2017). Parallel and Distributed Thompson Sampling for Large-scale Accelerated Exploration of Chemical Space. *ICML 2017*.
- 39 K. KANDASAMY et al. (2018). Parallelised Bayesian Optimisation via Thompson Sampling. *AISTATS 2018*.

A remarkable advantage of this policy is that the samples may be generated entirely in parallel, which allows linear scaling to arbitrarily large batch sizes. Batch Thompson sampling has delivered impressive performance in a real-world setting with batch sizes up to $b = 500$,³⁸ and is backed by theoretical guarantees on the asymptotic reduction of the simple regret (10.1).³⁹

11.4 ASYNCHRONOUS OBSERVATION WITH PENDING EXPERIMENTS

Some situations allow parallel observation with *asynchronous* execution. For example, when optimizing the result of a computational simulation, access to more than one CPU core (or even better, a cluster of machines) could enable many simulations to be run in parallel. To maximize throughput, we could immediately start a new simulation upon the termination of a previous job, without waiting for the other running processes to finish. An effective optimization policy for this setting must account for the pending experiments when designing each observation.

We may consider a general case where we wish to design a batch of experiments $\mathbf{x} \in \mathcal{X}^b$ when another batch of experiments $\mathbf{x}' \in \mathcal{X}^{b'}$ is under current evaluation, where the number of running and pending experiments may be arbitrary.⁴⁰ Here the action space for the current decision is $\mathcal{A} \in \mathcal{X}^b$, and we must make the decision under uncertainty both in the observations resulting from the chosen batch, \mathbf{y} , and the observations resulting from the pending experiments, \mathbf{y}' .

The one-step expected gain in utility from a set of proposed experiments and the pending experiments is

$$\beta_1(\mathbf{x}; \mathbf{x}', \mathcal{D}) = \mathbb{E}[u(\mathcal{D}_1) | \mathbf{x}, \mathbf{x}', \mathcal{D}] - u(\mathcal{D}),$$

where $\mathcal{D}_1 = \mathcal{D} \cup \{(\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}')\}$. This entails an expectation with respect to the unknown values \mathbf{y} and \mathbf{y}' :

$$\beta_1(\mathbf{x}; \mathcal{D}) + u(\mathcal{D}) = \iint u(\mathcal{D}_1) p(\mathbf{y}, \mathbf{y}' | \mathbf{x}, \mathbf{x}', \mathcal{D}) d\mathbf{y} d\mathbf{y}'.$$

reduction to synchronous case

This is simply the one-step marginal gain for the combined batch $\mathbf{x} \cup \mathbf{x}'$ from the synchronous case (11.7)! The only difference with respect to one-step lookahead is that we can only maximize this score with respect to \mathbf{x} as we are already committed to the pending experiments. Thus a one-step lookahead policy for the asynchronous case can be reduced to maximizing the corresponding score from the synchronous case with some batch members fixed. This reduction has been pointed out by numerous authors, and effectively every batch policy discussed above may be modified with little effort to work in the asynchronous case.

Moving beyond one-step lookahead may be extremely challenging, however, due to the implications of uncertainty in the order of termination for pending experiments. A full treatment would require a model for the time of termination and accounting for how the decision tree may branch after the present decision. Exact computation of even two-step

lookahead is likely intractable in most practical situations, but rollout might offer one path forward.

rollout: § 5.3, p. 102

11.5 MULTIFIDELITY OPTIMIZATION

In Bayesian optimization, we typically assume that observations of the objective function are expensive and should be made as sparingly as possible. However, some scenarios offer a potential shortcut: *indirect* inspection of the system of interest via a cheaper surrogate, such as the output of a computer simulation. In some cases, we may even have access to multiple surrogates of varying cost and fidelity. It is tempting to try to accelerate optimization using these surrogates to guide the search. This is the inspiration for *multipidelity optimization*, a cost-aware extension of optimization that has received significant attention.

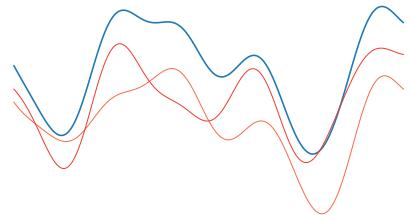
As a motivating example, consider a problem from materials science where we wish to optimize the properties of a material as a function of its composition, process parameters, etc. Materials scientists have several mechanisms available⁴¹ to investigate a proposed material, ranging from relatively inexpensive computer simulations (molecular dynamics, density functional theory, etc.) to extremely expensive synthesis and characterization in a laboratory. State-of-the-art materials discovery pipelines rely on these computational surrogates to winnow the search space for experimental campaigns.

Automated machine learning provides another motivating example. Suppose we wish to tune the hyperparameters of a model by minimizing validation error after training on a large training set. Although training on the full dataset may be costly, we can estimate performance by training on only a subset of the data⁴² or by terminating the training procedure early.⁴³ We may reasonably hope to accelerate hyperparameter tuning by exploiting these noisy, but cost-effective surrogates.

We will outline a decision-theoretic approach to multifidelity optimization below. The complexity of this setting will require readdressing every step of the procedure outlined at top of this chapter. We will focus on the first three steps, as deriving the optimal policy is mechanical once the model and decision problem are completely specified.

Formalization of problem and action space

Suppose that in addition to the objective function f , we have access to one-or-more surrogate functions $\{f_i\}: \mathcal{X} \rightarrow \mathbb{R}$, indexed by a parameter $i \in \mathcal{I}$. Most often we take the surrogate functions to form a discrete set, but in some cases we may wish to consider multidimensional and/or continuous surrogate spaces.⁴⁴ We denote the objective function itself with the special index $* \in \mathcal{I}$, writing $f = f_*$. We consider an optimization scenario where we may design each observation to be either of the objective or a surrogate as we see fit, by selecting a location $x \in \mathcal{X}$ for our next observation and an index $i \in \mathcal{I}$ specifying the desired surrogate. The action space for each such decision is $\mathcal{A} = \mathcal{X} \times \mathcal{I}$.



In multifidelity optimization, we wish to optimize an expensive objective function (blue) aided by access to cheaper – but still informative – surrogates (red).

materials science applications: p.309

⁴¹ “Relatively” should be stressed; these simulations can be quite expensive in absolute terms, but still much cheaper than synthesis.

⁴² A. KLEIN et al. (2015). Towards efficient Bayesian Optimization for Big Data. *Bayesian Optimization Workshop, NeurIPS 2015*.

⁴³ K. SWERSKY et al. (2014). Freeze–Thaw Bayesian Optimization. arXiv: 1406 . 3896 [stat.ML].

surrogate functions, $\{f_i\}$
surrogate index set, \mathcal{I}

objective function, f_*

⁴⁴ K. KANDASAMY et al. (2017). Multi-fidelity Bayesian Optimisation with Continuous Approximations. *ICML 2017*.

action space, $\mathcal{A} = \mathcal{X} \times \mathcal{I}$

joint prior process, $p(\{f_i\})$

observation model, $p(y | x, i, \phi)$
 posterior distribution, $p(\{f_i\} | \mathcal{D})$
 posterior predictive distribution,
 $p(y | x, i, \mathcal{D})$

joint Gaussian processes: § 2.4, p. 26

45 M. A. ÁLVAREZ et al. (2012). Kernels for Vector-Valued Functions: A Review. *Foundations and Trends in Machine Learning* 4(3):195–266.

46 K. ULRICH et al. (2015). GP Kernels for Cross-Spectrum Analysis. *NeurIPS 2015*.

shared domain covariance, $K_{\mathcal{X}}$
 cross-function covariance, $K_{\mathcal{I}}$

Modeling surrogate functions and observations

If surrogate observations are to be useful, they must provide information about the objective function, and the relationship between the objective and its surrogates is captured by a joint model over their values. We first design a joint prior process $p(\{f_i\})$ specifying the expected structure of each individual function and the nature of correlations between the functions. Next we must create an observation model linking the value y observed at a point $[x, i]$ to the underlying function value $\phi = f_i(x)$: $p(y | x, i, \phi)$. Now, given a set of observed data \mathcal{D} , we may derive the posterior belief over the functions, $p(\{f_i\} | \mathcal{D})$, and the posterior predictive distribution, $p(y | x, i, \mathcal{D})$, with which we can reason about proposed observations.

In practice, the joint prior process is usually a joint *Gaussian* process over the objective and its surrogates. The primary challenge in crafting a joint Gaussian process is in defining cross-covariance functions

$$K_{ij} = \text{cov}[f_i, f_j]$$

that adequately encode the correlations between the functions of interest, which can take some care to ensure the resulting joint covariance function over the collection $\{f_i\}$ is positive definite.

Fortunately, a great deal of effort has been invested in developing this model class into a flexible and expressive family.^{45,46} One simple construction offering some intuition is the *separable* covariance

$$K([x, i], [x', i']) = K_{\mathcal{X}}(x, x') K_{\mathcal{I}}(i, i'), \quad (11.18)$$

which decomposes the joint covariance into a covariance function on the domain $K_{\mathcal{X}}$ shared by each individual function and a covariance function between the functions, $K_{\mathcal{I}}$ (which would be a covariance *matrix* if \mathcal{I} is finite). In this construction, the marginal covariance and cross-covariance functions are all scaled versions of $K_{\mathcal{X}}$, with the $K_{\mathcal{I}}$ covariance scaling each marginal belief and encoding (constant) cross-correlations across functions as well. Figures 2.5 and 2.6 illustrate the behavior of this model for two highly correlated functions on a shared one-dimensional domain.

Defining a utility function

We have now addressed steps 2 and 3 of our general procedure for multifidelity optimization, by identifying the expanded action space implied by the problem and determining how to reason about the potential outcomes of these actions. Before we can proceed with deriving a policy, however, we must establish preferences over outcomes with a utility function $u(\mathcal{D})$. It is difficult to provide specific guidance for this choice, as these preferences are inherently bound to a given situation. One natural approach would be to choose a cost-aware utility function measuring optimization progress limited to the objective function alone, adjusted for the variable costs of each observation obtained. For example we

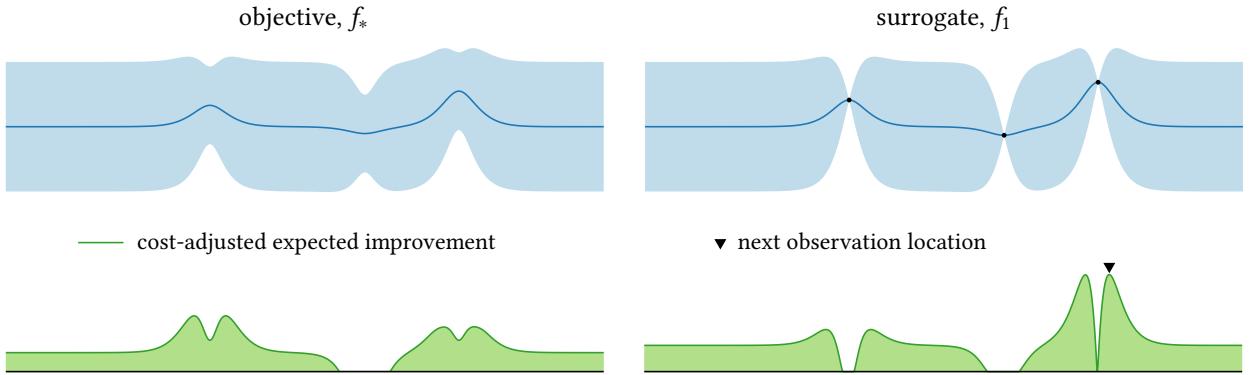


Figure 11.7: Cost-adjusted multifidelity expected improvement for a toy scenario. The objective (left) and its surrogate (right) are modeled as a joint Gaussian process with marginals equivalent to the example from figure 7.1 and constant cross-correlation 0.8. The cost of observation was assumed to be ten times greater for the objective than its surrogate. Maximizing the cost-adjusted multifidelity expected improvement elects to continue evaluating the surrogate.

might quantify the cost of observing the objective function or each of the surrogate functions with values $\{c_i\}$, and adjust a data utility $u'(\mathcal{D})$ appropriately, by defining

$$u(\mathcal{D}) = u'(\mathcal{D}) - \sum_{(x, i) \in \mathcal{D}} c_i.$$

With an appropriate utility function in hand, we can then proceed to derive the optimal policy as usual.

As an example, we might realize an analog of expected improvement (7.2) through a suitable redefinition of the simple reward (6.3). Suppose that at termination we wish to recommend a location visited during optimization, evaluated at *any* fidelity, using a risk-neutral utility. Given a multifidelity dataset $\mathcal{D} = ([\mathbf{x}, \mathbf{i}], \mathbf{y})$, the utility of this recommendation would be

$$u'(\mathcal{D}) = \max_{x' \in \mathbf{x}} \mu_{\mathcal{D}}([x', *]),$$

the maximum of the posterior mean for the objective function at the observed locations.

Figure 11.7 illustrates one-step lookahead policy with this utility function for a one-dimensional objective function f_* (left) and a surrogate f_1 (right). The marginal belief about each function is a Gaussian process identical to our running example from chapter 7; see figure 7.1. These are coupled together via the separable covariance (11.18) with $K_{\mathcal{I}}(*, *) = K_{\mathcal{I}}(1, 1) = 1$ and cross-correlation $K_{\mathcal{I}}(1, *) = 0.8$. We begin with three surrogate observations and compute the cost-adjusted expected improvement as described above, where the cost of observing the objective was set to ten times that of the surrogate. In this case, the one-step optimal decision is to continue evaluating the surrogate around the best-seen surrogate observation.

observation costs, $\{c_i\}$

multi-fidelity expected improvement

example and discussion

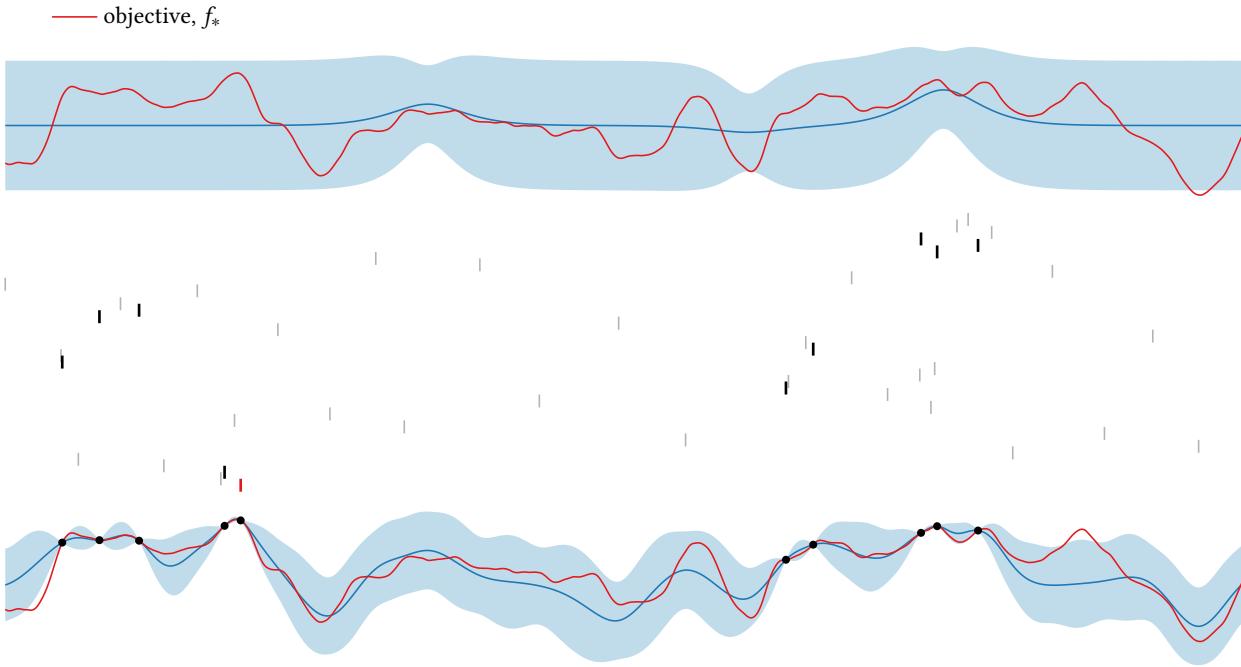


Figure 11.8: A simulation of optimization with the cost-adjusted multifidelity expected improvement starting from the scenario in figure 11.7. We simulate sequential observations of either the objective or the surrogate, illustrated using the running tick marks. The light gray marks correspond to surrogate observations and the heavy black marks objective observations. The optimum was found after 10 objective and 32 surrogate observations, marked in heavy red. The prior and posterior of the objective function conditioned on all observations are also shown.

Figure 11.8 simulates sequential multifidelity optimization using this policy; here the optimum was discovered after only 10 evaluations of the objective, guided by 32 observations of the cheaper surrogate. A remarkable feature we can see in the posterior is that *all* evaluations made of the objective function are above the prior mean, nearly all with z -scores of approximately $z = 1$ or greater. This can be ascribed not to extreme luck, but rather to efficient use of the surrogate to rule out regions unlikely to contain the optimum.

Multifidelity Bayesian optimization has enjoyed sustained interest from the research community, and numerous policies have been available. These include adaptations of the expected improvement,^{47,48} knowledge gradient,⁴⁹ upper confidence bound,^{50,51} and mutual information with x^* ^{52,53} and f^* ⁵⁴ acquisition functions, as well as novel approaches.⁵⁵

11.6 MULTITASK OPTIMIZATION

Multitask optimization addresses the sequential or simultaneous optimization of multiple objectives $\{f_i\}: \mathcal{X} \rightarrow \mathbb{R}$ representing performance

47 D. HUANG et al. (2006a). Sequential kriging optimization using multiple-fidelity evaluations. *Structural and Multidisciplinary Optimization* 32(5):369–382.

48 V. PICHENY et al. (2013a). Quantile-Based Optimization of Noisy Computer Experiments With Tunable Precision. *Technometrics* 55(1): 2–13.

49 J. WU et al. (2019). Practical Multi-fidelity Bayesian Optimization for Hyperparameter Tuning. *UAI 2019*.

50 K. KANDASAMY et al. (2016). Gaussian Process Bandit Optimisation with Multi-fidelity Evaluations. *NeurIPS 2016*.

51 K. KANDASAMY et al. (2017). Multi-fidelity Bayesian Optimisation with Continuous Approximations. *ICML 2017*.

52 K. SWERSKY et al. (2013). Multi-Task Bayesian Optimization. *NeurIPS 2013*.

on related tasks. Like multifidelity optimization, the underlying idea in multitask optimization is that if performance on the various tasks is *correlated* as a function of the input, we may accelerate optimization by transferring information between tasks.

As a motivating example of sequential multitask optimization, consider a web service wishing to retune the parameters of an ad placement algorithm on a regular basis to maximize revenue in the current climate. Here the revenue at each epoch represents the different tasks to be optimized, which are optimized individually one after another. Although we could treat each optimization problem separately, they are clearly related, and with some care we may be able to use past performance to provide a “warm start” to each new optimization problem rather than start from scratch.

We may also consider the simultaneous optimization of performance across tasks. For example, a machine learning practitioner may wish to tune model hyperparameters to maximize the average predictive performance on several related datasets.⁵⁶ A naïve approach would formulate the problem as maximizing a single objective defined to be the average performance across tasks, with each evaluation entailing retraining the model for each dataset. However, this would be potentially wasteful, as we may be able to eliminate poorly performing hyperparameters with high confidence after training on a fraction of the datasets. A multi-task approach would model each objective function separately (perhaps jointly) and consider evaluations of *single-task* performance to efficiently maximize the combined objective.

SWERSKY et al. described a particularly clever realization of this idea: selecting model hyperparameters via cross validation.⁵⁷ Here we recognize the predictive performance on each validation fold as being correlated due to shared training data across folds. If we can successfully share information across folds, we may potentially accelerate cross validation by iteratively selecting (hyperparameter, fold index) pairs rather than training proposed hyperparameters on every fold each time.

Formulation and approach

Let $\{f_i\}: \mathcal{X} \rightarrow \mathbb{R}$ be the set of objective functions we wish to consider, representing performance on the relevant tasks. As with multifidelity optimization, the key enabler of multitask optimization is a joint model $p(\{f_i\})$ over the tasks and a joint observation model $p(y | x, i, \phi)$ over evaluations thereof; this joint model allows us to share information between the tasks. This could, for example, take the form of a joint Gaussian process, as discussed previously.

Once this model has been chosen, we can turn to the problem of designing a multitask optimization policy. If each task is to be solved one at a time, a natural approach would be to design the utility function to ultimately evaluate performance on that task only. In this case, the problem devolves to single-objective optimization, and we may use any of the approaches discussed earlier in the book to derive a policy. The

53 Y. ZHANG et al. (2017). Information-Based Multi-Fidelity Bayesian Optimization. *Bayesian Optimization Workshop, NeurIPS 2017*.

54 S. TAKENO et al. (2020). Multi-fidelity Bayesian Optimization with Max-value Entropy Search and its Parallelization. *ICML 2020*.

55 J. SONG et al. (2019). A General Framework for Multi-fidelity Bayesian Optimization with Gaussian Processes. *AISTATS 2019*.

simultaneous tasks

56 This could also be formulated as a scalarized version of a multiobjective optimization problem, discussed in the next section.

57 K. SWERSKY et al. (2013). Multi-Task Bayesian Optimization. *NeurIPS 2013*.

modeling task objectives and observations

joint GPs for modeling multiple functions:
§ 11.5, p. 262

sequential tasks

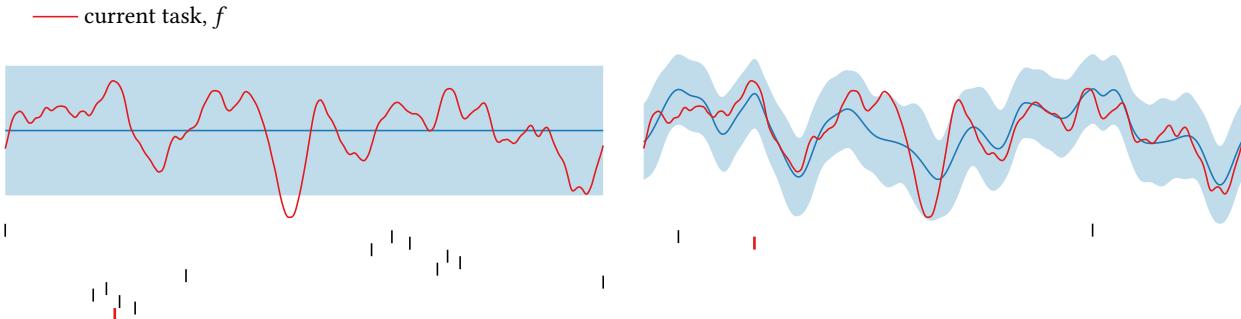


Figure 11.9: A demonstration of sequential multitask optimization. Left: a prior distribution over an objective function along with 13 observations selected by maximizing expected improvement, revealing the global maximum with the last evaluation. Right: the posterior distribution over the same objective conditioned on the observations of the two functions (now interpreted as related tasks) in figure 11.8. The global maximum is now found after three observations due to the better informed prior.

only difference is that the objective function model is now informed from our past experience with other tasks; as a result, our initial optimization decisions can be more targeted.

This procedure is illustrated in figure 11.9. Both panels illustrate the optimization of an objective function by sequentially maximizing expected improvement. The left panel begins with no information and locates the global optimum after 13 evaluations. The right panel begins the process instead with a *posterior* belief about the objective conditioned on the data obtained from the two functions in figure 11.8, modeled as related tasks with cross-correlation 0.8. Due to the better informed initial belief, we now find the global optimum after only three evaluations.

The case of simultaneous multitask optimization – where we may evaluate any task objective with each observation – requires somewhat more care. We must now design a utility function capturing our joint performance across the tasks and design each observation with respect to this utility. One simple option would be to select utility functions $\{u_i\}$ quantifying performance on each task separately and then take a weighted average:

$$u(\mathcal{D}) = \sum_i w_i u_i(\mathcal{D}).$$

This could be further adjusted for (perhaps task- and/or input-dependent) observation costs if needed. Now we may write the expected marginal gain in this combined utility as a weighted average of the expected marginal gain on each separate task. One-step lookahead would then maximize the weighted acquisition function

$$\alpha([x, i]; \mathcal{D}) = \sum_i w_i \alpha_i([x, i]; \mathcal{D})$$

over all possible observation location–task pairs $[x, i]$, where α_i is the one-step expected marginal gain in u_i . Note that observing a single task

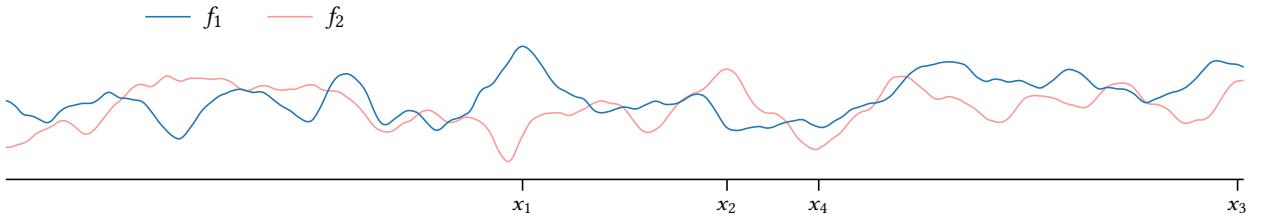


Figure 11.10: A simple multiobjective optimization example with two objectives $\{f_1, f_2\}$ on a one-dimensional domain. We compare four identified points: x_1 , the global optimum of f_1 , x_2 , the global optimum of f_2 , x_3 , a compromise with relatively high values of both objectives, and x_4 , a point with relatively low values of both objectives.

could in fact yield improvement in *all* utilities due to information sharing through the joint belief on task objectives.

11.7 MULTIOBJECTIVE OPTIMIZATION

Like multitask optimization, *multiobjective optimization* addresses the simultaneous optimization of multiple objectives $\{f_i\} : \mathcal{X} \rightarrow \mathbb{R}$. However, whereas in multitask optimization we usually seek to identify the global optimum of each function separately, in multiobjective optimization we seek to identify points *jointly* optimizing all of the objectives. Of course, this is not possible unequivocally unless all of the maxima happen to coincide, as we may need to sacrifice the value of one objective in order to increase another. Instead, we may consider the optimization of various *tradeoffs* between the objectives and rely on subjective preferences to determine which option is preferred in a given scenario. Multiobjective optimization may then be posed as the identification of one or more optimal tradeoffs among the objectives to support this analysis.

A classic example of multiobjective optimization can be found in finance, where we seek investment portfolios optimizing tradeoffs between *risk* (often captured by the standard deviation of return) and *reward* (often captured by the expected return). Generally, investments with higher risk yield higher reward, but the optimal investment strategy depends on the investor's risk tolerance – for example, when capital preservation is paramount, low-risk, low-reward investments are prudent. The set of investment portfolios maximizing reward for any given risk is known as the *Pareto frontier*,⁵⁸ which jointly span all rational solutions to the given problem. We generalize this concept below.

Pareto optimality

To illustrate the tradeoffs we may need to consider during multiobjective optimization, consider the two objectives in figure 11.10. The first objective f_1 has its global maximum at x_1 , which nearly coincides with the global minimum of the second objective f_2 . The reverse is true in the other direction: the global maximum of the second objective, x_2 ,

risk: standard deviation of return

reward: expected value of return

⁵⁸ In modern portfolio theory the term “efficient frontier” is more common for the equivalent concept.

Pareto frontier

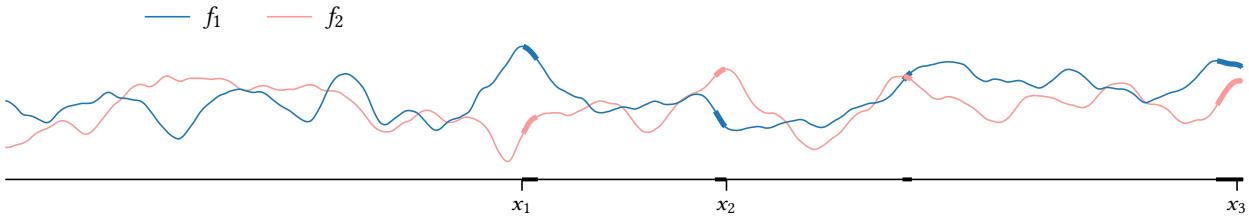
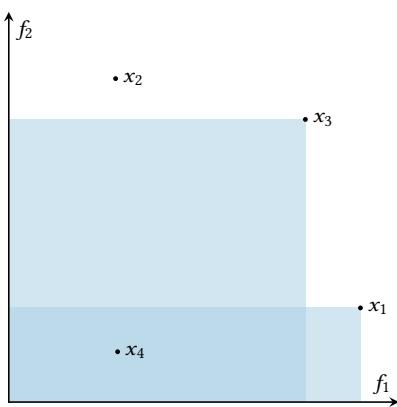
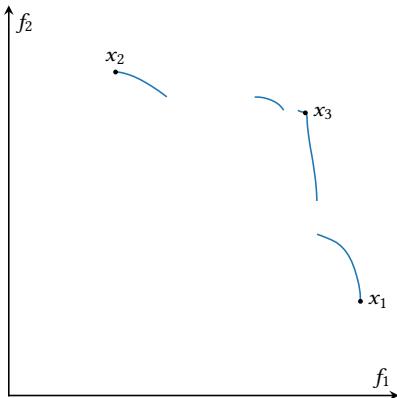


Figure 11.11: The objectives from figure 11.10 with the Pareto optimal solutions highlighted. All points along the intervals marked on the horizontal axis are Pareto optimal, with the highlighted corresponding objective values forming the Pareto frontier (see margin).



The regions dominated by points x_1 and x_3 in figure 11.10 are highlighted in blue. x_4 is dominated by both, and x_2 by neither.



The Pareto frontier for the scenario in figures 11.10–11.11. The four components correspond to the highlighted intervals in figure 11.11.

⁵⁹ K. M. MIETTINEN (1998). *Nonlinear Multiobjective Optimization*. Kluwer Academic Publishers.

achieves a relatively low value on the first objective. Neither point can be preferred over the other on any *objective* grounds, but a rational agent may have *subjective* preferences over these two locations depending on the relative importance of the two objectives. Some agents might even prefer a compromise location such as x_3 to either of these points, as it achieves relatively high – but suboptimal – values for both objectives.

It is clearly impossible to identify an unambiguously optimal location, even in this relatively simple example. We can, however, *eliminate* some locations as plainly subpar. For example, consider the point x_4 in figure 11.10. Assuming preferences are nondecreasing with each objective, no rational agent would prefer x_4 to x_3 as the latter point achieves higher value for *both* objectives. We may formalize this intuition by defining a partial order on potential solutions consistent with this reasoning.

We will say that a point x *dominates* another point x' , denoted $x' < x$, if no objective value is lower at x than at x' and if at least one objective value is higher at x than at x' . Assuming preferences are consistent with nondecreasing objective values, no agent could prefer a dominated point to any point dominating it. This concept is illustrated in the margin for the example from figure 11.10: all points in the blue regions are dominated, and in particular x_4 is dominated by both x_1 and x_3 . On the other hand, none of x_1 , x_2 , or x_3 is dominated by any of the other points.

A point $x \in \mathcal{X}$ that is *not* dominated by any other point is called *Pareto optimal*, and the image of all Pareto optimal points is called the *Pareto frontier*. The Pareto frontier is a central concept in multiobjective optimization – it represents the set of all possible solutions to the problem consistent with weakly monotone preferences for the objectives. Figure 11.11 shows the Pareto optimal points for our example from figure 11.10, which span four disconnected intervals. We may visualize the Pareto frontier by plotting the image of this set, as shown in the margin.

There are several approaches to multiobjective optimization that differ in when preferences among competing solutions are elicited.⁵⁹ So-called *a posteriori methods* seek to identify the entire Pareto frontier for a given problem, with preferences among possible solutions to be determined afterwards. In contrast, *a priori methods* assume that preferences are already predetermined, allowing us to seek a single Pareto

optimal solution consistent with those preferences. Bayesian realizations of both types of approaches have been realized, as we discuss below.

Formulation of decision problem and modeling objectives

Nearly all Bayesian multiobjective optimization procedures model each decision as choosing a location $x \in \mathcal{X}$, where we make an observation of every objective function. We could also consider a setting analogous to multitask or multifidelity optimization where we observe only one objective at a time, but this idea has not been sufficiently explored. For the following discussion, we will write y for the vector-valued observation resulting from an observation at a given point x , with y_i being associated with objective f_i .

As in the previous two sections, we build a joint model for the objectives $\{f_i\}$ and our observations of them via a prior process $p(\{f_i\})$ and an observation model $p(y | x, \{\phi_i\})$. The models are usually taken to be independent Gaussian processes on each objective combined with standard observation models. Joint Gaussian processes could also be used when appropriate, but a direct comparison of independent versus dependent models did not demonstrate improvement when modeling correlations between objectives, perhaps due to an increased burden in estimating model hyperparameters.⁶⁰

vector of observations at x , y

60 J. SVENSON and T. SANTNER (2016). Multiobjective optimization of expensive-to-evaluate deterministic computer simulator models. *Computational Statistics and Data Analysis* 94:250–264.

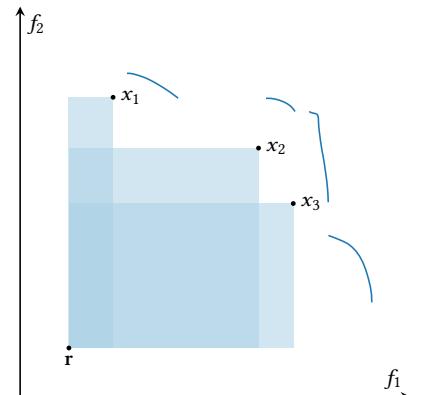
Expected hypervolume improvement

The majority of Bayesian multiobjective optimization approaches are a posteriori methods, seeking to identify the entire Pareto frontier or a representative portion of it. Many algorithms represent one-step lookahead for some utility function evaluating progress on this task.

One popular utility for multiobjective optimization is the volume under an estimate of the Pareto frontier, also known as the *S metric*.⁶¹ Namely, given observations of the objectives, we may build a natural statistical lower bound of the Pareto frontier by eliminating the outcomes dominated by the observations with high confidence. When observations are exact, we may simply enumerate the dominated regions and take their union; when observations are corrupted by noise, we may use a statistical lower bound of the underlying function values instead.⁶² This procedure is illustrated in the margin for our running example, where we have made exact observations of the objectives at three mutually nondominated locations; see figure 11.12. The upper-right boundary of the dominated region is a lower bound of the true Pareto frontier.

To evaluate progress on mapping out the Pareto frontier, we consider the volume of space dominated by the available observations and bounded below by an identified, clearly suboptimal reference point r ; see the blue shaded area in the margin. The reference point is necessary to ensure the dominated volume does not diverge to infinity. Assuming the reference point is chosen such that it will definitely be dominated, this utility is always positive and is maximized when the true Pareto fron-

61 E. ZITZLER (1999). Evolutionary Algorithms for Multiobjective Optimization: Methods and Applications. PhD thesis. Eidgenössische Technische Hochschule Zürich. [§ 3.1]



A lower bound of the Pareto frontier for our example given the data in figure 11.12.

62 M. EMMERICH and B. NAUJOKS (2004). Metamodel Assisted Multiobjective Optimisation Strategies and their Application in Airfoil Design. In: *Adaptive Computing in Design and Manufacture VI*.

Figure 11.12: The posterior belief about our example objectives from figure 11.10 given observations at the marked locations. The beliefs are separated vertically (by an arbitrary amount) for clarity, with the belief over the other function shown in gray for reference.

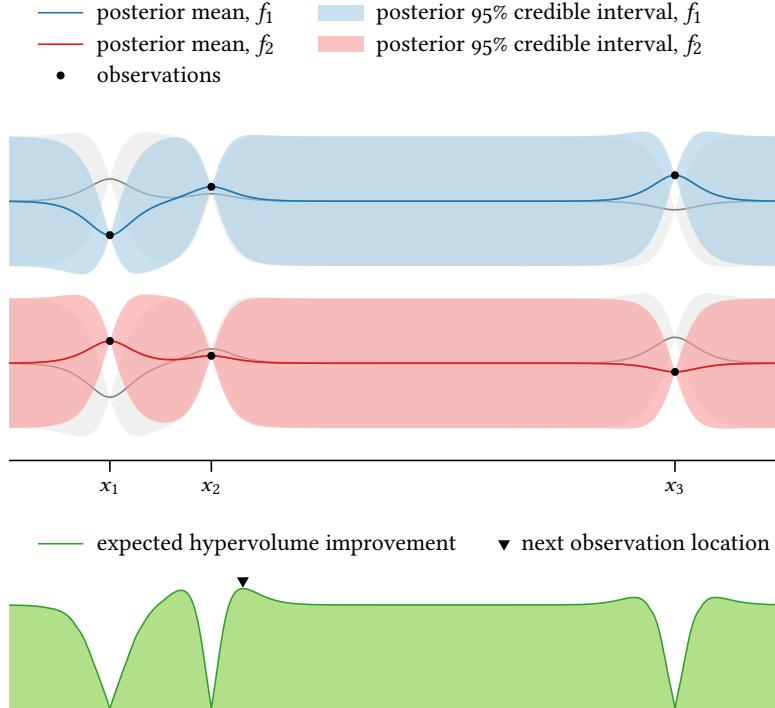


Figure 11.13: The expected hypervolume improvement acquisition function for the above example.

- 63 M. FLEISCHER (2003). The Measure of Pareto Optima: Applications to Multi-objective Metaheuristics. *EMO 2003*.
- 64 M. T. M. EMMERICH et al. (2006). Single- and Multiobjective Evolutionary Optimization Assisted by Gaussian Random Field Metamodels. *IEEE Transactions on Evolutionary Computation* 10(4):421–439.
- 65 W. PONWEISER et al. (2008). Multiobjective Optimization on a Limited Budget of Evaluations Using Model-Assisted S -Metric Selection. *PPSN X*.
- 66 K. YANG et al. (2019b). Multi-Objective Bayesian Global Optimization using expected hypervolume improvement gradient. *Swarm and Evolutionary Computation* 44:945–956.
- 67 K. YANG et al. (2019a). Efficient computation of expected hypervolume improvement using box decomposition algorithms. *Journal of Global Optimization* 75(1):3–34.

tier is revealed by the observed data.⁶³ Therefore it provides a sensible measure of progress for a posteriori multiobjective optimization.

The one-step marginal gain in this utility is known as *expected hypervolume improvement* (EHVI)^{64,65} and serves as a popular acquisition function. This score is shown in figure 11.13 for our example; the optimal decision attempts to refine the central portion of the Pareto frontier, but many alternatives are almost as favorable due to the roughness of the current estimate. Computation of EHVI is involved, and its cost grows considerably with the number of objectives. The primary difficulty is enumerating and integrating with respect to the lower bound to the Pareto front, which can become a complex region in higher dimensions. However, efficient algorithms are available for computing EVHI⁶⁶ and its gradient.⁶⁷

Information-theoretic a posteriori methods

Several popular information-theoretic policies for single-objective optimization have been adapted to a posteriori multiobjective optimization. The key idea behind these methods is to approximate the mutual information between a joint observation of the objectives and either the set of Pareto optimal points (in the domain), \mathcal{X}^* , or the Pareto frontier (in the codomain), \mathcal{F}^* . These approaches operate by maximizing the predictive form of mutual information (8.35) and largely follow the parallel single-objective cases in their approximation.

HERNÁNDEZ-LOBATO et al. proposed maximizing the mutual information between the observations y realized at a proposed observation location x and the set of Pareto optimal points \mathcal{X}^* :

$$\alpha_{\text{PESMO}}(x; \mathcal{D}) = H[y | x, \mathcal{D}] - \mathbb{E}_{\mathcal{X}^*}[H[y | x, \mathcal{D}, \mathcal{X}^*] | x, \mathcal{D}],$$

calling their policy *predictive entropy search for multiobjective optimization* (PESMO).⁶⁸ As in the single-objective case, for Gaussian process models with additive Gaussian noise, the first term of this expression can be computed exactly as the differential entropy of a multivariate normal distribution (A.16). However, the second term entails two computational barriers: computing an expectation with respect to \mathcal{X}^* and conditioning our objective function belief on this set. The authors provide approximations for each of these tasks for Gaussian process models based on Gaussian expectation propagation; these are reminiscent of the procedure used in predictive entropy search.

BELAKARIA et al. meanwhile proposed maximizing the mutual information with the Pareto frontier \mathcal{F}^* :

$$\alpha_{\text{MESMO}}(x; \mathcal{D}) = H[y | x, \mathcal{D}] - \mathbb{E}_{\mathcal{F}^*}[H[y | x, \mathcal{D}, \mathcal{F}^*] | x, \mathcal{D}].$$

The authors dubbed the resulting policy *max-value entropy search for multiobjective optimization* (MESMO).⁶⁹ This policy naturally shares many features with the PESMO policy. Again the chief difficulty is in addressing the expectation and conditioning with respect to \mathcal{F}^* in the second term of the mutual information; however, both of these tasks are rendered somewhat easier by working in the codomain. To approximate the expectation with respect to \mathcal{F}^* , the authors propose a Monte Carlo approach where posterior samples of the objectives are generated efficiently using a sparse-spectrum approximation, which are fed into an off-the-shelf, exhaustive multi-objective optimization routine. Conditioning y on a realization of the Pareto frontier now entails appropriate truncation of its multivariate normal belief.

A priori methods and scalarization

When preferences regarding tradeoffs among the objective functions can be sufficiently established prior to optimization, perhaps with input from a domain expert, we may reduce multiobjective optimization to a *single* objective problem by explicitly maximizing the desired criterion. This is called *scalarization* and is a prominent *a priori* method. Reframing in terms of a single objective offers the obvious benefit of allowing us to appeal to the expansive methodology for that purpose we have built up throughout the course of this book.

Scalarization is an important component of some a posteriori multiobjective optimization methods as well. The idea is to construct a family of exhaustive parametric scalarizations of the objectives such that the solution to any such problem is Pareto optimal, and that by spanning the parameter range we may reveal the entire Pareto frontier one point at a time.

68 D. HERNÁNDEZ-LOBATO et al. (2016a). Predictive Entropy Search for Multi-objective Bayesian Optimization. *ICML 2016*.

predictive entropy search: § 8.8, p. 180

69 S. BELAKARIA et al. (2019). Max-value Entropy Search for Multi-Objective Bayesian Optimization. *NeurIPS 2019*.

sparse spectrum approximation: § 8.7, p. 178

scalarization

a posteriori optimization via scalarization

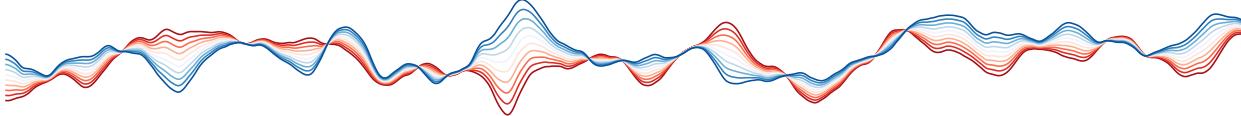
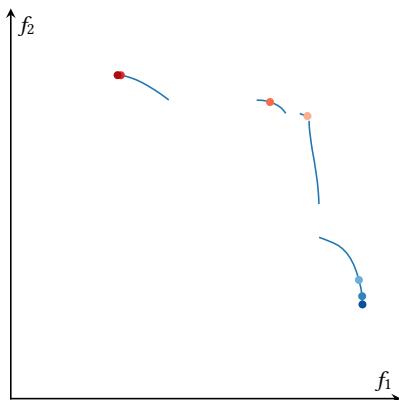


Figure 11.14: A series of linear scalarizations of the example objectives from figure 11.10. The “blue end” of the color spectrum is f_1 , and the “red end” f_2 , consistent with that plot.

vector of objective values at x , ϕ
scalarization function, g



The marked points represent the maxima of the similarly colored linear scalarizations in figure 11.14.

Let ϕ denote the vector of objective function values at an arbitrary location x , defining $\phi_i = f_i(x)$. A *scalarization function* $g: x \mapsto g(\phi) \in \mathbb{R}$ maps locations in the domain to scalars determined by their objective values. We may interpret the output as defining preferences over locations in a natural manner: namely, if $g(\phi) > g(\phi')$, then the outcomes at x are preferred to those at x' in the scalarization. With this interpretation, a scalarization function allows us to recast multiobjective optimization as a single-objective problem by maximizing g with respect to x .

A scalarization function can in principle be arbitrary, and a priori multiobjective optimization can be framed in terms of maximizing any such function. However, several tunable scalarization functions have been described in the literature that may be used in a general context.

A straightforward and intuitive example is the *linear scalarization* function:

$$g_{\text{LIN}}(x; \mathbf{w}) = \sum_i w_i \phi_i, \quad (11.19)$$

where each weight w_i is nonnegative; that is, we simply take a positive weighted sum of the objectives. A range of linear scalarizations for our running example is shown in figure 11.14, here constructed to smoothly interpolate between the two objectives. The maximum of a linear scalarization is guaranteed to lie on the Pareto frontier; however, not every Pareto optimal point can be recovered in this manner unless the frontier is strictly concave. That is the case for our example, illustrated in the marginal figure. If we model each objective with a (perhaps joint) Gaussian process, then the induced belief about any linear scalarization is conveniently also a Gaussian process, so no further modeling would be required for the scalarization function itself.

Another choice that has seen some use in Bayesian optimization is *augmented Chebyshev* scalarization, which augments the linear scalarization (11.19) with an additional, nonlinear term:

$$g_{\text{AC}}(x; \mathbf{w}, \rho) = \min_i [w_i(\phi_i - r_i)] + \rho g_{\text{LIN}}(x; \mathbf{w}). \quad (11.20)$$

Here r is a reference point as above and ρ is a small nonnegative constant; KNOWLES for example took $\rho = 0.05$.⁷⁰ The augmented Chebyshev scalarization function has the benefit that *all* points on the Pareto frontier can be realized by maximizing with respect to some corresponding setting of the weights, even if the frontier is nonconcave.

⁷⁰ J. KNOWLES (2005). PAREGO: A Hybrid Algorithm With On-Line Landscape Approximation for Expensive Multiobjective Optimization Problems. *IEEE Transactions on Evolutionary Computation* 10(1):50–66.

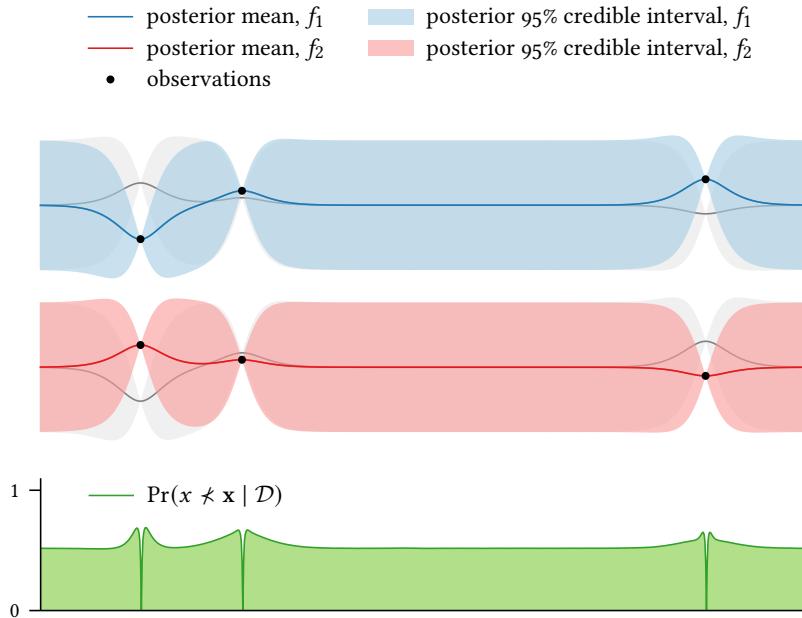


Figure 11.15: The probability of nondominance by the available data for our running example (above).

Several authors have derived Bayesian methods for a posteriori multi-objective optimization by solving a series of carefully constructed scalarized problems. KNOWLES for example proposed sampling random weights for the augmented Chebyshev scalarization (11.20) and optimizing the resulting objective by maximizing expected improvement, repeating this process until satisfied.⁷⁰ PARIA et al. proposed a similar approach incorporating a prior distribution over the parameters of a chosen scalarization function to allow the user to focus on an identified region of the Pareto frontier if desired.⁷¹ The procedure then proceeds by repeatedly sampling from that distribution and maximizing the resulting objective via Thompson sampling or maximizing an upper confidence bound. Ignoring the choice of policy for the scalarizations, this procedure is effectively the same as KNOWLES's, but the authors were able to establish theoretical regret bounds in their slightly different framework.

Other approaches

ZULUAGA et al. outlined an intriguing approach to a posteriori multi-objective optimization⁷² wherein the problem was recast as an *active learning* problem.⁷³ Namely, the authors considered the binary classification problem of predicting whether a given observation location was (approximately) Pareto optimal or not, then designed observations to maximize expected performance on this task. Their algorithm is supported by theoretical bounds on performance and performed admirably against the PareGO algorithm discussed above.⁷⁰

PICHENY proposed a spiritually similar approach also based on sequentially reducing a measure of uncertainty in the Pareto frontier.⁷⁴

expected improvement: § 7.3, p. 127

upper confidence bound, Thompson sampling: §§ 7.8–7.9, p. 145

⁷¹ B. PARIA et al. (2019). A Flexible Framework for Multi-Objective Bayesian Optimization using Random Scalarizations. *UAI 2019*.

⁷² M. ZULUAGA et al. (2016). ε -PAL: An Active Learning Approach to the Multi-Objective Optimization Problem. *Journal of Machine Learning Research* 17(104):1–32.

⁷³ B. SETTLES (2012). *Active Learning*. Morgan & Claypool.

⁷⁴ V. PICHENY (2015). Multiobjective optimization using Gaussian process emulators via stepwise uncertainty reduction. *Statistics and Computing* 25(6):1265–1280.

Given a dataset $\mathcal{D} = (\mathbf{x}, \mathbf{y})$, we consider the probability that a given point $x \in \mathcal{X}$ is *not* dominated by any point in the current dataset (that is, the probability x may lie on the Pareto frontier but not yet be discovered), $\Pr(x \not\prec \mathbf{x} \mid \mathcal{D})$. The integral of this score over the domain can be interpreted as a measure of uncertainty in the Pareto frontier as determined by the data, and negating this measure provides a plausible utility function for a posteriori multiobjective optimization:

$$u(\mathcal{D}) = - \int \Pr(x \not\prec \mathbf{x}) dx.$$

This probability is plotted for our running example in figure 11.15; here, there is a significant probability for many points to be nondominated by the rather sparse available data, indicating a significant degree of uncertainty in our understanding of the Pareto frontier. However, as the Pareto frontier is increasingly well determined by the data, this probability will vanish globally and the utility above will tend toward its maximal value of zero. After motivating this score, PICHENY proceeds to recommend designing observations via one-step lookahead.

11.8 GRADIENT OBSERVATIONS

Bayesian optimization is often described as a “derivative-free” approach to optimization, but this characterization is misleading. Although it is true that Bayesian optimization methods do not *require* the ability to observe derivatives, it is *certainly* not the case that we cannot make use of such observations when available. In fact, it is straightforward to condition a Gaussian process on derivative observations, even if corrupted by noise, and so from a modeling perspective we are already done.

conditioning a GP on derivative observations:
§ 2.6, p. 32

⁷⁵ This could be, for example, exact observation or additive Gaussian noise. Recall that for a GP on f , the joint distribution of $(\phi, \nabla\phi)$ is multivariate normal (2.28), so for these choices the predictive distribution of (y, g) is jointly Gaussian and exact inference is tractable.

Of course, ideally, our policy should also consider the acquisition of derivative information due to its influence on our belief and the utility of collected data, of which they now form a part. To do so is by now relatively simple. A fairly general scheme would assume that an observation at x yields a pair of measurements (y, g) respectively related to $(\phi, \nabla\phi)$ via a joint observation model.⁷⁵ We can then compute the one-step expected marginal gain in utility from observing these values:

$$\alpha_1(x; \mathcal{D}) = \iint [u(\mathcal{D}_1) - u(\mathcal{D})] p(y, g \mid x, \mathcal{D}) dy dg, \quad (11.21)$$

where the updated dataset \mathcal{D}_1 will reflect the entire observation (x, y, g) . Induction on the horizon gives the optimal policy as usual.

Figure 11.16 compares derivative-aware and derivative-unaware versions of the knowledge gradient (7.4) (assuming exact observation) for an example scenario. The derivative-aware version dominates the derivative-unaware one, as the acquisition of more information naturally leads to a greater expected marginal gain in utility. When derivative information is unavailable, the optimal decision is to evaluate nearby the previously best-seen point, effectively to estimate the derivative via finite differencing; when derivative information is available, an observation at this

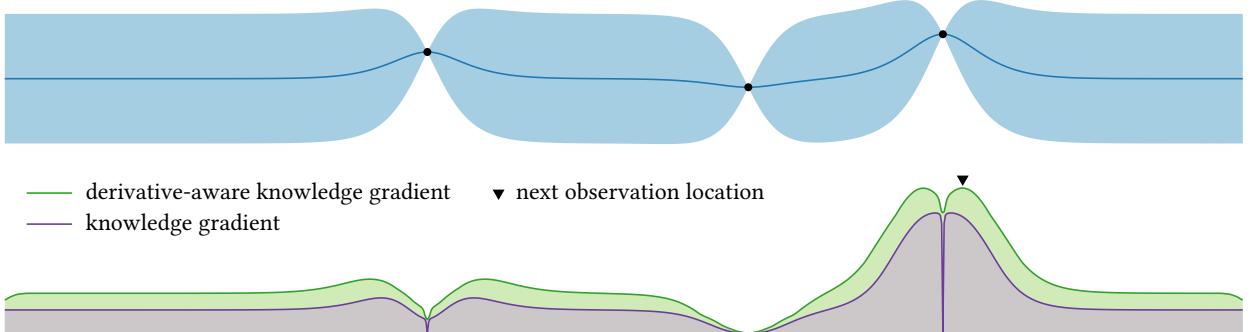


Figure 11.16: The knowledge gradient acquisition function for an example scenario reflecting the expected gain in global reward (6.5) provided exact observations of the objective function and its derivative. The vanilla knowledge gradient (7.4) based on an observation of the objective alone is shown for reference.

location yields effectively the same expected gain. However, we can fair even better in expectation by moving a bit farther astray, where the derivative information is less redundant.

When working with Gaussian process models in high dimension, it may not be wise to augment each observation of the objective with a full observation of the gradient due to the cubic scaling of inference. However, the scheme outlined above opens the door to consider the observation of *any* measurement related to the gradient. For example, we might condition on the value of a *directional* derivative, reducing the measurement to a scalar regardless of dimension and limiting computation. Such a scheme was promoted by Wu et al., who considered the acquisition of a single coordinate of the gradient; this could be extended to non-axis-aligned directional derivatives without major complication.⁷⁶ We could also consider a “multifidelity” extension where we weigh various possible gradient observations (including none at all) in light of their expected utility and the cost of acquisition/inference.

scaling of Gaussian process inference: § 9.1,
p. 201

⁷⁶ J. Wu et al. (2017). Bayesian Optimization with Gradients. *NeurIPS* 2017.

11.9 ENVIRONMENTAL VARIABLES

In many applications, the performance of a given system configuration depends on exogenous factors that are not under our control. For example, consider optimizing the parameters of a mobile robot’s gait to maximize some tradeoff of stability, efficiency, and speed. These objectives depend not only on the chosen parameters, but also on the nature of the environment, such as the composition and features of the surface to be traversed. These factors are not controllable by the robot at the time of performance.

In such a situation, we may seek to maximize the *expected* performance of a given configuration when presented with a random environment. Let us consider a function $g(x, \omega)$, where x as usual represents a configuration to be optimized, and ω represents the relevant parameters of the environment we wish to consider. In this context ω is called a *en-*

environmental parameter, ω

vironmental parameter or *environmental variable*. Now, if $p(\omega)$ encodes a distribution over environmental parameters we expect to encounter in practice, we may naturally seek to optimize the expected performance:

$$f(x) = \int g(x, \omega) p(\omega) d\omega. \quad (11.22)$$

Optimizing this objective presents a challenge: in general, the expectation with respect to ω cannot be evaluated directly but only estimated via repeated trials in different environments, and estimating $f(x)$ with some degree of precision may require numerous evaluations of $g(x, \omega)$. However, when g itself is expensive to evaluate – for example, if every trial requires manual manipulation of a robot and its environment before we can measure performance⁷⁷ – this may not be the most efficient approach, as we may waste significant resources shoring up our belief about suboptimal values.

Instead, when we can *control* the environment during optimization at will, we can gain some traction by designing a sequence of free parameter–environment pairs, potentially changing both configuration and environment in each iteration. The most direct way to design such an algorithm in our framework would be to model the environmental-conditional objective function g directly and define a utility function and policy with respect to this function, in light of the true environmental-marginal objective (11.22).

A Gaussian process model on g is particularly practical in this regard, as we may then use Bayesian quadrature to seamlessly estimate and quantify our uncertainty in (11.22) from any arbitrary set of observations. This approach was explored in depth by TOSCANO-PALMERIN and FRAZIER, who also provided an excellent review of the related literature.⁷⁸

11.10 INCREMENTAL OPTIMIZATION OF SEQUENTIAL PROCEDURES

In some applications, the objective function is determined by a *sequence* of dependent steps eventually producing its final value. If we have access to this sequential process and can model its progression, we may be able to accelerate optimization via shrewd “early stopping”: terminating evaluations still in progress when their final value can be forecasted with sufficient confidence.

Hyperparameter tuning presents one compelling example. Consider for example the optimization of neural network hyperparameters parameterized by θ .⁷⁹ The objective function in this setting is usually defined to be the value of some loss function ℓ (for example, validation error) after the network has been trained with the chosen hyperparameters. However, this training is an *iterative* procedure: if the network is parameterized by a vector of weights \mathbf{w} , then the objective function might be defined by

$$f(\theta) = \lim_{t \rightarrow \infty} \ell(\mathbf{w}_t; \theta), \quad (11.23)$$

the loss of the network after the weights have converged.⁸⁰ If we don’t

⁷⁷ M. TESCH et al. (2013). Expensive Function Optimization with Stochastic Binary Outcomes. *ICML 2013*.

Bayesian quadrature: § 2.6, p. 33

⁷⁸ S. TOSCANO-PALMERIN and P. I. FRAZIER (2018). Bayesian Optimization with Expensive Integrands. arXiv: 1803.08661 [stat.ML].

⁷⁹ We will use θ for the variable to be optimized here rather than x to be consistent with our previous discussion of Gaussian process hyperparameters in chapter 3.

⁸⁰ One could also consider early stopping in the training procedure as well, but this simple example is useful for exposition.

treat this objective function as a black box but take a peek inside, we may interpret it as the limiting value of the *learning curve* defined by the learning procedure’s loss at each stage of training.

The iterative nature of this objective function offers the opportunity for innovation in policy design. Learning curves typically exhibit fairly regular behavior, with the loss in each step of training generally falling over time until settling on its final value. This suggests we may be able to faithfully *extrapolate* the final value of a learning curve from the early stages of training; a toy example is presented in the margin. When possible, we may then be able to speed up optimization by not always training to convergence with every setting of the hyperparameters we explore. With this motivation in mind, several sophisticated methods for extrapolating learning curves have been developed, including carefully crafted parametric models^{81,82} and flexible Bayesian neural networks.⁸³

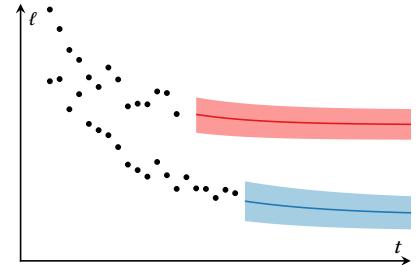
Exploiting the ability to extrapolate sequential objectives requires that we expand our action space (step 1) to allow fine-grained control over evaluation. One compelling scheme was proposed by SWERSKY et al.,⁸¹ who suggested maintaining a set of partial evaluations of the objective throughout optimization. In the context of our hyperparameter tuning example (11.23), we would maintain a set of thus-far investigated hyperparameters $\{\theta_i\}$, each accompanied by a sequence (of variable length) of thus-far evaluated weights $\{w_{t,i}\}$ and the associated losses $\{\ell_{t,i}\}$. Now, each action we design can either investigate a novel hyperparameter θ or extend an existing partial evaluation by one step. The authors dubbed this scheme *freeze–thaw Bayesian optimization*, as after each action we “freeze” the current evaluation, saving enough state such that we can “thaw” it later for further investigation if desired.

Regardless of modeling details, this scheme offers the potential for considerable savings when optimizing sequential objective functions. This idea of abandoning stragglers based on early progress is also the basis of the bandit-based *hyperband* algorithm.⁸⁴

11.11 NON-GAUSSIAN OBSERVATION MODELS AND ACTIVE SEARCH

Throughout this book, we have focused almost exclusively on the additive Gaussian noise observation model. There are good reasons for this: it is a reasonably faithful model of many systems and offers exact inference with Gaussian process models of the objective function. However, the assumption of Gaussian noise is not always warranted and may be fundamentally incompatible with some scenarios.

Fortunately, the decision-theoretic core of most Bayesian optimization approaches does not make any assumptions regarding our model of the objective function or our observations of it, and with some care the utility functions we developed for optimization can be adapted for virtually any scenario. Further, there are readily available pathways for incorporating non-Gaussian observation models into Gaussian process objective function models, so we do not need to abandon that rich model class in order to use alternatives.



Learning curve extrapolation from initial training results. At this point, we may already wish to abandon the red option.

- 81 K. SWERSKY et al. (2014). Freeze–Thaw Bayesian Optimization. arXiv: 1406 . 3896 [stat.ML].
- 82 T. DOMHAN et al. (2015). Speeding up Automatic Hyperparameter Optimization of Deep Neural Networks by Extrapolation of Learning Curves. *IJCAI 2015*.
- 83 A. KLEIN et al. (2017). Learning Curve Prediction with Bayesian Neural Networks. *ICLR 2017*.
- 84 L. LI et al. (2018b). Hyperband: A Novel Bandit-Based Approach to Hyperparameter Optimization. *Journal of Machine Learning Research* 18(185):1–52.

freeze–thaw Bayesian optimization

- decision theory for optimization: chapter 5, p. 87

utility functions for optimization: chapter 6, p. 109

GP inference with non-Gaussian observation models: § 2.8, p. 35

- 85 A. SHAH et al. (2014). Student-*t* Processes as Alternatives to Gaussian Processes. *AISTATS 2014*.
- 86 R. MARTINEZ-CANTIN et al. (2018). Practical Bayesian optimization in the presence of outliers. *AISTATS 2018*.
- 87 We explored this possibility at length in § 2.8.
- 88 M. TESCH et al. (2013). Expensive Function Optimization with Stochastic Binary Outcomes. *ICML 2013*.
- 89 Note that the simple reward utility function (6.3) we have been working with can be written in this exact form if we assume any additive observation noise has zero mean.
- 90 N. HOULSBY et al. (2012). Collaborative Gaussian Processes for Preference Learning. *NeurIPS 2012*.
- 91 E. BROCHU et al. (2015). Active Preference Learning with Discrete Choice Data. *NeurIPS 2007*.
- 92 J. GONZÁLEZ et al. (2017). Preferential Bayesian Optimization. *ICML 2017*.

- 93 R. GARNETT et al. (2012). Bayesian Optimal Active Search and Surveying. *ICML 2012*.

virtual screening: p. 310
modeling observations

Sequential optimization with non-Gaussian observation models

A decision-theoretic approach to optimization entails first selecting an objective function model $p(f)$ and observation model $p(y | x, \phi)$, which together are sufficient to derive the predictive distribution $p(y | x, \mathcal{D})$ relevant to every sequential decision made during optimization. After selecting a utility function $u(\mathcal{D})$, we may then follow the iterative procedure developed in chapter 5 to derive a policy.

This abstract approach has been realized in several specific settings. For example, both Student-*t* processes⁸⁵ and the Student-*t* observation model⁸⁶ have been explored to develop Bayesian optimization routines that are robust to the presence of outliers.⁸⁷

TESCH et al. explored the use of expected improvement for optimization from binary success/failure indicators, motivated by optimizing the probability of success of a robotic platform operating in an uncertain environment.⁸⁸ Here the utility function was taken to be this success probability maximized over the observed locations:

$$u(\mathcal{D}) = \max_x \mathbb{E}[y | x, \mathcal{D}] = \max_x \Pr(y = 1 | x, \mathcal{D}),$$

where for this expression we have assumed binary outcomes $y \in \{0, 1\}$ with $y = 1$ interpreted as indicating success.⁸⁹ The authors then derived a policy for this setting via one-step lookahead.

Another setting involving binary (or categorical) feedback is in optimizing human preferences, such as in A/B testing or user modeling. Here we might seek to optimize user preferences by repeatedly presenting a panel of options and asking for the most preferred item. HOULSBY et al. described a convenient reduction from preference learning to classification for Gaussian processes that allows the immediate use of standard policies such as expected improvement,^{88, 90, 91} although more sophisticated policies have also been proposed specifically for this setting.⁹²

Active search

GARNETT et al. introduced *active search* as a simple model of scientific discovery in a *discrete* domain $\mathcal{X} = \{x_i\}$.⁹³ In active search, we assume that among these points is hidden a rare, valuable subset exhibiting desirable properties for the task at hand. Given access to an oracle that can – at significant cost – determine whether an identified point belongs to the sought after class, the problem of active search is to design a sequence of experiments seeking to maximize the number of discoveries in a given budget. A motivating application is drug discovery, where the domain would represent a list of candidate molecules to search for those rare examples exhibiting significant binding activity with a chosen biological target. As the space of candidates is expansive and the cost of even virtual screening is nontrivial, intelligent experimental design has the potential to greatly improve the rate of discovery.

To derive an active search policy in our framework, we must first model the observation process and determine a suitable utility function.

The former requires consideration of the nuances of a given situation, but we may provide a barebones construction that is already sufficient to be of practical and theoretical interest. Given a discrete domain \mathcal{X} , we assume there is some identifiable subset $\mathcal{V} \subset \mathcal{X}$ of valuable points we wish to recover. We associate with each point $x \in \mathcal{X}$ a binary label $y = [x \in \mathcal{V}]$ indicating whether x is valuable ($y = 1$) or not ($y = 0$). A natural observation model is then to assume that selecting a point x for investigation reveals this binary label y in response.⁹⁴ Finally, we may define a natural utility function for active search by assuming that, all other things being held equal, we prefer a dataset containing more valuable points to one with fewer:

$$u(\mathcal{D}) = \sum_{x \in \mathcal{D}} y. \quad (11.24)$$

This is simply the cumulative reward utility (6.7), which here can be interpreted as counting the number of valuable points discovered.⁹⁵

To proceed with the Bayesian decision-theoretic approach, we must build a model for the uncertain elements inherent to each decision. Here the primary object of interest is the predictive posterior distribution $\Pr(y = 1 | x, \mathcal{D})$, the posterior probability that a given point x is valuable. We may build this model in any number of ways, for example by combining a Gaussian process prior on a latent function with an appropriate choice of observation model.⁹⁶

Equipped with a predictive model, deriving the optimal policy is a simple exercise. To begin, the one-step marginal gain in utility (5.8) is

$$\alpha_1(x; \mathcal{D}) = \Pr(y = 1 | x, \mathcal{D});$$

that is, the optimal one-step decision is to greedily maximize the probability of success. Although this is a simple (and somewhat obvious) policy that can perform well in practice, theoretical and empirical study on active search has established that massive gains can be had by adopting less myopic policies.

On the theoretical side, GARNETT et al. demonstrated by construction that the expected performance of *any* lookahead approximation can be exceeded by any arbitrary amount by extending the lookahead horizon even a single step.⁹³ This result was strengthened by JIANG et al., who showed – again by construction – that *no* policy that can be computed in time polynomial in $|\mathcal{X}|$ can approximate the performance of the optimal policy within any constant factor.⁹⁷ Thus the optimal active search policy is not only hard to compute, but also hard to approximate. These theoretical results, which rely on somewhat unnatural adversarial constructions, have been supported by empirical investigations on real-world data as well.^{93,97} For example, GARNETT et al. demonstrated that simply using two-step instead of one-step lookahead can significantly accelerate virtual screening for drug discovery across a broad range of biological targets.⁹⁸

This is perhaps a surprising state of affairs given the success of one-step lookahead – and the relative lack of less myopic alternatives –

⁹⁴ Other situations may call for other approaches; for example, if value is determined by thresholding a continuous measurement, we may wish to model that continuous observation process explicitly.

utility function

⁹⁵ The assumption of a discrete domain is to avoid repeatedly observing effectively the same point to trivially “max out” this score.

posterior predictive probability,
 $\Pr(y = 1 | x, \mathcal{D})$

⁹⁶ C. E. RASMUSSEN and C. K. I. WILLIAMS (2006). *Gaussian Processes for Machine Learning*. MIT Press. [chapter 3]

⁹⁷ S. JIANG et al. (2017). Efficient Nonmyopic Active Search. *ICML* 2017.

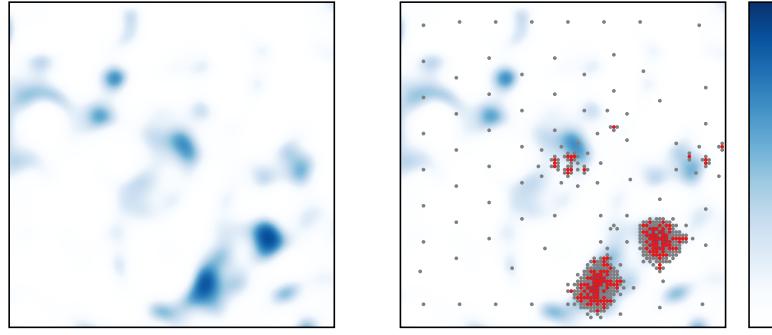
cost of computing optimal policy: § 5.3, p. 99

⁹⁸ R. GARNETT et al. (2015). Introducing the ‘active search’ method for iterative virtual screening. *Journal of Computer-Aided Molecular Design* 29(4):305–314.

Figure 11.17: A demonstration of active search. Left: a 2d domain colored according to the probability of a positive result. Right: 500 points chosen by a nonmyopic active search policy reflecting consideration of exploration versus exploitation. The observations colored red were positive and those colored gray were negative.

moving beyond one-step lookahead: § 7.10,
p. 150

batch rollout: § 5.3, p. 103



for traditional optimization. We can resolve this discrepancy by noting that utility functions used in that setting, such as the simple reward (6.3), inherently exhibit decreasing marginal gains as they are effectively bounded by the global maximum. Further, such a utility tends to remain relatively constant throughout optimization, punctuated by brief but significant increases when a new local maximum is discovered. On the other hand, for the cumulative reward (11.24), *every* observation has the potential to increase the utility by exactly one unit. As a result, every observation is on equal footing in terms of potential impact, and there is increased pressure to consider the entire search trajectory when designing each observation.

One thread of research on active search has focused on developing efficient, yet nonmyopic policies grounded in approximate dynamic programming, such as lookahead beyond one step.⁹³ JIANG et al. proposed an significantly less myopic alternative policy based on batch rollout.⁹⁷ The key observation is that we may construct the one-step optimal batch observation of size k by computing the posterior predictive probability $\Pr(y = 1 | x, \mathcal{D})$ for the unlabeled points, sorting, and taking the top k ; this is a consequence of linearity of expectation and utility (11.24). With this, we may realize an efficient batch rollout policy for horizon τ by maximizing the acquisition function

$$\Pr(y = 1 | x, \mathcal{D}) + \mathbb{E}_y \left[\sum'_{\tau=1} \Pr(y' = 1 | x', \mathcal{D}_1) | x, \mathcal{D} \right]. \quad (11.25)$$

Here the primed-sum notation $\sum'_{\tau=1}$ indicates the sum of the top- $(\tau - 1)$ values over the unlabeled data – the expected utility of the optimal batch observation consuming the remaining budget.⁹⁹ In experiments, this policy showed interesting emergent behavior: it *underperforms* lookahead policies in the early stages of search due to significant investment in early exploration.

Figure 11.17 illustrates active search in a 2d domain, which for the purposes of this demonstration was discretized into a 100×100 grid. The example is constructed so that a small number of discrete regions yield valuable items, which must be efficiently uncovered for a successful search. The right-hand panel shows a sequence of 500 observations designed by iteratively maximizing (11.25); their distribution clearly reflects

⁹⁹ The cost of computing this acquisition function is $\mathcal{O}(n^2 \log n)$, where $n = |\mathcal{X}|$; this is roughly the same order as the two-step expected marginal gain, $\mathcal{O}(n^2)$.

both exploration of the domain and exploitation of the most fruitful regions. The rate of discovery for the active search policy was approximately 3.8 times greater than would be expected from random search.

VANCHINATHAN et al. considered an expanded setting they dubbed *adaptive valuable item discovery* (AVID), for which active search is a special case. Here items may have nonnegative, *continuous* values and the cumulative reward utility (11.24) was augmented with a term encouraging diversity among the selected items.¹⁰⁰ The authors proposed an algorithm called GP-SELECT for this setting based on an acquisition function featuring two terms: a standard upper confidence bound score (8.25) and a term encouraging diversity related to *determinantal point processes*.¹⁰¹ The authors were further able to establish theoretical regret bounds for this algorithm under standard assumptions on the complexity of the value function.

¹⁰⁰ H. P. VANCHINATHAN et al. (2015). Discovering Valuable Items from Massive Data. *KDD 2015*.

¹⁰¹ A. KULESZA and B. TASKAR (2012). Determinantal Point Processes for Machine Learning. *Foundations and Trends in Machine Learning* 5(2–3):123–286.

12

A BRIEF HISTORY OF BAYESIAN OPTIMIZATION

In this chapter we provide a historical survey of the ideas underpinning Bayesian optimization, including important mathematical precedents. We also document the progression of major ideas in Bayesian optimization, from its first appearance in 1962 to the present day. A major goal will be to identify the point of introduction of prominent Bayesian optimization policies, as well as notable instances of subsequent reintroduction when ideas were forgotten and later rediscovered.

12.1 HISTORICAL PRECURSORS AND OPTIMAL DESIGN

The Bayesian approach to optimization is founded on a simple premise: experiments should be designed with purpose, both guided by our knowledge and aware of our ignorance. The optimization policies built on this principle can be understood as statistical manifestations of rational inquiry, where we design a sequence of experiments to systematically reveal the maximal value attained by the system of interest.

Statistical approaches to experimental design have a long history, with the earliest examples appearing over 200 years ago.^{1,2} A landmark early contribution was SMITH’s 1918 dissertation,³ which considered experimental design for polynomial regression models to minimize a measure of predictive uncertainty. Shortly thereafter, FISHER published a hugely influential guide to statistical experimental design based on his experience analyzing crop experiments at Rothamsted Experimental Station,⁴ which was instrumental in shaping modern statistical practice. These works served to establish the field of *optimal design*, an expansive subject which has now enjoyed a century of study. Numerous excellent references are available.^{5,6}

Early work in optimal design did not consider the possibility of adaptively designing a *sequence* of experiments, an essential feature of Bayesian optimization. Instead, the focus was optimizing fixed designs to minimize some measure of uncertainty when performing inference with the resulting data. This paradigm is practical when experiments are extremely time consuming but can easily run in parallel, such as the agricultural experiments studied extensively by FISHER. The most common goals considered in classical optimal design are accurate estimation of model parameters and confident prediction at unseen locations; these goals are usually formulated in terms of optimizing some statistical criterion as a function of the design.⁷ However, in 1941, HOTELLING notably studied experimental designs for estimating the location of the maximum of an unknown function in this nonadaptive setting.⁸ This was perhaps the first rigorous treatment of batch optimization.

12.2 SEQUENTIAL ANALYSIS AND BAYESIAN EXPERIMENTAL DESIGN

Concentrated study of sequential experiments began during World War II with WALD, who pioneered the field of *sequential analysis*. The seminal

¹ J. D. GERGONNE (1815). Application de la méthode des moindres quarrés à l’interpolation des suites. *Annales de Mathématiques pures et appliquées* 6:242–252.

² C. S. PEIRCE (1876). Note on the Theory of the Economy of Research. In: *Report of the Superintendent of the United States Coast Survey Showing the Progress of the Work for the Fiscal Year Ending with June, 1876*.

³ K. SMITH (1918). On the Standard Deviations of Adjusted and Interpolated Values of an Observed Polynomial Function and its Constants and the Guidance they Give Towards a Proper Choice of the Distribution of Observations. *Biometrika* 12(1–2):1–85.

⁴ R. A. FISHER (1935). *The Design of Experiments*. Oliver and Boyd.

⁵ G. E. P. BOX et al. (2005). *Statistics for Experimenters: Design, Innovation, and Discovery*. John Wiley & Sons.

⁶ D. C. MONTGOMERY (2019). *Design and Analysis of Experiments*. John Wiley & Sons.

⁷ These criteria often have alphabetic names: A-optimality, D-optimality, v-optimality, etc.

⁸ H. HOTELLING (1941). Experimental Determination of the Maximum of a Function. *The Annals of Mathematical Statistics* 12(1):20–45.

- 9 A. WALD (1945). Sequential Tests of Statistical Hypotheses. *The Annals of Mathematical Statistics* 16(2):117–186.
- 10 A. WALD (1947). *Sequential Analysis*. John Wiley & Sons.
- 11 M. FRIEDMAN and L. J. SAVAGE (1947). Planning Experiments Seeking Maxima. In: *Selected Techniques of Statistical Analysis for Scientific and Industrial Research, and Production and Management Engineering*.
- 12 H. HOTELLING (1941). Experimental Determination of the Maximum of a Function. *The Annals of Mathematical Statistics* 12(1):20–45.
- 13 G. E. P. BOX and K. B. WILSON (1951). On the Experimental Attainment of Optimum Conditions. *Journal of the Royal Statistical Society Series B (Methodological)* 13(1):1–45.
- 14 G. E. P. BOX (1954). The Exploration and Exploitation of Response Surfaces: Some General Considerations and Examples. *Biometrics* 10(1): 16–60.
- 15 G. E. P. BOX and P. V. YOULE (1954). The Exploration and Exploitation of Response Surfaces: An Example of the Link Between the Fitted Surface and the Basic Mechanism of the System. *Biometrics* 11(3):287–323.
- 16 H. CHERNOFF (1959). Sequential Design of Experiments. *The Annals of Mathematical Statistics* 30(3):755–770.
- 17 H. CHERNOFF (1972). *Sequential Analysis and Optimal Design*. Society for Industrial and Applied Mathematics.
- 18 J. BATHER (1996). A Conversation with Herman Chernoff. *Statistical Science* 11(4):335–350.

result was a hypothesis testing procedure for sequentially gathered data that can terminate dynamically once the result is known with sufficient confidence. Such sequential tests can be significantly more data efficient than tests requiring an a priori fixed sample size. The potential benefit of WALD's work to the war effort was immediately recognized, and his research was classified and only published after the war.^{9,10} The introduction of sequential analysis would kick off multiple parallel lines of investigation, leading to both Bayesian sequential experimental design and multi-armed bandits, and eventually all modern approaches to Bayesian optimization.

The success of sequential analysis lead numerous researchers to investigate *sequential experimental design* in the following years. Sequential experimental design for optimization was proposed by FRIEDMAN and SAVAGE as early as 1947.¹¹ They argued that nonadaptive optimization procedures, as considered earlier by HOTELLING,¹² can be wasteful as many experiments may be squandered needlessly exploring suboptimal regions. Instead, FRIEDMAN and SAVAGE suggested sequential optimization could be significantly more efficient, as poor regions of the domain can be quickly discarded while more time is spent exploiting more promising areas. Their proposed algorithm was a simple procedure entailing successive axis-aligned line searches, optimizing each input variable in turn while holding the others fixed – what would now be known as *cyclic coordinate descent*.

BOX greatly expounded on these ideas, working for years alongside a chemist (WILSON) experimentally optimizing chemical processes as a function of environmental and process parameters, for example, optimizing product yield as a function of reactant concentrations.^{13,14} BOX advocated the method of steepest ascent during early stages of optimization rather than the “one factor at a time” heuristic described by FRIEDMAN and SAVAGE, pointing out that the latter method is prone to becoming stuck in ridge-shaped features of the objective. BOX and YOULE also provided insightful commentary on how the process of optimization, and in particular geometric features of the objective function surface, may lead the experimenter to a greater fundamental understanding of underlying physical processes.¹⁵

In the following years, researchers developed general methods for sequential experimental design targeting a broad range of experimental goals. An early pioneer was CHERNOFF, a student of WALD's, who extended sequential analysis to the adaptive setting and provided asymptotically optimal procedures for sequential hypothesis testing.¹⁶ He also wrote a survey of this early work,¹⁷ and would eventually remark in an interview:¹⁸

Although I regard myself as non-Bayesian, I feel in sequential problems it is rather dangerous to play around with non-Bayesian procedures.

Another important contribution around this time was the reintroduction of multi-armed bandits by ROBBINS, which would quickly explode

into a massive body of literature. We will return to this line of work momentarily.

The Bayesian approach to sequential experimental design was formalized shortly after CHERNOFF's initial work. Authors such as RAIFFA and SCHLAIFER¹⁹ and LINDLEY²⁰ promoted a general approach based on Bayesian decision theory, wherein the experimenter selects a utility function reflecting their experimental goals and a model for reasoning about experimental outcomes given data, then designs each experiment to maximize the expected utility of the collected data. This is precisely the procedure we outlined in chapter 5. As we noted in our presentation, this framework yields theoretically *optimal* policies, but comes with an unwieldy computational burden.

12.3 THE RISE OF BAYESIAN OPTIMIZATION

By the early 1960s, the stage was set for Bayesian optimization, which could now be realized by appropriately adapting the now mature field of Bayesian experimental design.

KUSHNER was the first to seize the opportunity with a pair of papers on optimizing a one-dimensional objective observed with noise.^{21,22} All of the major ideas in modern Bayesian optimization were already in place in this initial work, including a Gaussian process model of the objective function and appealing to Bayesian decision theory to derive optimization policies. After dismissing the optimal policy (with respect to the global reward utility (6.5)) as “notoriously difficult”²¹ and “virtually impossible to compute,”²² KUSHNER suggests two alternatives “on the basis of heuristic or intuitive considerations”: maximizing an upper confidence bound (§ 7.8) and maximizing probability of improvement (§ 7.5), which share credit as the first Bayesian optimization policies to appear in the literature.²¹ The probability of improvement approach seems to have won his favor, and KUSHNER later provided extensive practical advice for realizing this method, including careful discussion of how the improvement threshold could be managed interactively throughout optimization by a human expert in the loop.²²

A significant body of literature on Bayesian optimization emerged in the Soviet Union following KUSHNER's seminal work. Many of these authors notably explored the promise of one-step lookahead for effective policy design, proposing and studying both expected improvement and the knowledge gradient for the first time. ŠALTENIS²³ was the first to introduce expected improvement (§ 7.3) in 1971.²⁴ This work contains an explicit formula for expected improvement for arbitrary Gaussian process models and the results of an impressive empirical investigation on a Soviet mainframe computer with objective functions in dimensions up to 32. ŠALTENIS concludes with the following observation:

The relatively large amounts of machine time spent in planning search and the complexity of the algorithm give us grounds to assume that the most effective sphere of application would be mul-

¹⁹ H. RAIFFA and R. SCHLAIFER (1961). *Applied Statistical Decision Theory*. Division of Research, Graduate School of Business Administration, Harvard University.

²⁰ D. V. LINDLEY (1972). *Bayesian Statistics, A Review*. Society for Industrial and Applied Mathematics.

²¹ H. J. KUSHNER (1962). A Versatile Stochastic Model of a Function of Unknown and Time Varying Form. *Journal of Mathematical Analysis and Applications* 5(1):150–167.

²² H. J. KUSHNER (1964). A New Method of Locating the Maximum Point of an Arbitrary Multi-peak Curve in the Presence of Noise. *Journal of Basic Engineering* 86(1):97–106.

²³ Also transliterated SHALTYANIS.

²⁴ V. R. ŠALTENIS (1971). One Method of Multiextremum Optimization. *Avtomatika i Vychislitel'naya Tekhnika (Automatic Control and Computer Sciences)* 5(3):33–38.

tiextremum target functions whose determination involves major computational difficulties.

This remains the target domain of Bayesian optimization today.

Another prominent early contributor was MOCKUS,²⁵ who wrote a series of papers on Bayesian optimization in the 1970s^{26,27,28} and has written two books on the subject.^{29,30} Like KUSHNER, MOCKUS begins his presentation in these papers by outlining the optimal policy for maximizing the global reward utility (6.5), but rejects it as computationally infeasible. As a practical alternative, MOCKUS instead promotes what he calls the “one-stage approach,” that is, one-step lookahead.²⁶ As he had chosen the global reward utility, the resulting optimization policy is to maximize the knowledge gradient (§ 7.4).

This claim may give some readers pause, as MOCKUS’s work is frequently cited as the origin of *expected improvement* instead. However, this is inaccurate for multiple reasons. Expected improvement had been introduced by ŠALTENIS in 1971, one year before MOCKUS’s first contribution on Bayesian optimization. Indeed, MOCKUS was aware of and cited ŠALTENIS’s work. Further, the acquisition function MOCKUS describes is defined with respect to the global reward utility underlying the knowledge gradient,³¹ not the simple reward utility underlying expected improvement.

That said, the situation is slightly more subtle. MOCKUS also discusses two convenient choices of models on the unit interval for which the knowledge gradient happens to equal expected improvement: the Wiener process and the Ornstein–Uhlenbeck (ou) process. Both are rare examples of *Gauss–Markov processes*, whose Markovian property renders sequential inference particularly convenient, and the early work on Bayesian optimization was dominated by these models due to the extreme computational limitations at the time. MOCKUS also points out these models have a “special propert[y]”,³² namely, their Markovian nature ensures that the posterior mean is always maximized *at an observed location*, and thus the simple and global reward utilities coincide!

12.4 LATER REDISCOVERY AND DEVELOPMENT

The expected improvement and the knowledge gradient acquisition strategies were both reintroduced decades later when computational power had increased to the point that Bayesian optimization could be a practical approach for real problems.

SCHONLAU³³ and JONES et al.,³⁴ working together, proposed maximizing expected improvement for efficient global optimization in the context of the *design and analysis of computer experiments* (DACE).³⁵ Here the objective function represents the output of a computational routine and is assumed to be observed without noise. In addition to promoting expected improvement as a policy, JONES et al. also provided extensive detail for practical implementation, including an insightful discussion on model validation and a branch-and-bound strategy for maximizing the

²⁵ Often transliterated MOČKUS in early work.

²⁶ J. MOCKUS (1972). Bayesian Methods of Search for an Extremum. *Avtomatika i Vychislitel'naya Tekhnika (Automatic Control and Computer Sciences)* 6(3):53–62.

²⁷ J. MOCKUS (1974). On Bayesian Methods for Seeking the Extremum. *Optimization Techniques IFIP Technical Conference*.

²⁸ J. MOCKUS et al. (1978). The Application of Bayesian Methods for Seeking the Extremeum. In: *Towards Global Optimization* 2.

²⁹ J. MOCKUS (1989). *Bayesian Approach to Global Optimization: Theory and Applications*. Kluwer Academic Publishers.

³⁰ J. MOCKUS et al. (2010). *Bayesian Heuristic Approach to Discrete and Global Optimization: Algorithms, Visualization, Software, and Applications*. Kluwer Academic Publishers.

³¹ See equation 8 in citation 26 above.

³² J. MOCKUS (1972). Bayesian Methods of Search for an Extremum. *Avtomatika i Vychislitel'naya Tekhnika (Automatic Control and Computer Sciences)* 6(3):53–62. [equations 36–37]

³³ M. SCHONLAU (1997). Computer Experiments and Global Optimization. Ph.D. thesis. University of Waterloo.

³⁴ D. R. JONES et al. (1998). Efficient Global Optimization of Expensive Black-Box Functions. *Journal of Global Optimization* 13(4):455–492.

³⁵ J. SACKS et al. (1989). Design and Analysis of Computer Experiments. *Statistical Science* 4(4): 409–435.

expected improvement acquisition function for a certain class of Gaussian process models. The knowledge gradient was picked up again and further developed by FRAZIER and POWELL,³⁶ who coined the name and studied the policy in the discrete and independent (bandit-like) setting, and SCOTT et al.,³⁷ who adopted the policy for continuous optimization.

Three years after reintroducing expected improvement, JONES wrote an extensive survey of then-current Bayesian optimization policies.³⁸ Despite the now-pervasive nature of expected improvement, it is striking that JONES actually promotes maximizing the probability of improvement as the most promising policy for Bayesian optimization throughout this survey. His concern with expected improvement was a potential lack of robustness if the objective function model is misspecified. His proposed alternative was maximizing the probability of improvement over a wide range of improvement targets and evaluating the objective function in parallel,³⁹ which he regarded as more robust when practical.

The concept of mutual information (§ 7.6) first appeared with SHANNON's introduction of information theory,⁴⁰ where it was called the *channel capacity* and served as a measure of the amount of information that could be transferred effectively over a noisy communication channel. LINDLEY later reinterpreted mutual information as the expected information gained by a proposed experiment and suggested maximizing this quantity as an effective means of general Bayesian experimental design.⁴¹

The application of this information-theoretic framework to Bayesian optimization was first proposed by VILLEMONTEIX et al.⁴² and later independently by HENNIG and SCHULER,⁴³ the latter of whom coined the now-prominent term *entropy search*. Both of these initial investigations considered maximizing the mutual information between a measurement and the location of the global optimum x^* (§ 7.6), using the formulation in (7.14). HERNÁNDEZ-LOBATO et al. later proposed a different set of approximations based on the equivalent formulation in (7.15) under the name *predictive entropy search*, as the key quantity is the expected reduction in entropy for the predictive distribution.⁴⁶ Both HOFFMAN and GHAHRAMANI⁴⁴ and WANG and JEJELKA⁴⁵ later pursued maximizing the mutual information with the value of the global optimum f^* (§ 7.6). These developments occurred contemporaneously and independently.

Prior to these algorithms designed to reduce the *entropy* of the value of the global maximum, there were occasional efforts to minimize some other measure of dispersion in this quantity. For example, CALVIN studied an algorithm for optimization on the unit interval $\mathcal{X} = [0, 1]$ wherein the *variance* of $p(f^* | \mathcal{D})$ was greedily minimized via one-step lookahead.⁴⁶ Here the model was again the Wiener process, which has the remarkable property that the distribution of $p(f^* | \mathcal{D})$ (and its variance) is analytically tractable.⁴⁷ The Wiener and OU processes are among the *only* nontrivial Gaussian processes with this property.⁴⁸

No history of Bayesian optimization would be complete without mentioning the role hyperparameter tuning has played in driving its recent development. With the advent of deep learning, the early 2010s

- 36 P. FRAZIER and W. POWELL (2007). The Knowledge Gradient Policy for Offline Learning with Independent Normal Rewards. *ADPRL 2007*.
- 37 W. SCOTT et al. (2011). The Correlated Knowledge Gradient for Simulation Optimization of Continuous Parameters Using Gaussian Process Regression. *SIAM Journal on Optimization* 21(3):996–1026.
- 38 D. R. JONES (2001). A Taxonomy of Global Optimization Methods Based on Response Surfaces. *Journal of Global Optimization* 21(4):345–383.
- 39 See § 7.5, p. 134 for a discussion of this proposal.
- 40 C. E. SHANNON (1948). A Mathematical Theory of Communication. *The Bell System Technical Journal* 27(3):379–423.
- 41 D. V. LINDLEY (1956). On a Measure of the Information Provided by an Experiment. *The Annals of Mathematical Statistics* 27(4):986–1005.
- 42 J. VILLEMONTEIX et al. (2009). An informational approach to the global optimization of expensive-to-evaluate functions. *Journal of Global Optimization* 44(4):509–534.
- 43 P. HENNIG and C. J. SCHULER (2012). Entropy Search for Information-Efficient Global Optimization. *Journal of Machine Learning Research* 13(Jun):1809–1837.
- 44 M. W. HOFFMAN and Z. GHAHRAMANI (2015). Output-Space Predictive Entropy Search for Flexible Global Optimization. *Bayesian Optimization Workshop, NeurIPS 2015*.
- 45 Z. WANG and S. JEJELKA (2017). Max-value Entropy Search for Efficient Bayesian Optimization. *ICML 2017*.
- 46 J. M. CALVIN (1993). Consistency of a Myopic Bayesian Algorithm for One-Dimensional Global Optimization. *Journal of Global Optimization* 3(2):223–232.
- 47 In fact, the joint density $p(x^*, f^* | \mathcal{D})$ is analytically tractable for the Wiener process with drift, which forms the posterior on each subinterval subdivided by the data:
- 48 L. A. SHEPP (1979). The Joint Density of the Maximum and its Location for a Wiener Process with Drift. *Journal of Applied Probability* 16(2): 423–427.
- 49 R. J. ADLER and J. E. TAYLOR (2007). *Random Fields and Geometry*. Springer-Verlag. [chapter 4, footnotes 1–2]

saw the rise of extraordinarily complex learning algorithms trained on extraordinarily large datasets. The great expense of training these models created unprecedented demand for efficient hyperparameter tuning to fuel the rapid development in the area. One could not ask for a more-perfect fit for Bayesian optimization!

Interest in Bayesian optimization for hyperparameter tuning was kicked off in earnest by SNOEK et al., who reported a dramatic improvement in performance when tuning a convolutional neural network via Bayesian optimization, even compared with carefully hand-tuned hyperparameters.⁴⁹ This served as a watershed moment for Bayesian optimization, leading to an explosion of interest and sustained effort from the machine learning in the following years – although this history began in 1815, over half of the works cited in this book are from after 2012!

12.5 MULTI-ARMED BANDITS TO INFINITE-ARMED BANDITS

We have now covered the evolution of decision-theoretic Bayesian optimization policies from WALD’s introduction of sequential analysis in 1945 to the state-of-the-art. Alongside these developments, a rich and expansive body of literature was forming on the multi-armed problem. A complete survey of this work would be out of this book’s scope; however, we can point the interested reader to comprehensive surveys.^{50,51} Both contain excellent bibliographic notes and combined serve as a indispensable guide to the literature on multi-armed bandits (§ 7.7). Our goal in the following will be to cover developments that directly influenced the evolution of Bayesian optimization.

THOMPSON was the first to seriously study the possibility of sequential experiments in the context of medical treatment;^{52,53} this work pre-dated WALD’s work by a decade. THOMPSON considered a model scenario where there are two possible treatments a doctor may prescribe for a disease, but it is not clear which should be preferred. To determine the better treatment, we must undertake a clinical trial, assigning each treatment to several patients and assessing the outcome. Traditionally, a single preliminary experiment would be conducted, after which the apparently better-performing treatment would be adopted and the worse-performing treatment discarded. However, THOMPSON argued that one should *never* eliminate either treatment at any stage of investigation, even in the face of overwhelming evidence. Rather, he proposed a *perpetual* clinical trial modeled as what we now call a two-armed bandit: the possible treatments represent alternatives available to the clinician, and patient outcomes determine the rewards. We can now consider assigning treatments for a sequence of patients guided by our evolving knowledge, hoping to efficiently and confidently determine the optimal treatment.

To conduct such a sequential clinical trial effectively, THOMPSON proposed maintaining a belief over which treatment is better in light of available evidence and always selecting the nominally better treatment according to its posterior probability of superiority. This is the eponymous Thompson sampling policy (§ 7.9), which elegantly addresses

⁴⁹ J. SNOEK et al. (2012). Practical Bayesian Optimization of Machine Learning Algorithms. *NeurIPS 2012*.

⁵⁰ D. A. BERRY and B. FRISTEDT (1985). *Bandit Problems: Sequential Allocation of Experiments*. Chapman & Hall.

⁵¹ T. LATTIMORE and C. SZEPESVÁRI (2020). *Bandit Algorithms*. Cambridge University Press.

⁵² W. R. THOMPSON (1933). On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of Two Samples. *Biometrika* 25(3–4):285–294.

⁵³ W. R. THOMPSON (1935). On the Theory of Apportionment. *American Journal of Mathematics* 57(2):450–456.

the exploration–exploitation dilemma. The more evidence we have in favor of one treatment, the more often we prescribe it, exploiting the apparently better choice. Meanwhile, we maintain a diminishing but nonzero probability of selecting the other treatment, forcing continual exploration until the better treatment becomes obvious. In the long term, we will eventually assign the correct treatment to new patients with probability approaching certainty.

THOMPSON’s work was in retrospect groundbreaking, but its significance was perhaps not fully realized at the time. However, Thompson sampling for multi-armed bandits has recently enjoyed an explosion of attention due to impressive empirical performance^{54,55} and strong theoretical guarantees.^{56,57,58} The first direct application of Thompson sampling to Bayesian optimization is due to SHAHRIARI et al.,⁵⁹ who adopted an efficient approximation first proposed by HERNÁNDEZ-LOBATO et al. in the context of entropy search.⁶⁰

Although THOMPSON had introduced bandits in the early 1930s, concerted effort began with ROBBINS’s landmark 1952 reintroduction and analysis of the problem.⁶¹ This work introduced the modern formulation of multi-armed bandits presented in § 7.7, where an arbitrary, unknown reward distribution is associated with each arm and the agent seeks a policy to maximize the expected cumulative reward. ROBBINS explores this problem in the special case of a two-armed bandit with Bernoulli rewards, demonstrating that a simple adaptive policy (“switch on lose, stay on win”) can achieve better performance than nonadaptive or random policies. He also presents a family of policies that can achieve asymptotically optimal behavior for any reward distributions. These policies are defined by a simple mechanism that explicitly forces continual but decreasingly frequent exploration of both arms according to a pregenerated schedule, effectively – if crudely – balancing exploration and exploitation.

The simple policies proposed by ROBBINS eventually achieve near-optimal cumulative reward, but they are not very *efficient* in the sense of achieving that behavior quickly. ROBBINS returned to this problem three decades later to further address the issue of efficiency.⁶² LAI and ROBBINS introduced policies that dynamically trade off exploration and exploitation by maximizing an upper confidence bound on the reward distributions of the arms (§ 7.8), and demonstrated that these policies are both asymptotically optimal and efficient. AUER et al. later proved that bandit policies based on maximizing upper confidence bounds can provide strong guarantees not only asymptotically but also in the finite-budget case.⁶³ Interestingly, KUSHNER had proposed optimization policies based on maximizing an upper confidence bound of the objective function in the continuous case two decades earlier than LAI and ROBBINS’s work,²¹ although he did not prove any performance guarantees for this procedure and abandoned the idea in favor of maximizing probability of improvement. This optimization policy would be rediscovered several times, including by COX and JOHN⁶⁴ and by SRINIVAS et al.,⁶⁵ who established theoretical guarantees on the rate of convergence of this policy in the Gaussian process case, as discussed in § 10.4.

- 54 O. CHAPELLE and L. LI (2011). An Empirical Evaluation of Thompson Sampling. *NeurIPS 2011*.
- 55 O.-C. GRANMO (2010). Solving two-armed Bernoulli bandit problems using a Bayesian learning automaton. *International Journal of Intelligent Computing and Cybernetics* 3(2): 207–234.
- 56 S. AGRAWAL and N. GOYAL (2012). Analysis of Thompson Sampling for the Multi-armed Bandit Problem. *COLT 2012*.
- 57 D. RUSSO and B. VAN ROY (2014). Learning to Optimize via Posterior Sampling. *Mathematics of Operations Research* 39(4):1221–1243.
- 58 D. RUSSO and B. VAN ROY (2016). An Information-Theoretic Analysis of Thompson Sampling. *Journal of Machine Learning Research* 17(68):1–30.
- 59 B. SHAHRIARI et al. (2014). An Entropy Search Portfolio for Bayesian Optimization. arXiv: 1406.4625 [stat.ML].
- 60 J. M. HERNÁNDEZ-LOBATO et al. (2014). Predictive Entropy Search for Efficient Global Optimization of Black-box Functions. *NeurIPS 2014*.
- 61 H. ROBBINS (1952). Some Aspects of the Sequential Design of Experiments. *Bulletin of the American Mathematical Society* 58(5):527–535.
- 62 T. L. LAI and H. ROBBINS (1985). Asymptotically Efficient Adaptive Allocation Rules. *Advances in Applied Mathematics* 6(1):4–22.
- 63 P. AUER et al. (2002). Finite-time Analysis of the Multiarmed Bandit Problem. *Machine Learning* 47(2–3):235–256.
- 64 D. D. COX and S. JOHN (1992). A Statistical Method for Global Optimization. *SMC 1992*.
- 65 N. SRINIVAS et al. (2010). Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design. *ICML 2010*.

12.6 WHAT'S NEXT?

The mathematical foundations of Bayesian optimization, as presented in this book and chronicled in this chapter, are by now well established. However, Bayesian optimization continues to deliver impressive performance on a wide variety of problems, and the field continues to develop at a rapid pace in light of this success. What big challenges remain on the horizon? We speculate on some potential opportunities below.

Gaussian process modeling in high dimension: § 3.5, p. 61

moving beyond one-step lookahead: § 7.10,
p. 150

utility functions: § 6, p. 109

- Effective modeling of objective functions remains a challenge, especially in high dimension. Meanwhile, methods such as stochastic gradient descent routinely (locally) optimize objectives in millions of dimensions, and are “unreasonably effective” at doing so – all with very weak guidance. Can we bridge this gap, either by extending or refining approaches for Gaussian process modeling, or by exploring another model class?
- *Nonmyopic* policies that reach beyond one-step lookahead have shown impressive empirical performance in initial investigations. However, these policies have not yet been widely adopted, presumably due to their (sometimes significantly) greater computational cost compared to myopic alternatives. The continued development of *efficient*, yet nonmyopic policies remains a promising avenue of research.
- A guiding philosophy in Bayesian optimization has been to “take the human out of the loop” and hand over complete control of experimental design to an algorithm.⁶⁶ This paradigm has demonstrated remarkable success on “black box” problems, where the user has little understanding of the system being optimized. However, in settings such as scientific discovery, the user has a *deep* understanding of and intuition for the mechanisms driving the system of interest, and we should perhaps consider how to “bring them back into the loop.” One could imagine an ecosystem of *cooperative* tools that enable Bayesian optimization algorithms to benefit from user knowledge while facilitating the presentation of experimental progress and evolving model beliefs back to the users.
- In the author’s experience, many consumers of Bayesian optimization have experimental goals that are not perfectly captured by any of the common utility functions used in Bayesian optimization – for example, users may want to ensure adequate coverage of the domain or to find a *diverse* set of locations with high objective values. Although the space of Bayesian optimization policies is fairly crowded, there is still room for innovation. If the ecosystem of cooperative tools envisioned above were realized, we might consider tuning the utility function “on the fly” via interactive preference elicitation.
- Bayesian optimization is not a particularly user-friendly approach, as effective optimization requires careful modeling of the system of interest. However, model construction remains something of a “black art,” even in the machine-learning community, and thus considerable machine-learning expertise can be required to get the most out of this approach. Can we build turnkey Bayesian optimization systems that achieve acceptable performance even in the absence of clear prior beliefs?

⁶⁶ B. SHAHRIARI et al. (2016). Taking the Human Out of the Loop: A Review of Bayesian Optimization. *Proceedings of the IEEE* 104(1):148–175.

A

THE GAUSSIAN DISTRIBUTION

The *Gaussian* (or *normal*) distribution is a fundamental probability distribution in probabilistic modeling and inference. Gaussian distributions are frequently encountered in Bayesian optimization as they serve as the foundation of *Gaussian processes*, an infinite-dimensional extension appropriate for reasoning about unknown objective functions. In this chapter we provide a brief introduction to finite-dimensional Gaussian distributions and establish important properties referenced throughout this book. We will begin with the univariate (one-dimensional) case, then construct of the multivariate (vector-valued) case via linear transformations of univariate Gaussians.

A.1 UNIVARIATE GAUSSIAN DISTRIBUTION

The univariate Gaussian distribution on a random variable $x \in \mathbb{R}$ has two scalar parameters corresponding to its first two moments: $\mu = \mathbb{E}[x]$ specifies the mean (also median and mode) and serves as a location parameter, and $\sigma^2 = \text{var}[x]$ specifies the variance and serves as a scale parameter.

Probability density function and degenerate case

When the variance is nonzero, the distribution has the probability density function

$$\mathcal{N}(x; \mu, \sigma^2 > 0) = Z^{-1} \exp\left(-\frac{1}{2}z^2\right), \quad (\text{A.1})$$

where $Z = \sqrt{2\pi}\sigma$ is a normalization constant and z is the familiar z -score of x :

$$z = \frac{x - \mu}{\sigma}. \quad (\text{A.2})$$

This PDF is illustrated in the margin. The probability density is rapidly decreasing with the magnitude of the z -score, with for example 99.7% of the density lying in the interval $|z| \leq 3$, or $x \in (\mu \pm 3\sigma)$.

In the degenerate case $\sigma^2 = 0$, the distribution collapses to a point mass at the mean and a probability density function does not exist. We can express this case with the Dirac delta distribution:

$$\mathcal{N}(x; \mu, \sigma^2 = 0) = \delta(x - \mu). \quad (\text{A.3})$$

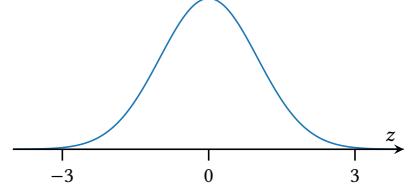
Standard normal distribution

The special case of zero mean and unit variance is called the *standard normal* distribution and enjoys a privileged role. We denote its density with the compact notation $\phi(x) = \mathcal{N}(x; 0, 1^2)$. The cumulative density function of the standard normal cannot be expressed in terms of elementary functions, but is such an important quantity that it also merits its

Gaussian processes: chapters 2–4, p. 15

mean, μ

variance, σ^2



A univariate Gaussian probability density function $\mathcal{N}(x; \mu, \sigma^2)$ as a function of the z -score.

degenerate case, $\sigma^2 = 0$

standard normal PDF, ϕ

standard normal CDF, Φ

own special notation:

$$\Phi(y) = \Pr(x < y) = \int_{-\infty}^y \phi(x) dx.$$

expressing arbitrary PDFs and CDFs in terms of the standard normal

We can write the PDF and CDF of an arbitrary nondegenerate Gaussian distribution $\mathcal{N}(x; \mu, \sigma^2)$ in terms of the standard normal PDF and CDF by appropriately rescaling and translating arguments to their z -scores (A.1):

$$p(x) = \frac{1}{\sigma} \phi\left(\frac{x - \mu}{\sigma}\right); \quad \Pr(x < y) = \Phi\left(\frac{y - \mu}{\sigma}\right),$$

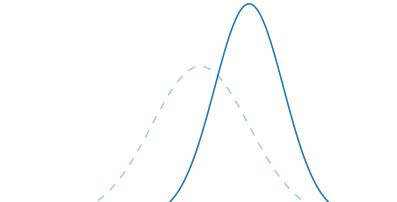
where the multiplicative factor in the PDF guarantees normalization.

Affine transformations

The family of univariate Gaussian distributions is closed under affine transformations, which simply translate and rescale the distribution and adjust its moments accordingly. If x has distribution $\mathcal{N}(x; \mu, \sigma^2)$, then the transformation $\xi = ax + b$ has distribution

$$p(\xi) = \mathcal{N}(\xi; a\mu + b, a^2\sigma^2).$$

The above results allow us to interpret *any* univariate Gaussian distribution as a translated and scaled copy of the standard normal after applying the transformation $x \mapsto \mu + \sigma x$. This process is illustrated in the margin, where a standard normal distribution is transformed via the mapping $x \mapsto 1 + x/\sqrt{2}$, resulting in a new Gaussian with increased mean $\mu = 1$ and decreased variance $\sigma^2 = \frac{1}{2}$. The PDF is appropriately translated and rescaled.



A standard normal PDF (dashed) and the PDF after applying the transformation $x \mapsto 1 + x/\sqrt{2}$ (solid).

A.2 MULTIVARIATE GAUSSIAN DISTRIBUTION

The multivariate Gaussian distribution extends the univariate case to an arbitrary random vector $\mathbf{x} \in \mathbb{R}^d$. We will provide an explicit construction of the multivariate Gaussian distribution as the result of applying an affine transformation to independent univariate standard normal random variables, extending the properties noted at the end of the previous section.

Standard multivariate normal and construction of general case

First we construct the *standard multivariate normal* distribution, represented by a random vector $\mathbf{z} \in \mathbb{R}^d$ whose entries are independent standard univariate normal random variables: $p(\mathbf{z}) = \prod_i \phi(z_i)$. It is clear from construction that the mean of this distribution is the zero vector and its covariance is the identity matrix:

$$\mathbb{E}[\mathbf{z}] = \mathbf{0}; \quad \text{cov}[\mathbf{z}] = \mathbf{I}, \tag{A.4}$$

and we will denote its density with $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$.

As before, this will serve as the basis of the general multivariate case by considering arbitrary affine transformations of this “standard” example. Suppose $\mathbf{x} \in \mathbb{R}^d$ is a vector-valued random variable and $\mathbf{z} \in \mathbb{R}^k$, $k \leq d$, is a k -dimensional standard multivariate normal vector. If we can write

$$\mathbf{x} = \boldsymbol{\mu} + \Lambda \mathbf{z} \quad (\text{A.5})$$

for some vector $\boldsymbol{\mu} \in \mathbb{R}^d$ and $d \times k$ matrix Λ , then \mathbf{x} has a multivariate normal distribution. We can compute its mean and covariance directly from (A.4–A.5):

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}; \quad \text{cov}[\mathbf{x}] = \Lambda \Lambda^\top = \Sigma.$$

This property completely characterizes the distribution. As in the univariate case, we can interpret every multivariate normal distribution as an affine transformation of a (possibly lower-dimensional) standard normal vector. We again parameterize this family by its first two moments: the mean vector $\boldsymbol{\mu}$ and the covariance matrix Σ . This covariance matrix is necessarily symmetric and *positive semidefinite*, which means all its eigenvalues are nonnegative. We can factor any such matrix as $\Sigma = \Lambda \Lambda^\top$ allowing us to recover the underlying transformation (A.5), although Λ need not be unique.

Probability density function and degenerate case

If Λ has full rank d , then the range of (A.5) is all of \mathbb{R}^d and a probability density function exists. This condition further implies that the covariance matrix Σ is *positive definite*; that is, its eigenvalues are strictly positive, implying its determinant is positive and the matrix is invertible. The distribution has a probability density function in this case analogous to the univariate PDF (A.1):

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = Z^{-1} \exp(-\frac{1}{2} \Delta^2). \quad (\text{A.6})$$

Here Z again represents a normalization constant:

$$Z = \sqrt{|2\pi\Sigma|} = (2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}, \quad (\text{A.7})$$

and Δ represents the *Mahalanobis distance*, a multivariate analog of the (absolute) z -score (A.2):

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}). \quad (\text{A.8})$$

It is easy to verify that these definitions are compatible with the univariate case when $d = 1$. Note in that case the condition of Σ being positive definite reduces to the previous condition for nondegeneracy, $\sigma^2 > 0$.

The dependence of the multivariate Gaussian density on \mathbf{x} is entirely through the value of the Mahalanobis distance Δ . To gain some geometric insight into the probability density, we can set this value to a constant and compute isoprobability contours. In the case of the *standard* multivariate

constructing general case via affine
transformations of standard normal vectors

parameterization in terms of first two
moments

mean vector and covariance matrix, $(\boldsymbol{\mu}, \Sigma)$
positive semidefinite

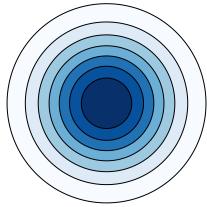
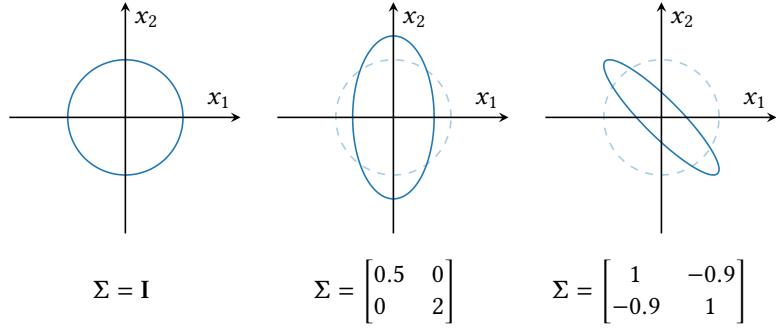
positive definite

normalization constant, Z

Mahalanobis distance, Δ

compatibility with univariate case

Figure A.1: Isoprobability contours $\Delta = 1$ for the standard bivariate normal distribution (left), a distribution with diagonal covariance, scaling the probability along the axes (middle), and a distribution with nonzero off-diagonal covariance, tilting the probability (right). The standard normal contour is shown for reference on the latter two examples.



Probability density and circular isoprobability contours of a standard bivariate normal distribution.

degenerate case, $|\Sigma| = 0$

- ¹ On this space, the PDF is similar to (A.6–A.8) but replaces the determinant with the pseudodeterminant and the inverse with the pseudoinverse. If $\Sigma = \mathbf{0}$, the distribution is a Dirac delta on μ .

Gaussian distribution, the Mahalanobis distance reduces to the normal Euclidean distance, and the set of points satisfying $\Delta = c > 0$ is then a sphere of radius c centered at the origin – see the illustration in the margin.

We can now understand the geometry of the general case $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$ purely in terms of the affine transformation (A.5) translating and warping the spherical distribution of the standard normal. Taking this view, the action of multiplying by Λ warps the isoprobability contours into ellipsoids, which are then translated to the new center $\boldsymbol{\mu}$. The standard theory of linear maps then gives further insight: the principal axes of the ellipsoids are given by the eigenvectors of Λ , and their axis semilengths are given by the eigenvalues scaled by c . See figure A.1 for an illustration.

The probability density function does not exist when Λ (and thus Σ) is rank-deficient, as the range of \mathbf{x} would then be restricted to the lower-dimensional affine subspace $\{\boldsymbol{\mu} + \Lambda \mathbf{z} \mid \mathbf{z} \in \mathbb{R}^d\}$ (A.5). However, it is still possible to define a probability density function in this degenerate case by restricting the support to this subspace.¹ This is analogous to the degenerate univariate case (A.3), where probability was restricted to the zero-dimensional subspace containing the mean only, $\{\boldsymbol{\mu}\}$.

Affine transformations

The multivariate Gaussian distribution has a number of convenient mathematical properties, many of which follow immediately from the characterization in (A.5). First it is obvious that any affine transformation of a multivariate normal distributed vector is also multivariate normal, as affine transformations are closed under composition. If $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$, then $\xi = \mathbf{Ax} + \mathbf{b}$ has distribution

$$p(\xi) = \mathcal{N}(\xi; \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\Sigma\mathbf{A}^\top). \quad (\text{A.9})$$

Further, if we apply this result with the transformation

$$\mathbf{x} \mapsto \begin{bmatrix} \mathbf{I} \\ \mathbf{A} \end{bmatrix} \mathbf{x} + \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} = \begin{bmatrix} \mathbf{x} \\ \xi \end{bmatrix},$$

we can see that \mathbf{x} and ξ in fact have a *joint* Gaussian distribution:

$$p(\mathbf{x}, \xi) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \xi \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma} & \boldsymbol{\Sigma}\mathbf{A}^\top \\ \mathbf{A}\boldsymbol{\Sigma} & \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top \end{bmatrix}\right). \quad (\text{A.10})$$

joint distribution with affine transformations

Sampling

The characterization of the multivariate normal in terms of affine transformations of standard normal random variables (A.5) also suggests a simple algorithm for drawing samples from the distribution. Given an arbitrary multivariate normal distribution $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$, we first factor the covariance as $\boldsymbol{\Sigma} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top$ where $\boldsymbol{\Lambda}$ has size $d \times k$; when $\boldsymbol{\Sigma}$ is positive definite, the *Cholesky decomposition* is the canonical choice. We now sample a k -dimensional standard normal vector \mathbf{z} by sampling each entry independently from a univariate standard normal; routines for this task are readily available. Finally, we transform this vector appropriately to provide a sample of \mathbf{x} : $\mathbf{z} \mapsto \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{z} = \mathbf{x}$. This procedure entails one-time $\mathcal{O}(d^3)$ work to compute $\boldsymbol{\Lambda}$ (which can be reused), followed by $\mathcal{O}(d^2)$ work to produce each sample.

Marginalization

Often we will have a vector \mathbf{x} with a multivariate Gaussian distribution, but only be interested in reasoning about a subset of its entries. Suppose $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ and partition the vector into two components:²

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}. \quad (\text{A.11})$$

We partition the mean vector and covariance matrix in the same way:

$$p(\mathbf{x}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}\right). \quad (\text{A.12})$$

Now writing the subvector \mathbf{x}_1 as $\mathbf{x}_1 = [\mathbf{I}, \mathbf{0}] \mathbf{x}$ and applying the affine property (A.9), we have:

$$p(\mathbf{x}_1) = \mathcal{N}(\mathbf{x}_1; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}). \quad (\text{A.13})$$

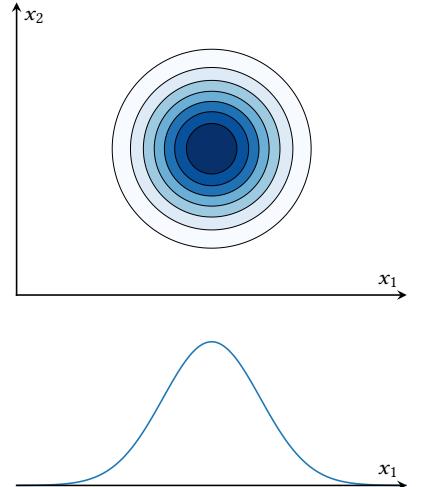
That is, to derive the marginal distribution of \mathbf{x}_1 we simply pick out the corresponding entries of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$.

Conditioning

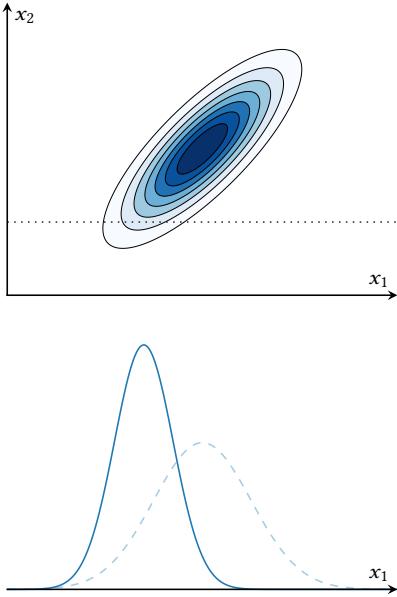
Multivariate Gaussian distributions are also closed under conditioning on the values of given entries. Suppose again that $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ and partition \mathbf{x} , $\boldsymbol{\mu}$, and $\boldsymbol{\Sigma}$ as before (A.11–A.12). Suppose now that we learn the exact value of the subvector \mathbf{x}_2 . The posterior on the remaining entries $p(\mathbf{x}_1 | \mathbf{x}_2)$ remains Gaussian, with distribution

$$p(\mathbf{x}_1 | \mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1; \boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{11|2}).$$

² We can first permute \mathbf{x} if required; as a linear transformation, this will simply permute the entries of $\boldsymbol{\mu}$ and the rows/columns of $\boldsymbol{\Sigma}$.



A bivariate Gaussian PDF $p(x_1, x_2)$ (top) and the Gaussian marginal $p(x_1)$ (bottom) (A.13).



The PDFs of a bivariate Gaussian $p(x_1, x_2)$ (top) and the conditional distribution $p(x_1 | x_2)$ given the value of x_2 marked by the dotted line (bottom) (A.14). The prior marginal distribution $p(x_1)$ is shown for reference; the observation decreased both the mean and standard deviation.

independent case

The posterior mean and covariance take the form of updates to the prior moments in light of the revealed information:

$$\boldsymbol{\mu}_{1|2} = \boldsymbol{\mu}_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2); \quad \Sigma_{11|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}. \quad (\text{A.14})$$

The mean is adjusted by an amount dependent on

1. the covariance between \mathbf{x}_1 and \mathbf{x}_2 , Σ_{12} ,
2. the uncertainty in \mathbf{x}_2 , Σ_{22} , and
3. the deviation of the observed values from the prior mean, $(\mathbf{x}_2 - \boldsymbol{\mu}_2)$.

Similarly, the uncertainty in \mathbf{x}_1 , Σ_{11} , is reduced by an amount dependent on factors 1–2. Notably, the correction to the covariance matrix does *not* depend on the observed values. Note that if \mathbf{x}_1 and \mathbf{x}_2 are independent, then $\Sigma_{12} = \mathbf{0}$, and conditioning does not alter the distribution of \mathbf{x}_1 .

Sums of normal vectors

Suppose \mathbf{x} and \mathbf{y} are d -dimensional random vectors with joint multivariate normal distribution

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{\nu} \end{bmatrix}, \begin{bmatrix} \Sigma & \mathbf{P} \\ \mathbf{P}^\top & \mathbf{T} \end{bmatrix}\right).$$

Then recognizing their sum $\mathbf{z} = \mathbf{x} + \mathbf{y} = [\mathbf{I}, \mathbf{I}][\mathbf{x}, \mathbf{y}]^\top$ as a linear transformation and applying (A.9), we have:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu} + \boldsymbol{\nu}, \Sigma + 2\mathbf{P} + \mathbf{T}).$$

When \mathbf{x} and \mathbf{y} are independent, $\mathbf{P} = \mathbf{0}$, and this simplifies to

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu} + \boldsymbol{\nu}, \Sigma + \mathbf{T}), \quad (\text{A.15})$$

where the moments simply add.

Differential entropy

The differential entropy of a multivariate normal random variable \mathbf{x} with distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$, expressed in nats, is

$$H[\mathbf{x}] = \frac{1}{2} \log |2\pi e \Sigma|. \quad (\text{A.16})$$

In the univariate case $p(x) = \mathcal{N}(x; \mu, \sigma^2)$, this reduces to

$$H[x] = \frac{1}{2} \log 2\pi e \sigma^2 \quad (\text{A.17})$$

Sequences of normal random variables

If $\{\mathbf{x}_i\}$ is a sequence of normal random variables with means $\{\boldsymbol{\mu}_i\}$ and covariances $\{\Sigma_i\}$ converging respectively to finite limits $\boldsymbol{\mu}_i \rightarrow \boldsymbol{\mu}$ and $\Sigma_i \rightarrow \Sigma$, then the sequence converges in distribution to a normal random variable \mathbf{x} with mean $\boldsymbol{\mu}$ and covariance Σ .

B

METHODS FOR APPROXIMATE BAYESIAN INFERENCE

In Bayesian optimization we occasionally face intractable posterior distributions that must be approximated before we can proceed. The *Laplace approximation* and *Gaussian expectation propagation* are two workhorses of approximate Bayesian inference, and at least one will suffice in most scenarios. Both result in Gaussian approximations to the posterior, especially convenient when working with Gaussian processes.

Approximate inference for GPS: § 2.8, p. 39

B.1 THE LAPLACE APPROXIMATION

Consider a vector-valued random variable $\mathbf{x} \in \mathbb{R}^d$ with arbitrary prior distribution $p(\mathbf{x})$. Suppose we obtain information \mathcal{D} , yielding an intractable posterior

$$p(\mathbf{x} | \mathcal{D}) = Z^{-1} p(\mathbf{x}) p(\mathcal{D} | \mathbf{x})$$

that we wish to approximate. The Laplace approximation is based on approximating the logarithm of the unnormalized posterior:

$$\Psi(\mathbf{x}) = \log p(\mathbf{x}) + \log p(\mathcal{D} | \mathbf{x})$$

with a Taylor expansion around its maximum:

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \Psi(\mathbf{x}).$$

unnormalized log posterior, Ψ

maximum a posteriori point, $\hat{\mathbf{x}}$

Taking a second-order Taylor expansion around this point yields:

$$\Psi(\mathbf{x}) \approx \Psi(\hat{\mathbf{x}}) - \frac{1}{2}(\mathbf{x} - \hat{\mathbf{x}})^T \mathbf{H}(\mathbf{x} - \hat{\mathbf{x}}),$$

where \mathbf{H} is the Hessian of the negative log posterior evaluated at $\hat{\mathbf{x}}$:

$$\mathbf{H} = -\frac{\partial^2 \Psi}{\partial \mathbf{x} \partial \mathbf{x}^\top}(\hat{\mathbf{x}}).$$

Hessian of negative log posterior, \mathbf{H}

Note the first-order term vanishes as we are expanding around a maximum. Exponentiating, we derive an approximation to the unnormalized posterior:

$$p(\mathbf{x} | \mathcal{D}) \propto \exp \Psi(\mathbf{x}) \approx \exp \Psi(\hat{\mathbf{x}}) \exp \left(-\frac{1}{2}(\mathbf{x} - \hat{\mathbf{x}})^T \mathbf{H}(\mathbf{x} - \hat{\mathbf{x}}) \right).$$

We recognize this as proportional to a Gaussian distribution, yielding a normal approximate posterior:

$$p(\mathbf{x} | \mathcal{D}) \approx q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \mathbf{H}^{-1}). \quad (\text{B.1})$$

Laplace approximation to posterior

Through some accounting when normalizing (B.1), the Laplace approximation also gives an approximation to the normalizing constant Z :

$$Z \approx \hat{Z}_{\text{LA}} = (2\pi)^{\frac{d}{2}} |\mathbf{H}|^{-\frac{1}{2}} \exp \Psi(\hat{\mathbf{x}}). \quad (\text{B.2})$$

Laplace approximation to normalizing constant, \hat{Z}_{LA}

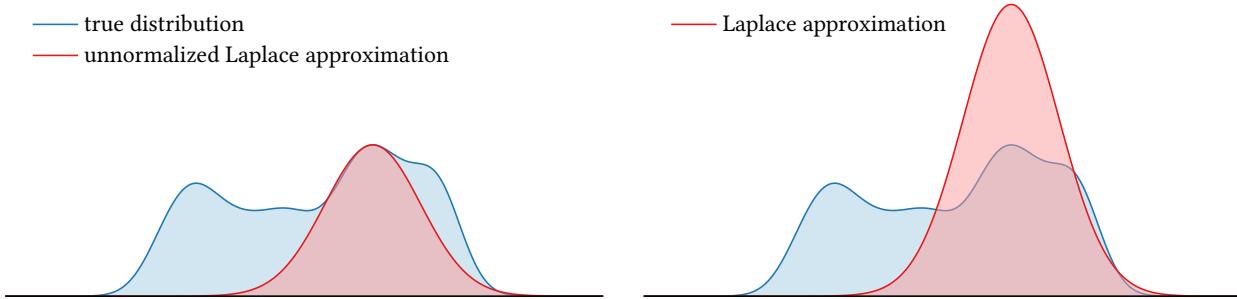


Figure B.1: A Laplace approximation to a one-dimensional posterior distribution. The left panel shows the Laplace approximation before normalization, and the right panel afterwards.

The Laplace approximation procedure is illustrated in figure B.1, where we show the approximate posterior both before and after normalization. The posterior density is an excellent local approximation around the maximum but is not a great global fit as a significant fraction of the true posterior mass is ignored. However, the Laplace approximation is remarkably simple and general and is sometimes the only viable approximation scheme.

B.2 GAUSSIAN EXPECTATION PROPAGATION

Expectation propagation (EP) is a technique for approximate Bayesian inference that enjoys some use in Bayesian optimization. We will give a brief and incomplete introduction that should nonetheless suffice for common applications in this context. A complete introduction can be found in MINKA’s thesis,¹ and CUNNINGHAM et al. provide in-depth advice regarding efficient and stable computation for the rank-one case we consider here.²

Consider a multivariate Gaussian random variable ξ with distribution

$$p(\xi) = \mathcal{N}(\xi; \mu_0, \Sigma_0).$$

Suppose we obtain information \mathcal{D} about ξ in the form of a collection of *factors*, each of which specifies the likelihood of a scalar product $x = \mathbf{a}^\top \xi$ associated with that factor.³ We consider the posterior distribution

$$p(\xi | \mathcal{D}) = Z^{-1} p(\xi) \prod_i t_i(x_i), \quad (\text{B.3})$$

where i th factor t_i informs our belief about $x_i = \mathbf{a}_i^\top \xi$. Unfortunately, this posterior is intractable except the notable case when all factors are Gaussian.

Gaussian expectation propagation proceeds by replacing each of the factors with an unnormalized Gaussian distribution:

$$t_i(x_i) \approx \tilde{t}_i(x_i) = \tilde{Z}_i \mathcal{N}(x_i; \tilde{\mu}_i, \tilde{\sigma}_i^2).$$

Here $(\tilde{Z}_i, \tilde{\mu}_i, \tilde{\sigma}_i^2)$ are called *site parameters* for the approximate factor \tilde{t}_i , which we may design to optimize the fit. We will consider this issue further shortly. Given arbitrary site parameters for each factor, the resulting approximation to (B.3) is

$$p(\xi | \mathcal{D}) \approx q(\xi) = \hat{Z}_{\text{EP}}^{-1} p(\xi) \prod_i \tilde{t}_i(x_i). \quad (\text{B.4})$$

As a product of Gaussians, the approximate posterior is also Gaussian:

$$q(\xi) = \mathcal{N}(\xi; \boldsymbol{\mu}, \Sigma), \quad (\text{B.5})$$

with parameters:⁴

$$\boldsymbol{\mu} = \Sigma \left(\Sigma_0^{-1} \boldsymbol{\mu}_0 + \sum_i \frac{\tilde{\mu}_i}{\tilde{\sigma}_i^2} \mathbf{a}_i \right); \quad \Sigma = \left(\Sigma_0^{-1} + \sum_i \frac{1}{\tilde{\sigma}_i^2} \mathbf{a}_i \mathbf{a}_i^\top \right)^{-1}.$$

Gaussian EP also yields an approximation of the normalizing constant Z , if desired:

$$Z \approx \hat{Z}_{\text{EP}} = \frac{\mathcal{N}(\mathbf{0}; \boldsymbol{\mu}, \Sigma)}{\mathcal{N}(\mathbf{0}; \boldsymbol{\mu}_0, \Sigma_0)} \prod_i \tilde{Z}_i \mathcal{N}(0; \tilde{\mu}_i, \tilde{\sigma}_i^2). \quad (\text{B.6})$$

What remains to be determined is an effective means of choosing the site parameters to maximize the approximation fidelity. One reasonable goal would be to minimize the Kullback–Leibler (KL) divergence between the true and approximate distributions; for our Gaussian approximation (B.5), this is achieved through moment matching.² Unfortunately, determining the moments of the true posterior (B.3) may be difficult, so expectation propagation instead matches the *marginal* moments for each of the $\{x_i\}$, approximately minimizing the KL divergence. This is accomplished through an iterative procedure where we repeatedly sweep over each of the approximate factors and refine its parameters until convergence.

We initialize all site parameters to $(\tilde{Z}, \tilde{\mu}, \tilde{\sigma}^2) = (1, 0, \infty)$; with these choices the approximate factors drop away, and our initial approximation is simply the prior: $(\boldsymbol{\mu}, \Sigma) = (\boldsymbol{\mu}_0, \Sigma_0)$. Now we perform a series of updates for each of the approximate factors in turn. These updates take a convenient general form, and we will drop factor index subscripts below to simplify notation.

Let $\tilde{t}(x) = \tilde{Z} \mathcal{N}(x; \tilde{\mu}, \tilde{\sigma}^2)$ be an arbitrary factor in our approximation (B.4). The idea behind expectation propagation is to drop this factor from the approximation entirely, forming the *cavity distribution*:

$$\bar{q}(\xi) = \frac{q(\xi)}{\tilde{t}(x)},$$

and replace it with the true factor $t(x)$, forming the *tilted distribution* $\bar{q}(\xi) t(x)$. The tilted distribution is closer to the true posterior (B.3) as the factor in question is no longer approximated. We now adjust the site

site parameters, $(\tilde{Z}, \tilde{\mu}, \tilde{\sigma}^2)$

Gaussian EP approximate posterior, $q(\xi)$

⁴ The updated covariance incorporates only a series of rank-one updates, which can be applied using the Sherman–Morrison formula.

Gaussian EP parameters, $(\boldsymbol{\mu}, \Sigma)$

Gaussian EP approximation to normalizing constant, \hat{Z}_{EP}

setting site parameters

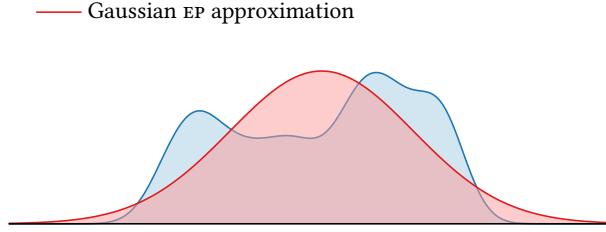
expectation propagation approximately minimizes KL divergence

site parameter initialization

cavity distribution, \bar{q}

tilted distribution

Figure B.2: A Gaussian EP approximation to the distribution in figure B.1.



parameters to minimize the KL divergence between the tilted distribution and the new approximation $q'(\xi) = \tilde{q}(\xi) \tilde{t}'(x)$:

$$(\tilde{Z}, \tilde{\mu}, \tilde{\sigma}^2) = \arg \min D_{\text{KL}} [\tilde{q}(\xi) t(x) \parallel q'(\xi)]$$

by matching zeroth, first, and second moments.

Because Gaussian distributions are closed under marginalization, we can simplify this procedure by manipulating only *marginal* distributions for x rather than the full joint distribution.⁵ The marginal belief about $x = \mathbf{a}^\top \xi$ in our current approximation (B.5) is:

$$q(x) = \mathcal{N}(x; \mu, \sigma^2); \quad \mu = \mathbf{a}_i^\top \boldsymbol{\mu}; \quad \sigma^2 = \mathbf{a}_i^\top \Sigma \mathbf{a}_i.$$

By dividing by the approximate factor $\tilde{t}(\xi)$, we arrive at the marginal cavity distribution, which is Gaussian:

$$\tilde{q}(x) = \mathcal{N}(x; \bar{\mu}, \bar{\sigma}^2); \quad \bar{\mu} = \tilde{\sigma}^2 (\mu \sigma^{-2} - \tilde{\mu} \tilde{\sigma}^{-2}); \quad \bar{\sigma}^2 = (\sigma^{-2} - \tilde{\sigma}^{-2})^{-1} \quad (\text{B.7})$$

Consider the zeroth moment of the marginal tilted distribution:

$$Z = \int t(x) \mathcal{N}(x; \bar{\mu}, \bar{\sigma}^2) dx; \quad (\text{B.8})$$

this quantity clearly depends on the cavity parameters $(\bar{\mu}, \bar{\sigma}^2)$. If we define

$$\alpha = \frac{\partial \log Z}{\partial \bar{\mu}}; \quad \beta = \frac{\partial \log Z}{\partial \bar{\sigma}^2}; \quad (\text{B.9})$$

and an auxiliary variable $\gamma = (\alpha^2 - 2\beta)^{-1}$ then we may achieve the desired moment matching by updating the site parameters to:⁶

$$\tilde{\mu} = \bar{\mu} + \alpha \gamma; \quad \tilde{\sigma}^2 = \gamma - \bar{\sigma}^2; \quad \tilde{Z} = Z \sqrt{2\pi} \sqrt{\bar{\sigma}^2 + \tilde{\sigma}^2} \exp\left(\frac{1}{2}\alpha^2 \gamma\right). \quad (\text{B.10})$$

This completes our update for the chosen factor; the full EP procedure repeatedly updates each factor in this manner until convergence. The result of Gaussian expectation propagation for the distribution from figure B.1 is shown in figure B.2. The fit is good and reflects the more global nature of the expectation propagation scheme achieved through moment matching rather than merely maximizing the posterior.

A convenient aspect expectation propagation is that incorporating a new factor only requires computing the zeroth moment against an arbitrary normal distribution (B.8) and the partial derivatives in (B.9). We provide these computations for several useful factor types below.

⁵ M. SEAGER (2008). *Expectation Propagation for Exponential Families*. Technical report. University of California, Berkeley.

⁶ T. MINKA (2008). EP: A quick reference.

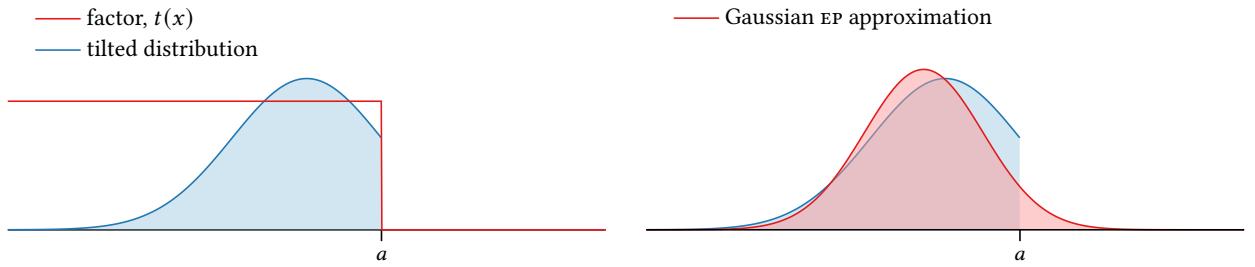


Figure B.3: A Gaussian EP approximation to a one-dimensional normal distribution truncated at a .

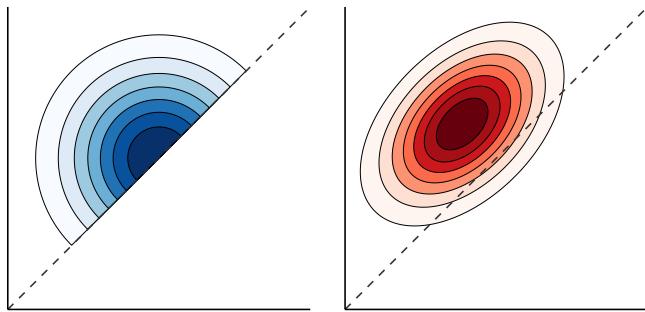


Figure B.4: A Gaussian EP approximation to a joint normal distribution conditioned on one coordinate being less than the other, corresponding to the dashed boundary. Contours of the tilted distribution are shown in the left panel, and contours of the approximate posterior in the right panel.

Truncating a variable

A common use of expectation propagation in Bayesian optimization is to approximately constrain a Gaussian random variable x to be less than a threshold a . We may capture this information by a single factor $t(x) = [x < a]$. In the context of expectation propagation, we must consider the normalizing constant of the tilted distribution, which is a truncated normal:

$$Z = \int [x < a] \mathcal{N}(x; \bar{\mu}, \bar{\sigma}^2) dx = \Phi(z); \quad z = \frac{a - \bar{\mu}}{\bar{\sigma}}. \quad (\text{B.11})$$

The required quantities for an expectation propagation update are now (B.10):

$$\alpha = -\frac{\phi(z)}{\Phi(z)\bar{\sigma}}; \quad \beta = \frac{z\alpha}{2\bar{\sigma}}; \quad \gamma = -\frac{\bar{\sigma}}{\alpha} \left(\frac{\phi(z)}{\Phi(z)} + z \right)^{-1}. \quad (\text{B.12})$$

A Gaussian EP approximation to a truncated normal distribution is illustrated in figure (B.3). The fit is good, but not perfect: approximately 5% of its mass exceeds the threshold. This inaccuracy is the price of approximation.

We may also apply this approach to approximately condition our belief on ξ on one entry being dominated by another: $\xi_i < \xi_j$. Consider the vector \mathbf{a} where $a_i = 1$, $a_j = -1$ and all other entries are zero. Then $x = \mathbf{a}^\top \xi = \xi_i - \xi_j$. The condition is now equivalent to $[x < 0]$, and we can proceed as outlined above. This approximation is illustrated for a bivariate normal in figure B.4; again the fit appears reasonable.

conditioning on one variable being greater than another

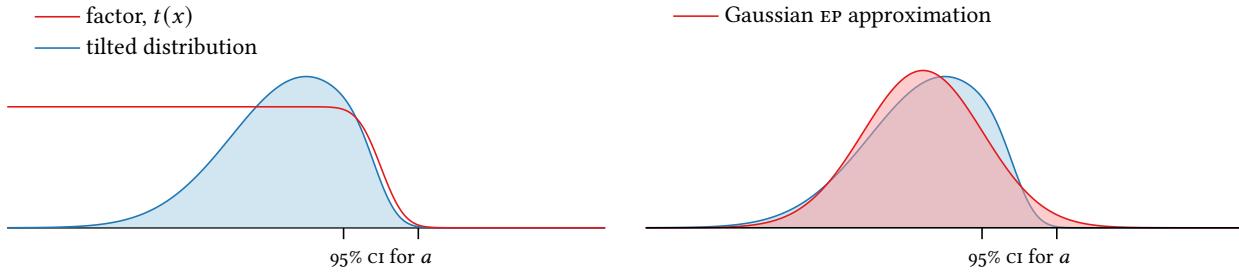


Figure B.5: A Gaussian EP approximation to a one-dimensional normal distribution truncated at an unknown threshold a with the marked 95% credible interval.

Truncation at an uncertain threshold

A sometimes useful extension of the above is to consider truncation at an *uncertain* threshold a . Suppose we have a Gaussian belief about a :

$$p(a) = \mathcal{N}(a; \mu, \sigma^2).$$

Integrating the hard truncation factor $[x < a]$ against this belief yields the following “soft truncation” factor:

$$t(x) = \int [x < a] p(a) da = \Phi\left(\frac{\mu - x}{\sigma}\right). \quad (\text{B.13})$$

We consider again the normalizing constant of the tilted distribution:

$$Z = \int \Phi\left(\frac{\mu - x}{\sigma}\right) \mathcal{N}(x; \bar{\mu}, \bar{\sigma}^2) dx = \Phi(z); \quad z = \frac{\mu - \bar{\mu}}{\sqrt{\sigma^2 + \bar{\sigma}^2}},$$

Defining $s = \sqrt{\sigma^2 + \bar{\sigma}^2}$, we may compute:

$$\alpha = -\frac{\phi(z)}{\Phi(z)s}; \quad \beta = \frac{z\alpha}{2s}; \quad \gamma = -\frac{s}{\alpha} \left(\frac{\phi(z)}{\Phi(z)} + z \right)^{-1}.$$

The hard truncation formulas above may be interpreted as a special case of this result by setting $(\mu, \sigma^2) = (a, 0)$. This procedure is illustrated in figure B.5, where we softly truncate a one-dimensional Gaussian distribution.

C

GRADIENTS

Under mild continuity assumptions, for Gaussian processes conditioned with exact inference, we may compute the gradient of both the log marginal likelihood with respect to model hyperparameters and of the posterior predictive moments with respect to observation location. The former aids in maximizing or sampling from the model posterior, and the latter in maximizing acquisition functions derived from the predictive distribution. In modern software these gradients are often computed via automatic differentiation, but we present their functional forms here to offer insight into their behavior.

We will consider a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ with distribution $\mathcal{GP}(f; m, K)$, observed with independent (but possibly heteroskedastic) additive Gaussian noise:

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

where the noise scale σ_n may optionally depend on x . We will notate the prior moment and noise scale functions with:

$$m(x; \theta); \quad K(x, x'; \theta); \quad \sigma_n(x; \theta),$$

and assume these are differentiable with respect to observation location and any parameters they may have. In a slight abuse of notation, we will use θ to indicate a vector collating the values of *all* hyperparameters of this model.

We will consider an arbitrary set of observations $\mathcal{D} = (\mathbf{x}, \mathbf{y})$. We will write the prior moments of $\phi = f(\mathbf{x})$ and the noise covariance associated with these observations as:

$$\boldsymbol{\mu} = m(\mathbf{x}; \theta); \quad \Sigma = K(\mathbf{x}, \mathbf{x}; \theta); \quad \mathbf{N} = \text{diag } \sigma_n^2(\mathbf{x}; \theta),$$

all of which have implicit dependence on the hyperparameters. It will also be useful to introduce notation for two repeating quantities:

$$\mathbf{V} = \text{cov}[\mathbf{y} | \mathbf{x}, \theta] = \Sigma + \mathbf{N}; \quad \boldsymbol{\alpha} = \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}).$$

C.1 GRADIENT OF LOG MARGINAL LIKELIHOOD

The log marginal likelihood of the data is (4.8):

$$\mathcal{L}(\theta) = \log p(\mathbf{y} | \mathbf{x}, \theta) = -\frac{1}{2} [\boldsymbol{\alpha}^\top (\mathbf{y} - \boldsymbol{\mu}) + \log |\mathbf{V}| + n \log 2\pi].$$

The partial derivatives with respect to mean function parameters have the form:

$$\frac{\partial \mathcal{L}}{\partial \theta} = \boldsymbol{\alpha}^\top \frac{\partial \boldsymbol{\mu}}{\partial \theta},$$

and partial derivatives with respect to covariance function and likelihood parameters (that is, the parameters of \mathbf{V}) take the form:

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{1}{2} \left[\boldsymbol{\alpha}^\top \frac{\partial \mathbf{V}}{\partial \theta} \boldsymbol{\alpha} - \text{tr} \left[\mathbf{V}^{-1} \frac{\partial \mathbf{V}}{\partial \theta} \right] \right].$$

C.2 GRADIENT OF PREDICTIVE DISTRIBUTION WITH RESPECT TO LOCATION

For a given observation location x , let us define the vectors

$$\mathbf{k} = K(\mathbf{x}, x); \quad \boldsymbol{\beta} = \mathbf{V}^{-1}\mathbf{k}.$$

The posterior moments of $\phi = f(x)$ are (2.19):

$$\mu = m(x) + \boldsymbol{\alpha}^\top \mathbf{k}; \quad \sigma^2 = K(x, x) - \boldsymbol{\beta}^\top \mathbf{k},$$

and the partial derivatives of these moments with respect to observation location are:

$$\frac{\partial \mu}{\partial x} = \frac{\partial m}{\partial x} - \boldsymbol{\alpha}^\top \frac{\partial \mathbf{k}}{\partial x}; \quad \frac{\partial \sigma^2}{\partial x} = \frac{\partial K(x, x)}{\partial x} - 2\boldsymbol{\beta}^\top \frac{\partial \mathbf{k}}{\partial x}.$$

The predictive distribution for a noisy observation y at x is (8.5):

$$p(y | x, \mathcal{D}, \sigma_n^2) = \mathcal{N}(y; \mu, s^2); \quad s^2 = \sigma^2 + \sigma_n^2,$$

and the partial derivative of the predictive variance and standard deviation with respect to x are:

$$\frac{\partial s^2}{\partial x} = \frac{\partial \sigma^2}{\partial x} + \frac{\partial \sigma_n^2}{\partial x}; \quad \frac{\partial s}{\partial x} = \frac{1}{2s} \frac{\partial s^2}{\partial x}.$$

C.3 GRADIENTS OF COMMON ACQUISITION FUNCTIONS

Gradient of noisy expected improvement

To aid in the optimization of expected improvement, we may also compute the gradient of g ,¹ and thus the gradient of expected improvement with respect to the proposed observation location x . This will require several applications of the chain rule due to numerous dependencies among the variables involved.

First we note that the $\{c_i\}$ values defining the intervals of dominance (8.15) depend on the vectors \mathbf{a} and \mathbf{b} . In particular, for $2 \leq i \leq n$, c_i is the z -value where the $(i-1)$ th and i th lines intersect, which occurs at

$$c_i = \frac{a_i - a_{i-1}}{b_{i-1} - b_i},$$

making the dependence explicit. For $2 \leq i \leq n$, we have:

$$\frac{\partial c_i}{\partial a_i} = \frac{1}{b_{i-1} - b_i}; \quad \frac{\partial c_i}{\partial a_{i-1}} = -\frac{\partial c_i}{\partial a_i}; \quad \frac{\partial c_i}{\partial b_i} = \frac{a_i - a_{i-1}}{(b_i - b_{i-1})^2}; \quad \frac{\partial c_i}{\partial b_{i-1}} = -\frac{\partial c_i}{\partial b_i},$$

and at the fixed endpoints at infinity $i \in \{1, n+1\}$, we have

$$\frac{\partial c_i}{\partial \mathbf{a}} = \frac{\partial c_i}{\partial \mathbf{b}} = \mathbf{0}^\top.$$

Now we may compute the gradient of $g(\mathbf{a}, \mathbf{b})$ with respect to its inputs, accounting for the implicit dependence of interval endpoints on these

gradient of g with respect to \mathbf{a} and \mathbf{b}

values:

$$\begin{aligned}\frac{\partial g}{\partial a_i} &= [\Phi(c_{i+1}) - \Phi(c_i)] \\ &\quad + \frac{\partial c_{i+1}}{\partial a_i} [a_i + b_i c_{i+1} - [i \leq n] (a_{i+1} + b_{i+1} c_{i+1})] \phi(c_{i+1}) \\ &\quad - \frac{\partial c_i}{\partial a_i} [a_i + b_i c_i - [i > 1] (a_{i-1} + b_{i-1} c_i)] \phi(c_i); \\ \frac{\partial g}{\partial b_i} &= [\phi(c_i) - \phi(c_{i+1})] \\ &\quad + \frac{\partial c_{i+1}}{\partial b_i} [a_i + b_i c_{i+1} - [i \leq n] (a_{i+1} + b_{i+1} c_{i+1})] \phi(c_{i+1}) \\ &\quad - \frac{\partial c_i}{\partial b_i} [a_i + b_i c_i - [i > 1] (a_{i-1} + b_{i-1} c_i)] \phi(c_i).\end{aligned}$$

Here $[i > 1]$ and $[i \leq n]$ represent the Iverson bracket. We will also require the gradient of the \mathbf{a} and \mathbf{b} vectors with respect to x :

$$\frac{\partial \mathbf{a}}{\partial x} = \left[\frac{\partial a_u}{\partial x} \right]; \quad \frac{\partial \mathbf{b}}{\partial x} = \frac{1}{s} \left[\frac{\partial K_D(\mathbf{x}', x)}{\partial x} - \mathbf{b} \frac{\partial s}{\partial x} \right]. \quad (\text{c.1})$$

Note the dependence on the gradient of the predictive parameters. Finally, if we preprocess the inputs by identifying an appropriate transformation matrix \mathbf{P} such that $g(\mathbf{a}, \mathbf{b}) = g(\mathbf{Pa}, \mathbf{Pb}) = g(\boldsymbol{\alpha}, \boldsymbol{\beta})$ (8.14), then the desired gradient of expected improvement is:

$$\frac{\partial \alpha_{\text{EI}}}{\partial x} = \frac{\partial g}{\partial \boldsymbol{\alpha}} \mathbf{P} \frac{\partial \mathbf{a}}{\partial x} + \frac{\partial g}{\partial \boldsymbol{\beta}} \mathbf{P} \frac{\partial \mathbf{b}}{\partial x}.$$

Gradient of noisy probability of improvement

We may compute the gradient of (8.23) via several applications of the chain rule, under the (mild) assumption that the endpoints ℓ and u correspond to *unique* lines (a_ℓ, b_ℓ) and (a_u, b_u) .² Then we have

$$\frac{\partial \alpha_{\text{PI}}}{\partial a_\ell} = -\frac{\phi(\ell)}{b_\ell}; \quad \frac{\partial \alpha_{\text{PI}}}{\partial b_\ell} = -\frac{\ell \phi(\ell)}{b_\ell}; \quad \frac{\partial \alpha_{\text{PI}}}{\partial a_u} = \frac{\phi(-u)}{b_u}; \quad \frac{\partial \alpha_{\text{PI}}}{\partial b_u} = \frac{u \phi(-u)}{b_u},$$

and we may compute the gradient with respect to the proposed observation location as

$$\frac{\partial \alpha_{\text{PI}}}{\partial x} = \frac{\partial \alpha_{\text{PI}}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial x} + \frac{\partial \alpha_{\text{PI}}}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial x}.$$

The gradient with respect to \mathbf{a} and \mathbf{b} was computed previously (c.1).

gradient of \mathbf{a} and \mathbf{b} with respect to x

this picks up the discussion from § 8.3, p. 170

² If not, probability of improvement is not differentiable as moving one of the lines at the shared intersection will alter the probability of improvement only in favorable directions.

Approximating gradient via Gauss–Hermite quadrature

We may also use numerical techniques to estimate the gradient of the acquisition function with respect to the proposed observation location:

$$\frac{\partial \alpha}{\partial x} = \frac{\partial}{\partial x} \int \Delta(x, y) \mathcal{N}(y; \mu, s^2) dy \quad (\text{c.2})$$

this picks up the discussion from § 8.5, p. 172

GRADIENTS

- 3 It is sufficient for example that both Δ and $\frac{\partial\Delta}{\partial x}$ be continuous in x and y , and that the moments of the predictive distribution of y be continuous with respect to x .

approximation with Gauss–Hermite quadrature

accounting for dependence of $\{y_i\}$ on predictive parameters

this picks up the discussion from § 8.6, p. 175

If we assume sufficient regularity,³ then we may swap the order of expectation and differentiation:

$$\frac{\partial\alpha}{\partial x} = \int \frac{\partial}{\partial x} \Delta(x, y) \mathcal{N}(y; \mu, s^2) dy,$$

reducing the computation of the gradient to Gaussian expectation.

Gauss–Hermite quadrature remains a natural choice, resulting in the approximation

$$\frac{\partial\alpha}{\partial x} \approx \sum_{i=1}^n \bar{w}_i \frac{\partial\Delta}{\partial x}(x, y_i), \quad (\text{c.3})$$

which is simply the derivative of the estimator in (8.29). When using this estimate we must remember that the $\{y_i\}$ samples depend on x through the parameters of predictive distribution (8.29). Accounting for this dependence, the required gradient of the marginal gain at the i th integration node z_i is:

$$\frac{\partial\Delta}{\partial x}(x, y_i) = \left[\frac{\partial\Delta}{\partial x} \right]_y + \frac{\partial\Delta}{\partial y} \left[\frac{\partial\mu}{\partial x} + \frac{y_i - \mu}{s} \frac{\partial s}{\partial x} \right], \quad (\text{c.4})$$

where $\left[\frac{\partial\Delta}{\partial x} \right]_y$ indicates the partial derivative of the updated utility when the observed value y is held constant.

Approximating the gradient of knowledge gradient

Special care is needed to estimate the gradient of the knowledge gradient using numerical techniques. Inspecting (c.3), we must compute the gradient of the marginal gain

$$\frac{\partial\Delta}{\partial x} = \frac{\partial}{\partial x} \max_{x' \in \mathcal{X}} \mu_{\mathcal{D}'}(x').$$

It is not clear how we can differentiate through the max operator in this expression. However, under mild assumptions we may appeal to the *envelope theorem* for arbitrary choice sets to proceed.⁴ Consider a point x in the interior of the domain $\mathcal{X} \subset \mathbb{R}^n$ and fix an arbitrary associated observation value y .⁵ Take any point maximizing the updated posterior mean:

$$x^* \in \arg \max_{x' \in \mathcal{X}} \mu_{\mathcal{D}'}(x');$$

this point need not be unique. Assuming the updated posterior mean is differentiable,⁶ the envelope theorem states that the gradient of the global reward utility is equal to the gradient of the updated posterior mean evaluated at x^* :

$$\frac{\partial\Delta}{\partial x} = \frac{\partial}{\partial x} \mu_{\mathcal{D}'}(x^*).$$

For the knowledge gradient, and accounting for the dependence on x in the locations of the y samples (c.4), we may compute:

$$\frac{\partial\Delta}{\partial x}(x, y_i) = \frac{y_i - \mu}{s^3} \left[s \frac{\partial K_{\mathcal{D}}}{\partial x}(x, x_i^*) - \frac{\partial s}{\partial x} K_{\mathcal{D}}(x, x_i^*) \right],$$

where x_i^* maximizes the updated posterior mean corresponding to y_i . Combining this result with (c.3) yields a Gauss–Hermite approximation to the gradient of the acquisition function. This approach requires maximizing the updated posterior mean for each sample y_i . We may appeal to standard techniques for this task, making use of the gradient

$$\frac{\partial \mu_{\mathcal{D}'}^*}{\partial x}.$$

Gradient of predictive entropy search acquisition function

The explicit formula (8.52) and results above allow us to compute – after investing considerable tedium – the gradient of the predictive entropy search approximation to the mutual information (8.41). We begin by differentiating that estimate:

$$\frac{\partial \alpha_{\text{PES}}}{\partial x} = \frac{1}{s} \frac{\partial s}{\partial x} + \frac{1}{2n} \sum_{i=1}^n \frac{1}{s_{*i}^2} \frac{\partial s_{*i}^2}{\partial x}.$$

Given a fixed sample x^* and dropping the subscript, we may compute (8.39):

$$\frac{\partial s_*^2}{\partial x} = \frac{\partial \sigma_*^2}{\partial x} + \frac{\partial \sigma_n^2}{\partial x}.$$

We proceed by differentiating the approximate latent predictive variance (8.52):

$$\frac{\partial \sigma_*^2}{\partial x} = \frac{\partial \zeta^2}{\partial x} - \frac{2\zeta^2 - 2\rho}{\gamma} \left[\frac{\partial \zeta^2}{\partial x} - \frac{\partial \rho}{\partial x} \right] + \frac{(\zeta^2 - \rho)^2}{\gamma^2} \frac{\partial \gamma}{\partial x},$$

which depends on the derivatives of the expectation propagation update terms:

$$\begin{aligned} \frac{\partial \gamma}{\partial x} &= \frac{1}{\alpha} \left[1 - \frac{(\Phi(z) - \phi(z))^2}{(z\Phi(z) + \phi(z))^2} \right] \left[\frac{\partial m}{\partial x} + z \frac{\partial \bar{\sigma}}{\partial x} \right] + \frac{2\gamma}{\bar{\sigma}} \frac{\partial \bar{\sigma}}{\partial x}; \\ \frac{\partial \bar{\sigma}}{\partial x} &= \frac{1}{2\bar{\sigma}} \left[\frac{\partial \zeta^2}{\partial x} - 2 \frac{\partial \rho}{\partial x} \right]; \end{aligned}$$

and the predictive parameters (8.51):

$$\begin{aligned} \frac{\partial m}{\partial x} &= \frac{\partial \mu}{\partial x} + (\boldsymbol{\alpha}^*)^\top \frac{\partial \mathbf{k}}{\partial x}; \\ \frac{\partial \zeta^2}{\partial x} &= \frac{\partial \sigma^2}{\partial x} - 2\mathbf{k}^\top \mathbf{V}_*^{-1} \frac{\partial \mathbf{k}}{\partial x}; \quad \frac{\partial \rho}{\partial x} = \frac{\partial k_1}{\partial x} - \mathbf{k}_*^\top \mathbf{V}_*^{-1} \frac{\partial \mathbf{k}}{\partial x}. \end{aligned}$$

Gradient of OPES acquisition function

To compute the gradient of the OPES approximation, we differentiate (8.63):

$$\frac{\partial \alpha_{\text{OPES}}}{\partial x} = \frac{1}{s} \frac{\partial s}{\partial x} + \frac{1}{2} \sum_i \frac{w_i}{s_{*i}^2} \frac{\partial s_{*i}^2}{\partial x}.$$

this picks up the discussion from § 8.8, p. 187

Given a fixed sample f^* and dropping the subscript, we may compute (8.63):

$$\frac{\partial s_*^2}{\partial x} = \frac{\partial \sigma_*^2}{\partial x} + \frac{\partial \sigma_n^2}{\partial x}.$$

Finally, we differentiate (8.61):

$$\begin{aligned} \frac{\partial \sigma_*^2}{\partial x} &= \sigma \frac{\phi(z)}{\Phi(z)} \left[\frac{\partial \mu}{\partial x} + z \frac{\partial \sigma}{\partial x} \right] \left[1 - 2 \frac{\phi(z)^2}{\Phi(z)^2} - 3z \frac{\phi(z)}{\Phi(z)} - z^2 \right] \\ &\quad + \frac{\partial \sigma^2}{\partial x} \left[1 - z \frac{\phi(z)}{\Phi(z)} - \frac{\phi(z)^2}{\Phi(z)^2} \right]. \end{aligned}$$

D

ANNOTATED BIBLIOGRAPHY OF APPLICATIONS

Countless settings across science, engineering, and beyond involve free parameters that can be tuned at will to achieve some objective. However, in many cases the evaluation of a given parameter setting can be extremely costly, stymieing exhaustive exploration of the design space.

Of course, such a situation is a perfect use case for Bayesian optimization. Through careful modeling and intelligent policy design, Bayesian optimization algorithms can deliver impressive optimization performance even with small observation budgets. This capability has been demonstrated in hundreds of studies across a wide range of domains.

Here we provide a brief survey of some notable applications of Bayesian optimization. The selected references are not intended to be exhaustive (that would be impossible given the size of the source material!), but rather to be representative, diverse, and good starting points for further investigation.

CHEMISTRY AND MATERIALS SCIENCE

At a high level, the synthesis of everything from small molecules to bulk materials proceeds in the same manner: initial materials are combined and subjected to suitable conditions such that they transform into a final product. Although this seems like a simple recipe, we face massive challenges in its realization:

- We usually wish that the resulting products be *useful*, that is, that they exhibit desirable properties. In drug discovery, for example, we seek molecules exhibiting binding activity against an identified biological target. In other settings we might seek products exhibiting favorable optical, electronic, mechanical, thermal, and/or other properties.
- The space of possible products can be enormous and the sought after properties exceptionally rare. For example, there are an estimated 10^{60} pharmacologically active molecules, only a tiny fraction of which might exhibit binding activity against any given biological target.
- Determining the properties of a candidate molecule or material ultimately requires synthesis and characterization in a laboratory, which can be complicated, costly, and slow.

For these reasons, exploration of molecular or material spaces can devolve into a cumbersome trial-and-error search for “needles in a haystack.”

Over the past few decades, the fields of *computational chemistry* and *computational materials science* have developed sophisticated techniques for estimating chemical and material properties from simulation. As accurate simulation often requires consideration of quantum mechanical interactions, these surrogates can still be quite costly, but are nonetheless cheaper and more easily parallelized than laboratory experiments. This has enabled computer-guided exploration of large molecular and materials spaces at relatively little cost, but in many cases exhaustive

computational chemistry, computational materials science

exploration remains untenable, and we must appeal to methods such as Bayesian optimization for guidance.

Virtual screening of molecular spaces

The ability to estimate molecular properties using computational methods has enabled so-called *virtual screening*, where early stages of discovery can be performed *in silico*. Bayesian optimization and active search can dramatically increase the rate of discovery in this setting.

chemoinformatics
molecular fingerprints

neural fingerprints

One challenge here is constructing predictive models for the properties of molecules, which are complex structured objects. The field of *chemoinformatics* has developed a number of *molecular fingerprints* intended to serve as useful feature representations for molecules when predicting chemical properties. Traditionally, these fingerprints were developed based on chemical intuition; however, numerous *neural* fingerprints have emerged in recent years, designed via deep representation learning on huge molecular databases.

GRAFF, DAVID E. et al. (2021). Accelerating high-throughput virtual screening through molecular pool-based active learning. *Chemical Science* 12(22):7866–7881.

HERNÁNDEZ-LOBATO, JOSÉ MIGUEL et al. (2017). Parallel and Distributed Thompson Sampling for Large-scale Accelerated Exploration of Chemical Space. *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*.

JIANG, SHALI et al. (2017). Efficient Nonmyopic Active Search. *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*.

De novo design

A recent trend in molecular discovery has been to learn deep *generative* models for molecular structures and perform optimization in the (continuous) latent space of such a model. This enables *de novo* design, where the optimization routine can propose entirely novel structures for evaluation by feeding selected points in the latent space through the generational procedure.

the question of synthesizability

Although appealing, a major complication with this scheme is identifying a synthetic route – if one even exists! – for proposed structures. This issue deserves careful consideration when designing a system that can effectively transition from virtual screening to the laboratory.

GAO, WENHAO et al. (2020). The Synthesizability of Molecules Proposed by Generative Models. *Journal of Chemical Information and Modeling* 60(12):5714–5723.

GÓMEZ-BOMBARELLI, RAFAEL et al. (2018). Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules. *ACS Central Science* 4(2):268–276.

- GRIFFITHS, RYAN-RHYS et al. (2020). Constrained Bayesian optimization for automatic chemical design using variational autoencoders. *Chemical Science* 11(2):577–586.
- KOROVINA, KSENIA et al. (2020). Chembo: Bayesian Optimization of Small Organic Molecules with Synthesizable Recommendations. *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS 2020)*.

Reaction optimization

Not all applications of Bayesian optimization in chemistry take the form of optimizing chemical properties as a function of molecular structure. For example, even once a useful molecule has been identified, there may remain numerous parameters of its synthesis – reaction environment, processing parameters, etc. – that can be further optimized, seeking for example to reduce the cost and/or increase the yield of production.

- SHIELDS, BENJAMIN J. et al. (2021). Bayesian reaction optimization as a tool for chemical synthesis. *Nature* 590:89–96.

Conformational search

Conformational search seeks to identify the configuration(s) of a molecule with the lowest potential energy. Even for molecules of moderate size, the space of possible configurations can be enormous and the potential energy surface can exhibit numerous local minima. Interaction with other structures such as a surface (*adsorption*) can further complicate the computation of potential energy, rendering the search even more difficult.

adsorption

- CARR, SHANE F. et al. (2017). Accelerating the Search for Global Minima on Potential Energy Surfaces using Machine Learning. *Journal of Chemical Physics* 145(15):154106.
- FANG, LINCAN et al. (2021). Efficient Amino Acid Conformer Search with Bayesian Optimization. *Journal of Chemical Theory and Computation* 17(3):1955–1966.
- PACKWOOD, DANIEL (2017). *Bayesian Optimization for Materials Science*. Springer–Verlag.

We may also use related Bayesian methods to map out minimal-energy pathways between neighboring minima on a potential energy surface in order to understand the intermediate geometry of a molecule as it transforms from one energetically favorable state to another.

mapping minimal-energy pathways

- KOISTINEN, OLLI-PEKKA et al. (2017). Nudged elastic band calculations accelerated with Gaussian process regression. *Journal of Chemical Physics* 147(15):152720.

Optimization of material properties and performance

Bayesian optimization has demonstrated remarkable success in accelerating materials design. As in molecular design, in many cases we can perform early screening *in silico*, using computational methods to approximate properties of interest. The literature in this space is vast, and the studies below are representative examples targeting a range of material types and properties.

- ATTIA, PETER M. et al. (2020). Closed-loop optimization of fast-charging protocols for batteries with machine learning. *Nature* 578:397–402.
- FUKAZAWA, TARO et al. (2019). Bayesian optimization of chemical composition: A comprehensive framework and its application to $R\text{Fe}_{12}$ -type magnet compounds. *Physical Review Materials* 3(5):053807.
- HAGHANIFAR, SAJAD et al. (2020). Discovering high-performance broadband and broad angle antireflection surfaces by machine learning. *Optica* 7(7):784–789.
- HERBOL, HENRY C. et al. (2018). Efficient search of compositional space for hybrid organic–inorganic perovskites via Bayesian optimization. *npj Computational Materials* 4(51).
- JU, SHENGHONG et al. (2017). Designing Nanostructures for Phonon Transport via Bayesian Optimization. *Physical Review X* 7:021024.
- MIYAGAWA, SHINSUKE et al. (2021). Application of Bayesian optimization for improved passivation performance in $\text{TiO}_x/\text{SiO}_y/\text{c-Si}$ heterostructure by hydrogen plasma treatment. *Applied Physics Express* 14(2):025503.
- NAKAMURA, KENSAKU et al. (2021). Multi-objective Bayesian optimization of optical glass compositions. *Ceramics International* 47(11):15819–15824.
- NUGRAHA, ASEP SUGIH et al. (2020). Mesoporous trimetallic PtPdAu alloy films toward enhanced electrocatalytic activity in methanol oxidation: unexpected chemical compositions discovered by Bayesian optimization. *Journal of Materials Chemistry A* 8(27):13532–13540.
- OSADA, KEIICHI et al. (2020). Adaptive Bayesian optimization for epitaxial growth of Si thin films under various constraints. *Materials Today Communications* 25:1015382.
- SEKO, ATSUTO et al. (2015). Prediction of Low-Thermal-Conductivity Compounds with First-Principles Anharmonic Lattice-Dynamics Calculations and Bayesian Optimization. *Physical Review Letters* 115(20):205901.

Structural search

A fundamental question in materials science is how the structure of a material gives rise to its material properties. However, predicting the likely structure of a material – for example by evaluating the potential energy of plausible structures and minimizing – can be extraordinarily difficult, as the number of possible configurations can be astronomical.

- KIYOHARA, SHIN et al. (2016). Acceleration of stable interface structure searching using a kriging approach. *Japanese Journal of Applied Physics* 55(4):045502.
- OKAMOTO, YASUHARU (2017). Applying Bayesian Approach to Combinatorial Problem in Chemistry. *Journal of Physical Chemistry A* 121(17):3299–3304.
- TODOROVIĆ, MILICA et al. (2019). Bayesian inference of atomistic structure in functional materials. *npj Computational Materials* 5(35).

Software for Bayesian optimization in chemistry and materials science

- HÄSE, FLORIAN et al. (2018). Phoenics: A Bayesian Optimizer for Chemistry. *ACS Central Science* 4(9):1134–1145.
- UENO, TSUYOSHI et al. (2016). COMBO: An efficient Bayesian optimization library for materials science. *Materials Discovery* 4:18–21.

PHYSICS

Modern physics is driven by experiments of massive scope, which have enabled the refinement of theories of the Universe on all scales. The complexity of these experiments – from data acquisition to the following analysis and inference – offers many opportunities for optimization, but the same complexity often renders optimization difficult due to large parameter spaces and/or the expense of evaluating a particular setting.

Experimental physics

Complex physical instruments such as particle accelerators offer the capability of extraordinarily fine tuning through careful setting of their control parameters. In some cases, these parameters can be altered on-the-fly during operation, resulting in a huge space of possible configurations. Bayesian optimization can help accelerate the tuning process.

- DURIS, J. et al. (2020). Bayesian Optimization of a Free-Electron Laser. *Physical Review Letters* 124(12):124801.
- ROUSSEL, RYAN et al. (2021). Multiobjective Bayesian optimization for online accelerator tuning. *Physical Review Accelerators and Beams* 24(6):062801.
- SHALLOO, R. J. et al. (2020). Automation and control of laser wakefield accelerators using Bayesian optimization. *Nature Communications* 11:6355.
- WIGLEY, P. B. et al. (2016). Fast machine-learning online optimization of ultra-cold-atom experiments. *Scientific Reports* 6:25890.

Inverse problems in physics

The goal of a *inverse problem* is to determine the free parameters of a generative model¹ from observations of its output. Inverse problems

¹ This model may not be probabilistic but rather a complex simulation procedure.

forward map	are pervasive in physics – many physical models feature parameters (physical constants) that cannot be determined from first principles. However, we can <i>infer</i> these parameters experimentally by comparing the predictions of the model with the behavior of relevant observations.
Bayesian analog	In this context, determining the predictions of the model given a setting of its parameters is called the <i>forward map</i> . The inverse problem is then the task of “inverting” this map: searching the parameter space for the settings in the greatest accord with observed data.
	We can draw parallels here with Bayesian inference, where the forward map is characterized by the likelihood, which determines the distribution of observed data given the parameters. The Bayesian answer to the inverse problem is then encapsulated by the posterior distribution of the parameters given the data, which identifies the most plausible parameters in light of the observations.

For complex physical models, even the forward map can be exceptionally expensive to compute, making it difficult to completely explore its parameter space. For example, the forward map of a cosmological model may require simulating the evolution of an entire universe at a fine enough resolution to observe its large scale structure. Exhaustively traversing the space of cosmological parameters and comparing with observed structure would be infeasible, but progress may be possible with careful guidance. Bayesian optimization has proven useful to this end on a range of difficult inverse problems.

ILTEN, P. et al. (2017). Event generator tuning using Bayesian optimization. *Journal of Instrumentation* 12(4):Po4028.

LECLERCQ, FLORENT (2018). Bayesian optimization for likelihood-free cosmological inference. *Physical Review D* 98(6):063511.

ROGERS, KEIR K. et al. (2019). Bayesian emulator optimisation for cosmology: application to the Lyman-alpha forest. *Journal of Cosmology and Astroparticle Physics* 2019(2):031.

VARGAS-HERNÁNDEZ, R. A. et al. (2019). Bayesian optimization for the inverse scattering problem in quantum reaction dynamics. *New Journal of Physics* 21(2):022001.

BIOLOGICAL SCIENCES AND ENGINEERING

Biological systems are extraordinarily complex. Obtaining experimental measurements to shed light on the behavior of these systems can be difficult, slow, and expensive, and the resulting data can be corrupted with significant noise. Efficient experimental design is thus critical to make progress, and Bayesian optimization is a natural tool to consider.

LI, YAN et al. (2018). A Knowledge Gradient Policy for Sequencing Experiments to Identify the Structure of RNA Molecules Using a Sparse Additive Belief Model. *INFORMS Journal on Computing* 30(4):625–786.

- LORENZ, ROMY et al. (2018). Dissociating frontoparietal brain networks with neuroadaptive Bayesian optimization. *Nature Communications* 9:1227.
- NIKITIN, ARTYOM et al. (2019). Bayesian optimization for seed germination. *Plant Methods* 15:43.

Inverse problems in the biological sciences

Challenging inverse problems (see above) are pervasive in biology. Effective modeling of biological systems often requires complicated, nonlinear models with numerous free parameters. Bayesian optimization can help guide the search for the parameters offering the best explanation of observed data.

- DOKOOHAKI, HAMZE et al. (2018). Use of inverse modelling and Bayesian optimization for investigating the T effect of biochar on soil hydrological properties. *Agricultural Water Management* 2018: 268–274.
- THOMAS, MARCUS et al. (2018). A method for efficient Bayesian optimization of self-assembly systems from scattering data. *BMC Systems Biology* 12:65.
- ULMASOV, DONIYOR et al. (2016). Bayesian Optimization with Dimension Scheduling: Application to Biological Systems. *Computer Aided Chemical Engineering* 38:1051–1056.

Gene and protein design

The tools of modern biology enable the custom design of genetic sequences and even entire proteins, which we can – at least theoretically – tailor as we see fit to achieve a particular purpose. It is natural to pose both gene and protein design in terms of optimizing some figure of merit over a space of alternatives. However, in either case we face an immediate combinatorial explosion in the number of possible genetic or amino acid sequences we might consider. Bayesian optimization has shown promise for overcoming this obstacle through careful experimental design.

- GONZÁLEZ, JAVIER et al. (2015). Bayesian Optimization for Synthetic Gene Design. *Bayesian Optimization: Scalability and Flexibility Workshop (BayesOpt 2015), Conference on Neural Information Processing Systems (neurIPS 2015)*.
- HIE, BRIAN L. et al. (2021). Adaptive machine learning for protein engineering. arXiv: 2106.06466 [q-bio.QM].
- YANG, KEVIN K. et al. (2019). Machine-learning-guided directed evolution for protein engineering. *Nature Methods* 16:687–694.

A related problem is faced in modern plant breeding, where the goal is to develop plant varieties with desirable characteristics: yield, drought or pest resistance, etc. Plant phenotyping is an inherently slow

plant breeding

process, as we must wait for planted seeds to germinate and grow until sufficiently mature that traits can be measured. This slow turnover rate makes it impossible to fully explore the space of possible genotypes, the genetic information that (in combination with other factors including environment and management) gives rise to the phenotypes we seek to optimize. Modern plant breeding uses genetic sequencing to guide the breeding process by building models to predict phenotype from genotype and using these models in combination with large gene banks to inform the breeding process. A challenge in this approach is that the space of possible genotypes is huge, but Bayesian optimization can help accelerate the search.

TANAKA, RYOKEI et al. (2018). Bayesian optimization for genomic selection: a method for discovering the best genotype among a large number of candidates. *Theoretical and Applied Genetics* 131(1):93–105.

Biomedical engineering

Difficult optimization problems are pervasive in biomedical engineering due to the complexity of the systems involved and the often considerable cost of gathering data, whether through studies with human subjects, complicated laboratory testing, and/or nontrivial simulation.

COOPY, GLEN WRIGHT et al. (2018). Bayesian Optimization of Personalized Models for Patient Vital-Sign Monitoring. *IEEE Journal of Biomedical and Health Informatics* 22(2):301–310.

GHASSEMI, MOHAMMAD et al. (2014). Global Optimization Approaches for Parameter Tuning in Biomedical Signal Processing: A Focus of Multi-scale Entropy. *Computing in Cardiology* 41(12–2):993–996.

KIM, GILHWAN et al. (2021). Using Bayesian Optimization to Identify Optimal Exoskeleton Parameters Targeting Propulsion Mechanics: A Simulation Study. *bioRxiv*: 2021.01.14.426703.

KIM, MYUNGHEE et al. (2017). Human-in-the-loop Bayesian optimization of wearable device parameters. *PLOS ONE* 12(9):e0184054.

OLOFSSON, SIMON et al. (2019). Bayesian Multiobjective Optimisation With Mixed Analytical and Black-Box Functions: Application to Tissue Engineering. *IEEE Transactions on Biomedical Engineering* 66(3):727–739.

ROBOTICS

Robotics is fraught with difficult optimization problems. A robotic platform may have numerous tunable parameters influencing its behavior, and the dependence of its performance on these parameters may be highly complex. Further, empirical evaluation can be difficult – real-world experiments must proceed in real time, and there may be a considerable set-up time between experiments.

- BANSAL, SOMIL et al. (2017). Goal-Driven Dynamics Learning via Bayesian Optimization. *Proceedings of the 56th Annual IEEE Conference on Decision and Control (CDC 2017)*.
- CALANDRA, ROBERTO et al. (2016). Bayesian optimization for learning gaits under uncertainty: An experimental comparison on a dynamic bipedal walker. *Annals of Mathematics and Artificial Intelligence* 76(1–2):5–23.
- JUNGE, KAI et al. (2020). Improving Robotic Cooking Using Batch Bayesian Optimization. *IEEE Robotics and Automation Letters* 5(2):760–765.
- MARTINEZ-CANTIN, RUBEN et al. (2009). A Bayesian exploration-exploitation approach for optimal online sensing and planning with a visually guided mobile robot. *Autonomous Robotics* 27(2):93–103.

In some cases we may be able to accelerate optimization by augmenting real-world evaluation with simulation in a multifidelity setup.

multifidelity optimization: § 11.5, p. 261

- MARCO, ALONSO et al. (2017). Virtual vs. Real: Trading Off Simulations and Physical Experiments in Reinforcement Learning with Bayesian Optimization. *Proceedings of the 2017 IEEE International Conference on Robotics and Automation (ICRA 2017)*.
- RAI, AKSHARA et al. (n.d.). Using Simulation to Improve Sample-Efficiency of Bayesian Optimization for Bipedal Robots. *Journal of Machine Learning Research* 20(49).

Modeling for robotics

Modeling robot performance as a function of its parameters can be difficult due to nominally high-dimensional parameter spaces and the potential for context-dependent nonstationarity. This difficulty has motivated sophisticated modeling approaches for addressing these issues.

- JAQUIER, NOÉMIE et al. (2019). Bayesian Optimization Meets Riemannian Manifolds in Robot Learning. *Proceedings of the 3rd Conference on Robot Learning (CORL 2019)*.
- MARTINEZ-CANTIN, RUBEN (2017). Bayesian Optimization with Adaptive Kernels for Robot Control. *Proceedings of the 2017 IEEE International Conference on Robotics and Automation (ICRA 2017)*.
- YUAN, KAI et al. (2019). Bayesian Optimization for Whole-Body Control of High-Degree-of-Freedom Robots Through Reduction of Dimensionality. *IEEE Robotics and Automation Letters* 4(3):2268–2275.

Safe and robust optimization

A complication faced in some robotic optimization settings is in ensuring that the evaluated parameters are both *robust* (that is, that performance is not overly sensitive to minor perturbations in the parameters) and *safe*

utility functions: chapter 6, p. 109
 constrained optimization: § 11.2, p. 247

(that is, that there is no chance of catastrophic failure in the robotic platform). We may address these concerns in the design of the optimization policy. For example, we might realize robustness by redefining the utility function in terms of the expected performance of perturbed parameters, and we might realize safety by incorporating appropriate constraints.

- BERKENKAMP, FELIX et al. (2021). Bayesian optimization with safety constraints: safe and automatic parameter tuning in robotics. *Machine Learning Special Issue on Robust Machine Learning*.
- GARCÍA-BARCOS, JAVIER et al. (2021). Robust policy search for robot navigation. *IEEE Robotics and Automation Letters* 6(2):2389–2396.
- NOGUEIRA, JOSÉ et al. (2016). Unscented Bayesian Optimization for Safe Robot Grasping. *Proceedings of the 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2016)*.

Adversarial attacks

A specific safety issue to consider in robotic platforms incorporating deep neural networks is the possibility of adversarial attacks that may be able to alter the environment so as to induce unsafe behavior. Bayesian optimization has been applied to the efficient construction of adversarial attacks; this capability can in turn be used during the design phase seeking to build robotic controllers that are robust to such attack.

- BOLOOR, ADITH et al. (2020). Attacking vision-based perception in end-to-end autonomous driving models. *Journal of Systems Architecture* 110:101766.
- GHOSH, SHROMONA et al. (2018). Verifying Controllers Against Adversarial Examples with Bayesian Optimization. *Proceedings of the 2018 IEEE International Conference on Robotics and Automation (ICRA 2018)*.

CIVIL ENGINEERING

The optimization of large-scale critical systems such as power, transportation, water distribution, and sensor networks can be difficult due to complex dynamics and the considerable expense of reconfiguration. Optimization can sometimes be aided via computer simulation, but the configuration spaces involved can nonetheless be huge, precluding exhaustive search.

- BAHERI, ALI et al. (2017). Altitude Optimization of Airborne Wind Energy Systems: A Bayesian Optimization Approach. *Proceedings of the 2017 American Control Conference (ACC 2017)*.
- CORNEJO-BUENO, L. et al. (2018). Bayesian optimization of a hybrid system for robust ocean wave features prediction. *Neurocomputing* 275: 818–828.

- GARNETT, R. et al. (2010). Bayesian Optimization for Sensor Set Selection. *Proceedings of the 9th ACM/IEEE International Conference on Information Processing in Sensor Networks (IPSN 2010)*.
- GRAMACY, ROBERT B. et al. (2016). Modeling an Augmented Lagrangian for Blackbox Constrained Optimization. *Technometrics* 58(1):1–11.
- HICKISH, BOB et al. (2020). Investigating Bayesian Optimization for rail network optimization. *International Journal of Rail Transporation* 8(4):307–323.
- KOPSIAFTIS, GEORGE et al. (2019). Gaussian Process Regression Tuned by Bayesian Optimization for Seawater Intrusion Prediction. *Computational Intelligence and Neuroscience* 2019:2859429.
- MARCHANT, ROMAN et al. (2012). Bayesian Optimisation for Intelligent Environmental Monitoring. *Proceedings of the 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2012)*.

Structural engineering

The following study applied Bayesian optimization in a structural engineering setting: tuning the hyperparameters of a real-time predictive model for the typhoon-induced response of a ~1000 m real-world bridge.

- ZHANG, YI-MING et al. (2021). Probabilistic Framework with Bayesian Optimization for Predicting Typhoon-Induced Dynamic Responses of a Long-Span Bridge. *Journal of Structural Engineering* 147(1): 04020297.

ELECTRICAL ENGINEERING

Over the past few decades, sophisticated tools have enabled the automation of many aspects of digital circuit design, even for extremely large circuits. However, the design of *analog* circuits remains largely a manual process. As circuits grow increasingly complex, the optimization of analog circuits (for example, to minimize power consumption subject to performance constraints) is becoming increasingly more difficult. Even computer simulation of complex analog circuits can entail significant cost, so careful experimental design is imperative for exploring the design space. Bayesian optimization has proven effective in this regard.

- CHEN, PENG (2015). Bayesian Optimization for Broadband High-Efficiency Power Amplifier Designs. *IEEE Transactions on Microwave Theory and Techniques* 63(12):4263–4272.
- FANG, YAORAN et al. (2018). A Bayesian Optimization and Partial Element Equivalent Circuit Approach to Coil Design in Inductive Power Transfer Systems. *Proceedings of the 2018 IEEE PELS Workshop on Emerging Technologies: Wireless Power Transfer (wow 2018)*.
- LIU, MINGJIE et al. (2020). Closing the Design Loop: Bayesian Optimization Assisted Hierarchical Analog Layout Synthesis. *Proceedings of the 57th ACM/IEEE Design Automation Conference (DAC 2020)*.

- LYU, WENLONG et al. (2018). An Efficient Bayesian Optimization Approach for Automated Optimization of Analog Circuits. *IEEE Transactions on Circuits and Systems—I: Regular Papers* 65(6):1954–1967.
- TORUN, HAKKI MERT et al. (2018). A Global Bayesian Optimization Algorithm and Its Application to Integrated System Design. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 26(4):792–802.

MECHANICAL ENGINEERING

The following study applied Bayesian optimization in a mechanical engineering setting: tuning the parameters of a welding process (via slow and expensive real-world experiments) to maximize weld quality.

- STERLING, DILLON et al. (2015). Welding Parameter Optimization Based on Gaussian Process Regression Bayesian Optimization Algorithm. *Proceedings of the 2015 IEEE International Conference on Automation Science and Engineering (CASE 2015)*.

Aerospace engineering

Nuances in airfoil design can have significant impact on aerodynamic performance, improvements in which can in turn lead to significant cost savings from increased fuel efficiency. Airfoil optimization, however, is challenging due to the large design spaces involved and the nontrivial cost of evaluating a proposed configuration. Empirical measurement requires constructing an airfoil and testing its performance in a wind chamber. This process is too slow to explore the configuration space effectively, but we can *simulate* the process with high fidelity via methods from computational fluid dynamics. These computational surrogates are still fairly costly due to the need to numerically solve nonlinear partial differential equations (the Navier–Stokes equations) at a sufficiently fine resolution, but they are nonetheless cheaper and more easily parallelized than empirical evaluation.²

Bayesian optimization can accelerate airfoil optimization via careful and cost-aware experimental design. One important idea here that can lead to significant computational savings is multifidelity modeling and optimization. It is relatively easy to control the cost–fidelity tradeoff in computational fluid dynamics simulations by altering its resolution accordingly. This allows to rapidly explore the design space with cheap-but-rough simulations, then progressively refine the most promising regions discovered.

- CHAITANYA, PARUCHURI et al. (2020). Bayesian optimisation for low-noise aerofoil design with aerodynamic constraints. *International Journal of Aeroacoustics* 20(1–2):109–129.
- FORRESTER ALEXANDER, I. J. et al. (2009). Recent advances in surrogate-based optimization. *Progress in Aerospace Sciences* 45(1–3):50–79.

² There are many parallels with this situation and that faced in (computational) chemistry and materials science, where quantum-mechanical simulation also requires numerically solving a partial differential equation (the Schrödinger equation).

multifidelity optimization: § 11.5, p. 261

- HEBBAL, ALI et al. (2019). Multi-objective optimization using Deep Gaussian Processes: Application to Aerospace Vehicle Design.
- LAM, REMI R. et al. (2018). Advances in Bayesian Optimization with Applications in Aerospace Engineering.
- PRIEM, RÉMY et al. (2020). An efficient application of Bayesian optimization to an industrial MDO framework for aircraft design. *Proceedings of the 2020 AIAA Aviation Forum*.
- REISENTHEL, PATRICK H. et al. (2011). A Numerical Experiment on Allocating Resources Between Design of Experiment Samples and Surrogate-Based Optimization Infills. *Proceedings of the 2011 AIAA/ASME/ASCE/AHS/ ASC Structures, Structural Dynamics and Materials Conference*.
- ZHENG, HONGYU et al. (2020). Multifidelity kinematic parameter optimization of a flapping airfoil. *Physical Review E* 101(1):013107.

Automobile engineering

Bayesian optimization has also proven useful in automobile engineering. Automobile components and subsystems can have numerous tunable parameters affecting performance, and evaluating a given configuration can be complicated due to complex interactions among vehicle components. Bayesian optimization has shown success in this setting.

- LIESSNER, ROMAN et al. (2019). Simultaneous Electric Powertrain Hardware and Energy Management Optimization of a Hybrid Electric Vehicle Using Deep Reinforcement Learning and Bayesian Optimization. *Proceedings of the 2019 IEEE Vehicle Power and Propulsion Conference (VPPC 2019)*.
- NEUMANN-BROSIG, MATTHIAS et al. (2020). Data-Efficient Autotuning With Bayesian Optimization: An Industrial Control Study. *IEEE Transactions on Control Systems Technology* 28(3):730–740.
- THOMAS, SINNU SUSAN et al. (2019). Designing MacPherson Suspension Architectures using Bayesian Optimization. *Proceedings of the 28th Belgian Dutch Conference on Machine Learning (Benelearn 2019)*.

ALGORITHM CONFIGURATION, HYPERPARAMETER TUNING, AND AUTOML

Complex algorithms and software pipelines can have numerous parameters influencing their performance, and determining the optimal configuration for a given situation can require a significant trial-and-error. Bayesian optimization has demonstrated remarkable success on a range of problems under the umbrella of *algorithm configuration*.

In the most general setting, we may simply model algorithm performance as a black-box function. A challenge here is dealing with high-dimensional parameter spaces that may have complex structure, such as the presence of conditional parameters that only become relevant depending on the settings of other parameters. This can pose a challenge

algorithm configuration

alternative models: § 8.11, p. 196

for Gaussian process models, but alternatives such as random forests can perform admirably.

- DALIBARD, VALENTIN et al. (2017). BOAT: Building Auto-Tuners with Structured Bayesian Optimization. *Proceedings of the 26th International Conference on World Wide Web (www 2017)*.
- GONZALVEZ, JOAN et al. (2019). Financial Applications of Gaussian Processes and Bayesian Optimization. arXiv: 1903.04841 [q-fin.PM].
- HOOS, HOLGER H. (2012). Programming by Optimization. *Communications of the ACM* 55(2):70–80.
- HUTTER, FRANK et al. (2011). Sequential Model-Based Optimization for General Algorithm Configuration. *Proceedings of the 5th Learning and Intelligent Optimization Conference (LION 5)*.
- KUNJIR, MAYURESH (2019). Guided Bayesian Optimization to AutoTune Memory-based Analytics. *Proceedings of the 35th IEEE International Conference on Data Engineering Workshops (ICDEW 2019)*.
- SŁOWIK, AGNIESZKA et al. (2019). Bayesian Optimisation for Heuristic Configuration in Automated Theorem Proving. *Proceedings of the 5th and 6th Vampire Workshops (Vampire 2019)*.
- VARGAS-HERNÁNDEZ, R. A. (2020). Bayesian Optimization for Calibrating and Selecting Hybrid-Density Functional Models. *Journal of Physical Chemistry A* 124(20):4053–4061.

Hyperparameter tuning of machine learning algorithms

hyperparameter tuning

The task of configuring machine learning algorithms in particular is known as *hyperparameter tuning*. Hyperparameter tuning is especially challenging in the era of deep learning, due to the often considerable cost of training and validating a proposed configuration. However, Bayesian optimization has proven effective at this task, and the results of a recent black-box optimization competition (run by TURNER et al. below) focusing on optimization problems from machine learning soundly established its superiority to alternatives such as random search.

- QUITADAMO, ANDREW et al. (2017). Bayesian Hyperparameter Optimization for Machine Learning Based eQTL Analysis. *Proceedings of the 8th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics (BCB 2017)*.
- SNOEK, JASPER et al. (2012). Practical Bayesian Optimization of Machine Learning Algorithms. *Advances in Neural Information Processing Systems 25 (NeurIPS 2012)*.
- TURNER, RYAN et al. (2021). Bayesian Optimization is Superior to Random Search for Machine Learning Hyperparameter Tuning: Analysis of the Black-Box Optimization Challenge 2020. *Proceedings of the NeurIPS 2020 Competition and Demonstration Track*.
- YOGATAMA, DANI et al. (2015). Bayesian Optimization of Text Representations. *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP 2015)*.

In addition to tuning machine learning algorithms to maximize predictive performance, Bayesian optimization can also tune the parameters of other algorithms for learning and inference such as Monte Carlo samplers.

tuning samplers, etc.

HAMZE, FIRAS et al. (2013). Self-Avoiding Random Dynamics on Integer Complex Systems. *ACM Transactions on Modeling and Computer Simulation* 23(1):9.

MAHENDRAN, NIMALAN et al. (2010). Adaptive MCMC with Bayesian Optimization. *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS 2010)*.

One recent high-profile use of Bayesian optimization was in tuning the hyperparameters of DeepMind’s AlphaGo agent. Bayesian optimization was able to improve AlphaGo’s self-play win rate from one-half of games to nearly two-thirds, and the version tuned with Bayesian optimization was used in the final match against Lee Sedol.

CHEN, YUTIAN et al. (2018). Bayesian Optimization in AlphaGo. arXiv: 1812.06855 [cs.LG].

Automated machine learning

The goal of *automated machine learning* (autoML) is to develop automated procedures for tasks related to machine learning in order to boost efficiency and open the power of machine learning to a wider audience. Hyperparameter tuning is one particular instance of this overall vision, but we can also consider the automation of other aspects of machine learning.

automated machine learning (autoML)

In *model selection*, for example, we seek to tune not only the hyperparameters of a machine learning model but also the structure of the model itself. One notable special case is *neural architecture search*, the optimization of neural network architecture. Such problems are particularly difficult due to the discrete, structured nature of the search spaces involved. However, with careful modeling and acquisition strategies, Bayesian optimization becomes a feasible solution.

model assessment, selection, and averaging:
chapter 4, p. 67

neural architecture search

JIN, HAIFENG et al. (2019). Auto-Keras: An Efficient Neural Architecture Search System. *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD 2019)*.

KANDASAMY, KIRTHEVASAN et al. (2018). Neural Architecture Search with Bayesian Optimisation and Optimal Transport. *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*.

MALKOMES, GUSTAVO et al. (2016). Bayesian optimization for automated model selection. *Advances in Neural Information Processing Systems 29 (NeurIPS 2016)*.

WHITE, COLIN et al. (2021). BANANAS: Bayesian Optimization with Neural Architectures for Neural Architecture Search. *Proceedings of 35th AAAI Conference on Artificial Intelligence (AAAI 2021)*.

Researchers have even considered using Bayesian optimization for dynamic model selection during Bayesian optimization itself to realize fully autonomous and robust optimization routines.

MALKOMES, GUSTAVO et al. (2018). Automating Bayesian optimization with Bayesian optimization. *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*.

Optimization of entire machine learning pipelines

We may also consider the joint optimization of entire machine learning pipelines, starting from raw data and automatically constructing a bespoke machine learning system. The space of possible pipelines is absolutely enormous, as we must consider preprocessing steps such as feature selection and engineering in addition to model selection and hyperparameter tuning. Nonetheless, Bayesian optimization has proven up to the task.

Building a working automl system of this scope entails a number of subtle design questions: the space of pipelines to consider, the allocation of training resources to proposed pipelines (which may vary enormously in their complexity), and the modeling of performance as a function of dataset and pipeline, to name a few. All of these questions have received careful consideration in the literature, and the following reference is an excellent entry point to that body of work.

FEURER, MATTHIAS et al. (2020). Auto-Sklearn 2.0: The Next Generation. arXiv: 2007.04074 [cs.LG].

ADAPTIVE HUMAN-COMPUTER INTERFACES

Adaptive human-computer interfaces seek to tailor themselves on-the-fly to suit the preferences of the user and/or the system provider. For example:

- a content provider may wish to learn user preferences to ensure recommended content is relevant,
- a website depending on ad revenue may wish to tune their interface and advertising placement algorithms to maximize ad revenue,
- a computer game may seek to adjust its difficulty dynamically to keep the player engaged, or
- a data-visualization system may seek to infer the user's goals in interacting with a dataset and customize the presentation of data accordingly.

The design of such a system can be challenging, as the space of possible user preferences and/or the space of possible algorithmic settings can be large, and the optimal configuration may change over time or depend on other context. Further, we may only assess the utility of a given interface configuration through user interaction, which is a slow,

cumbersome, and noisy channel from which to glean information. Finally, we face the additional challenge that if we are not careful, the user may become annoyed and simply abandon the platform altogether! Nonetheless, Bayesian optimization has shown success in tuning adaptive interfaces and in related problems such as preference optimization, A/B testing, etc.

- BROCHU, ERIC et al. (2010). A Bayesian Interactive Optimization Approach to Procedural Animation Design. *Proceedings of the 2010 ACM SIGGRAPH/Eurographics Symposium on Computer Animation (SCA 2010)*.
- BROCHU, ERIC et al. (2015). Active Preference Learning with Discrete Choice Data. *Advances in Neural Information Processing Systems 20 (NeurIPS 2007)*.
- GONZÁLEZ, JAVIER et al. (2017). Preferential Bayesian Optimization. *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*.
- KHAJAH, MOHAMMAD M. et al. (2016). Designing Engaging Games Using Bayesian Optimization. *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems (CHI 2016)*.
- LETHAM, BENJAMIN et al. (2019). Constrained Bayesian Optimization with Noisy Experiments. *Bayesian Analysis* 14(2):495–519.
- MONADJEMI, SHAYAN et al. (2020). Active Visual Analytics: Assisted Data Discovery in Interactive Visualizations via Active Search. arXiv: 2010.08155 [cs.HC].

BIBLIOGRAPHY

- ABBASI-YADKORI, YASIN (2012). Online Learning for Linearly Parameterized Control Problems. PhD thesis. University of Alberta.
- ACERBI, LUIGI (2018). Variational Bayesian Monte Carlo. *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*, pp. 8213–8223.
- ACERBI, LUIGI and WEI JI MA (2017). Practical Bayesian Optimization for Model Fitting with Bayesian Adaptive Direct Search. *Advances in Neural Information Processing Systems 30 (NeurIPS 2017)*, pp. 1836–1846.
- ADAMS, RYAN PRESCOTT, IAIN MURRAY, and DAVID J. C. MACKAY (2009). Tractable Nonparametric Bayesian Inference in Poisson Processes with Gaussian Process Intensities. *Proceedings of the 26th International Conference on Machine Learning (ICML 2009)*, pp. 9–16.
- ADLER, ROBERT J. and JONATHAN E. TAYLOR (2007). *Random Fields and Geometry*. Springer Monographs in Mathematics. Springer-Verlag.
- AGRAWAL, RAJEEV (1995). The Continuum-Armed Bandit Problem. *SIAM Journal on Control and Optimization* 33(6):1926–1951.
- AGRAWAL, SHIPRA and NAVIN GOYAL (2012). Analysis of Thompson Sampling for the Multi-armed Bandit Problem. *Proceedings of the 25th Annual Conference on Learning Theory (COLT 2012)*. Vol. 23. Proceedings of Machine Learning Research, pp. 39.1–39.26.
- ÁLVAREZ, MAURICIO A., LORENZO ROSASCO, and NEIL D. LAWRENCE (2012). Kernels for Vector-Valued Functions: A Review. *Foundations and Trends in Machine Learning* 4(3):195–266.
- ARCONES, MIGUEL A. (1992). On the arg max of a Gaussian process. *Statistics & Probability Letters* 15(5):373–374.
- AUER, PETER, NICOLÒ CESÀ-BIANCHI, and PAUL FISCHER (2002). Finite-time Analysis of the Multiarmed Bandit Problem. *Machine Learning* 47(2–3):235–256.
- BAPTISTA, RICARDO and MATTHIAS POLOCZEK (2018). Bayesian Optimization of Combinatorial Structures. *Proceedings of the 35th International Conference on Machine Learning (ICML 2018)*. Vol. 80. Proceedings of Machine Learning Research, pp. 462–471.
- BATHER, JOHN (1996). A Conversation with Herman Chernoff. *Statistical Science* 11(4):335–350.
- BELAKARIA, SYRINE, ARYAN DESHWAL, and JANARDHAN RAO DOPPA (2019). Max-value Entropy Search for Multi-Objective Bayesian Optimization. *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, pp. 7825–7835.
- BELLMAN, RICHARD (1952). On the Theory of Dynamic Programming. *Proceedings of the National Academy of Sciences* 38(8):716–719.
- BELLMAN, RICHARD (1957). *Dynamic Programming*. Princeton University Press.
- BERGER, JAMES O. (1985). *Statistical Decision Theory and Bayesian Analysis*. 2nd ed. Springer Series in Statistics. Springer-Verlag.
- BERGSTRA, JAMES, RÉMI BARDET, YOSHUA BENGIO, and BALÁZS KÉGL (2011). Algorithms for Hyper-Parameter Optimization. *Advances in Neural Information Processing Systems 24 (NeurIPS 2011)*, pp. 2546–2554.
- BERGSTRA, JAMES and YOSHUA BENGIO (2012). Random Search for Hyper-Parameter Optimization. *Journal of Machine Learning Research* 13:281–305.
- BERKENKAMP, FELIX, ANGELA P. SCHOELLING, and ANDREAS KRAUSE (2019). No-Regret Bayesian Optimization with Unknown Hyperparameters. *Journal of Machine Learning Research* 20(50):1–24.

BIBLIOGRAPHY

- BERRY, DONALD A. and BERT FRISTEDT (1985). *Bandit Problems: Sequential Allocation of Experiments*. Monographs on Statistics and Applied Probability. Chapman & Hall.
- BERTSEKAS, DIMITRI P. (2017). *Dynamic Programming and Optimal Control*. 4th ed. Vol. 1. Athena Scientific.
- BOCHNER, S (1933). Monotone Funktionen, Stieltjessche Integrale und harmonische Analyse. *Mathematische Annalen* 108:378–410.
- BOX, G. E. P. (1954). The Exploration and Exploitation of Response Surfaces: Some General Considerations and Examples. *Biometrics* 10(1):16–60.
- BOX, GEORGE E. P., J. STUART HUNTER, and WILLIAM G. HUNTER (2005). *Statistics for Experimenters: Design, Innovation, and Discovery*. 2nd ed. Wiley Series in Probability and Statistics. John Wiley & Sons.
- BOX, G. E. P. and K. B. WILSON (1951). On the Experimental Attainment of Optimum Conditions. *Journal of the Royal Statistical Society Series B (Methodological)* 13(1):1–45.
- BOX, G. E. P. and P. V. YOULE (1954). The Exploration and Exploitation of Response Surfaces: An Example of the Link Between the Fitted Surface and the Basic Mechanism of the System. *Biometrics* 11(3):287–323.
- BREIMAN, LEO (2001). Random Forests. *Machine Learning* 45(1):5–32.
- BRENT, RICHARD P. (1973). *Algorithms for Minimization Without Derivatives*. Prentice-Hall Series in Automatic Computation. Prentice-Hall.
- BROCHU, ERIC, VLAD M. CORA, and NANDO DE FREITAS (2010). A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning. arXiv: 1012.2599 [cs.LG].
- BROCHU, ERIC, NANDO DE FREITAS, and ABHIJEET GHOSH (2015). Active Preference Learning with Discrete Choice Data. *Advances in Neural Information Processing Systems 20 (NeurIPS 2007)*, pp. 409–416.
- BROOKS, STEVE, ANDREW GELMAN, GALIN L. JONES, and XIAO-LI MENG, eds. (2011). *Handbook of Markov Chain Monte Carlo*. Handbooks of Modern Statistical Methods. Chapman & Hall.
- BUBECK, SÉBASTIEN and NICOLÒ CESA-BIANCHI (2012). Regret Analysis of Stochastic and Nonstochastic Multi-armed Bandit Problems. *Foundations and Trends in Machine Learning* 5(1):1–122.
- BUBECK, SÉBASTIEN, RÉMI MUNOS, and GILLES STOLTZ (2009). Pure Exploration in Multi-armed Bandits Problems. *Proceedings of the 20th International Conference on Algorithmic Learning Theory (ALT 2009)*. Vol. 5809. Lecture Notes in Computer Science. Springer-Verlag, pp. 23–37.
- BULL, ADAM D. (2011). Convergence Rates of Efficient Global Optimization Algorithms. *Journal of Machine Learning Research* 12(88):2879–2904.
- CAFLISCH, RUSSEL E. (1998). Monte Carlo and quasi-Monte Carlo methods. *Acta Numerica* 7: 1–49.
- CALANDRA, ROBERTO, JAN PETERS, CARL EDWARD RASMUSSEN, and MARC PETER DEISENROTH (2016). Manifold Gaussian Processes for Regression. *Proceedings of the 2016 International Joint Conference on Neural Networks (IJCNN 2016)*, pp. 3338–3345.
- CALVIN, J. and A. ŽILINSKAS (1999). On the Convergence of the P-Algorithm for One-Dimensional Global Optimization of Smooth Functions. *Journal of Optimization Theory and Applications* 102(3):479–495.
- CALVIN, JAMES M. (1993). Consistency of a Myopic Bayesian Algorithm for One-Dimensional Global Optimization. *Journal of Global Optimization* 3(2):223–232.

- CALVIN, JAMES M. (2000). Convergence Rate of the P-Algorithm for Optimization of Continuous Functions. In: *Approximation and Complexity in Numerical Optimization: Continuous and Discrete Problems*. Ed. by PANOS M. PARDALOS. Vol. 42. Nonconvex Optimization and Its Applications. Springer–Verlag, pp. 116–129.
- CALVIN, JAMES M. and ANTANAS ŽILINSKAS (2001). On Convergence of a P-Algorithm Based on a Statistical Model of Continuously Differentiable Functions. *Journal of Global Optimization* 19(3):229–245.
- CHAPELLE, OLIVIER and LIHONG LI (2011). An Empirical Evaluation of Thompson Sampling. *Advances in Neural Information Processing Systems 24 (neurips 2011)*, pp. 2249–2257.
- CHERNOFF, HERMAN (1959). Sequential Design of Experiments. *The Annals of Mathematical Statistics* 30(3):755–770.
- CHERNOFF, HERMAN (1972). *Sequential Analysis and Optimal Design*. CBMS–NSF Regional Conference Series in Applied Mathematics. Society for Industrial and Applied Mathematics.
- CHEVALIER, CLÉMENT and DAVID GINSBOURGER (2013). Fast Computation of the Multi-points Expected Improvement with Applications in Batch Selection. *Proceedings of the 7th Learning and Intelligent Optimization Conference (LION 7)*. Vol. 7997. Lecture Notes in Computer Science. Springer–Verlag, pp. 59–69.
- CHOWDHURY, SAYAK RAY and ADITYA GOPALAN (2017). On Kernelized Multi-armed Bandits. *Proceedings of the 34th International Conference on Machine Learning (icml 2017)*. Vol. 70. Proceedings of Machine Learning Research, pp. 844–853.
- CLARK, CHARLES E. (1961). The Greatest of a Finite Set of Random Variables. *Operations Research* 9(2):145–162.
- CONTAL, EMILE, DAVID BUFFONI, ALEXANDRE ROBICQUET, and NICOLAS VAYATIS (2013). Parallel Gaussian Process Optimization with Upper Confidence Bound and Pure Exploration. *Proceedings of the 2013 European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ecml pkdd 2013)*. Vol. 8188. Lecture Notes in Computer Science. Springer–Verlag, pp. 225–240.
- COVER, THOMAS M. and JOY A. THOMAS (2006). *Elements of Information Theory*. 2nd ed. John Wiley & Sons.
- COX, DENNIS D. and SUSAN JOHN (1992). A Statistical Method for Global Optimization. *Proceedings of the 1992 IEEE International Conference on Systems, Man, and Cybernetics (smc 1992)*, pp. 1241–1246.
- CUNNINGHAM, JOHN P., PHILIPP HENNIG, and SIMON LACOSTE-JULIEN (2011). Gaussian Probabilities and Expectation Propagation. arXiv: 1111.6832 [stat.ML].
- CUTAJAR, KURT, MICHAEL A. OSBORNE, JOHN P. CUNNINGHAM, and MAURIZIO FILIPPONE (2016). Preconditioning Kernel Matrices. *Proceedings of the 33rd International Conference on Machine Learning (icml 2016)*. Vol. 48. Proceedings of Machine Learning Research, pp. 2529–2538.
- DAI, ZHONGXIANG, HAIBIN YU, BRYAN KIAN HSIANG LOW, and PATRICK JAILLET (2019). Bayesian Optimization Meets Bayesian Optimal Stopping. *Proceedings of the 36th International Conference on Machine Learning (icml 2019)*. Vol. 97. Proceedings of Machine Learning Research, pp. 1496–1506.
- DALIBARD, VALENTIN, MICHAEL SCHAARSCHMIDT, and EIKO YONEKI (2017). BOAT: Building Auto-Tuners with Structured Bayesian Optimization. *Proceedings of the 26th International Conference on World Wide Web (www 2017)*, pp. 479–488.

BIBLIOGRAPHY

- DANI, VARSHA, THOMAS P. HAYES, and SHAM M. KAKADE (2008). Stochastic Linear Optimization Under Bandit Feedback. *Proceedings of the 21st Conference on Learning Theory (COLT 2008)*, pp. 355–366.
- DAVIS, PHILIP J. and PHILIP RABINOWITZ (1984). *Methods of Numerical Integration*. 2nd ed. Computer Science and Applied Mathematics. Academic Press.
- DE ATH, GEORGE, JONATHAN E. FIELDSEND, and RICHARD M. EVERSON (2020). What do you Mean? The Role of the Mean Function in Bayesian Optimization. *Proceedings of the 2020 Genetic and Evolutionary Computation Conference (GECCO 2020)*, pp. 1623–1631.
- DE FREITAS, NANDO, ALEX J. SMOLA, and MASROUR ZOGHI (2012a). Regret Bounds for Deterministic Gaussian Process Bandits. arXiv: 1203.2177 [cs.LG].
- DE FREITAS, NANDO, ALEX J. SMOLA, and MASROUR ZOGHI (2012b). Exponential Regret Bounds for Gaussian Process Bandits with Deterministic Observations. *Proceedings of the 29th International Conference on Machine Learning (ICML 2012)*, pp. 955–962.
- DEGROOT, MORRIS H. (1970). *Optimal Statistical Decisions*. McGraw-Hill.
- DESAUTELS, THOMAS, ANDREAS KRAUSE, and JOEL W. BURDICK (2014). Parallelizing Exploration-Exploitation Tradeoffs in Gaussian Process Bandit Optimization. *Journal of Machine Learning Research* 15(119):4053–4103.
- DIACONIS, PERSI (1988). Bayesian Numerical Analysis. In: *Statistical Decision Theory and Related Topics IV*. Ed. by SHANTI S. GUPTA and JAMES O. BERGER. Vol. 1, pp. 163–175.
- DJOLONGA, JOSIP, ANDREAS KRAUSE, and VOLKAN CEVHER (2013). High-Dimensional Gaussian Process Bandits. *Advances in Neural Information Processing Systems 26 (NeurIPS 2013)*, pp. 1025–1033.
- DOMHAN, TOBIAS, JOST TOBIAS SPRINGENBERG, and FRANK HUTTER (2015). Speeding up Automatic Hyperparameter Optimization of Deep Neural Networks by Extrapolation of Learning Curves. *Proceedings of the 24th International Conference on Artificial Intelligence (IJCAI 2015)*, pp. 3460–3468.
- DUVENAUD, DAVID, JAMES ROBERT LLOYD, ROGER GROSSE, JOSHUA B. TENENBAUM, and ZOUBIN GHAHRAMANI (2013). Structure Discovery in Nonparametric Regression through Compositional Kernel Search. *Proceedings of the 30th International Conference on Machine Learning (ICML 2013)*. Vol. 28. Proceedings of Machine Learning Research, pp. 1166–1174.
- EMMERICH, MICHAEL and BORIS NAUJOKS (2004). Metamodel Assisted Multiobjective Optimisation Strategies and their Application in Airfoil Design. In: *Adaptive Computing in Design and Manufacture VI*. Ed. by I. C. PARMEE, pp. 249–260.
- EMMERICH, MICHAEL T. M., KYRIAKOS C. GIANNAKOGLOU, and BORIS NAUJOKS (2006). Single- and Multiobjective Evolutionary Optimization Assisted by Gaussian Random Field Metamodels. *IEEE Transactions on Evolutionary Computation* 10(4):421–439.
- FERNÁNDEZ-DELGADO, MANUEL, EVA CERNADAS, SENÉN BARRO, and DINANI AMORIM (2014). Do we Need Hundreds of Classifiers to Solve Real World Classification Problems? *Journal of Machine Learning Research* 15(90):3133–3181.
- FEURER, MATTHIAS, JOST TOBIAS SPRINGENBERG, and FRANK HUTTER (2015). Initializing Bayesian Hyperparameter Optimization via Meta-Learning. *Proceedings of 29th AAAI Conference on Artificial Intelligence (AAAI 2015)*, pp. 1128–1135.
- FISHER, RONALD A. (1935). *The Design of Experiments*. Oliver and Boyd.
- FLEISCHER, M. (2003). The Measure of Pareto Optima: Applications to Multi-objective Metaheuristics. *Proceedings of the 2nd International Conference on Evolutionary Multi-Criterion Optimization (EMO 2003)*. Vol. 2632. Lecture Notes in Computer Science. Springer-Verlag, pp. 519–533.

- FORRESTER, ALEXANDER I. J., ANDY J. KEANE, and NEIL W. BRESSLOFF (2006). Design and Analysis of “Noisy” Computer Experiments. *AIAA Journal* 44(10):2331–2339.
- FRAZIER, PETER and WARREN POWELL (2007). The Knowledge Gradient Policy for Offline Learning with Independent Normal Rewards. *Proceedings of the 2007 IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning (ADPRL 2007)*, pp. 143–150.
- FRAZIER, PETER, WARREN POWELL, and SAVAS DAYANIK (2009). The Knowledge-Gradient Policy for Correlated Normal Beliefs. *INFORMS Journal on Computing* 21(4):599–613.
- FRIEDMAN, MILTON and L. J. SAVAGE (1947). Planning Experiments Seeking Maxima. In: *Selected Techniques of Statistical Analysis for Scientific and Industrial Research, and Production and Management Engineering*. Ed. by CHURCHILL EISENHART, MILLARD W. HARTAY, and W. ALLEN WALLIS. McGraw-Hill, pp. 363–372.
- GARDNER, JACOB R., CHUAN GUO, KILIAN Q. WEINBERGER, ROMAN GARNETT, and ROGER GROSSE (2017). Discovering and Exploiting Additive Structure for Bayesian Optimization. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS 2017)*. Vol. 54. Proceedings of Machine Learning Research, pp. 1311–1319.
- GARDNER, JACOB R., MATT J. KUSNER, ZHIXIANG (EDDIE) XU, KILIAN Q. WEINBERGER, and JOHN P. CUNNINGHAM (2014). Bayesian Optimization with Inequality Constraints. *Proceedings of the 31st International Conference on Machine Learning (ICML 2014)*. Vol. 32. Proceedings of Machine Learning Research, pp. 937–945.
- GARDNER, JACOB R., GEOFF PLEISS, DAVID BINDEL, KILIAN Q. WEINBERGER, and ANDREW GORDON WILSON (2018). GPyTorch: Blackbox Matrix–Matrix Gaussian Process Inference with GPU Acceleration. *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*, pp. 7576–7586.
- GARNETT, ROMAN, THOMAS GÄRTNER, MARTIN VOGT, and JÜRGEN BAJORATH (2015). Introducing the ‘active search’ method for iterative virtual screening. *Journal of Computer-Aided Molecular Design* 29(4):305–314.
- GARNETT, ROMAN, YAMUNA KRISHNAMURTHY, XUEHAN XIONG, JEFF SCHNEIDER, and RICHARD MANN (2012). Bayesian Optimal Active Search and Surveying. *Proceedings of the 29th International Conference on Machine Learning (ICML 2012)*, pp. 1239–1246.
- GARNETT, R., M. A. OSBORNE, and S. J. ROBERTS (2010). Bayesian Optimization for Sensor Set Selection. *Proceedings of the 9th ACM/IEEE International Conference on Information Processing in Sensor Networks (IPSN 2010)*, pp. 209–219.
- GARNETT, ROMAN, MICHAEL A. OSBORNE, and PHILIPP HENNIG (2014). Active Learning of Linear Embeddings for Gaussian Processes. *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence (UAI 2014)*, pp. 230–239.
- GELBART, MICHAEL A., JASPER SNOEK, and RYAN P. ADAMS (2014). Bayesian Optimization with Unknown Constraints. *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence (UAI 2014)*, pp. 250–259.
- GELMAN, ANDREW and AKI VEHTARI (2021). What are the most important statistical ideas of the past 50 years? *Journal of the American Statistical Association*.
- GERGONNE, JOSEPH DIEZ (1815). Application de la méthode des moindres quarrés à l’interpolation des suites. *Annales de Mathématiques pures et appliquées* 6:242–252.
- GHOSAL, SUBHASHIS and ANINDYA ROY (2006). Posterior Consistency of Gaussian Process Prior for Nonparametric Binary Regression. *The Annals of Statistics* 34(5):2413–2429.
- GIBBS, MARK N. (1997). Bayesian Gaussian Processes for Regression and Classification. PhD thesis. University of Cambridge.

BIBLIOGRAPHY

- GILBOA, ELAD, YUNUS SAATÇI, and JOHN P. CUNNINGHAM (2013). Scaling Multidimensional Gaussian Processes using Projected Additive Approximations. *Proceedings of the 30th International Conference on Machine Learning (ICML 2013)*. Vol. 28. Proceedings of Machine Learning Research, pp. 454–461.
- GINSBOURGER, DAVID and RODOLPHE LE RICHE (2010). Towards Gaussian Process-based Optimization with Finite Time Horizon. *Proceedings of the 9th International Workshop in Model-Oriented Design and Analysis (MODA 9)*. Contributions to Statistics. Springer-Verlag, pp. 89–96.
- GINSBOURGER, DAVID, RODOLPHE LE RICHE, and LAURENT CARRARO (2010). Kriging is well-suited to parallelize optimization. In: *Computational Intelligence in Expensive Optimization Problems*. Ed. by YOEL YENNE and CHI-KEONG GO. Adaptation Learning and Optimization. Springer-Verlag, pp. 131–162.
- GOLUB, GENE H. and CHARLES F. VAN LOAN (2013). *Matrix Computations*. 4th ed. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press.
- GÓMEZ-BOMBARELLI, RAFAEL, JENNIFER N. WEI, DAVID DUVENAUD, JOSÉ MIGUEL HERNÁNDEZ-LOBATO, BENJAMÍN SÁNCHEZ-LENGELING, DNNIS SHEBERLA, et al. (2018). Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules. *ACS Central Science* 4(2):268–276.
- GONZÁLEZ, JAVIER, ZHENWEN DAI, ANDREAS DAMIANOU, and NEIL D. LAWRENCE (2017). Preferential Bayesian Optimization. *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*. Vol. 70. Proceedings of Machine Learning Research, pp. 1282–1291.
- GONZÁLEZ, JAVIER, ZHENWEN DAI, PHILIPP HENNIG, and NEIL LAWRENCE (2016a). Batch Bayesian Optimization via Local Penalization. *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS 2016)*. Vol. 51. Proceedings of Machine Learning Research, pp. 648–657.
- GONZÁLEZ, JAVIER, MICHAEL OSBORNE, and NEIL D. LAWRENCE (2016b). GLASSES: Relieving The Myopia Of Bayesian Optimization. *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS 2016)*. Vol. 51. Proceedings of Machine Learning Research, pp. 790–799.
- GRAMACY, ROBERT B. and HERBERT K. H. LEE (2011). Optimization Under Unknown Constraints. In: *Bayesian Statistics 9*. Ed. by J. M. BERNARDO, M. J. BAYARRI, J. O. BERGER, A. P. DAWID, D. HECKERMAN, A. F. M. SMITH, et al. Oxford University Press, pp. 229–256.
- GRANMO, OLE-CHRISTOFFER (2010). Solving two-armed Bernoulli bandit problems using a Bayesian learning automaton. *International Journal of Intelligent Computing and Cybernetics* 3(2):207–234.
- GRÜNEWÄLDER, STEFFEN, JEAN-YVES AUDIBERT, MANFRED OPPER, and JOHN SHawe-TAYLOR (2010). Regret Bounds for Gaussian Process Bandit Problems. *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS 2010)*. Vol. 9. Proceedings of Machine Learning Research, pp. 273–280.
- HANSEN, NIKOLAUS (2016). The cma Evolution Strategy: A Tutorial. arXiv: 1604 . 00772 [cs.LG].
- HASTIE, TREVOR and ROBERT TIBSHIRANI (1986). Generalized Additive Models. *Statistical Science* 1(3):297–318.
- HENNIG, PHILIPP, MICHAEL A. OSBORNE, and MARK GIROLAMI (2015). Probabilistic numerics and uncertainty in computations. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 471(2179):20150142.

- HENNIG, PHILIPP, MICHAEL A. OSBORNE, and HANS KERSTING (2022). *Probabilistic Numerics: Computation as Machine Learning*. Cambridge University Press.
- HENNIG, PHILIPP and CHRISTIAN J. SCHULER (2012). Entropy Search for Information-Efficient Global Optimization. *Journal of Machine Learning Research* 13(Jun):1809–1837.
- HENSMAN, JAMES, ALEXANDER G. DE G. MATTHEWS, MAURIZIO FILIPPONE, and ZOUBIN GHAHRAMANI (2015). MCMC for Variationally Sparse Gaussian Processes. *Advances in Neural Information Processing Systems* 28 (*NeurIPS 2015*), pp. 1648–1656.
- HERNÁNDEZ-LOBATO, DANIEL, JOSÉ MIGUEL HERNÁNDEZ-LOBATO, AMAR SHAH, and RYAN P. ADAMS (2016a). Predictive Entropy Search for Multi-objective Bayesian Optimization. *Proceedings of the 33rd International Conference on Machine Learning (ICML 2016)*. Vol. 48. Proceedings of Machine Learning Research, pp. 1492–1501.
- HERNÁNDEZ-LOBATO, JOSÉ MIGUEL, MICHAEL A. GELBART, RYAN P. ADAMS, MATTHEW W. HOFFMAN, and ZOUBIN GHAHRAMANI (2016b). A General Framework for Constrained Bayesian Optimization using Information-based Search. *Journal of Machine Learning Research* 17:1–53.
- HERNÁNDEZ-LOBATO, JOSÉ MIGUEL, MATTHEW W. HOFFMAN, and ZOUBIN GHAHRAMANI (2014). Predictive Entropy Search for Efficient Global Optimization of Black-box Functions. *Advances in Neural Information Processing Systems* 27 (*NeurIPS 2014*), pp. 918–926.
- HERNÁNDEZ-LOBATO, JOSÉ MIGUEL, JAMES REQUEIMA, EDWARD O. PYZER-KNAPP, and ALÁN ASPURU-GUZIK (2017). Parallel and Distributed Thompson Sampling for Large-scale Accelerated Exploration of Chemical Space. *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*. Vol. 70. Proceedings of Machine Learning Research, pp. 1470–1479.
- HESTENES, MAGNUS R. and EDUARD STIEFEL (1952). Methods of Conjugate Gradients for Solving Linear Systems. *Journal of Research of the National Bureau of Standards* 49(6): 409–436.
- HOFFMAN, MATTHEW D. and ANDREW GELMAN (2014). The No-U-turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research* 15(4):1593–1623.
- HOFFMAN, MATTHEW W. and ZOUBIN GHAHRAMANI (2015). Output-Space Predictive Entropy Search for Flexible Global Optimization. *Bayesian Optimization: Scalability and Flexibility Workshop (BayesOpt 2015)*, Conference on Neural Information Processing Systems (*NeurIPS 2015*).
- HOFFMAN, MATTHEW W. and BOBAK SHAHRIARI (2014). Modular mechanisms for Bayesian optimization. *Bayesian Optimization in Academia and Industry (BayesOpt 2014)*, Conference on Neural Information Processing Systems (*NeurIPS 2014*).
- HOTELLING, HAROLD (1941). Experimental Determination of the Maximum of a Function. *The Annals of Mathematical Statistics* 12(1):20–45.
- HOULSBY, NEIL, JOSÉ MIGUEL HERNÁNDEZ-LOBATO, FERENC HUSZÁR, and ZOUBIN GHAHRAMANI (2012). Collaborative Gaussian Processes for Preference Learning. *Advances in Neural Information Processing Systems* 25 (*NeurIPS 2012*), pp. 2096–2104.
- HUANG, D., T. T. ALLEN, W. I. NOTZ, and R. A. MILLER (2006a). Sequential kriging optimization using multiple-fidelity evaluations. *Structural and Multidisciplinary Optimization* 32(5):369–382.
- HUANG, D., T. T. ALLEN, W. I. NOTZ, and N. ZENG (2006b). Global Optimization of Stochastic Black-Box Systems via Sequential Kriging Meta-Models. *Journal of Global Optimization* 34(3):441–466.

BIBLIOGRAPHY

- HUTTER, FRANK, HOLGER H. HOOS, and KEVIN LEYTON-BROWN (2011). Sequential Model-Based Optimization for General Algorithm Configuration. *Proceedings of the 5th Learning and Intelligent Optimization Conference (LION 5)*. Vol. 6683. Lecture Notes in Computer Science. Springer-Verlag, pp. 507–523.
- HUTTER, FRANK, LIN XU, HOLGER H. HOOS, and KEVIN LEYTON-BROWN (2014). Algorithm runtime prediction: Methods & evaluation. *Artificial Intelligence* 206:79–111.
- INGERSOLL JR., JONATHAN E. (1987). *Theory of Financial Decision Making*. Rowman & Littlefield Studies in Financial Economics. Rowman & Littlefield.
- JANZ, DAVID, DAVID R. BURT, and JAVIER GONZÁLEZ (2020). Bandit optimisation of functions in the Matérn kernel RKHS. *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS 2020)*. Vol. 108. Proceedings of Machine Learning Research, pp. 2486–2495.
- JIANG, SHALI, HENRY CHAI, JAVIER GONZÁLEZ, and ROMAN GARNETT (2020a). BINOCULARS for Efficient, Nonmyopic Sequential Experimental Design. *Proceedings of the 37th International Conference on Machine Learning (ICML 2020)*. Vol. 119. Proceedings of Machine Learning Research, pp. 4794–4803.
- JIANG, SHALI, DANIEL R. JIANG, MAXIMILIAN BALANDAT, BRIAN KARRER, JACOB R. GARDNER, and ROMAN GARNETT (2020b). Efficient Nonmyopic Bayesian Optimization via One-Shot Multi-Step Trees. *Advances in Neural Information Processing Systems 33 (NeurIPS 2020)*, pp. 18039–18049.
- JIANG, SHALI, GUSTAVO MALKOMES, GEOFF CONVERSE, ALYSSA SHOFNER, BENJAMIN MOSELEY, and ROMAN GARNETT (2017). Efficient Nonmyopic Active Search. *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*. Vol. 70. Proceedings of Machine Learning Research, pp. 1714–1723.
- JONES, D. R., C. D. PERTTUNEN, and B. E. STUCKMAN (1993). Lipschitzian Optimization Without the Lipschitz Constant. *Journal of Optimization Theory and Application* 79(1):157–181.
- JONES, DONALD R. (2001). A Taxonomy of Global Optimization Methods Based on Response Surfaces. *Journal of Global Optimization* 21(4):345–383.
- JONES, DONALD R., MATTHIAS SCHONLAU, and WILLIAM J. WELCH (1998). Efficient Global Optimization of Expensive Black-Box Functions. *Journal of Global Optimization* 13(4): 455–492.
- JYLÄNKI, PASI, JARNO VANHATALO, and AKI VEHTARI (2011). Robust Gaussian Process Regression with a Student-*t* Likelihood. *Journal of Machine Learning Research* 12(99): 3227–3257.
- KANAGAWA, MOTONOBU, PHILIPP HENNIG, DINO SEJDINOVIC, and BHARATH K. SRIPERUM-BUDUR (2018). Gaussian Processes and Kernel Methods: A Review on Connections and Equivalences. arXiv: 1807.02582 [stat.ML].
- KANDASAMY, KIRTHEVANAN, GAUTAM DASARATHY, JUNIER OLIVA, JEFF SCHNEIDER, and BARNABÁS PÓCZOS (2016). Gaussian Process Bandit Optimisation with Multi-fidelity Evaluations. *Advances in Neural Information Processing Systems 29 (NeurIPS 2016)*, pp. 992–1000.
- KANDASAMY, KIRTHEVANAN, GAUTAM DASARATHY, JEFF SCHNEIDER, and BARNABÁS PÓCZOS (2017). Multi-fidelity Bayesian Optimisation with Continuous Approximations. *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*. Vol. 70. Proceedings of Machine Learning Research, pp. 1799–1808.
- KANDASAMY, KIRTHEVANAN, AKSHAY KRISHNAMURTHY, JEFF SCHNEIDER, and BARNABÁS PÓCZOS (2018). Parallelised Bayesian Optimisation via Thompson Sampling. *Proceedings*

- of the 21st International Conference on Artificial Intelligence and Statistics (AISTATS 2018).* Vol. 84. Proceedings of Machine Learning Research, pp. 133–142.
- KANDASAMY, KIRTHEVASAN, JEFF SCHNEIDER, and BARNABÁS PÓCZOS (2015). High Dimensional Bayesian Optimisation and Bandits via Additive Models. *Proceedings of the 32nd International Conference on Machine Learning (ICML 2015)*. Vol. 37. Proceedings of Machine Learning Research, pp. 295–304.
- KATHURIA, TARUN, AMIT DESHPANDE, and PUSHMEET KOHLI (2016). Batched Gaussian Process Bandit Optimization via Determinantal Point Processes. *Advances in Neural Information Processing Systems 29 (NeurIPS 2016)*, pp. 4206–4214.
- KIM, JEANKYUNG and DAVID POLLARD (1990). Cube Root Asymptotics. *The Annals of Statistics* 18(1):191–219.
- KIM, JUNGTAEK, MICHAEL MCCOURT, TACKGEUN YOU, SAEHOON KIM, and SEUNGJIN CHOI (2021). Bayesian optimization with approximate set kernels. *Machine Learning* 110(5):857–879.
- KLEIN, AARON, SIMON BARTELS, STEFAN FALKNER, PHILIPP HENNIG, and FRANK HUTTER (2015). Towards efficient Bayesian Optimization for Big Data. *Bayesian Optimization: Scalability and Flexibility Workshop (BayesOpt 2015), Conference on Neural Information Processing Systems (NeurIPS 2015)*.
- KLEIN, AARON, STEFAN FALKNER, JOST TOBIAS SPRINGENBERG, and FRANK HUTTER (2017). Learning Curve Prediction with Bayesian Neural Networks. *Proceedings of the 5th International Conference on Learning Representations (ICLR 2017)*.
- KNOWLES, JOSHUA (2005). PareGO: A Hybrid Algorithm With On-Line Landscape Approximation for Expensive Multiobjective Optimization Problems. *IEEE Transactions on Evolutionary Computation* 10(1):50–66.
- KO, CHUN-WA, JON LEE, and MAURICE QUEYRANNE (1995). An Exact Algorithm for Maximum Entropy Sampling. *Operations Research* 43(4):684–691.
- KONISHI, SADANORI and GENSHIRO KITAGAWA (2008). *Information Criteria and Statistical Modeling*. Springer Series in Statistics. Springer-Verlag.
- KSCHISCHANG, FRANK R., BRENDAN J. FREY, and HANS-ANDREA LEOLIGER (2001). Factor Graphs and the Sum-Product Algorithm. *IEEE Transactions on Information Theory* 47(2):498–519.
- KULESZA, ALEX and BEN TASKAR (2012). Determinantal Point Processes for Machine Learning. *Foundations and Trends in Machine Learning* 5(2–3):123–286.
- KUSHNER, HAROLD J. (1962). A Versatile Stochastic Model of a Function of Unknown and Time Varying Form. *Journal of Mathematical Analysis and Applications* 5(1):150–167.
- KUSHNER, H. J. (1964). A New Method of Locating the Maximum Point of an Arbitrary Multipeak Curve in the Presence of Noise. *Journal of Basic Engineering* 86(1):97–106.
- KUSS, MALTE (2006). Gaussian Process Models for Robust Regression, Classification, and Reinforcement Learning. PhD thesis. Technische Universität Darmstadt.
- LAI, T. L. and HERBERT ROBBINS (1985). Asymptotically Efficient Adaptive Allocation Rules. *Advances in Applied Mathematics* 6(1):4–22.
- LAM, REMI R., KAREN E. WILCOX, and DAVID H. WOLPERT (2016). Bayesian Optimization with a Finite Budget: An Approximate Dynamic Programming Approach. *Advances in Neural Information Processing Systems 29 (NeurIPS 2016)*, pp. 883–891.
- LANGE, KENNETH L., RODERICK J. A. LITTLE, and JEREMY M. G. TAYLOR (1989). Robust Statistical Modeling Using the t Distribution. *Journal of the American Statistical Association* 84(408):881–896.

BIBLIOGRAPHY

- LATTIMORE, TOR and CSABA SZEPESVÁRI (2020). *Bandit Algorithms*. Cambridge University Press.
- LÁZARO-GREDILLA, MIGUEL, JOAQUIN QUIÑONERO-CANDELA, CARL EDWARD RASMUSSEN, and ANÍBAL R. FIGUEIRAS-VIDAL (2010). Sparse Spectrum Gaussian Process Regression. *Journal of Machine Learning Research* 11(Jun):1865–1881.
- LETHAM, BENJAMIN, BRIAN KARRER, GUILHERME OTTONI, and EYTAN BAKSHY (2019). Constrained Bayesian Optimization with Noisy Experiments. *Bayesian Analysis* 14(2): 495–519.
- LEVINA, ELIZAVETA and PETER J. BICKEL (2004). Maximum Likelihood Estimation of Intrinsic Dimension. *Advances in Neural Information Processing Systems 17 (NeurIPS 2004)*, pp. 777–784.
- LÉVY, PAUL (1948). *Processus stochastiques et mouvement brownien*. Gauthier–Villars.
- LI, CHUN-LIANG, KIRTHEVASAN KANDASAMY, BARNABÁS PÓCZOS, and JEFF SCHNEIDER (2016). High Dimensional Bayesian Optimization via Restricted Projection Pursuit Models. *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS 2016)*. Vol. 51. Proceedings of Machine Learning Research, pp. 884–892.
- LI, CHUNYUAN, HEERAD FARKHOOR, ROSANNE LIU, and JASON YOSINSKI (2018a). Measuring the Intrinsic Dimension of Objective Landscapes. *Proceedings of the 6th International Conference on Learning Representations (ICLR 2018)*.
- LI, LISHA, KEVIN JAMIESON, GIULIA DESALVO, AFSHIN ROSTAMIZADEH, and AMEET TALWALKAR (2018b). Hyperband: A Novel Bandit-Based Approach to Hyperparameter Optimization. *Journal of Machine Learning Research* 18(185):1–52.
- LINDLEY, D. V. (1956). On a Measure of the Information Provided by an Experiment. *The Annals of Mathematical Statistics* 27(4):986–1005.
- LINDLEY, D. V. (1972). *Bayesian Statistics, A Review*. CBMS–NSF Regional Conference Series in Applied Mathematics. Society for Industrial and Applied Mathematics.
- LOCATELLI, M. (1997). Bayesian Algorithms for One-Dimensional Global Optimization. *Journal of Global Optimization* 10(1):57–76.
- LÖWNER, KARL (1934). Über monotone Matrixfunktionen. *Mathematische Zeitschrift* 38:177–216.
- LUKIĆ, MILAN N. and JAY H. BEDER (2001). Stochastic Processes with Sample Paths in Reproducing Kernel Hilbert Spaces. *Transactions of the American Mathematical Society* 353(10):3945–3969.
- MACKAY, DAVID J. C. (1998). Introduction to Gaussian Processes. *Neural Networks and Machine Learning*. Ed. by CHRISTOPHER M. BISHOP. Vol. 168. NATO ASI Series F: Computer and Systems Sciences. Springer–Verlag.
- MACKAY, DAVID J. C. (2003). *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press.
- MALKOMES, GUSTAVO and ROMAN GARNETT (2018). Automating Bayesian optimization with Bayesian optimization. *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*, pp. 5984–5994.
- MALKOMES, GUSTAVO, CHIP SCHAFF, and ROMAN GARNETT (2016). Bayesian optimization for automated model selection. *Advances in Neural Information Processing Systems 29 (NeurIPS 2016)*, pp. 2900–2908.
- MARMIN, SÉBASTIEN, CLÉMENT CHEVALIER, and DAVID GINSBOURGER (2015). Differentiating the Multipoint Expected Improvement for Optimal Batch Design. *Proceedings of the 1st International Workshop on Machine Learning, Optimization, and Big Data (MOD 2015)*. Vol. 9432. Lecture Notes in Computer Science. Springer–Verlag, pp. 37–48.

- MARSCHAK, JACOB and ROY RADNER (1972). *Economic Theory of Teams*. Yale University Press.
- MARTINEZ-CANTIN, RUBEN, KEVIN TEE, and MICHAEL MCCOURT (2018). Practical Bayesian optimization in the presence of outliers. *Proceedings of the 21st International Conference on Artificial Intelligence and Statistics (AISTATS 2018)*. Vol. 84. Proceedings of Machine Learning Research, pp. 1722–1731.
- MASSART, PASCAL (2007). *Concentration Inequalities and Model Selection: Ecole d'Eté de Probabilités de Saint-Flour XXXIII – 2003*. Vol. 1896. Lecture Notes in Mathematics. Springer-Verlag.
- MCCULLAGH, P. and J. A. NELDER (1989). *Generalized Linear Models*. 2nd ed. Monographs on Statistics and Applied Probability. Chapman & Hall.
- MEINSHAUSEN, NICOLAI (2006). Quantile Regression Forests. *Journal of Machine Learning Research* 7(35):983–999.
- MIETTINEN, KAISA M. (1998). *Nonlinear Multiobjective Optimization*. International Series in Operations Research & Management Science. Kluwer Academic Publishers.
- MILGROM, PAUL and ILYA SEGAL (2002). Envelope Theorems for Arbitrary Choice Sets. *Econometrica* 70(2):583–601.
- MINKA, THOMAS (2008). EP: A quick reference. URL: https://tminka.github.io/papers/ep_minka-ep-quickref.pdf.
- MINKA, THOMAS P. (2001). A family of algorithms for approximate Bayesian inference. Ph.D. thesis. Massachusetts Institute of Technology.
- MOCKUS, JONAS (1972). Bayesian Methods of Search for an Extremum. *Avtomatika i Vychislitel'naya Tekhnika (Automatic Control and Computer Sciences)* 6(3):53–62.
- MOCKUS, JONAS (1974). On Bayesian Methods for Seeking the Extremum. *Optimization Techniques IFIP Technical Conference*. Vol. 27. Lecture Notes in Computer Science. Springer-Verlag, pp. 400–404.
- MOCKUS, JONAS (1989). *Bayesian Approach to Global Optimization: Theory and Applications*. Mathematics and its Applications. Kluwer Academic Publishers.
- MOCKUS, JONAS, WILLIAM EDDY, AUDRIS MOCKUS, LINAS MOCKUS, and GINTARAS REKLAITAS (2010). *Bayesian Heuristic Approach to Discrete and Global Optimization: Algorithms, Visualization, Software, and Applications*. Nonconvex Optimization and its Applications. Kluwer Academic Publishers.
- MOCKUS, J., V. TIEŠIS, and A. ŽILINSKAS (1978). The Application of Bayesian Methods for Seeking the Extremum. In: *Towards Global Optimization 2*. Ed. by L. C. W. DIXON and G. P. SZEGÖ. North-Holland, pp. 117–129.
- MØLLER, JESPER, ANNE RANDI SYVERSVEEN, and RASMUS PLENGE WAAGEPETERSEN (1998). Log Gaussian Cox Processes. *Scandinavian Journal of Statistics* 25(3):451–482.
- MONTGOMERY, DOUGLAS C. (2019). *Design and Analysis of Experiments*. 10th ed. John Wiley & Sons.
- MOORE, ANDREW W. and CHRISTOPHER G. ATKESON (1993). Memory-based Reinforcement Learning: Efficient Computation with Prioritized Sweeping. *Advances in Neural Information Processing Systems 5 (NeurIPS 1992)*, pp. 263–270.
- MORICONI, RICCARDO, MARC PETER DEISENROTH, and K. S. SESH KUMAR (2020). High-dimensional Bayesian optimization using low-dimensional feature spaces. *Machine Learning* 109(9–10):1925–1943.
- MURRAY, IAIN (2016). Differentiation of the Cholesky decomposition. arXiv: 1602 . 07527 [stat.CO].

BIBLIOGRAPHY

- MURRAY, IAIN, RYAN PRESCOTT ADAMS, and DAVID J. C. MACKAY (2010). Elliptical slice sampling. *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS 2010)*. Vol. 9. Proceedings of Machine Learning Research, pp. 541–548.
- MUTNÝ, MOJMÍR and ANDREAS KRAUSE (2018). Efficient High Dimensional Bayesian Optimization with Additivity and Quadrature Fourier Features. *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*, pp. 9005–9016.
- NEAL, RADFORD M. (1997). *Monte Carlo Implementation of Gaussian Process Models for Bayesian Regression and Classification*. Technical report (9702). Department of Statistics, University of Toronto.
- NEAL, RADFORD M. (1998). Regression and Classification Using Gaussian Process Priors. In: *Bayesian Statistics 6*. Ed. by J. M. BERNARDO, J. O. BERGER, A. P. DAWID, and A. F. M. SMITH. Oxford University Press, pp. 475–490.
- NGUYEN, VU, SUNIL GUPTA, SANTU RANA, CHENG LI, and SVETHA VENKATESH (2017). Regret for Expected Improvement over the Best-Observed Value and Stopping Condition. *Proceedings of the 9th Asian Conference on Machine Learning (ACML 2017)*. Vol. 77. Proceedings of Machine Learning Research, pp. 279–294.
- NICKISCH, HANNES and CARL EDWARD RASMUSSEN (2008). Approximations for Binary Gaussian Process Classification. *Journal of Machine Learning Research* 9(Oct):2035–2078.
- O'HAGAN, A. (1978). Curve Fitting and Optimal Design for Prediction. *Journal of the Royal Statistical Society Series B (Methodological)* 40(1):1–42.
- O'HAGAN, A. (1991). Bayes-Hermite quadrature. *Journal of Statistical Planning and Inference* 29(3):245–260.
- O'HAGAN, ANTHONY and JONATHAN FORSTER (2004). *Kendall's Advanced Theory of Statistics*. 2nd ed. Vol. 2B: Bayesian Inference. Arnold.
- OH, CHANGYONG, JAKUB M. TOMCZAK, EFSTRATIOS GAVVES, and MAX WELLING (2019). Combinatorial Bayesian Optimization using the Graph Cartesian Product. *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, pp. 2914–2924.
- ØKSENDAL, BERNT (2013). *Stochastic Differential Equations: An Introduction with Applications*. 6th ed. Universitext. Springer-Verlag.
- OSBORNE, MICHAEL A., DAVID DUVENAUD, ROMAN GARNETT, CARL E. RASMUSSEN, STEPHEN J. ROBERTS, and ZOUBIN GHAHRAMANI (2012). Active Learning of Model Evidence Using Bayesian Quadrature. *Advances in Neural Information Processing Systems 25 (NeurIPS 2012)*, pp. 46–54.
- OSBORNE, MICHAEL A., ROMAN GARNETT, and STEPHEN J. ROBERTS (2009). Gaussian Processes for Global Optimization. *Proceedings of the 3rd Learning and Intelligent Optimization Conference (LION 3)*.
- PARIA, BISWAJIT, KIRTHEVASAN KANDASAMY, and BARNABÁS PÓCZOS (2019). A Flexible Framework for Multi-Objective Bayesian Optimization using Random Scalarizations. *Proceedings of the 35th Conference on Uncertainty in Artificial Intelligence (UAI 2019)*. Vol. 115. Proceedings of Machine Learning Research, pp. 766–776.
- PEIRCE, C. S. (1876). Note on the Theory of the Economy of Research. In: *Report of the Superintendent of the United States Coast Survey Showing the Progress of the Work for the Fiscal Year Ending with June, 1876*. Government Printing Office, pp. 197–201.
- PICHENY, VICTOR (2014). A Stepwise uncertainty reduction approach to constrained global optimization. *Proceedings of the 17th International Conference on Artificial Intelligence and Statistics (AISTATS 2014)*. Vol. 33. Proceedings of Machine Learning Research, pp. 787–795.

- PICHENY, VICTOR (2015). Multiobjective optimization using Gaussian process emulators via stepwise uncertainty reduction. *Statistics and Computing* 25(6):1265–1280.
- PICHENY, VICTOR, DAVID GINSBOURGER, YANN RICHET, and GREGORY CAPLIN (2013a). Quantile-Based Optimization of Noisy Computer Experiments With Tunable Precision. *Technometrics* 55(1):2–13.
- PICHENY, VICTOR, TOBIAS WAGNER, and DAVID GINSBOURGER (2013b). A benchmark of kriging-based infill criteria for noisy optimization. *Structural and Multidisciplinary Optimization* 48(3):607–626.
- PLEISS, GEOFF, MARTIN JANKOWIAK, DAVID ERIKSSON, ANIL DAMLE, and JACOB R. GARDNER (2020). Fast Matrix Square Roots with Applications to Gaussian Processes and Bayesian Optimization. *Advances in Neural Information Processing Systems 33 (NeurIPS 2020)*, pp. 22268–22281.
- POINCARÉ, HENRI (1912). *Calcul des probabilités*. 2nd ed. Gauthier–Villars.
- PONWEISER, WOLFGANG, TOBIAS WAGNER, DIRK BIERMANN, and MARKUS VINCZE (2008). Multiobjective Optimization on a Limited Budget of Evaluations Using Model-Assisted S -Metric Selection. *Proceedings of the 10th International Conference on Parallel Problem Solving from Nature (PPSN X)*. Vol. 5199. Lecture Notes in Computer Science. Springer–Verlag, pp. 784–794.
- POWELL, WARREN B. (2011). *Approximate Dynamic Programming: Solving the Curses of Dimensionality*. 2nd ed. Wiley Series in Probability and Statistics. John Wiley & Sons.
- RAIFFA, HOWARD and ROBERT SCHLAIFER (1961). *Applied Statistical Decision Theory*. Division of Research, Graduate School of Business Administration, Harvard University.
- RASMUSSEN, CARL EDWARD and ZOUBIN GHAHRAMANI (2002). Bayesian Monte Carlo. *Advances in Neural Information Processing Systems 15 (NeurIPS 2002)*, pp. 505–512.
- RASMUSSEN, CARL EDWARD and CHRISTOPHER K. I. WILLIAMS (2006). *Gaussian Processes for Machine Learning*. Adaptive Computation and Machine Learning. MIT Press.
- ROBBINS, HERBERT (1952). Some Aspects of the Sequential Design of Experiments. *Bulletin of the American Mathematical Society* 58(5):527–535.
- ROLLAND, PAUL, JONATHAN SCARLETT, ILIJA BOGUNOVIC, and VOLKAN CEVHER (2018). High-Dimensional Bayesian Optimization via Additive Models with Overlapping Groups. *Proceedings of the 21st International Conference on Artificial Intelligence and Statistics (AISTATS 2018)*. Vol. 84. Proceedings of Machine Learning Research, pp. 298–307.
- ROSS, ANDREW M. (2010). Computing Bounds on the Expected Maximum of Correlated Normal Variables. *Methodology and Computing in Applied Probability* 12(1):111–138.
- RUDIN, WALTER (1976). *Principles of Mathematical Analysis*. 3rd ed. International Series in Pure and Applied Mathematics. McGraw–Hill.
- RUE, HÅVARD, SARA MARTINO, and NICOLAS CHOPIN (2009). Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society Series B (Methodological)* 71(2):319–392.
- RUSSO, DANIEL and BENJAMIN VAN ROY (2014). Learning to Optimize via Posterior Sampling. *Mathematics of Operations Research* 39(4):1221–1243.
- RUSSO, DANIEL and BENJAMIN VAN ROY (2016). An Information-Theoretic Analysis of Thompson Sampling. *Journal of Machine Learning Research* 17(68):1–30.
- SACKS, JEROME, WILLIAM J. WELCH, TOBY J. MITCHELL, and HENRY P. WYNN (1989). Design and Analysis of Computer Experiments. *Statistical Science* 4(4):409–435.
- ŠALTENIS, VYDŪNAS R. (1971). One Method of Multiextremum Optimization. *Avtomatika i Vychislitel'naya Tekhnika (Automatic Control and Computer Sciences)* 5(3):33–38.

BIBLIOGRAPHY

- SANCHEZ, SUSAN M. and PAUL J. SANCHEZ (2005). Very Large Fractional Factorial and Central Composite Designs. *ACM Transactions on Modeling and Computer Simulation* 15(4): 362–377.
- SCARLETT, JONATHAN (2018). Tight Regret Bounds for Bayesian Optimization in One Dimension. *Proceedings of the 35th International Conference on Machine Learning (ICML 2018)*. Vol. 80. Proceedings of Machine Learning Research, pp. 4500–4508.
- SCARLETT, JONATHAN, ILIJA BOGUNOVIC, and VOLKAN CEVHER (2017). Lower Bounds on Regret for Noisy Gaussian Process Bandit Optimization. *Proceedings of the 2017 Conference on Learning Theory (COLT 2017)*. Vol. 65. Proceedings of Machine Learning Research, pp. 1723–1742.
- SCARLETT, JONATHAN and VOLKAN CEVHAR (2021). An Introductory Guide to Fano’s Inequality with Applications in Statistical Estimation. In: *Information-Theoretic Methods in Data Science*. Ed. by MIGUEL R. D. RODRIGUES and YONINA C. ELDAR. Cambridge University Press.
- SCHONLAU, MATTHIAS, WILLIAM J. WELCH, and DONALD R. JONES (1998). Global versus Local Search in Constrained Optimization of Computer Models. In: *New Developments and Applications in Experimental Design*. Vol. 34. Lecture Notes – Monograph Series. Institute of Mathematical Statistics, pp. 11–25.
- SCHONLAU, MATTHIAS (1997). Computer Experiments and Global Optimization. Ph.D. thesis. University of Waterloo.
- SCOTT, WARREN, PETER FRAZIER, and WARREN POWELL (2011). The Correlated Knowledge Gradient for Simulation Optimization of Continuous Parameters Using Gaussian Process Regression. *SIAM Journal on Optimization* 21(3):996–1026.
- SEEGER, MATTHIAS (2008). *Expectation Propagation for Exponential Families*. Technical report. University of California, Berkeley.
- SETTLES, BURR (2012). *Active Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool.
- SHAH, AMAR and ZOUBIN GHAHRAMANI (2015). Parallel Predictive Entropy Search for Batch Global Optimization of Expensive Objective Functions. *Advances in Neural Information Processing Systems* 28 (*NeurIPS 2015*), pp. 3330–3338.
- SHAH, AMAR, ANDREW GORDON WILSON, and ZOUBIN GHAHRAMANI (2014). Student-*t* Processes as Alternatives to Gaussian Processes. *Proceedings of the 17th International Conference on Artificial Intelligence and Statistics (AISTATS 2014)*. Vol. 33. Proceedings of Machine Learning Research, pp. 877–885.
- SHAHRIARI, BOBAK, KEVIN SWERSKY, ZIYU WANG, RYAN P. ADAMS, and NANDO DE FREITAS (2016). Taking the Human Out of the Loop: A Review of Bayesian Optimization. *Proceedings of the IEEE* 104(1):148–175.
- SHAHRIARI, BOBAK, ZIYU WANG, MATTHEW W. HOFFMAN, ALEXANDRE BOUCHARD-CÔTÉ, and NANDO DE FREITAS (2014). An Entropy Search Portfolio for Bayesian Optimization. arXiv: 1406.4625 [stat.ML].
- SHANNON, C. E. (1948). A Mathematical Theory of Communication. *The Bell System Technical Journal* 27(3):379–423.
- SHAO, T. S., T. C. CHEN, and R. M. FRANK (1964). Tables of Zeros and Gaussian Weights of Certain Associated Laguerre Polynomials and the Related Generalized Hermite Polynomials. *Mathematics of Computation* 18(88):598–616.
- SHEPP, L. A. (1979). The Joint Density of the Maximum and its Location for a Wiener Process with Drift. *Journal of Applied Probability* 16(2):423–427.

- SILVERMAN, B. W. (1986). *Density Estimation for Statistics and Data Analysis*. Monographs on Statistics and Applied Probability. Chapman & Hall.
- SLEPIAN, DAVID (1962). The One-Sided Barrier Problem for Gaussian Noise. *The Bell System Technical Journal* 41(2):463–501.
- SMITH, KIRSTINE (1918). On the Standard Deviations of Adjusted and Interpolated Values of an Observed *Polynomial Function* and its Constants and the Guidance they Give Towards a Proper Choice of the Distribution of Observations. *Biometrika* 12(1–2):1–85.
- SMOLA, ALEX J. and BERNHARD SCHÖLKOPF (2000). Sparse Greedy Matrix Approximation for Machine Learning. *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 911–918.
- SNELSON, EDWARD and ZOUBIN GHAHRAMANI (2005). Sparse Gaussian Processes using Pseudo-inputs. *Advances in Neural Information Processing Systems 18 (NeurIPS 2005)*, pp. 1257–1264.
- SNOEK, JASPER, HUGO LAROCHELLE, and RYAN P. ADAMS (2012). Practical Bayesian Optimization of Machine Learning Algorithms. *Advances in Neural Information Processing Systems 25 (NeurIPS 2012)*, pp. 2951–2959.
- SNOEK, JASPER, OREN RIPPEL, KEVIN SWERSKY, RYAN KIROS, NADATHUR SATISH, NARAYANAN SUNDARAM, et al. (2015). Scalable Bayesian Optimization Using Deep Neural Networks. *Proceedings of the 32nd International Conference on Machine Learning (ICML 2015)*. Vol. 37. Proceedings of Machine Learning Research, pp. 2171–2180.
- SNOEK, JASPER, KEVIN SWERSKY, RICHARD ZEMEL, and RYAN P. ADAMS (2014). Input Warping for Bayesian Optimization of Non-Stationary Functions. *Proceedings of the 31st International Conference on Machine Learning (ICML 2014)*. Vol. 32. Proceedings of Machine Learning Research, pp. 1674–1682.
- SONG, JIALIN, YUXIN CHEN, and YISONG YUE (2019). A General Framework for Multi-fidelity Bayesian Optimization with Gaussian Processes. *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS 2019)*. Vol. 89. Proceedings of Machine Learning Research, pp. 3158–3167.
- SPRINGENBERG, JOST TOBIAS, AARON KLEIN, STEFAN FALKNER, and FRANK HUTTER (2016). Bayesian Optimization with Robust Bayesian Neural Networks. *Advances in Neural Information Processing Systems 29 (NeurIPS 2016)*, pp. 4134–4142.
- SRINIVAS, NIRANJAN, ANDREAS KRAUSE, SHAM KAKADE, and MATTHIAS SEEGER (2010). Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design. *Proceedings of the 27th International Conference on Machine Learning (ICML 2010)*, pp. 1015–1022.
- STEIN, MICHAEL L. (1999). *Interpolation of Spatial Data: Some Theory for Kriging*. Springer Series in Statistics. Springer-Verlag.
- STERLING, TEAGUE and JOHN J. IRWIN (2015). ZINC 15 – Ligand Discovery for Everyone. *Journal of Chemical Information and Modeling* 55(11):2324–2337.
- STRELTSOV, SIMON and PIROOZ VAKILI (1999). A Non-myopic Utility Function for Statistical Global Optimization Algorithms. *Journal of Global Optimization* 14(3):283–298.
- SUTTON, RICHARD S. (1990). Integrated Architectures for Learning, Planning, and Reacting Based on Approximating Dynamic Programming. *Proceedings of the 7th International Conference on Machine Learning (ICML 1990)*, pp. 216–224.
- SVENSON, JOSHUA and THOMAS SANTNER (2016). Multiobjective optimization of expensive-to-evaluate deterministic computer simulator models. *Computational Statistics and Data Analysis* 94:250–264.

BIBLIOGRAPHY

- SWERSKY, KEVIN, JASPER SNOEK, and RYAN P. ADAMS (2013). Multi-Task Bayesian Optimization. *Advances in Neural Information Processing Systems 26 (NeurIPS 2013)*, pp. 2004–2012.
- SWERSKY, KEVIN, JASPER SNOEK, and RYAN P. ADAMS (2014). Freeze–Thaw Bayesian Optimization. arXiv: 1406.3896 [stat.ML].
- TAKENO, SHION, HITOSHI FUKUOKA, YUHKI TSUKADA, TOSHIYUKI KOYAMA, MOTOKI SHIGA, ICHIRO TAKEUCHI, et al. (2020). Multi-fidelity Bayesian Optimization with Max-value Entropy Search and its Parallelization. *Proceedings of the 37th International Conference on Machine Learning (ICML 2020)*. Vol. 119. Proceedings of Machine Learning Research, pp. 9334–9345.
- TALLIS, G. M. (1961). The Moment Generating Function of the Truncated Multi-normal Distribution. *Journal of the Royal Statistical Society Series B (Methodological)* 23(1):223–229.
- TESCH, MATTHEW, JEFF SCHNEIDER, and HOWIE CHOSET (2013). Expensive Function Optimization with Stochastic Binary Outcomes. *Proceedings of the 30th International Conference on Machine Learning (ICML 2013)*. Vol. 28. Proceedings of Machine Learning Research, pp. 1283–1291.
- THOMPSON, WILLIAM R. (1933). On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of Two Samples. *Biometrika* 25(3–4):285–294.
- THOMPSON, WILLIAM R. (1935). On the Theory of Apportionment. *American Journal of Mathematics* 57(2):450–456.
- TITSIAS, MICHALIS (2009). Variational Learning of Inducing Variables in Sparse Gaussian Processes. *Proceedings of the 12th International Conference on Artificial Intelligence and Statistics (AISTATS 2009)*. Vol. 5. Proceedings of Machine Learning Research, pp. 567–574.
- TOSCANO-PALMERIN, SAUL and PETER I. FRAZIER (2018). Bayesian Optimization with Expensive Integrands. arXiv: 1803.08661 [stat.ML].
- TURBAN, SEBASTIEN (2010). *Convolution of a truncated normal and a centered normal variable*. Technical report. Columbia University.
- TURNER, RYAN, DAVID ERIKSSON, MICHAEL MCCOURT, JUHA KIILI, EERO LAAKSONEN, ZHEN XU, et al. (2021). Bayesian Optimization is Superior to Random Search for Machine Learning Hyperparameter Tuning: Analysis of the Black-Box Optimization Challenge 2020. *Proceedings of the NeurIPS 2020 Competition and Demonstration Track*. Vol. 133. Proceedings of Machine Learning Research, pp. 3–26.
- ULRICH, KYLE, DAVID E. CARLSON, KAFUI DZIRASA, and LAWRENCE CARIN (2015). GP Kernels for Cross-Spectrum Analysis. *Advances in Neural Information Processing Systems 28 (NeurIPS 2015)*, pp. 1999–2007.
- VAKILI, SATTAR, KIA KHEZELI, and VICTOR PICHENY (2021). On Information Gain and Regret Bounds in Gaussian Process Bandits. *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS 2021)*. Vol. 130. Proceedings of Machine Learning Research, pp. 82–90.
- VAKILI, SATTAR, VICTOR PICHENY, and NICOLAS DURRANDE (2020). Regret Bounds for Noise-Free Bayesian Optimization. arXiv: 2002.05096 [stat.ML].
- VALKO, MICHAL, NATHAN KORDA, RÉMI MUNOS, ILIAS FLAOUNAS, and NELLO CRISTIANINI (2013). Finite-Time Analysis of Kernelised Contextual Bandits. *Proceedings of the 29th Conference on Uncertainty in Artificial Intelligence (UAI 2013)*, pp. 654–663.
- VAN DE CORPUT, J. G. (1935). Verteilungsfunktionen: Erste Mitteilung. *Proceedings of the Koninklijke Nederlandse Akademie van Wetenschappen* 38:813–821.

- VAN DER VAART, AAD W. and JON A. WELLNER (1996). *Weak Convergence and Empirical Processes with Applications to Statistics*. Springer Series in Statistics. Springer-Verlag.
- VANCHINATHAN, HASTAGIRI P., ANDREAS MARFURT, CHARLES-ANTOINE ROBELIN, DONALD KOSSMANN, and ANDREAS KRAUSE (2015). Discovering Valuable Items from Massive Data. *Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD 2015)*, pp. 1195–1204.
- VAZQUEZ, EMMANUEL, JULIEN VILLEMONTEIX, MARYAN SIDORKIEWICZ, and ÉRIC WALTER (2008). Global optimization based on noisy evaluations: an empirical study of two statistical approaches. *Proceedings of the 6th International Conference on Inverse Problems in Engineering: Theory and Practice (ICIPE 2008)*. Vol. 135. Journal of Physics: Conference Series, paper number 012100.
- VEHTARI, AKI and JANNE OJANEN (2012). A survey of Bayesian predictive methods for model assessment, selection and comparison. *Statistics Surveys* 6:142–228.
- VILLEMONTEIX, JULIEN, EMMANUEL VAZQUEZ, and ERIC WALTER (2009). An informational approach to the global optimization of expensive-to-evaluate functions. *Journal of Global Optimization* 44(4):509–534.
- VIVARELLI, FRANCESCO and CHRISTOPHER K. I WILLIAMS (1998). Discovering hidden features with Gaussian process regression. *Advances in Neural Information Processing Systems 11 (NeurIPS 1998)*, pp. 613–619.
- VON NEUMANN, JOHN and OSKAR MORGENTERN (1944). *Theory of Games and Economic Behavior*. Princeton University Press.
- VONDRAK, JAN (2005). Probabilistic Methods in Combinatorial and Stochastic Optimization. Ph.D. thesis. Massachusetts Institute of Technology.
- WALD, A. (1945). Sequential Tests of Statistical Hypotheses. *The Annals of Mathematical Statistics* 16(2):117–186.
- WALD, ABRAHAM (1947). *Sequential Analysis*. Wiley Mathematical Statistics Series. John Wiley & Sons.
- WANG, JIALEI, SCOTT C. CLARK, ERIC LIU, and PETER I. FRAZIER (2020a). Parallel Bayesian Global Optimization of Expensive Functions. *Operations Research* 68(6):1850–1865.
- WANG, ZEXIN, VINCENT Y. F. TAN, and JONATHAN SCARLETT (2020b). Tight Regret Bounds for Noisy Optimization of a Brownian Motion. arXiv: 2001.09327 [cs.LG].
- WANG, ZI and STEFANIE JEGELKA (2017). Max-value Entropy Search for Efficient Bayesian Optimization. *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*. Vol. 70. Proceedings of Machine Learning Research, pp. 3627–3635.
- WANG, ZI, BOLEI ZHOU, and STEFANIE JEGELKA (2016a). Optimization as Estimation with Gaussian Processes in Bandit Settings. *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS 2016)*. Vol. 51. Proceedings of Machine Learning Research, pp. 1022–1031.
- WANG, ZIYU and NANDO DE FREITAS (2014). Theoretical Analysis of Bayesian Optimization with Unknown Gaussian Process Hyper-Parameters. arXiv: 1406.7758 [stat.ML].
- WANG, ZIYU, FRANK HUTTER, MASROUR ZOGHI, DAVID MATHESON, and NANDO DE FREITAS (2016b). Bayesian Optimization in a Billion Dimensions via Random Embeddings. *Journal of Artificial Intelligence Research* 55:361–387.
- WENDLAND, HOLGER (2004). *Scattered Data Approximation*. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press.
- WHITTLE, PETER (1982). *Optimization Over Time: Dynamic Programming and Stochastic Control*. Vol. 1. Wiley Series in Probability and Mathematical Statistics. John Wiley & Sons.

BIBLIOGRAPHY

- WHITTLE, P. (1988). Restless Bandits: Activity Allocation in a Changing World. *Journal of Applied Probability* 25(A):287–298.
- WILLIAMS, CHRISTOPHER K. I. and MATTHIAS SEEGER (2000). Using the Nyström Method to Speed Up Kernel Machines. *Advances in Neural Information Processing Systems 13 (neurIPS 2000)*, pp. 682–688.
- WILSON, ANDREW GORDON and RYAN PRESCOTT ADAMS (2013). Gaussian Process Kernels for Pattern Discovery and Extrapolation. *Proceedings of the 30th International Conference on Machine Learning (ICML 2013)*. Vol. 28. Proceedings of Machine Learning Research, pp. 1067–1075.
- WILSON, ANDREW GORDON, ZHITING HU, RUSLAN SALAKHUTDINOV, and ERIC P. XING (2016). Deep Kernel Learning. *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS 2016)*. Vol. 51. Proceedings of Machine Learning Research, pp. 370–378.
- WU, JIAN and PETER I. FRAZIER (2016). The Parallel Knowledge Gradient Method for Batch Bayesian Optimization. *Advances in Neural Information Processing Systems 29 (NeurIPS 2016)*, pp. 3126–3134.
- WU, JIAN and PETER I. FRAZIER (2019). Practical Two-Step Look-Ahead Bayesian Optimization. *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, pp. 9813–9823.
- WU, JIAN, MATTHIAS POLOCZEK, ANDREW GORDON WILSON, and PETER I. FRAZIER (2017). Bayesian Optimization with Gradients. *Advances in Neural Information Processing Systems 30 (NeurIPS 2017)*, pp. 5267–5278.
- WU, JIAN, SAUL TOSCANO-PALMERIN, PETER I. FRAZIER, and ANDREW GORDON WILSON (2019). Practical Multi-fidelity Bayesian Optimization for Hyperparameter Tuning. *Proceedings of the 35th Conference on Uncertainty in Artificial Intelligence (UAI 2019)*. Vol. 115. Proceedings of Machine Learning Research, pp. 788–798.
- YANG, KAIFENG, MICHAEL EMMERICH, ANDRÉ DEUTZ, and THOMAS BÄCK (2019a). Efficient computation of expected hypervolume improvement using box decomposition algorithms. *Journal of Global Optimization* 75(1):3–34.
- YANG, KAIFENG, MICHAEL EMMERICH, ANDRÉ DEUTZ, and THOMAS BÄCK (2019b). Multi-Objective Bayesian Global Optimization using expected hypervolume improvement gradient. *Swarm and Evolutionary Computation* 44:945–956.
- YUE, XUBO and RAED AL KONTAR (2020). Why Non-myopic Bayesian Optimization is Promising and How Far Should We Look-ahead? A Study via Rollout. *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS 2020)*. Vol. 108. Proceedings of Machine Learning Research, pp. 2808–2818.
- ZHANG, YEHONG, TRONG NGHIA HOANG, BRYAN KIAN HSIANG LOW, and MOHAN KANKAN-HALLI (2017). Information-Based Multi-Fidelity Bayesian Optimization. *Bayesian Optimization for Science and Engineering Workshop (BayesOpt 2017), Conference on Neural Information Processing Systems (NeurIPS 2017)*.
- ZILBERSTEIN, SCHLOMO (1996). Using Anytime Algorithms in Intelligent Systems. *AI Magazine* 17(3):73–83.
- ŽILINSKAS, ANTANAS G. (1975). Single-Step Bayesian Search Method for an Extremum of Functions of a Single Variable. *Kibernetika (Cybernetics)* 11(1):160–166.
- ZITZLER, ECKART (1999). Evolutionary Algorithms for Multiobjective Optimization: Methods and Applications. PhD thesis. Eidgenössische Technische Hochschule Zürich.
- ZULUAGA, MARCELA, ANDREAS KRAUSE, and MARKUS PÜSCHEL (2016). ε -PAL: An Active Learning Approach to the Multi-Objective Optimization Problem. *Journal of Machine Learning Research* 17(104):1–32.

INDEX

Primary references and definitions are indicated in bold.

A

- a posteriori methods, 268, 270, 272
- a priori methods, 268, 271
- acquisition function, $\alpha(\mathbf{x}; \mathcal{D})$, 11, **88**, 94, 96, 98, 150, 245, *see also* expected improvement; knowledge gradient; mutual information; probability of improvement; upper confidence bound
- batch, $\beta(\mathbf{x}; \mathcal{D})$, 250
- gradient, 158, 208
- optimization, 207
- action space, \mathcal{A} , **90**, 243
 - for batch observations, 250
 - for dynamic termination, 104
 - for multifidelity optimization, 261
 - for optimization with fixed budget, 91
 - for sequential procedures, 277
 - for terminal recommendation, 110, 248
- active learning, 136, 273
- active search, 115, **278**
- adaptivity gap, 252
- additive decomposition, 63, 83
- aerospace engineering, applications in, 320
- algorithm configuration, 244, 321
- anytime algorithm, 208, 246
- approximate dynamic programming, **100**, 150, 280
- augmented Chebyshev scalarization, 272
- automated machine learning (automl), 261, 323, *see also* hyperparameter tuning
- automatic relevance determination (ARD), 57, 62, 239
- automobile engineering, applications in, 321

B

- bandits, *see* multi-armed bandits
- batch observations, 250, 260
 - connection to sequential observations, 252
- batch rollout, **103**, 152, 154, 280, *see also* rollout
- Bayes' theorem, 7
- Bayesian decision theory, 10, **89**, 124
 - isolated decisions, 90
 - sequential decisions
 - dynamic termination, 103

fixed budget, 91

multi-armed bandits, 143

Bayesian inference

introduction to, 6

of objective function, 8

Bayesian information criterion, 79, 83

Bayesian neural networks, 198, 277

Bayesian Occam's razor, 71

Bayesian quadrature, 33, 276

Bayesian regret, $\mathbb{E}[r_\tau]$, $\mathbb{E}[R_\tau]$, 218, 224, 236

Bellman equation, *see* Bellman optimality

Bellman optimality, 94, 99

beta warping, 58

BINOCULARS algorithm, 153

biology, applications in, 314

biomedical engineering, applications in, 316

Bochner's theorem, 50

branch and bound, 231, 237

C

central composite design, 76, 81

certainty equivalent, 111

characteristic length scale, *see* length scale

chemistry, applications in, 284, 309

chemoinformatics, 310

Cholesky decomposition, 201, 256

low-rank updates, 202

civil engineering, applications in, 318

combinatorial optimization, 209

compactness of domain, 34

conditional entropy, $H[\omega | \mathcal{D}]$, 115, *see also* entropy

conditional mutual information, $I(\omega; \psi | \mathcal{D})$, 137, *see also* mutual information

confidence ellipsoid, 225, 232, 235, 238

conformational search, 311

conjugate gradients, 203

constant liar heuristic, 254

constrained optimization, 247

constraint functions, 247

unknown, 247

constraints, 247

constraints on objective function, 36, 39, 56

continuity in mean square, 29, *see also* sample path continuity

continuous differentiability, 31, 221

cost-aware optimization, 103, 243, 251, 263

- covariance function, *see* prior covariance function
- cross-covariance function, 19, 23, 24, 27, 30, 201, 262
- cumulative regret, R_t , 145, 214
- cumulative reward, 114, 142, 155, 215, 279
- curse of dimensionality, 61, 208
- D**
- de novo* design, 310
- decoupled constraint observations, 250
- deep kernel, 59, 61
- deep neural networks, vii, 1, 59, 61, 288
- design and analysis of computer experiments (DACE), 286
- determinantal point process, 259, 281
- differentiability in mean square, 30, *see also* continuous differentiability
- differential entropy, $H[\omega]$, *see* entropy
- dilation, 56
- disjoint union, 27
- drug discovery, *see* molecular design
- dynamic termination, *see* termination decisions
- E**
- early stopping, 210, 276
- electrical engineering, applications in, 319
- elliptical slice sampling, 38
- embedding, *see* linear embedding, neural embedding
- entropy search, *see* mutual information
- entropy, $H[\omega]$, 115, *see also* conditional entropy
- environmental variables, 275
- expectation propagation, 39, 182, 190, 271, 298
- expected gain per unit cost, 246
- expected hypervolume improvement (EHVI), 269
- expected improvement, α_{EI} , 81, 95, 113, 117, 127, 151, 158, 193, 196, 197, 199, 263, 264, 266, *see also* simple reward augmented, 166
- batch, 257
- comparison with probability of improvement, 132
- computation with noise, 160
- alternative formulations, 163
- gradient, 304
- computation without noise, 159
- gradient, 159
- convergence, 217
- modified, 154
- origin, 285
- worst-case regret with noise, 238, 241
- expected utility, 90, 93
- exploration bonus, 145
- exploration vs. exploitation dilemma, 11, 83, 123, 128, 131, 133, 143, 145, 146, 148, 154, 159, 214, 289
- exponential covariance function, $K_{M^{1/2}}$, 52, 219
- extreme value theorem, 34
- F**
- factor graph, 37
- Fano's inequality, 230
- feasible region, \mathcal{F} , 247
- figure of merit, *see* acquisition function
- fill distance, δ_X , 235
- Fourier transform, 50, 53, 234
- freeze–thaw Bayesian optimization, 277
- fully independent training conditional (FITC) approximation, 207
- G**
- Gauss–Hermite quadrature, 172, 256
- gradient, 305
- Gaussian process (GP), $\mathcal{GP}(f; \mu, K)$, 8, 16, 95, 124, *see also* prior mean function; prior covariance function
- approximate inference, 35, *see also* sparse spectrum approximation; sparse approximation
- classification, 41, 278
- computation of policies with, 157
- continuity, 28
- credible intervals, 18
- differentiability, 30
- exact inference, 19
- additive Gaussian noise, 20, 23
- computation, 201
- derivative observations, 32
- exact observation, 22
- interpretation of posterior moments, 21
- joint, *see* joint Gaussian process
- marginal likelihood $p(\mathbf{y} | \mathbf{x}, \theta)$, 72, 202, 220
- gradient, 303
- maxima, existence and uniqueness, 33
- mixture, 38, 75, 193
- model assessment, selection, and averaging, 67
- modeling, 45

- posterior predictive distribution, 25, 157, 208
 gradient, 303
 sampling, 18
- gene design, 315
 generalized additive models, 63
 generalized linear models, 41
 GLASSES algorithm, 152
 global reward, 113, 117, 129, 172, *see also*
 knowledge gradient
 GP-SELECT algorithm, 281
 Gram matrix, $K(\mathbf{x}, \mathbf{x})$, 17, 49
 grid search, 3, 236
- H**
 Hamiltonian Monte Carlo (HMC), 75
 Heine–Borel theorem, 30
 heteroskedastic noise, 4, 23, 166
 Hölder continuity, 35, 221, 236
 horizon, 93, 125, 151, 243, 252
 hubris of youth, vii
 human–computer interfaces, 324
 hyperband, 277
 hyperparameter tuning, 1, 61, 109, 261, 265, 287, 322
 hyperparameters, θ , 68, *see also* length scale; output scale
 unknown, effect of convergence, 239
- I**
 ill-conditioning, 203
 incumbent value, ϕ^* , 128, 159, 249
 inducing values, v , 205
 infill function, *see* acquisition function
 information capacity, γ_τ , 221, 225, 229, 232, 235
 bounds, 222
 information gain, 115, 117, 135, 180, 187, *see also* mutual information
 initialization, 210
 instantaneous regret, ρ_τ , 214, 237
 integrated expected conditional improvement, 249
 intrinsic dimensionality, 61
 inverse problems, 313, 315
 isotropy, 50, 54
 iterative numerical methods, 203
 Iverson bracket, xi
- J**
 joint Gaussian process, 27, 265, *see also*
 cross-covariance function
 between function and gradient, 30
- exact inference, 28
 for multifidelity modeling, 262
 marginals, 27
- K**
 kernel, *see* prior covariance function
 kernel kernel, 83
 knowledge gradient, α_{KG} , 113, 129, 172, 193, 264, 274, *see also* global reward
 batch, 257
 computation, 172
 discrete domain, 173
 gradient, 306
 KGCP approximation, 175
 origin, 286
 Kolmogorov extension theorem, 16
 kriging believer heuristic, 254
 Kullback–Leibler divergence, $D_{KL}[p \parallel q]$, 115, 137, 206
- L**
 Laplace approximation, 39, 41, 76, 79, 83, 297
 learning curve, 277
 length scale, 54, 56, 69, *see also* automatic relevance determination
 likelihood, 7, *see also* observation model
 marginal, *see* marginal likelihood
 limited lookahead, *see* lookahead
 linear covariance function, K_{LIN} , 54
 linear embedding, 58, 62, 209
 linear scalarization, 272
 linear transformations
 of domain, 55, 56, 62
 of Gaussian processes, 33, *see also*
 Bayesian quadrature; joint
 Gaussian process: between
 function and gradient
 Lipschitz continuity, 226, 236, 255
 local optimum, conditioning a Gaussian process on a, 36, 182
 local penalization, 255
 lookahead, 101, 125, 150, 279, *see also* one-step lookahead; two-step lookahead
 low-dimensional structure, 58, 61, 62, 209, 290
 low-discrepancy sequence, 177, 210
 lower regret bounds
 Bayesian regret with noise, 229
 Bayesian regret without noise, 236
 worst-case regret with noise, 234
 worst-case regret without noise, 237
 Löwner order, 22, 48

- M**
- Mahalanobis distance, 48, 57, 293
 - manifold Gaussian process, 59, 61
 - marginal likelihood, $p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})$, 71
 - materials science, applications in, 261, 309
 - Matérn covariance function, K_M , 51, 124, 221, 236, 239
 - matrix calculus convention, xi
 - maximum a posteriori inference, *see* model selection
 - maximum information gain, *see* information capacity
 - max-value entropy search, α_{MES} , 187
 - for multiobjective optimization, 271
 - gradient, 191
 - mean function, *see* prior mean function
 - mechanical engineering, applications in, 320
 - model assessment, 45, 67, 70, *see also* model posterior
 - model averaging, 74, 116
 - in acquisition function, 192
 - multiple model structures, 79
 - model evidence, *see* marginal likelihood
 - model posterior, $p(\boldsymbol{\theta} \mid \mathcal{D})$, 71
 - model prior, $p(\boldsymbol{\theta})$, 70
 - model selection, 73, 324
 - multiple model structures, 79
 - model space, *see* model structure
 - model structure, \mathcal{M} , 68, *see also*
 - hyperparameters
 - posterior, $\Pr(\mathcal{M} \mid \mathcal{D})$, 79
 - prior, $\Pr(\mathcal{M})$, 78
 - search, 81
 - model, $p(\mathbf{y} \mid \mathbf{x})$, 68
 - molecular design, 62, 209, 278, 309
 - molecular fingerprint, 310
 - Monte Carlo sampling, 37, 75, 84, 181, 187, 256, 271
 - multi-armed bandits, 141
 - infinite-armed, 144
 - optimal policy, 143
 - multifidelity optimization, 261, 317
 - multiobjective optimization, 267
 - multitask optimization, 264
 - multifidelity optimization, 26
 - multiobjective optimization, 26
 - multit-armed bandits
 - origin, 288
 - mutual information, $I(\omega; \psi)$, 116, 135, 287, *see also*
 - information gain; conditional mutual information
- with f^* , α_{f^*} , 140, 187, 193, 264, 287, *see also*
 - output-space predictive entropy search; max-value entropy search
 - with x^* , α_{x^*} , 139, 180, 193, 264, 287, *see also*
 - predictive entropy search
- myopic approximation, *see* lookahead
- N**
- nats, xi
 - needle in a haystack analogy, 216, 234, 236, 238
 - neural architecture search, 323
 - neural embedding, 59, 61, 209
 - no U-turn sampler (NUTS), 76
 - nonmyopic policies, 150, 279, 290
 - no-regret property, 145, 215
 - Nyström method, 207
- O**
- objective function posterior, $p(f \mid \mathcal{D})$, 9, 74, 92, *see also* Gaussian process: exact inference, approximate inference
 - model-marginal, *see* model averaging
 - objective function prior, $p(f)$, 8, *see also* Gaussian process
 - observation costs, 104
 - unknown, 243
 - observation model, $p(y \mid x, \phi)$, 4
 - additive Gaussian noise, 4, 23, 69, 78, 157
 - additive Student- t noise, 36, 38, 41, 278
 - exact observation, 4, 22
 - for unknown costs, 244
 - observation noise scale, σ_n , 4, 23, 69, 78, 157, 203, 222
 - one-step lookahead, 94, 102, 126, 171, 243, 245, 250, 279, 285
 - cost-aware optimization, 106
 - with cumulative reward, 155
 - with global reward, *see* knowledge gradient
 - with information gain, *see* mutual information
 - with simple reward, *see* expected improvement
 - one-step lookehead, 249
 - optimal design, 283
 - optimal policy
 - batch observations, 251
 - computational cost, 99, 125, 279
 - generic, 243
 - multi-armed bandits, 143
 - sequential optimization, 98

- optimism in the face of uncertainty, 145
 optimization policy, 3, *see also* acquisition function; grid search; optimal policy; random search; Thompson sampling
 optimal, *see* optimal policy
 Ornstein–Uhlenbeck (ou) process, 52, 174, 286
 output scale, 55, 69, 239
 output-space predictive entropy search, α_{OPES} , 187
 gradient, 307
- P**
- PareGO algorithm, 273
 Pareto dominance, 268
 Pareto frontier, 153, 267
 Pareto optimality, *see* Pareto frontier
 Parzen estimation, 196
 periodic covariance function, 35, 58, 60
 physics, applications in, 313
 plant breeding, 1, 315
 posterior distribution, 7
 posterior predictive distribution, $p(y | x, \mathcal{D})$, 8, 25, 74, 92, 157, 208, 279
 for multifidelity modeling, 262
 for unknown costs, 244
 model-marginal, *see* model averaging
 predictive entropy search, α_{PES} , 180
 batch, 258
 for multiobjective optimization, 271
 gradient, 307
 preference optimization, 278, 325
 prior covariance function, $K(x, x')$, 17, 48, 67,
 see also exponential, linear, Matérn, spectral mixture, squared exponential covariance functions; automatic relevance determination
 addition, 59, 63
 multiplication, 59
 scaling, 55
 warping, 56
 prior distribution, 6
 prior mean function, $\mu(x)$, 16, 46, 67
 concave quadratic, 48
 constant, 47, 69, 124, 207
 marginalization of parameter, 47, 69
 impact on posterior mean, 46
 linear combination of bases, 48
 probability of improvement, α_{PI} , 131, 167, 193, 196, 199
 batch, 258
- comparison with expected improvement, 132
 computation with noise, 169
 gradient, 305
 computation without noise, 167
 gradient, 169
 convergence, 217
 correspondence with upper confidence bound, 170
 origin, 285
 selection of improvement target, 133, 285
 protein design, 315
 pseudopoints, *see* inducing values
- Q**
- quantile function, $q(\pi)$, 145, 165, 170
- R**
- random embedding, 63
 random forests, 196
 random search, 3
 reaction optimization, 284, 311
 regret, 213, *see also* simple regret; cumulative regret; Bayesian regret; worst-case regret
 representer points, 180, 188
 reproducing kernel Hilbert space (RKHS), \mathcal{H}_K , 219, 224, 231, 240
 RKHS ball, $\mathcal{H}_K[B]$, 220, 224, 231
 RKHS norm, $\|f\|_{\mathcal{H}_K}$, 220
 risk neutrality, 111, 127, 248
 risk tolerance, 111
 risk vs. reward tradeoff, 112, 267
 robotics, applications in, 276, 278, 316
 robust optimization, 317
 rolling horizon, 101
 rollout, 102, 151, 261
 batch, *see* batch rollout
- S**
- \mathcal{S} metric, 269
 safe optimization, 317
 sample path continuity, 30, 34, 218, 221
 scalarization, 271
 second-derivative test, 36, 40
 separable covariance function, 262
 sequential analysis, 283
 sequential experimental design, 284
 sequential simulation, 253
 signal variance, *see* output scale
 simple regret, r_τ , 214, 215, 216, 229, 235

- simple reward, 95, 112, 117, 127, 158, 165, 249,
 280, *see also* expected
 improvement
- small data, 68
- sparse approximation, 204
- sparse spectrum approximation, 51, 178
- spectral density, κ , 51, 53
- spectral measure, ν , 51, 53
- spectral mixture covariance function, K_{SM} , 51,
 53
- spectral points, 178
- squared exponential covariance function, K_{SE} ,
 17, 52, 221
- stationarity, 50, 56, 58, 178, 207, 234
- stochastic process, 8, 16, *see also* Gaussian
 process
- structural search, 312
- Student- t process, 278
- sub-Gaussian distribution, 232
- T**
- terminal recommendation, 118
- terminal recommendations, 90, 109
- termination decisions, 5, 251
 optimal, 103
 practical, 211
- termination option, \emptyset , 104
- Thompson sampling, 148, 176, 181, 187, 195, 257
 acquisition function view, 148
 batch, 259
 computation, 176
 origin, 288
 regret bounds
 Bayesian regret with noise, 228
 worst-case regret with noise, 232
 worst-case regret without noise, 239
- truncated normal distribution, $\mathcal{T}\mathcal{N}(\phi; \mu, \sigma^2 I)$,
 40, 159, 189, 301
- two-step lookahead, 96, 150, 249
- U**
- upper confidence bound, α_{UCB} , 145, 170, 195,
 264
- batch, 259
- computation, 170
- correspondence with probability of
 improvement, 170
- gradient, 170
- origin, 285
- regret bounds
 Bayesian regret with noise, 225
 worst-case regret with noise, 232, 241
 worst-case regret without noise, 238
 selecting confidence parameter, 147
- utility function
 for active search, 279
 for constrained optimization, 248
 for cost-aware optimization, 104
 for isolated decisions, $u(a, \psi, \mathcal{D})$, 90
 for multifidelity optimization, 262
 for multitask optimization, 266
 for optimization, $u(\mathcal{D})$, 93, 109, 243, *see also*
 cumulative reward; global
 reward; simple reward;
 information gain
 for terminal recommendations, $v(\phi)$, 111
- V**
- value of data, α_t^* , 95, 101
- value of sample information, 126
- variational inference, 39, 206
- virtual screening, 310
- von Neumann–Morgenstern theorem, 90, 120
- W**
- weighted Euclidean distance, 57, *see also*
 automatic relevance determination
- Wiener process, 174, 217, 231, 239, 286
- wiggliness, 56, 198, *see also* length scale
- worst-case regret, \bar{r}_t, \bar{R}_t , 218, 231, 237