

Probabilistic Artificial Intelligence

Andreas Krause, Jonas Hübotter

ETH zürich

Institute for Machine Learning
Department of Computer Science

Compiled on **February 11, 2025**.

This manuscript is based on the course PROBABILISTIC ARTIFICIAL INTELLIGENCE (263-5210-00L) at ETH Zürich.



This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

© 2025 ETH Zürich. All rights reserved.

Preface

Artificial intelligence commonly refers to the science and engineering of artificial systems that can carry out tasks generally associated with requiring aspects of human intelligence, such as playing games, translating languages, and driving cars. In recent years, there have been exciting advances in learning-based, data-driven approaches towards AI, and machine learning and deep learning have enabled computer systems to perceive the world in unprecedented ways. Reinforcement learning has enabled breakthroughs in complex games such as Go and challenging robotics tasks such as quadrupedal locomotion.

A key aspect of intelligence is to not only make predictions, but reason about the *uncertainty* in these predictions, and to consider this uncertainty when making decisions. This is what “Probabilistic Artificial Intelligence” is about. The first part covers probabilistic approaches to machine learning. We discuss the differentiation between “epistemic” uncertainty due to lack of data and “aleatoric” uncertainty, which is irreducible and stems, e.g., from noisy observations and outcomes. We discuss concrete approaches towards probabilistic inference, such as Bayesian linear regression, Gaussian process models and Bayesian neural networks. Often, inference and making predictions with such models is intractable, and we discuss modern approaches to efficient approximate inference.

The second part of the manuscript is about taking uncertainty into account in sequential decision tasks. We consider active learning and Bayesian optimization — approaches that collect data by proposing experiments that are informative for reducing the epistemic uncertainty. We then consider reinforcement learning, a rich formalism for modeling agents that learn to act in uncertain environments. After covering the basic formalism of Markov Decision Processes, we consider modern deep RL approaches that use neural network function approximation. We close by discussing modern approaches in model-based RL, which harness epistemic and aleatoric uncertainty to guide exploration, while also reasoning about safety.

Guide to the Reader

The material covered in this manuscript may support a one semester graduate introduction to probabilistic machine learning and sequential decision-making. We welcome readers from all backgrounds. However, we assume familiarity with basic concepts in probability, calculus, linear algebra, and machine learning (e.g., neural networks) as covered in a typical introductory course to machine learning. In Chapter 1, we give a gentle introduction to probabilistic inference, which serves as the foundation for the rest of the manuscript. As part of this first chapter, we also review key concepts from probability theory. We provide a chapter reviewing key concepts of further mathematical background in the back of the manuscript.

Throughout the manuscript, we focus on key concepts and ideas rather than their historical development. We encourage you to consult the provided references for further reading and historical context to delve deeper into the covered topics.

Finally, we have included a set of exercises at the end of each chapter. When we highlight an exercise throughout the text, we use this question mark: (?) — so don't be surprised when you stumble upon it. You will find solutions to all exercises in the back of the manuscript.

Problem 1.1

We hope you will find this resource useful.

Contributing

We encourage you to raise issues and suggest fixes for anything you think can be improved. We are thankful for any such feedback!

CONTACT: pai-script@lists.inf.ethz.ch

Acknowledgements

We are grateful to Sebastian Curi for creating the original Jupyter notebooks that accompany the course at ETH Zürich and which were instrumental in the creation of many figures. We thank Hado van Hasselt for kindly contributing Figure 12.1, and thank Tuomas Haarnoja (Haarnoja et al., 2018a) and Roberto Calandra (Chua et al., 2018) for kindly agreeing to have their figures included in this manuscript. Furthermore, many of the exercises in these notes are adapted from iterations of the course at ETH Zürich. Special thanks to all instructors that contributed to the course material over the years. We also thank all students of the course in the Fall of 2022, 2023, and 2024 who provided valuable feedback on various iterations of this manuscript and corrected many

mistakes. Finally, we thank Zhiyuan Hu, Shyam Sundhar Ramesh, Leander Diaz-Bone, Nicolas Menet, and Ido Hakimi for proofreading parts of various drafts of this text.

Contents

1	<i>Fundamentals of Inference</i>	1
1.1	<i>Probability</i>	2
1.2	<i>Probabilistic Inference</i>	15
1.3	<i>Supervised Learning and Point Estimates</i>	22
1.4	<i>Outlook: Decision Theory</i>	29
<i>I Probabilistic Machine Learning</i>		35
2	<i>Linear Regression</i>	39
2.1	<i>Weight-space View</i>	40
2.2	<i>Aleatoric and Epistemic Uncertainty</i>	44
2.3	<i>Non-linear Regression</i>	45
2.4	<i>Function-space View</i>	45
3	<i>Filtering</i>	51
3.1	<i>Conditioning and Prediction</i>	53
3.2	<i>Kalman Filters</i>	54
4	<i>Gaussian Processes</i>	59
4.1	<i>Learning and Inference</i>	60
4.2	<i>Sampling</i>	61
4.3	<i>Kernel Functions</i>	62
4.4	<i>Model Selection</i>	67
4.5	<i>Approximations</i>	70

5	<i>Variational Inference</i>	83
	5.1 <i>Laplace Approximation</i>	83
	5.2 <i>Predictions with a Variational Posterior</i>	87
	5.3 <i>Blueprint of Variational Inference</i>	88
	5.4 <i>Information Theoretic Aspects of Uncertainty</i>	89
	5.5 <i>Evidence Lower Bound</i>	100
6	<i>Markov Chain Monte Carlo Methods</i>	113
	6.1 <i>Markov Chains</i>	114
	6.2 <i>Elementary Sampling Methods</i>	121
	6.3 <i>Sampling using Gradients</i>	124
7	<i>Deep Learning</i>	139
	7.1 <i>Artificial Neural Networks</i>	139
	7.2 <i>Bayesian Neural Networks</i>	142
	7.3 <i>Approximate Probabilistic Inference</i>	144
	7.4 <i>Calibration</i>	152
	<i>II Sequential Decision-Making</i>	157
8	<i>Active Learning</i>	161
	8.1 <i>Conditional Entropy</i>	161
	8.2 <i>Mutual Information</i>	163
	8.3 <i>Submodularity of Mutual Information</i>	166
	8.4 <i>Maximizing Mutual Information</i>	168
	8.5 <i>Learning Locally: Transductive Active Learning</i>	172
9	<i>Bayesian Optimization</i>	177
	9.1 <i>Exploration-Exploitation Dilemma</i>	177
	9.2 <i>Online Learning and Bandits</i>	178
	9.3 <i>Acquisition Functions</i>	180
10	<i>Markov Decision Processes</i>	197
	10.1 <i>Bellman Expectation Equation</i>	199
	10.2 <i>Policy Evaluation</i>	201

10.3 Policy Optimization	203
10.4 Partial Observability	209
11 Tabular Reinforcement Learning	217
11.1 The Reinforcement Learning Problem	217
11.2 Model-based Approaches	219
11.3 Balancing Exploration and Exploitation	220
11.4 Model-free Approaches	224
12 Model-free Reinforcement Learning	233
12.1 Tabular Reinforcement Learning as Optimization	233
12.2 Value Function Approximation	235
12.3 Policy Approximation	238
12.4 On-policy Actor-Critics	244
12.5 Off-policy Actor-Critics	251
12.6 Maximum Entropy Reinforcement Learning	256
12.7 Learning from Preferences	260
13 Model-based Reinforcement Learning	273
13.1 Planning	274
13.2 Learning	281
13.3 Exploration	287
A Mathematical Background	299
A.1 Probability	299
A.2 Quadratic Forms and Gaussians	301
A.3 Parameter Estimation	302
A.4 Optimization	313
A.5 Useful Matrix Identities and Inequalities	319
B Solutions	321
<i>Bibliography</i>	385
<i>Summary of Notation</i>	393

Acronyms 399

Index 401

1

Fundamentals of Inference

Boolean logic is the algebra of statements which are either true or false. Consider, for example, the statements

“If it is raining, the ground is wet.” and “It is raining.”

A quite remarkable property of Boolean logic is that we can combine these premises to draw logical inferences which are *new* (true) statements. In the above example, we can conclude that the ground must be wet. This is an example of logical reasoning which is commonly referred to as *logical inference*, and the study of artificial systems that are able to perform logical inference is known as *symbolic artificial intelligence*.

But is it really raining? Perhaps it is hard to tell by looking out of the window. Or we have seen it rain earlier, but some time has passed since we have last looked out of the window. And is it really true that if it rains, the ground is wet? Perhaps the rain is just light enough that it is absorbed quickly, and therefore the ground still appears dry.

This goes to show that in our experience, the real world is rarely black and white. We are frequently (if not usually) uncertain about the truth of statements, and yet we are able to reason about the world and make predictions. We will see that the principles of Boolean logic can be extended to reason in the face of uncertainty. The mathematical framework that allows us to do this is probability theory, which — as we will find in this first chapter — can be seen as a natural extension of Boolean logic from the domain of certainty to the domain of uncertainty. In fact, in the 20th century, Richard Cox and Edwin Thompson Jaynes have done early work to formalize probability theory as the “logic under uncertainty” (Cox, 1961; Jaynes, 2002).

In this first chapter, we will briefly recall the fundamentals of probability theory, and we will see how *probabilistic inference* can be used

to reason about the world. In the remaining chapters, we will then discuss how probabilistic inference can be performed efficiently given limited computational resources and limited time, which is the key challenge in *probabilistic artificial intelligence*.

1.1 Probability

Probability is commonly interpreted in two different ways. In the frequentist interpretation, one interprets the probability of an event (say a coin coming up “heads” when flipping it) as the limit of relative frequencies in repeated independent experiments. That is,

$$\text{Probability} = \lim_{N \rightarrow \infty} \frac{\# \text{ events happening in } N \text{ trials}}{N}.$$

This interpretation is natural, but has a few issues. It is not very difficult to conceive of settings where repeated experiments do not make sense. Consider the outcome:

“Person X will live for at least 80 years.”

There is no way in which we could conduct multiple independent experiments in this case. Still, this statement is going to turn out either true or false, as humans we are just not able to determine its truth value beforehand. Nevertheless, humans commonly have *beliefs* about statements of this kind. We also commonly reason about statements such as

“The Beatles were more groundbreaking than The Monkees.”

This statement does not even have an objective truth value, and yet we as humans tend to have opinions about it.

While it is natural to consider the relative frequency of the outcome in repeated experiments as our belief, if we are not able to conduct repeated experiments, our notion of probability is simply a subjective measure of uncertainty about outcomes. In the early 20th century, Bruno De Finetti has done foundational work to formalize this notion which is commonly called *Bayesian reasoning* or the Bayesian interpretation of probability (De Finetti, 1970).

We will see that modern approaches to probabilistic inference often lend themselves to a Bayesian interpretation, even if such an interpretation is not strictly necessary. For our purposes, probabilities will be a means to an end: the end usually being solving some task. This task may be to make a prediction or to take an action with an uncertain outcome, and we can evaluate methods according to how well they perform on this task. No matter the interpretation, the mathematical

framework of probability theory which we will formally introduce in the following is the same.

1.1.1 Probability Spaces

A probability space is a mathematical model for a random experiment. The set of all possible outcomes of the experiment Ω is called *sample space*. An *event* $A \subseteq \Omega$ of interest may be any combination of possible outcomes. The set of all events $\mathcal{A} \subseteq \mathcal{P}(\Omega)$ that we are interested in is often called the *event space* of the experiment.¹ This set of events is required to be a σ -algebra over the sample space.

Definition 1.1 (σ -algebra). Given the set Ω , the set $\mathcal{A} \subseteq \mathcal{P}(\Omega)$ is a σ -algebra over Ω if the following properties are satisfied:

1. $\Omega \in \mathcal{A}$;
2. if $A \in \mathcal{A}$, then $\bar{A} \in \mathcal{A}$ (*closedness under complements*); and
3. if we have $A_i \in \mathcal{A}$ for all i , then $\bigcup_{i=1}^{\infty} A_i \in \mathcal{A}$ (*closedness under countable unions*).

Note that the three properties of σ -algebras correspond to characteristics we universally expect when working with random experiments. Namely, that we are able to reason about the event Ω that any of the possible outcomes occur, that we are able to reason about an event not occurring, and that we are able to reason about events that are composed of multiple (smaller) events.

Example 1.2: Event space of throwing a die

The event space \mathcal{A} can also be thought of as “how much information is available about the experiment”. For example, if the experiment is a throw of a die and Ω is the set of possible values on the die: $\Omega = \{1, \dots, 6\}$, then the following \mathcal{A} implies that the observer cannot distinguish between 1 and 3:

$$\mathcal{A} \doteq \{\emptyset, \Omega, \{1, 3, 5\}, \{2, 4, 6\}\}.$$

Intuitively, the observer only understands the parity of the face of the die.

Definition 1.3 (Probability measure). Given the set Ω and the σ -algebra \mathcal{A} over Ω , the function

$$\mathbb{P} : \mathcal{A} \rightarrow \mathbb{R}$$

is a *probability measure* on \mathcal{A} if the *Kolmogorov axioms* are satisfied:

1. $0 \leq \mathbb{P}(A) \leq 1$ for any $A \in \mathcal{A}$;
2. $\mathbb{P}(\Omega) = 1$; and

¹ We use $\mathcal{P}(\Omega)$ to denote the *power set* (set of all subsets) of Ω .

3. $\mathbb{P}(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} \mathbb{P}(A_i)$ for any countable set of mutually disjoint events $\{A_i \in \mathcal{A}\}_i$.²

Remarkably, all further statements about probability follow from these three natural axioms. For an event $A \in \mathcal{A}$, we call $\mathbb{P}(A)$ the *probability* of A . We are now ready to define a probability space.

Definition 1.4 (Probability space). A *probability space* is a triple $(\Omega, \mathcal{A}, \mathbb{P})$ where

- Ω is a sample space,
- \mathcal{A} is a σ -algebra over Ω , and
- \mathbb{P} is a probability measure on \mathcal{A} .

Example 1.5: Borel σ -algebra over \mathbb{R}

In our context, we often have that Ω is the set of real numbers \mathbb{R} or a compact subset of it. In this case, a natural event space is the σ -algebra generated by the set of events

$$A_x \doteq \{x' \in \Omega : x' \leq x\}.$$

The smallest σ -algebra \mathcal{A} containing all sets A_x is called the *Borel σ -algebra*. \mathcal{A} contains all “reasonable” subsets of Ω (except for some pathological examples). For example, \mathcal{A} includes all singleton sets $\{x\}$, as well as all countable unions of intervals.

In the case of discrete Ω , in fact $\mathcal{A} = \mathcal{P}(\Omega)$, i.e., the Borel σ -algebra contains *all* subsets of Ω .

1.1.2 Random Variables

The set Ω is often rather complex. For example, take Ω to be the set of all possible graphs on n vertices. Then the outcome of our experiment is a graph. Usually, we are not interested in a specific graph but rather a property such as the number of edges, which is shared by many graphs. A function that maps a graph to its number of edges is a random variable.

Definition 1.6 (Random variable). A *random variable* X is a function

$$X : \Omega \rightarrow \mathcal{T}$$

where \mathcal{T} is called *target space* of the random variable,³ and where X respects the information available in the σ -algebra \mathcal{A} . That is,⁴

$$\forall S \subseteq \mathcal{T} : \quad \{\omega \in \Omega : X(\omega) \in S\} \in \mathcal{A}. \quad (1.1)$$

² We say that a set of sets $\{A_i\}_i$ is disjoint if for all $i \neq j$ we have $A_i \cap A_j = \emptyset$.

³ For a random variable that maps a graph to its number of edges, $\mathcal{T} = \mathbb{N}_0$. For our purposes, you can generally assume $\mathcal{T} \subseteq \mathbb{R}$.

⁴ In our example of throwing a die, X should assign the same value to the outcomes 1, 3, 5.

Concrete values x of a random variable X are often referred to as *states* or *realizations* of X . The probability that X takes on a value in $S \subseteq \mathcal{T}$ is

$$\mathbb{P}(X \in S) = \mathbb{P}(\{\omega \in \Omega : X(\omega) \in S\}). \quad (1.2)$$

1.1.3 Distributions

Consider a random variable X on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$, where Ω is a compact subset of \mathbb{R} , and \mathcal{A} the Borel σ -algebra.

In this case, we can refer to the probability that X assumes a particular state or set of states by writing

$$p_X(x) \doteq \mathbb{P}(X = x) \quad (\text{in the discrete setting}), \quad (1.3)$$

$$P_X(x) \doteq \mathbb{P}(X \leq x). \quad (1.4)$$

Note that " $X = x$ " and " $X \leq x$ " are merely events (that is, they characterize subsets of the sample space Ω satisfying this condition) which are in the Borel σ -algebra, and hence their probability is well-defined.

Hereby, p_X and P_X are referred to as the probability mass function (PMF) and cumulative distribution function (CDF) of X , respectively. Note that we can also *implicitly* define probability spaces through random variables and their associated PMF/CDF, which is often very convenient.

We list some common examples of discrete distributions in Appendix A.1.1. Further, note that for continuous variables, $\mathbb{P}(X = x) = 0$. Here, instead we typically use the probability density function (PDF), to which we (with slight abuse of notation) also refer with p_X . We discuss densities in greater detail in Section 1.1.4.

We call the subset $S \subseteq \mathcal{T}$ of the domain of a PMF or PDF p_X such that all elements $x \in S$ have positive probability, $p_X(x) > 0$, the *support* of the distribution p_X . This quantity is denoted by $X(\Omega)$.

1.1.4 Continuous Distributions

As mentioned, a continuous random variable can be characterized by its *probability density function* (PDF). But what is a density? We can derive some intuition from physics.

Let M be a (non-homogeneous) physical object, e.g., a rock. We commonly use $m(M)$ and $\text{vol}(M)$ to refer to its mass and volume, respectively. Now, consider for a point $x \in M$ and a ball $B_r(x)$ around x with radius r the following quantities:

$$\lim_{r \rightarrow 0} \text{vol}(B_r(x)) = 0 \quad \lim_{r \rightarrow 0} m(B_r(x)) = 0.$$

They appear utterly uninteresting at first, yet, if we divide them, we get what is called the *density* of M at x .

$$\lim_{r \rightarrow 0} \frac{m(B_r(x))}{\text{vol}(B_r(x))} \doteq \rho(x).$$

We know that the relationship between density and mass is described by the following formula:

$$m(M) = \int_M \rho(x) dx.$$

In other words, the density is to be integrated. For a small region I around x , we can approximate $m(I) \approx \rho(x) \cdot \text{vol}(I)$.

Crucially, observe that even though the mass of any particular point x is zero, i.e., $m(\{x\}) = 0$, assigning a density $\rho(x)$ to x is useful for integration and approximation. The same idea applies to continuous random variables, only that volume corresponds to intervals on the real line and mass to probability. Recall that probability density functions are normalized such that their probability mass across the entire real line integrates to one.

Example 1.7: Normal distribution / Gaussian

A famous example of a continuous distribution is the *normal distribution*, also called *Gaussian*. We say, a random variable X is *normally distributed*, $X \sim \mathcal{N}(\mu, \sigma^2)$, if its PDF is

$$\mathcal{N}(x; \mu, \sigma^2) \doteq \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right). \quad (1.5)$$

We have $\mathbb{E}[X] = \mu$ and $\text{Var}[X] = \sigma^2$. If $\mu = 0$ and $\sigma^2 = 1$, this distribution is called the *standard normal distribution*. The Gaussian CDF cannot be expressed in closed-form.

Note that the mean of a Gaussian distribution coincides with the maximizer of its PDF, also called *mode* of a distribution.

We will focus in the remainder of this chapter on continuous distributions, but the concepts we discuss extend mostly to discrete distributions simply by “replacing integrals by sums”.

1.1.5 Joint Probability

A joint probability (as opposed to a marginal probability) is the probability of two or more events occurring simultaneously:

$$\mathbb{P}(A, B) \doteq \mathbb{P}(A \cap B). \quad (1.6)$$

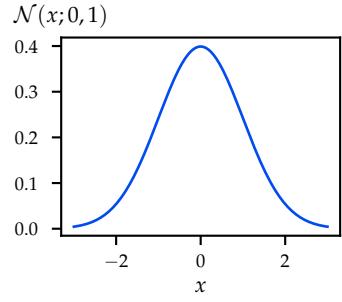


Figure 1.1: PDF of the standard normal distribution. Observe that the PDF is symmetric around the mode.

In terms of random variables, this concept extends to joint distributions. Instead of characterizing a single random variable, a *joint distribution* is a function $p_{\mathbf{X}} : \mathbb{R}^n \rightarrow \mathbb{R}$, characterizing a *random vector* $\mathbf{X} \doteq [X_1 \cdots X_n]^\top$. For example, if the X_i are discrete, the joint distribution characterizes joint probabilities of the form

$$\mathbb{P}(\mathbf{X} = [x_1, \dots, x_n]) = \mathbb{P}(X_1 = x_1, \dots, X_n = x_n),$$

and hence describes the relationship among all variables X_i . For this reason, a joint distribution is also called a *generative model*. We use $X_{i:j}$ to denote the random vector $[X_i \cdots X_j]^\top$.

We can “sum out” (respectively “integrate out”) variables from a joint distribution in a process called “marginalization”:

Fact 1.8 (Sum rule). *We have that*

$$p(x_{1:i-1}, x_{i+1:n}) = \int_{X_i(\Omega)} p(x_{1:i-1}, x_i, x_{i+1:n}) dx_i. \quad (1.7)$$

1.1.6 Conditional Probability

Conditional probability updates the probability of an event A given some new information, for example, after observing the event B .

Definition 1.9 (Conditional probability). Given two events A and B such that $\mathbb{P}(B) > 0$, the probability of A conditioned on B is given as

$$\mathbb{P}(A | B) \doteq \frac{\mathbb{P}(A, B)}{\mathbb{P}(B)}. \quad (1.8)$$

Simply rearranging the terms yields,

$$\mathbb{P}(A, B) = \mathbb{P}(A | B) \cdot \mathbb{P}(B) = \mathbb{P}(B | A) \cdot \mathbb{P}(A). \quad (1.9)$$

Thus, the probability that both A and B occur can be calculated by multiplying the probability of event A and the probability of B conditional on A occurring.

We say $\mathbf{Z} \sim \mathbf{X} | \mathbf{Y} = \mathbf{y}$ (or simply $\mathbf{Z} \sim \mathbf{X} | \mathbf{y}$) if \mathbf{Z} follows the *conditional distribution*

$$p_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}) \doteq \frac{p_{\mathbf{X},\mathbf{Y}}(\mathbf{x}, \mathbf{y})}{p_{\mathbf{Y}}(\mathbf{y})}. \quad (1.10)$$

If \mathbf{X} and \mathbf{Y} are discrete, we have that $p_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}) = \mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = \mathbf{y})$ as one would naturally expect.

Extending Equation (1.9) to arbitrary random vectors yields the product rule (also called the *chain rule of probability*):

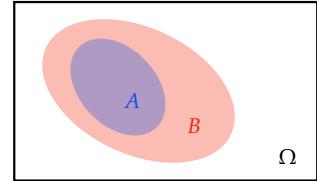


Figure 1.2: Conditioning an event A on another event B can be understood as replacing the universe of all possible outcomes Ω by the observed outcomes B . Then, the conditional probability is simply expressing the likelihood of A given that B occurred.

Fact 1.10 (Product rule). *Given random variables $X_{1:n}$,*

$$p(x_{1:n}) = p(x_1) \cdot \prod_{i=2}^n p(x_i | x_{1:i-1}). \quad (1.11)$$

Combining sum rule and product rule, we can compute marginal probabilities too:

$$p(\mathbf{x}) = \int_{\mathbf{Y}(\Omega)} p(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \int_{\mathbf{Y}(\Omega)} p(\mathbf{x} | \mathbf{y}) \cdot p(\mathbf{y}) d\mathbf{y} \quad (1.12)$$

This is called the *law of total probability* (LOTP), which is colloquially often referred to as *conditioning* on \mathbf{Y} . If it is difficult to compute $p(\mathbf{x})$ directly, conditioning can be a useful technique when \mathbf{Y} is chosen such that the densities $p(\mathbf{x} | \mathbf{y})$ and $p(\mathbf{y})$ are straightforward to understand.

first using the sum rule (1.7) then the product rule (1.11)

1.1.7 Independence

Two random vectors \mathbf{X} and \mathbf{Y} are *independent* (denoted $\mathbf{X} \perp \mathbf{Y}$) if and only if knowledge about the state of one random vector does not affect the distribution of the other random vector, namely if their conditional CDF (or in case they have a joint density, their conditional PDF) simplifies to

$$P_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}) = P_{\mathbf{X}}(\mathbf{x}), \quad p_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}) = p_{\mathbf{X}}(\mathbf{x}). \quad (1.13)$$

For the conditional probabilities to be well-defined, we need to assume that $p_{\mathbf{Y}}(\mathbf{y}) > 0$.

The more general characterization of independence is that \mathbf{X} and \mathbf{Y} are independent if and only if their joint CDF (or in case they have a joint density, their joint PDF) can be decomposed as follows:

$$P_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}) = P_{\mathbf{X}}(\mathbf{x}) \cdot P_{\mathbf{Y}}(\mathbf{y}), \quad p_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}) = p_{\mathbf{X}}(\mathbf{x}) \cdot p_{\mathbf{Y}}(\mathbf{y}). \quad (1.14)$$

The equivalence of the two characterizations (when $p_{\mathbf{Y}}(\mathbf{y}) > 0$) is easily proven using the product rule: $p_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}) = p_{\mathbf{Y}}(\mathbf{y}) \cdot p_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y})$.

A “weaker” notion of independence is conditional independence.⁵ Two random vectors \mathbf{X} and \mathbf{Y} are *conditionally independent* given a random vector \mathbf{Z} (denoted $\mathbf{X} \perp \mathbf{Y} | \mathbf{Z}$) iff, given \mathbf{Z} , knowledge about the value of one random vector \mathbf{Y} does not affect the distribution of the other random vector \mathbf{X} , namely if

$$P_{\mathbf{X}|\mathbf{Y}, \mathbf{Z}}(\mathbf{x} | \mathbf{y}, \mathbf{z}) = P_{\mathbf{X}|\mathbf{Z}}(\mathbf{x} | \mathbf{z}), \quad (1.15a)$$

$$p_{\mathbf{X}|\mathbf{Y}, \mathbf{Z}}(\mathbf{x} | \mathbf{y}, \mathbf{z}) = p_{\mathbf{X}|\mathbf{Z}}(\mathbf{x} | \mathbf{z}). \quad (1.15b)$$

⁵ We discuss in Remark 1.11 how “weaker” is to be interpreted in this context.

Similarly to independence, we have that \mathbf{X} and \mathbf{Y} are conditionally independent given \mathbf{Z} if and only if their joint CDF or joint PDF can be decomposed as follows:

$$P_{\mathbf{X}, \mathbf{Y} | \mathbf{Z}}(\mathbf{x}, \mathbf{y} | \mathbf{z}) = P_{\mathbf{X} | \mathbf{Z}}(\mathbf{x} | \mathbf{z}) \cdot P_{\mathbf{Y} | \mathbf{Z}}(\mathbf{y} | \mathbf{z}), \quad (1.16a)$$

$$p_{\mathbf{X}, \mathbf{Y} | \mathbf{Z}}(\mathbf{x}, \mathbf{y} | \mathbf{z}) = p_{\mathbf{X} | \mathbf{Z}}(\mathbf{x} | \mathbf{z}) \cdot p_{\mathbf{Y} | \mathbf{Z}}(\mathbf{y} | \mathbf{z}). \quad (1.16b)$$

Remark 1.11: Common causes

How can conditional independence be understood as a “weaker” notion of independence? Clearly, conditional independence does not imply independence: a trivial example is $X \perp X | X \not\Rightarrow X \perp X$.⁶ Neither does independence imply conditional independence: for example, $X \perp Y \not\Rightarrow X \perp Y | X + Y$.⁷

When we say that conditional independence is a weaker notion we mean to emphasize that X and Y can be “made” (conditionally) independent by conditioning on the “right” Z even if X and Y are dependent. This is known as *Reichenbach’s common cause principle* which says that for any two random variables $X \not\perp Y$ there exists a random variable Z (which may be X or Y) that causally influences both X and Y , and which is such that $X \perp Y | Z$.

⁶ $X \perp X | X$ is true trivially.

⁷ Knowing X and $X + Y$ already implies the value of Y , and hence, $X \not\perp Y | X + Y$.

1.1.8 Directed Graphical Models

Directed graphical models (also called *Bayesian networks*) are often used to visually denote the (conditional) independence relationships of a large number of random variables. They are a schematic representation of the factorization of the generative model into a product of conditional distributions as a directed acyclic graph. Given the sequence of random variables $\{X_i\}_{i=1}^n$, their generative model can be expressed as

$$p(x_{1:n}) = \prod_{i=1}^n p(x_i | \text{parents}(x_i)) \quad (1.17)$$

where $\text{parents}(x_i)$ is the set of parents of the vertex X_i in the directed graphical model. In other words, the parenthood relationship encodes a conditional independence of a random variable X with a random variable Y given their parents:⁸

$$X \perp Y | \text{parents}(X), \text{parents}(Y). \quad (1.18)$$

Equation (1.17) simply uses the product rule and the conditional independence relationships to factorize the generative model. This can greatly reduce the model’s complexity, i.e., the length of the product.

⁸ More generally, vertices u and v are conditionally independent given a set of vertices Z if Z *d-separates* u and v , which we will not cover in depth here.

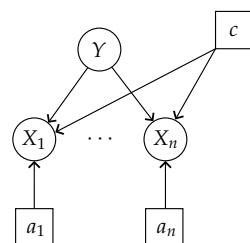


Figure 1.3: Example of a directed graphical model. The random variables X_1, \dots, X_n are mutually independent given the random variable Y . The squared rectangular nodes are used to represent dependencies on parameters c, a_1, \dots, a_n .

An example of a directed graphical model is given in Figure 1.3. Circular vertices represent random quantities (i.e., random variables). In contrast, square vertices are commonly used to represent deterministic quantities (i.e., parameters that the distributions depend on). In the given example, we have that X_i is conditionally independent of all other X_j given Y . *Plate notation* is a condensed notation used to represent repeated variables of a graphical model. An example is given in Figure 1.4.

1.1.9 Expectation

The *expected value* or *mean* $\mathbb{E}[\mathbf{X}]$ of a random vector \mathbf{X} is the (asymptotic) arithmetic mean of an arbitrarily increasing number of independent realizations of \mathbf{X} . That is,⁹

$$\mathbb{E}[\mathbf{X}] \doteq \int_{\mathbf{X}(\Omega)} \mathbf{x} \cdot p(\mathbf{x}) d\mathbf{x} \quad (1.19)$$

A very special and often used property of expectations is their *linearity*, namely that for any random vectors \mathbf{X} and \mathbf{Y} in \mathbb{R}^n and any $\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^m$ it holds that

$$\mathbb{E}[\mathbf{AX} + \mathbf{b}] = \mathbf{A}\mathbb{E}[\mathbf{X}] + \mathbf{b} \quad \text{and} \quad \mathbb{E}[\mathbf{X} + \mathbf{Y}] = \mathbb{E}[\mathbf{X}] + \mathbb{E}[\mathbf{Y}]. \quad (1.20)$$

Note that \mathbf{X} and \mathbf{Y} do not necessarily have to be independent! Further, if \mathbf{X} and \mathbf{Y} are independent then

$$\mathbb{E}[\mathbf{XY}^\top] = \mathbb{E}[\mathbf{X}] \cdot \mathbb{E}[\mathbf{Y}]^\top. \quad (1.21)$$

The following intuitive lemma can be used to compute expectations of transformed random variables.

Fact 1.12 (Law of the unconscious statistician, LOTUS).

$$\mathbb{E}[g(\mathbf{X})] = \int_{\mathbf{X}(\Omega)} g(\mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} \quad (1.22)$$

where $g : \mathbf{X}(\Omega) \rightarrow \mathbb{R}^n$ is a “nice” function¹⁰ and \mathbf{X} is a continuous random vector. The analogous statement with a sum replacing the integral holds for discrete random variables.

This is a nontrivial fact that can be proven using the change of variables formula which we discuss in Section 1.1.11.

Similarly to conditional probability, we can also define conditional expectations. The expectation of a continuous random vector \mathbf{X} given that $\mathbf{Y} = \mathbf{y}$ is defined as

$$\mathbb{E}[\mathbf{X} | \mathbf{Y} = \mathbf{y}] \doteq \int_{\mathbf{X}(\Omega)} \mathbf{x} \cdot p_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}) d\mathbf{x}. \quad (1.23)$$

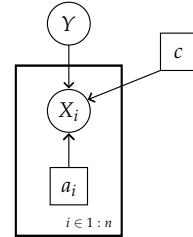


Figure 1.4: The same directed graphical model as in Figure 1.3 using plate notation.

⁹ In infinite probability spaces, absolute convergence of $\mathbb{E}[\mathbf{X}]$ is necessary for the existence of $\mathbb{E}[\mathbf{X}]$.

¹⁰ g being a continuous function, which is either bounded or absolutely integrable (i.e., $\int |g(\mathbf{x})| p(\mathbf{x}) d\mathbf{x} < \infty$), is sufficient. This is satisfied in most cases.

Observe that $\mathbb{E}[\mathbf{X} | \mathbf{Y} = \cdot]$ defines a deterministic mapping from \mathbf{y} to $\mathbb{E}[\mathbf{X} | \mathbf{Y} = \mathbf{y}]$. Therefore, $\mathbb{E}[\mathbf{X} | \mathbf{Y}]$ is itself a random vector:

$$\mathbb{E}[\mathbf{X} | \mathbf{Y}](\omega) = \mathbb{E}[\mathbf{X} | \mathbf{Y} = \mathbf{Y}(\omega)] \quad (1.24)$$

where $\omega \in \Omega$. This random vector $\mathbb{E}[\mathbf{X} | \mathbf{Y}]$ is called the *conditional expectation* of \mathbf{X} given \mathbf{Y} .

Analogously to the law of total probability (1.12), one can condition an expectation on another random vector. This is known as the *tower rule* or the *law of total expectation* (LOTE):

Theorem 1.13 (Tower rule). *Given random vectors \mathbf{X} and \mathbf{Y} , we have*

$$\mathbb{E}_{\mathbf{Y}}[\mathbb{E}_{\mathbf{X}}[\mathbf{X} | \mathbf{Y}]] = \mathbb{E}[\mathbf{X}]. \quad (1.25)$$

Proof sketch. We only prove the case where \mathbf{X} and \mathbf{Y} have a joint density. We have

$$\begin{aligned} \mathbb{E}[\mathbb{E}[\mathbf{X} | \mathbf{Y}]] &= \int \left(\int \mathbf{x} \cdot p(\mathbf{x} | \mathbf{y}) d\mathbf{x} \right) p(\mathbf{y}) d\mathbf{y} \\ &= \int \int \mathbf{x} \cdot p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} && \text{by definition of conditional densities} \\ &= \int \mathbf{x} \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} d\mathbf{x} && \text{(1.10)} \\ &= \int \mathbf{x} \cdot p(\mathbf{x}) d\mathbf{x} && \text{by Fubini's theorem} \\ &= \mathbb{E}[\mathbf{X}]. && \text{using the sum rule (1.7)} \end{aligned} \quad \square$$

1.1.10 Covariance and Variance

Given two random vectors \mathbf{X} in \mathbb{R}^n and \mathbf{Y} in \mathbb{R}^m , their *covariance* is defined as

$$\text{Cov}[\mathbf{X}, \mathbf{Y}] \doteq \mathbb{E}\left[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{Y} - \mathbb{E}[\mathbf{Y}])^\top\right] \quad (1.26)$$

$$= \mathbb{E}\left[\mathbf{XY}^\top\right] - \mathbb{E}[\mathbf{X}] \cdot \mathbb{E}[\mathbf{Y}]^\top \quad (1.27)$$

$$= \text{Cov}[\mathbf{Y}, \mathbf{X}]^\top \in \mathbb{R}^{n \times m}. \quad (1.28)$$

Covariance measures the linear dependence between two random vectors since a direct consequence of its definition (1.26) is that given linear maps $A \in \mathbb{R}^{n' \times n}$, $B \in \mathbb{R}^{m' \times m}$, vectors $\mathbf{c} \in \mathbb{R}^{n'}$, $\mathbf{d} \in \mathbb{R}^{m'}$ and random vectors \mathbf{X} in \mathbb{R}^n and \mathbf{Y} in \mathbb{R}^m , we have that

$$\text{Cov}[A\mathbf{X} + \mathbf{c}, B\mathbf{Y} + \mathbf{d}] = A\text{Cov}[\mathbf{X}, \mathbf{Y}]B^\top. \quad (1.29)$$

Two random vectors \mathbf{X} and \mathbf{Y} are said to be *uncorrelated* if and only if $\text{Cov}[\mathbf{X}, \mathbf{Y}] = \mathbf{0}$. Note that if \mathbf{X} and \mathbf{Y} are independent, then Equation (1.21) implies that \mathbf{X} and \mathbf{Y} are uncorrelated. The reverse does not hold in general.

Remark 1.14: Correlation

The *correlation* of the random vectors \mathbf{X} and \mathbf{Y} is a normalized covariance,

$$\text{Cor}[\mathbf{X}, \mathbf{Y}](i, j) \doteq \frac{\text{Cov}[X_i, Y_j]}{\sqrt{\text{Var}[X_i]\text{Var}[Y_j]}} \in [-1, 1]. \quad (1.30)$$

Two random vectors \mathbf{X} and \mathbf{Y} are therefore uncorrelated if and only if $\text{Cor}[\mathbf{X}, \mathbf{Y}] = 0$.

There is also a nice geometric interpretation of covariance and correlation. For zero mean random variables X and Y , $\text{Cov}[X, Y]$ is an inner product.¹¹

The cosine of the angle θ between X and Y (that are not deterministic) coincides with their correlation,

$$\cos \theta = \frac{\text{Cov}[X, Y]}{\|X\| \|Y\|} = \text{Cor}[X, Y]. \quad (1.31)$$

$\cos \theta$ is also called a *cosine similarity*. Thus,

$$\theta = \arccos \text{Cor}[X, Y]. \quad (1.32)$$

For example, if X and Y are uncorrelated, then they are orthogonal in the inner product space. If $\text{Cor}[X, Y] = -1$ then $\theta \equiv \pi$ (that is, X and Y “point in opposite directions”), whereas if $\text{Cor}[X, Y] = 1$ then $\theta \equiv 0$ (that is, X and Y “point in the same direction”).

The covariance of a random vector \mathbf{X} in \mathbb{R}^n with itself is called its *variance*:

$$\text{Var}[\mathbf{X}] \doteq \text{Cov}[\mathbf{X}, \mathbf{X}] \quad (1.33)$$

$$= \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top] \quad (1.34)$$

$$= \mathbb{E}[\mathbf{X}\mathbf{X}^\top] - \mathbb{E}[\mathbf{X}] \cdot \mathbb{E}[\mathbf{X}]^\top \quad (1.35)$$

$$= \begin{bmatrix} \text{Cov}[X_1, X_1] & \cdots & \text{Cov}[X_1, X_n] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_n, X_1] & \cdots & \text{Cov}[X_n, X_n] \end{bmatrix}. \quad (1.36)$$

The scalar variance $\text{Var}[X]$ of a random variable X is a measure of uncertainty about the value of X since it measures the average squared deviation from $\mathbb{E}[X]$. We will see that the eigenvalue spectrum of a covariance matrix can serve as a measure of uncertainty in the multivariate setting.¹²

¹¹ That is,

- $\text{Cov}[X, Y]$ is symmetric,
- $\text{Cov}[X, Y]$ is linear (here we use $\mathbb{E}X = \mathbb{E}Y = 0$), and
- $\text{Cov}[X, X] \geq 0$.

using the Euclidean inner product formula, $\text{Cov}[X, Y] = \|X\| \|Y\| \cos \theta$

¹² The *multivariate* setting (as opposed to the *univariate* setting) studies the joint distribution of multiple random variables.

Remark 1.15: Standard deviation

The length of a random variable X in the inner product space described in Remark 1.14 is called its *standard deviation*,

$$\|X\| = \sqrt{\text{Cov}[X, X]} = \sqrt{\text{Var}[X]} \doteq \sigma[X]. \quad (1.37)$$

That is, the longer a random variable is in the inner product space, the more “uncertain” we are about its value. If a random variable has length 0, then it is deterministic.

The variance of a random vector \mathbf{X} is also called the *covariance matrix* of \mathbf{X} and denoted by $\Sigma_{\mathbf{X}}$ (or Σ if the correspondence to \mathbf{X} is clear from context). A covariance matrix is symmetric by definition due to the symmetry of covariance, and is always positive semi-definite $\textcircled{?}$.

Problem 1.4

Two useful properties of variance are the following:

- It follows from Equation (1.29) that for any linear map $A \in \mathbb{R}^{m \times n}$ and vector $\mathbf{b} \in \mathbb{R}^m$,

$$\text{Var}[A\mathbf{X} + \mathbf{b}] = A\text{Var}[\mathbf{X}]A^\top. \quad (1.38)$$

In particular, $\text{Var}[-\mathbf{X}] = \text{Var}[\mathbf{X}]$.

- It follows from the definition of variance (1.34) that for any two random vectors \mathbf{X} and \mathbf{Y} ,

$$\text{Var}[\mathbf{X} + \mathbf{Y}] = \text{Var}[\mathbf{X}] + \text{Var}[\mathbf{Y}] + 2\text{Cov}[\mathbf{X}, \mathbf{Y}]. \quad (1.39)$$

In particular, if \mathbf{X} and \mathbf{Y} are independent then the covariance term vanishes and $\text{Var}[\mathbf{X} + \mathbf{Y}] = \text{Var}[\mathbf{X}] + \text{Var}[\mathbf{Y}]$.

Analogously to conditional probability and conditional expectation, we can also define conditional variance. The *conditional variance* of a random vector \mathbf{X} given another random vector \mathbf{Y} is the random vector

$$\text{Var}[\mathbf{X} | \mathbf{Y}] \doteq \mathbb{E}\left[(\mathbf{X} - \mathbb{E}[\mathbf{X} | \mathbf{Y}])(\mathbf{X} - \mathbb{E}[\mathbf{X} | \mathbf{Y}])^\top \mid \mathbf{Y}\right]. \quad (1.40)$$

Intuitively, the conditional variance is the remaining variance when we use $\mathbb{E}[\mathbf{X} | \mathbf{Y}]$ to predict \mathbf{X} rather than if we used $\mathbb{E}[\mathbf{X}]$. One can also condition a variance on another random vector, analogously to the laws of total probability (1.12) and expectation (1.25).

Theorem 1.16 (Law of total variance, LOTV).

$$\text{Var}[\mathbf{X}] = \mathbb{E}_{\mathbf{Y}}[\text{Var}_{\mathbf{X}}[\mathbf{X} | \mathbf{Y}]] + \text{Var}_{\mathbf{Y}}[\mathbb{E}_{\mathbf{X}}[\mathbf{X} | \mathbf{Y}]]. \quad (1.41)$$

Here, the first term measures the average deviation from the mean of \mathbf{X} across realizations of \mathbf{Y} and the second term measures the uncertainty

in the mean of \mathbf{X} across realizations of \mathbf{Y} . In Section 2.2, we will see that both terms have a meaningful characterization in the context of probabilistic inference.

Proof sketch of LOTV. To simplify the notation, we present only a proof for the univariate setting.

$$\begin{aligned}
 \text{Var}[X] &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 \\
 &= \mathbb{E}[\mathbb{E}[X^2 | Y]] - \mathbb{E}[\mathbb{E}[X | Y]]^2 && \text{by the tower rule (1.25)} \\
 &= \mathbb{E}[\text{Var}[X | Y] + \mathbb{E}[X | Y]^2] - \mathbb{E}[\mathbb{E}[X | Y]]^2 \\
 &= \mathbb{E}[\text{Var}[X | Y]] + (\mathbb{E}[\mathbb{E}[X | Y]^2] - \mathbb{E}[\mathbb{E}[X | Y]]^2) \\
 &= \mathbb{E}[\text{Var}[X | Y]] + \text{Var}[\mathbb{E}[X | Y]]. && \square \quad \text{by the definition of variance (1.35)}
 \end{aligned}$$

1.1.11 Change of Variables

It is often useful to understand the distribution of a transformed random variable $Y = g(X)$ that is defined in terms of a random variable X , whose distribution is known. Let us first consider the univariate setting. We would like to express the distribution of Y in terms of the distribution of X , that is, we would like to find

$$P_Y(y) = \mathbb{P}(Y \leq y) = \mathbb{P}(g(X) \leq y) = \mathbb{P}\left(X \leq g^{-1}(y)\right). \quad (1.42)$$

When the random variables are continuous, this probability can be expressed as an integration over the domain of X . We can then use the substitution rule of integration to “change the variables” to an integration over the domain of Y . Taking the derivative yields the density p_Y .¹³ There is an analogous change of variables formula for the multivariate setting.

Fact 1.17 (Change of variables formula). *Let \mathbf{X} be a random vector in \mathbb{R}^n with density $p_{\mathbf{X}}$ and let $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a differentiable and invertible function. Then $\mathbf{Y} = g(\mathbf{X})$ is another random variable, whose density can be computed based on $p_{\mathbf{X}}$ and g as follows:*

$$p_{\mathbf{Y}}(\mathbf{y}) = p_{\mathbf{X}}(g^{-1}(\mathbf{y})) \cdot |\det(Dg^{-1}(\mathbf{y}))| \quad (1.43)$$

where $Dg^{-1}(\mathbf{y})$ is the Jacobian of g^{-1} evaluated at \mathbf{y} .

¹³ The full proof of the change of variables formula in the univariate setting can be found in section 6.7.2 of “Mathematics for machine learning” (Deisenroth et al., 2020).

Here, the term $|\det(Dg^{-1}(\mathbf{y}))|$ measures how much a unit volume changes when applying g . Intuitively, the change of variables swaps the coordinate system over which we integrate. The factor $|\det(Dg^{-1}(\mathbf{y}))|$ corrects for the change in volume that is caused by this change in coordinates.

Intuitively, you can think of the vector field g as a perturbation to \mathbf{X} , “pushing” the probability mass around. The perturbation of a density $p_{\mathbf{X}}$ by g is commonly denoted by the *pushforward*

$$g_* p_{\mathbf{X}} \doteq p_{\mathbf{Y}} \quad \text{where } \mathbf{Y} = g(\mathbf{X}). \quad (1.44)$$

This concludes our quick tour of probability theory, and we are well-prepared to return to the topic of probabilistic inference.

1.2 Probabilistic Inference

Recall the logical implication “If it is raining, the ground is wet.” from the beginning of this chapter. Suppose that we look outside a window and see that it is not raining: will the ground be dry? Logical reasoning does not permit drawing an inference of this kind, as there might be reasons other than rain for which the ground could be wet (e.g., sprinklers). However, intuitively, by observing that it is not raining, we have just excluded the possibility that the ground is wet because of rain, and therefore we would deem it “more likely” that the ground is dry than before. In other words, if we were to walk outside now and the ground was wet, we would be more surprised than we would have been if we had not looked outside the window before.

As humans, we are constantly making such “plausible” inferences of our beliefs: be it about the weather, the outcomes of our daily decisions, or the behavior of others. *Probabilistic inference* is the process of updating such a prior belief $\mathbb{P}(\bar{W})$ to a posterior belief $\mathbb{P}(\bar{W} | \bar{R})$ upon observing \bar{R} where — to reduce clutter — we write W for “The ground is wet” and R for “It is raining”.

The central principle of probabilistic inference is Bayes’ rule:

Theorem 1.18 (Bayes’ rule). *Given random vectors \mathbf{X} in \mathbb{R}^n and \mathbf{Y} in \mathbb{R}^m , we have for any $x \in \mathbb{R}^n, y \in \mathbb{R}^m$ that*

$$p(x | y) = \frac{p(y | x) \cdot p(x)}{p(y)}. \quad (1.45)$$

Proof. Bayes’ rule is a direct consequence of the definition of conditional densities (1.10) and the product rule (1.11). \square

Let us consider the meaning of each term separately:

- the *prior* $p(x)$ is the initial belief about x ,
- the (*conditional*) *likelihood* $p(y | x)$ describes how likely the observations y are under a given value x ,

- the *posterior* $p(x | y)$ is the updated belief about x after observing y ,
- the *joint likelihood* $p(x, y) = p(y | x)p(x)$ combines prior and likelihood,
- the *marginal likelihood* $p(y)$ describes how likely the observations y are across all values of x .

The marginal likelihood can be computed using the sum rule (1.7) or the law of total probability (1.12),

$$p(y) = \int_{X(\Omega)} p(y | x) \cdot p(x) dx. \quad (1.46)$$

Note, however, that the marginal likelihood is simply normalizing the conditional distribution to integrate to one, and therefore a constant with respect to x . For this reason, $p(y)$ is commonly called the *normalizing constant*.

Example 1.19: Plausible inferences

Let us confirm our intuition from the above example. The logical implication “If it is raining, the ground is wet.” (denoted $R \rightarrow W$) can be succinctly expressed as $\mathbb{P}(W | R) = 1$. Since $\mathbb{P}(W) \leq 1$, we know that

$$\mathbb{P}(R | W) = \frac{\mathbb{P}(W | R) \cdot \mathbb{P}(R)}{\mathbb{P}(W)} = \frac{\mathbb{P}(R)}{\mathbb{P}(W)} \geq \mathbb{P}(R).$$

That is, observing that the ground is wet makes it more likely to be raining. From $\mathbb{P}(R | W) \geq \mathbb{P}(R)$ we know $\mathbb{P}(\bar{R} | W) \leq \mathbb{P}(\bar{R})$,¹⁴ which leads us to follow that

$$\mathbb{P}(W | \bar{R}) = \frac{\mathbb{P}(\bar{R} | W) \cdot \mathbb{P}(W)}{\mathbb{P}(\bar{R})} \leq \mathbb{P}(W),$$

that is, having observed it not to be raining made the ground less likely to be wet.

¹⁴ since $\mathbb{P}(\bar{X}) = 1 - \mathbb{P}(X)$

Example 1.19 is called a *plausible inference* because the observation of \bar{R} does not completely determine the truth value of \bar{W} , and hence, does not permit logical inference. In the case, however, that logical inference is permitted, it coincides with probabilistic inference.

Example 1.20: Logical inferences

For example, if we were to observe that the ground is not wet, then logical inference implies that it must not be raining: $\bar{W} \rightarrow \bar{R}$. This is called the *contrapositive* of $R \rightarrow W$.

Indeed, by probabilistic inference, we obtain analogously

$$\mathbb{P}(R \mid \overline{W}) = \frac{\mathbb{P}(\overline{W} \mid R) \cdot \mathbb{P}(R)}{\mathbb{P}(\overline{W})} = \frac{(1 - \mathbb{P}(W \mid R)) \cdot \mathbb{P}(R)}{\mathbb{P}(\overline{W})} = 0.$$

as $\mathbb{P}(W \mid R) = 1$

Observe that a logical inference does not depend on the prior $\mathbb{P}(R)$: Even if the prior was $\mathbb{P}(R) = 1$ in Example 1.20, after observing that the ground is not wet, we are forced to conclude that it is not raining to maintain logical consistency. The examples highlight that while *logical* inference does not require the notion of a prior, plausible (*probabilistic!*) inference does.

1.2.1 Where do priors come from?

Bayes' rule necessitates the specification of a prior $p(x)$. Different priors can lead to the deduction of dramatically different posteriors, as one can easily see by considering the extreme cases of a prior that is a point density at $x = x_0$ and a prior that is “uniform” over \mathbb{R}^n .¹⁵ In the former case, the posterior will be a point density at x_0 regardless of the likelihood. In other words, no evidence can alter the “prior belief” the learner ascribed to x . In the latter case, the learner has “no prior belief”, and therefore the posterior will be proportional to the likelihood. Both steps of probabilistic inference are perfectly valid, though one might debate which prior is more reasonable.

Someone who follows the Bayesian interpretation of probability might argue that everything is conditional, meaning that the prior is simply a posterior of all former observations. While this might seem natural (“my world view from today is the combination of my world view from yesterday and the observations I made today”), this lacks an explanation for “the first day”. Someone else who is more inclined towards the frequentist interpretation might also object to the existence of a prior belief altogether, arguing that a prior is *subjective* and therefore not a valid or desirable input to a learning algorithm. Put differently, a frequentist “has the belief not to have any belief”. This is perfectly compatible with probabilistic inference, as long as the prior is chosen to be *noninformative*:

$$p(x) \propto \text{const.} \quad (1.47)$$

Choosing a noninformative prior in the absence of any evidence is known as the *principle of indifference* or the *principle of insufficient reason*, which dates back to the famous mathematician Pierre-Simon Laplace.

¹⁵ The latter is not a valid probability distribution, but we can still derive meaning from the posterior as we discuss in Remark 1.22.

Example 1.21: Why be indifferent?

Consider a criminal trial with three suspects, A, B, and C. The collected evidence shows that suspect C can not have committed the crime, however it does not yield any information about suspects A and B. Clearly, any distribution respecting the data must assign zero probability of having committed the crime to suspect C. However, any distribution interpolating between $(1, 0, 0)$ and $(0, 1, 0)$ respects the data. The principle of indifference suggests that the desired distribution is $(\frac{1}{2}, \frac{1}{2}, 0)$, and indeed, any alternative distribution seems unreasonable.

Remark 1.22: Noninformative and improper priors

It is not necessarily required that the prior $p(x)$ is a valid distribution (i.e., integrates to 1). Consider for example, the noninformative prior $p(x) \propto \mathbb{1}\{x \in I\}$ where $I \subseteq \mathbb{R}^n$ is an infinitely large interval. Such a prior which is not a valid distribution is called an *improper prior*. We can still derive meaning from the posterior of a given likelihood and (improper) prior as long as the posterior is a valid distribution.

Laplace's principle of indifference can be generalized to cases where *some* evidence is available. The *maximum entropy principle*, originally proposed by Jaynes (1968), states that one should choose as prior from all possible distributions that are *consistent* with prior knowledge, the one that makes the *least* "additional assumptions", i.e., is the least "informative". In philosophy, this principle is known as *Occam's razor* or the *principle of parsimony*. The "informativeness" of a distribution p is quantified by its *entropy* which is defined as

$$H[p] \doteq \mathbb{E}_{x \sim p}[-\log p(x)]. \quad (1.48)$$

The more concentrated p is, the less is its entropy; the more diffuse p is, the greater is its entropy.¹⁶

In the absence of any prior knowledge, the uniform distribution has the highest entropy,¹⁷ and hence, the maximum entropy principle suggests a noninformative prior (as does Laplace's principle of indifference). In contrast, if the evidence perfectly determines the value of x , then the only consistent explanation is the point density at x . The maximum entropy principle characterizes a reasonable choice of prior for these two extreme cases and all cases in between. Bayes' rule can in fact be derived as a consequence of the maximum entropy principle in the sense that the posterior is the least "informative" distribution among all distributions that are consistent with the prior and the ob-

¹⁶ We give a thorough introduction to entropy in Section 5.4.

¹⁷ This only holds true when the set of possible outcomes of x finite (or a bounded continuous interval), as in this case, the noninformative prior is a proper distribution — the uniform distribution. In the "infinite case", there is no uniform distribution and the noninformative prior can be attained from the maximum entropy principle as the limiting solution as the number of possible outcomes of x is increased.

servations (??).

Problem 5.7

1.2.2 Conjugate Priors

If the prior $p(x)$ and posterior $p(x | y)$ are of the same family of distributions, the prior is called a *conjugate prior* to the likelihood $p(y | x)$. This is a very desirable property, as it allows us to recursively apply the same learning algorithm implementing probabilistic inference. We will see in Chapter 2 that under some conditions the Gaussian is *self-conjugate*. That is, if we have a Gaussian prior and a Gaussian likelihood then our posterior will also be Gaussian. This will provide us with the first *efficient* implementation of probabilistic inference.

Example 1.23: Conjugacy of beta and binomial distribution

As an example for conjugacy, we will show that the beta distribution is a conjugate prior to a binomial likelihood. Recall the PMF of the binomial distribution

$$\text{Bin}(k; n, \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k} \quad (1.49)$$

and the PDF of the beta distribution,

$$\text{Beta}(\theta; \alpha, \beta) \propto \theta^{\alpha-1} (1 - \theta)^{\beta-1}, \quad (1.50)$$

We assume the prior $\theta \sim \text{Beta}(\alpha, \beta)$ and likelihood $k | \theta \sim \text{Bin}(n, \theta)$. Let $n_H = k$ be the number of heads and $n_T = n - k$ the number of tails in the binomial trial k . Then,

$$\begin{aligned} p(\theta | k) &\propto p(k | \theta)p(\theta) \\ &\propto \theta^{n_H} (1 - \theta)^{n_T} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \\ &= \theta^{\alpha+n_H-1} (1 - \theta)^{\beta+n_T-1}. \end{aligned}$$

using Bayes' rule (1.45)

Thus, $\theta | k \sim \text{Beta}(\alpha + n_H, \beta + n_T)$.

This same conjugacy can be shown for the multivariate generalization of the beta distribution, the *Dirichlet distribution*, and the multivariate generalization of the binomial distribution, the *multinomial distribution*.

1.2.3 Tractable Inference with the Normal Distribution

Using arbitrary distributions for learning and inference is computationally very expensive when the number of dimensions is large — even in the discrete setting. For example, computing marginal distributions using the sum rule yields an exponentially long sum in the

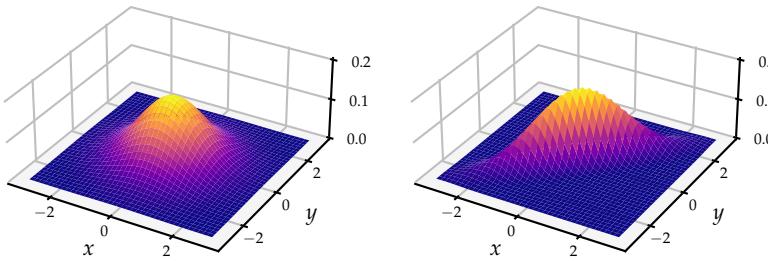
size of the random vector. Similarly, the normalizing constant of the conditional distribution is a sum of exponential length. Even to represent any discrete joint probability distribution requires space that is exponential in the number of dimensions (cf. Figure 1.5).

One strategy to get around this computational blowup is to restrict the class of distributions. Gaussians are a popular choice for this purpose since they have extremely useful properties: they have a compact representation and — as we will see in Chapter 2 — they allow for closed-form probabilistic inference.

In Equation (1.5), we have already seen the PDF of the univariate Gaussian distribution. A random vector \mathbf{X} in \mathbb{R}^n is *normally distributed*, $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, if its PDF is

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \doteq \frac{1}{\sqrt{\det(2\pi\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (1.51)$$

where $\boldsymbol{\mu} \in \mathbb{R}^n$ is the mean vector and $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$ the covariance matrix [\(?\)](#). We call $\boldsymbol{\Lambda} \doteq \boldsymbol{\Sigma}^{-1}$ the *precision matrix*. \mathbf{X} is also called a *Gaussian random vector* (GRV). $\mathcal{N}(\mathbf{0}, \mathbf{I})$ is the multivariate *standard normal distribution*. We call a Gaussian *isotropic* if its covariance matrix is of the form $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$ for some $\sigma^2 \in \mathbb{R}$. In this case, the sublevel sets of the PDF are perfect spheres as can be seen in Figure 1.6.



Note that a Gaussian can be represented using only $O(n^2)$ parameters. In the case of a diagonal covariance matrix, which corresponds to n independent univariate Gaussians [\(?\)](#), we just need $O(n)$ parameters.

In Equation (1.51), we assume that the covariance matrix $\boldsymbol{\Sigma}$ is invertible, i.e., does not have the eigenvalue 0. This is not a restriction since it can be shown that a covariance matrix has a zero eigenvalue if and only if there exists a deterministic linear relationship between some variables in the joint distribution [\(?\)](#). As we have already seen that a

X_1	\dots	X_{n-1}	X_n	$\mathbb{P}(X_{1:n})$
0	\dots	0	0	0.01
0	\dots	0	1	0.001
0	\dots	1	0	0.213
\vdots	\vdots	\vdots	\vdots	\vdots
1	\dots	1	1	0.0003

Figure 1.5: A table representing a joint distribution of n binary random variables. The table has 2^n rows. The number of parameters is $2^n - 1$ since the final probability is determined by all other probabilities as they must sum to one.

Problem 1.11

Figure 1.6: Shown are the PDFs of two-dimensional Gaussians with mean $\mathbf{0}$ and covariance matrices

$$\boldsymbol{\Sigma}_1 \doteq \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{\Sigma}_2 \doteq \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$

respectively.

Problem 1.8

Problem 1.6

covariance matrix does not have negative eigenvalues $\textcircled{?}$, this ensures that Σ and Λ are positive definite.¹⁸

An important property of the normal distribution is that it is closed under marginalization and conditioning.

Theorem 1.24 (Marginal and conditional distribution). $\textcircled{?}$ Consider the Gaussian random vector \mathbf{X} and fix index sets $A \subseteq [n]$ and $B \subseteq [n]$. Then, we have that for any such marginal distribution,

$$\mathbf{X}_A \sim \mathcal{N}(\boldsymbol{\mu}_A, \boldsymbol{\Sigma}_{AA}), \quad (1.52)$$

and that for any such conditional distribution,

$$\mathbf{X}_A | \mathbf{X}_B = \mathbf{x}_B \sim \mathcal{N}(\boldsymbol{\mu}_{A|B}, \boldsymbol{\Sigma}_{A|B}) \quad \text{where} \quad (1.53a)$$

$$\boldsymbol{\mu}_{A|B} \doteq \boldsymbol{\mu}_A + \boldsymbol{\Sigma}_{AB}\boldsymbol{\Sigma}_{BB}^{-1}(\mathbf{x}_B - \boldsymbol{\mu}_B), \quad (1.53b)$$

$$\boldsymbol{\Sigma}_{A|B} \doteq \boldsymbol{\Sigma}_{AA} - \boldsymbol{\Sigma}_{AB}\boldsymbol{\Sigma}_{BB}^{-1}\boldsymbol{\Sigma}_{BA}. \quad (1.53c)$$

Theorem 1.24 provides a closed-form characterization of probabilistic inference for the case that random variables are jointly Gaussian. We will discuss in Chapter 2, how this can be turned into an efficient inference algorithm.

Observe that upon inference, the variance can only shrink! Moreover, how much the variance is reduced depends purely on *where* the observations are made (i.e., the choice of B) but not on *what* the observations are. In contrast, the posterior mean $\boldsymbol{\mu}_{A|B}$ depends affinely on $\boldsymbol{\mu}_B$. These are special properties of the Gaussian and do not generally hold true for other distributions.

It can be shown that Gaussians are additive and closed under affine transformations $\textcircled{?}$. The closedness under affine transformations (1.78) implies that a Gaussian $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is equivalently characterized as

$$\mathbf{X} = \boldsymbol{\Sigma}^{1/2}\mathbf{Y} + \boldsymbol{\mu}. \quad (1.54)$$

where $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\boldsymbol{\Sigma}^{1/2}$ is the square root of $\boldsymbol{\Sigma}$.¹⁹ Importantly, this implies together with Theorem 1.24 and additivity (1.79) that:

Any affine transformation of a Gaussian random vector is a Gaussian random vector.

A consequence of this is that given any jointly Gaussian random vectors \mathbf{X}_A and \mathbf{X}_B , \mathbf{X}_A can be expressed as an affine function of \mathbf{X}_B with added independent Gaussian noise. Formally, we define

$$\mathbf{X}_A \doteq \mathbf{A}\mathbf{X}_B + \mathbf{b} + \boldsymbol{\varepsilon} \quad \text{where} \quad (1.55a)$$

Problem 1.4

¹⁸ The inverse of a positive definite matrix is also positive definite.

Problem 1.9

By $\boldsymbol{\mu}_A$ we denote $[\mu_{i_1}, \dots, \mu_{i_k}]$ where $A = \{i_1, \dots, i_k\}$. $\boldsymbol{\Sigma}_{AA}$ is defined analogously.

Here, $\boldsymbol{\mu}_A$ characterizes the prior belief and $\boldsymbol{\Sigma}_{AB}\boldsymbol{\Sigma}_{BB}^{-1}(\mathbf{x}_B - \boldsymbol{\mu}_B)$ represents “how different” \mathbf{x}_B is from what was expected.

Problem 1.10

¹⁹ More details on the square root of a symmetric and positive definite matrix can be found in Appendix A.2.

$$\mathbf{A} \doteq \boldsymbol{\Sigma}_{AB}\boldsymbol{\Sigma}_{BB}^{-1}, \quad (1.55b)$$

$$\mathbf{b} \doteq \boldsymbol{\mu}_A - \boldsymbol{\Sigma}_{AB}\boldsymbol{\Sigma}_{BB}^{-1}\boldsymbol{\mu}_B, \quad (1.55c)$$

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{A|B}). \quad (1.55d)$$

It directly follows from the closedness of Gaussians under affine transformations (1.78) that the characterization of \mathbf{X}_A via Equation (1.55) is equivalent to $\mathbf{X}_A \sim \mathcal{N}(\boldsymbol{\mu}_A, \boldsymbol{\Sigma}_{AA})$, and hence, *any* Gaussian \mathbf{X}_A can be modeled as a so-called *conditional linear Gaussian*, i.e., an affine function of another Gaussian \mathbf{X}_B with additional independent Gaussian noise. We will use this fact frequently to represent Gaussians in a compact form.

1.3 Supervised Learning and Point Estimates

Throughout the first part of this manuscript, we will focus mostly on the *supervised learning* problem where we want to learn a function

$$f^* : \mathcal{X} \rightarrow \mathcal{Y}$$

from labeled training data. That is, we are given a collection of labeled examples, $\mathcal{D}_n \doteq \{(x_i, y_i)\}_{i=1}^n$, where the $x_i \in \mathcal{X}$ are *inputs* and the $y_i \in \mathcal{Y}$ are *outputs* (called *labels*), and we want to find a function \hat{f} that best-approximates f^* . It is common to choose \hat{f} from a parameterized *function class* $\mathcal{F}(\Theta)$, where each function f_θ is described by some parameters $\theta \in \Theta$.

Remark 1.25: What this manuscript is about and not about

As illustrated in Figure 1.7, the restriction to a function class leads to two sources of error: the *estimation error* of having “incorrectly” determined \hat{f} within the function class, and the *approximation error* of the function class itself. Choosing a “good” function class / architecture with small approximation error is therefore critical for any practical application of machine learning. We will discuss various function classes, from linear models to deep neural networks, however, determining the “right” function class will not be the focus of this manuscript. To keep the exposition simple, we will assume in the following that $f^* \in \mathcal{F}(\Theta)$ with parameters $\theta^* \in \Theta$.

Instead, we will focus on the problem of estimation/inference within a given function class. We will see that inference in smaller function classes is often more computationally efficient since the search space is smaller or — in the case of Gaussians — has a known tractable structure. On the other hand, larger function

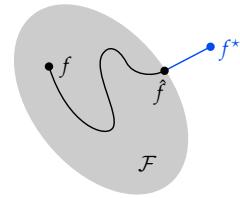


Figure 1.7: Illustration of **estimation error** and **approximation error**. f^* denotes the true function and \hat{f} is the best approximation from the function class \mathcal{F} . We do not specify here, how one could quantify “error”. For more details, see Appendix A.3.5.

classes are more expressive and therefore can typically better approximate the ground truth f^* .

We differentiate between the task of *regression* where $\mathcal{Y} \doteq \mathbb{R}^k$ ²⁰ and the task of *classification* where $\mathcal{Y} \doteq \mathcal{C}$ and \mathcal{C} is an m -element set of classes. In other words, regression is the task of predicting a continuous label, whereas classification is the task of predicting a discrete class label. These two tasks are intimately related: in fact, we can think of classification tasks as a regression problem where we learn a probability distribution over class labels. In this regression problem, $\mathcal{Y} \doteq \Delta^{\mathcal{C}}$ where $\Delta^{\mathcal{C}}$ denotes the set of all probability distributions over the set of classes \mathcal{C} which is an $(m - 1)$ -dimensional convex polytope in the m -dimensional space of probabilities $[0, 1]^m$ (cf. Appendix A.1.2).

For now, let us stick to the regression setting. We will assume that the observations are noisy, that is, $y_i \stackrel{\text{iid}}{\sim} p(\cdot | x_i, \theta^*)$ for some *known* conditional distribution $p(\cdot | x_i, \theta)$ but *unknown* parameter θ^* .²¹ Our assumption can equivalently be formulated as

$$y_i = \underbrace{f_{\theta}(x_i)}_{\text{signal}} + \underbrace{\varepsilon_i(x_i)}_{\text{noise}} \quad (1.56)$$

where $f_{\theta}(x_i)$ is the mean of $p(\cdot | x_i, \theta)$ and $\varepsilon_i(x_i) = y_i - f_{\theta}(x_i)$ is some independent zero-mean noise, for example (but not necessarily) Gaussian.²² When the noise distribution may depend on x_i , the noise is said to be *heteroscedastic* and otherwise the noise is called *homoscedastic*.

1.3.1 Maximum Likelihood Estimation

A common approach to finding \hat{f} is to select the model $f \in \mathcal{F}(\Theta)$ under which the training data is most likely. This is called the *maximum likelihood estimate* (or MLE):

$$\begin{aligned} \hat{\theta}_{\text{MLE}} &\doteq \arg \max_{\theta \in \Theta} p(y_{1:n} | x_{1:n}, \theta) \\ &= \arg \max_{\theta \in \Theta} \prod_{i=1}^n p(y_i | x_i, \theta). \end{aligned} \quad (1.57)$$

Such products of probabilities are often numerically unstable, which is why one typically takes the logarithm:

$$= \arg \max_{\theta \in \Theta} \sum_{i=1}^n \underbrace{\log p(y_i | x_i, \theta)}_{\text{log-likelihood}}. \quad (1.58)$$

We will denote the *negative log-likelihood* by $\ell_{\text{nll}}(\theta; \mathcal{D}_n)$.

²⁰The labels are usually scalar, so $k = 1$.

²¹The case where the labels are deterministic is the special case of $p(\cdot | x_i, \theta^*)$ being a point density at $f^*(x_i)$.

²²It is crucial that the assumed noise distribution accurately reflects the noise of the data. For example, using a (light-tailed) Gaussian noise model in the presence of heavy-tailed noise will fail! We discuss the distinction between light and heavy tails in Appendix A.3.2.

using the independence of the training data (1.56)

The MLE is often used in practice due to its desirable asymptotic properties as the sample size n increases. We give a brief summary here and provide additional background and definitions in Appendix A.3. To give any guarantees on the convergence of the MLE, we necessarily need to assume that θ^* is identifiable.²³ If additionally, ℓ_{nll} is “well-behaved” then standard results say that the MLE is *consistent* and *asymptotically normal* (Van der Vaart, 2000):

$$\hat{\theta}_{\text{MLE}} \xrightarrow{\mathbb{P}} \theta^* \quad \text{and} \quad \hat{\theta}_{\text{MLE}} \xrightarrow{\mathcal{D}} \mathcal{N}(\theta^*, S_n) \quad \text{as } n \rightarrow \infty. \quad (1.59)$$

Here, we denote by S_n the asymptotic variance of the MLE which can be understood as measuring the “quality” of the estimate.²⁴ This implies in some sense that the MLE is asymptotically unbiased. Moreover, the MLE can be shown to be *asymptotically efficient* which is to say that there exists no other consistent estimator with a “smaller” asymptotic variance.²⁵

The situation is quite different in the finite sample regime. Here, the MLE need not be unbiased, and it is susceptible to *overfitting* to the (finite) training data as we discuss in more detail in Appendix A.3.5.

1.3.2 Using Priors: Maximum a Posteriori Estimation

We can incorporate prior assumptions about the parameters θ^* into the estimation procedure. One approach of this kind is to find the mode of the posterior distribution, called the *maximum a posteriori estimate* (or MAP estimate):

$$\hat{\theta}_{\text{MAP}} \doteq \arg \max_{\theta \in \Theta} p(\theta | \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) \quad (1.60)$$

$$= \arg \max_{\theta \in \Theta} p(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}, \theta) \cdot p(\theta) \quad (1.61) \quad \text{by Bayes' rule (1.45)}$$

$$= \arg \max_{\theta \in \Theta} \log p(\theta) + \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \theta) \quad (1.62) \quad \text{taking the logarithm}$$

$$= \arg \min_{\theta \in \Theta} \underbrace{-\log p(\theta)}_{\text{regularization}} + \underbrace{\ell_{\text{nll}}(\theta; \mathcal{D}_n)}_{\text{quality of fit}}. \quad (1.63)$$

Here, the *log-prior* $\log p(\theta)$ acts as a regularizer. Common regularizers are given, for example, by

- $p(\theta) = \mathcal{N}(\theta; \mathbf{0}, (2\lambda)^{-1} \mathbf{I})$ which yields $-\log p(\theta) = \lambda \|\theta\|_2^2 + \text{const}$,
- $p(\theta) = \text{Laplace}(\theta; \mathbf{0}, \lambda^{-1})$ which yields $-\log p(\theta) = \lambda \|\theta\|_1 + \text{const}$,
- a uniform prior (cf. Section 1.2.1) for which the MAP is equivalent to the MLE. In other words, the MLE is merely the mode of the posterior distribution under a uniform prior.

The Gaussian and Laplace regularizers act as simplicity biases, preferring simpler models over more complex ones, which empirically

²³ That is, $\theta^* \neq \theta \implies f^* \neq f_\theta$ for any $\theta \in \Theta$. In words, there is no other parameter θ that yields the same function f_θ as θ^* .

²⁴ A “smaller” variance means that we can be more confident that the MLE is close to the true parameter.

²⁵ see Appendix A.3.4.

tends to reduce the risk of overfitting. However, one may also encode more nuanced information about the (assumed) structure of θ^* into the prior.

An alternative way of encoding a prior is by restricting the function class to some $\tilde{\Theta} \subset \Theta$, for example to rotation- and translation-invariant models as often done when the inputs are images. This effectively sets $p(\theta) = 0$ for all $\theta \in \Theta \setminus \tilde{\Theta}$ but is better suited for numerical optimization than to impose this constraint directly on the prior.

Encoding prior assumptions into the function class or into the parameter estimation can accelerate learning and improve generalization performance dramatically, yet importantly, incorporating a prior can also inhibit learning in case the prior is “wrong”. For example, when the learning task is to differentiate images of cats from images of dogs, consider the (stupid) prior that only permits models that exclusively use the upper-left pixel for prediction. No such model will be able to solve the task, and therefore starting from this prior makes the learning problem effectively unsolvable which illustrates that priors have to be chosen with care.

1.3.3 When does the prior matter?

We have seen that the MLE has desirable asymptotic properties, and that MAP estimation can be seen as a regularized MLE where the type of regularization is encoded by the prior. Is it possible to derive similar asymptotic results for the MAP estimate?

To answer this question, we will look at the asymptotic effect of the prior on the posterior more generally. *Doob’s consistency theorem* states that assuming parameters are identifiable,²⁶ there exists $\tilde{\Theta} \subseteq \Theta$ with $p(\tilde{\Theta}) = 1$ such that the posterior is consistent for any $\theta^* \in \tilde{\Theta}$ (Doob, 1949; Miller, 2016, 2018):

$$\theta \mid \mathcal{D}_n \xrightarrow{\mathbb{P}} \theta^* \quad \text{as } n \rightarrow \infty. \quad (1.64)$$

In words, Doob’s consistency theorem tells us that for *any* prior distribution, the posterior is guaranteed to converge to a point density in the (small) neighborhood $\theta^* \in B$ of the true parameter as long as $p(B) > 0$.²⁷ We call such a prior a *well-specified prior*.

Remark 1.26: Cromwell’s rule

In the case where $|\Theta|$ is finite, Doob’s consistency theorem strongly suggests that the prior should not assign 0 probability (or probability 1 for that matter) to any individual parameter $\theta \in \Theta$, unless we know with certainty that $\theta^* \neq \theta$. This is called *Cromwell’s rule*,

²⁶This is akin to the assumption required for consistency of the MLE, cf. Equation (1.59).

²⁷ B can for example be a ball of radius ϵ around θ^* (with respect to some geometry of Θ).

and a prior obeying by this rule is always well-specified.

Under the same assumption that the prior is well-specified (and regularity conditions²⁸), the *Bernstein-von Mises theorem*, which was first discovered by Pierre-Simon Laplace in the early 19th century, establishes the asymptotic normality of the posterior distribution (Van der Vaart, 2000; Miller, 2016):

$$\theta \mid \mathcal{D}_n \xrightarrow{\mathcal{D}} \mathcal{N}(\theta^*, S_n) \quad \text{as } n \rightarrow \infty \quad (1.65)$$

and where S_n is the same as the asymptotic variance of the MLE.²⁹

These results link probabilistic inference to maximum likelihood estimation in the asymptotic limit of infinite data. Intuitively, in the limit of infinite data, the prior is “overwhelmed” by the observations and the posterior becomes equivalent to the limiting distribution of the MLE.³⁰ One can interpret the regime of infinite data as the regime where computational resources and time are unlimited and plausible inferences evolve into logical inferences. This transition signifies a shift from the realm of uncertainty to that of certainty. The importance of the prior surfaces precisely in the non-asymptotic regime where plausible inferences are necessary due to *limited computational resources and limited time*.

1.3.4 Estimation vs Inference

You can interpret a single parameter vector $\theta \in \Theta$ as “one possible explanation” of the data. Maximum likelihood and maximum a posteriori estimation are examples of *estimation* algorithms which return a one such parameter vector — called a *point estimate*. That is, given the training set \mathcal{D}_n , they return a single parameter vector $\hat{\theta}_n$. We give a more detailed account of estimation in Appendix A.3.

Example 1.27: Point estimates and invalid logical inferences

To see why point estimates can be problematic, recall Example 1.19.

We have seen that the logical implication

“If it is raining, the ground is wet.”

can be expressed as $\mathbb{P}(W \mid R) = 1$. Observing that “The ground is wet.” does not permit logical inference, yet, the maximum likelihood estimate of R is $\hat{R}_{\text{MLE}} = 1$. This is logically inconsistent since there might be other *explanations* for the ground to be wet, such as a sprinkler! With only a finite sample (say independently observing n times that the ground is wet), we cannot rule out with

²⁸ These regularity conditions are akin to the assumptions required for asymptotic normality of the MLE, cf. Equation (1.59).

²⁹ This has also been called the “Bayesian central limit theorem”, which is a bit of a misnomer since the theorem also applies to likelihoods (when the prior is noninformative) which are often used in frequentist statistics.

³⁰ More examples and discussion can be found in section 17.8 of Le Cam (1986), chapter 8 of Le Cam and Yang (2000), chapter 10 of Van der Vaart (2000), and in Tanner (1991).

certainty that the ground is wet for other reasons than rain.

In practice, we never observe an infinite amount of data. Example 1.27 demonstrates that on a finite sample, point estimates may perform invalid logical inferences, and can therefore lure us into a false sense of certainty.

Remark 1.28: MLE and MAP are approximations of inference

The MLE and MAP estimate can be seen as a naive approximation of probabilistic inference, represented by a point density which “collapses” all probability mass at the mode of the posterior distribution. This can be a relatively decent — even if overly simple — approximation when the distribution is unimodal, symmetric, and light-tailed as in Figure 1.8, but is usually a very poor approximation for practical posteriors that are complex and multimodal.

In this manuscript, we will focus mainly on algorithms for *probabilistic inference* which compute or approximate the distribution $p(\theta | \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$ over parameters. Returning a distribution over parameters is natural since this acknowledges that given a finite sample with noisy observations, more than one parameter vector can explain the data.

1.3.5 Probabilistic Inference and Prediction

The prior distribution $p(\theta)$ can be interpreted as the degree of our belief that the model parameterized by θ “describes the (previously seen) data best”. The likelihood captures how likely the training data is under a particular model:

$$p(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}, \theta) = \prod_{i=1}^n p(y_i | \mathbf{x}_i, \theta). \quad (1.66)$$

The posterior then represents our belief about the best model after seeing the training data. Using Bayes’ rule (1.45), we can write it as³¹

$$p(\theta | \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \frac{1}{Z} p(\theta) \prod_{i=1}^n p(y_i | \mathbf{x}_i, \theta) \quad \text{where} \quad (1.67a)$$

$$Z \doteq \int_{\Theta} p(\theta) \prod_{i=1}^n p(y_i | \mathbf{x}_i, \theta) d\theta \quad (1.67b)$$

is the *normalizing constant*. We refer to this process of learning a model from data as *learning*. We can then use our learned model for *prediction* at a new input \mathbf{x}^* by conditioning on θ ,

$$p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \int_{\Theta} p(y^*, \theta | \mathbf{x}^*, \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) d\theta$$

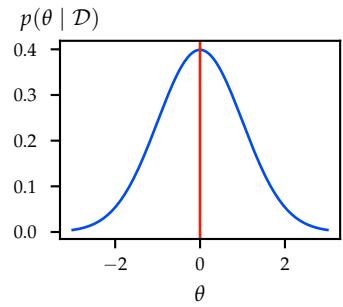


Figure 1.8: The MLE/MAP are point estimates at the mode $\hat{\theta}$ of the posterior distribution $p(\theta | \mathcal{D})$.

³¹ We generally assume that

$$p(\theta | \mathbf{x}_{1:n}) = p(\theta).$$

For our purposes, you can think of the inputs $\mathbf{x}_{1:n}$ as fixed deterministic parameters, but one can also consider inputs drawn from a distribution over \mathcal{X} .

by the sum rule (1.7)

$$= \int_{\Theta} p(y^* | x^*, \theta) \cdot p(\theta | x_{1:n}, y_{1:n}) d\theta. \quad (1.68)$$

Here, the distribution over models $p(\theta | x_{1:n}, y_{1:n})$ is called the *posterior* and the distribution over predictions $p(y^* | x^*, x_{1:n}, y_{1:n})$ is called the *predictive posterior*. The predictive posterior quantifies our posterior uncertainty about the “prediction” y^* , however, since this is typically a complex distribution, it is difficult to communicate this uncertainty to a human. One statistic that can be used for this purpose is the smallest set $\mathcal{C}_\delta(x^*) \subseteq \mathbb{R}$ for a fixed $\delta \in (0, 1)$ such that

$$\mathbb{P}(y^* \in \mathcal{C}_\delta(x^*) | x^*, x_{1:n}, y_{1:n}) \geq 1 - \delta. \quad (1.69)$$

That is, we believe with “confidence” at least $1 - \delta$ that the true value of y^* lies in $\mathcal{C}_\delta(x^*)$. Such a set $\mathcal{C}_\delta(x^*)$ is called a *credible set*.

We have seen here that the tasks of learning and prediction are intimately related. Indeed, “prediction” can be seen in many ways as a natural by-product of “reasoning” (i.e., probabilistic inference), where we evaluate the likelihood of outcomes given our learned explanations for the world. This intuition can be read off directly from Equation (1.68) where $p(y^* | x^*, \theta)$ corresponds to the likelihood of an outcome given the explanation θ and $p(\theta | x_{1:n}, y_{1:n})$ corresponds to our inferred belief about the world. We will see many more examples of this link between probabilistic inference and prediction throughout this manuscript.

The high-dimensional integrals of Equations (1.67b) and (1.68) are typically intractable, and represent the main computational challenge in probabilistic inference. Throughout the first part of this manuscript, we will describe settings where exact inference is tractable, as well as modern approximate inference algorithms that can be used when exact inference is intractable.

1.3.6 Recursive Probabilistic Inference and Memory

We have already alluded to the fact that probabilistic inference has a recursive structure, which lends itself to continual learning and which often leads to efficient algorithms. Let us denote by

$$p^{(t)}(\theta) \doteq p(\theta | x_{1:t}, y_{1:t}) \quad (1.70)$$

the posterior after the first t observations with $p^{(0)}(\theta) = p(\theta)$. Now, suppose that we have already computed $p^{(t)}(\theta)$ and observe y_{t+1} . We can recursively update the posterior as follows,

$$\begin{aligned} p^{(t+1)}(\theta) &= p(\theta | y_{1:t+1}) \\ &\propto p(\theta | y_{1:t}) \cdot p(y_{t+1} | \theta, y_{1:t}) \end{aligned}$$

by the product rule (1.11) and
 $y^* \perp x_{1:n}, y_{1:n} | \theta$

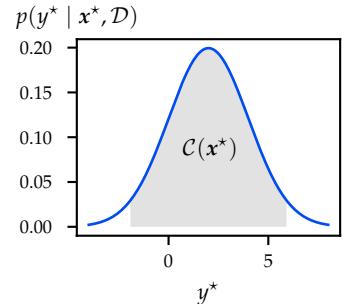


Figure 1.9: Example of a 95% credible set at x^* where the predictive posterior is Gaussian with mean $\mu(x^*)$ and standard deviation $\sigma(x^*)$. In this case, the gray area integrates to ≈ 0.95 for

$$\mathcal{C}_{0.05}(x^*) = [\mu(x^*) \pm 1.96\sigma(x^*)].$$

using Bayes' rule (1.45)

$$= p^{(t)}(\theta) \cdot p(y_{t+1} | \theta). \quad (1.71) \quad \text{using } y_{t+1} \perp y_{1:t} | \theta, \text{ see Figure 2.3}$$

Intuitively, the posterior distribution at time t “absorbs” or “summarizes” all seen data.

By unrolling the recursion of Equation (1.71), we see that regardless of the philosophical interpretation of probability, probabilistic inference is a fundamental mechanism of learning. Even the MLE which performs naive approximate inference without a prior (i.e., a uniform prior), is based on $p^{(n)}(\theta) \propto p(y_{1:n} | x_{1:n}, \theta)$ which is the result of n individual plausible inferences, where the $(t + 1)$ -st inference uses the posterior of the t -th inference as its prior.

So far we have been considering the supervised learning setting, where all data is available a-priori. However, by sequentially obtaining the new posterior and replacing our prior, we can also perform probabilistic inference as data arrives online (i.e., in “real-time”). This is analogous to recursive logical inference where derived consequences are repeatedly added to the set of propositions to derive new consequences. This also highlights the intimate connection between “reasoning” and “memory”. Indeed, the posterior distribution $p^{(t)}(\theta)$ can be seen as a form of memory that evolves with time t .

1.4 Outlook: Decision Theory

How can we use our predictions to make concrete decisions under uncertainty? We will study this question extensively in Part II of this manuscript, but briefly introduce some fundamental concepts here. Making decisions using a probabilistic model $p(y | x)$ of output $y \in \mathcal{Y}$ given input $x \in \mathcal{X}$, such as the ones we have discussed in the previous section, is commonly formalized by

- a set of possible actions \mathcal{A} , and
- a reward function $r(y, a) \in \mathbb{R}$ that computes the reward or utility of taking action $a \in \mathcal{A}$, assuming the true output is $y \in \mathcal{Y}$.

Standard decision theory recommends picking the action with the largest expected utility:

$$a^*(x) \doteq \arg \max_{a \in \mathcal{A}} \mathbb{E}_{y|x}[r(y, a)]. \quad (1.72)$$

Here, a^* is called the *optimal decision rule* because, under the given probabilistic model, no other rule can yield a higher expected utility.

Let us consider some examples of reward functions and their corresponding optimal decisions:

Example 1.29: Reward functions

Under the decision rule from Equation (1.72), different reward functions r can lead to different decisions. Let us examine two reward functions for the case where $\mathcal{Y} = \mathcal{A} = \mathbb{R}$ (?)

- Alternatively to considering r as a reward function, we can interpret $-r$ as the loss of taking action a when the true output is y . If our goal is for our actions a to “mimic” the output y , a natural choice is the squared loss, $-r(y, a) = (y - a)^2$. It turns out that under the squared loss, the optimal decision is simply the mean: $a^*(x) = \mathbb{E}[y | x]$.
- To contrast this, we consider the asymmetric loss,

$$-r(y, a) = c_1 \underbrace{\max\{y - a, 0\}}_{\text{underestimation error}} + c_2 \underbrace{\max\{a - y, 0\}}_{\text{overestimation error}},$$

which penalizes underestimation and overestimation differently.

When $y | x \sim \mathcal{N}(\mu_x, \sigma_x^2)$ then the optimal decision is

$$a^*(x) = \mu_x + \sigma_x \cdot \underbrace{\Phi^{-1}\left(\frac{c_1}{c_1 + c_2}\right)}_{\text{pessimism / optimism}}$$

where Φ is the CDF of the standard normal distribution.³² Note that if $c_1 = c_2$, then the second term vanishes and the optimal decision is the same as under the squared loss. If $c_1 > c_2$, the second term is positive (i.e., *optimistic*) to avoid underestimation, and if $c_1 < c_2$, the second term is negative (i.e., *pessimistic*) to avoid overestimation. We will find these notions of optimism and pessimism to be useful in many decision-making scenarios.

Problem 1.13

While Equation (1.72) describes how to make optimal decisions given a (posterior) probabilistic model, it does not tell us how to learn or improve this model in the first place. That is, these decisions are only optimal under the assumption that we cannot use their outcomes and our resulting observations to update our model and inform future decisions. When we start to consider the effect of our decisions on future data and future posteriors, answering “how do I make optimal decisions?” becomes more complex, and we will study this in Part II on sequential decision-making.

Discussion

In this chapter, we have learned about the fundamental concepts of probabilistic inference. We have seen that probabilistic inference is the

³² Recall that the CDF Φ of the standard normal distribution is a sigmoid with its inverse satisfying

$$\Phi^{-1}(u) \begin{cases} < 0 & \text{if } u < 0.5, \\ = 0 & \text{if } u = 0.5, \\ > 0 & \text{if } u > 0.5. \end{cases}$$

natural extension of logical reasoning to domains with uncertainty. We have also derived the central principle of probabilistic inference, Bayes' rule, which is simple to state but often computationally challenging. In the next part of this manuscript, we will explore settings where exact inference is tractable, as well as modern approaches to approximate probabilistic inference.

Overview of Mathematical Background

We have included brief summaries of the fundamentals of **parameter estimation** (mean estimation in particular) and **optimization** in Appendices A.3 and A.4, respectively, which we will refer back to throughout the manuscript. Appendix A.2 discusses the **correspondence of Gaussians and quadratic forms**. Appendix A.5 comprises a list of useful matrix identities and inequalities.

Problems

1.1. Properties of probability.

Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Derive the following properties of probability from the Kolmogorov axioms:

1. For any $A, B \in \mathcal{A}$, if $A \subseteq B$ then $\mathbb{P}(A) \leq \mathbb{P}(B)$.
2. For any $A \in \mathcal{A}$, $\mathbb{P}(\overline{A}) = 1 - \mathbb{P}(A)$.
3. For any countable set of events $\{A_i \in \mathcal{A}\}_i$,

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) \leq \sum_{i=1}^{\infty} \mathbb{P}(A_i). \quad (1.73)$$

which is called a *union bound*.

1.2. Random walks on graphs.

Let G be a simple connected finite graph. We start at a vertex u of G . At every step, we move to one of the neighbors of the current vertex uniformly at random, e.g., if the vertex has 3 neighbors, we move to one of them, each with probability $1/3$. What is the probability that the walk visits a given vertex v eventually?

1.3. Law of total expectation.

Show that if $\{A_i\}_{i=1}^k$ are a partition of Ω and \mathbf{X} is a random vector,

$$\mathbb{E}[\mathbf{X}] = \sum_{i=1}^k \mathbb{E}[\mathbf{X} | A_i] \cdot \mathbb{P}(A_i). \quad (1.74)$$

1.4. Covariance matrices are positive semi-definite.

Prove that a covariance matrix Σ is always positive semi-definite. That is, all of its eigenvalues are greater or equal to zero, or equivalently, $x^\top \Sigma x \geq 0$ for any $x \in \mathbb{R}^n$.

1.5. Probabilistic inference.

As a result of a medical screening, one of the tests revealed a serious disease in a person. The test has a high accuracy of 99% (the probability of a positive response in the presence of a disease is 99% and the probability of a negative response in the absence of a disease is also 99%). However, the disease is quite rare and occurs only in one person per 10 000. Calculate the probability of the examined person having the identified disease.

1.6. Zero eigenvalues of covariance matrices.

We say that a random vector X in \mathbb{R}^n is not linearly independent if for some $\alpha \in \mathbb{R}^n \setminus \{0\}$, $\alpha^\top X = 0$.

1. Show that if X is not linearly independent, then $\text{Var}[X]$ has a zero eigenvalue.
2. Show that if $\text{Var}[X]$ has a zero eigenvalue, then X is not linearly independent.

Hint: Consider the variance of $\lambda^\top X$ where λ is the eigenvector corresponding to the zero eigenvalue.

Thus, we have shown that $\text{Var}[X]$ has a zero eigenvalue if and only if X is not linearly independent.

1.7. Product of Gaussian PDFs.

Let $\mu_1, \mu_2 \in \mathbb{R}^n$ be mean vectors and $\Sigma_1, \Sigma_2 \in \mathbb{R}^{n \times n}$ be covariance matrices. Prove that

$$\mathcal{N}(x; \mu, \Sigma) \propto \mathcal{N}(x; \mu_1, \Sigma_1) \cdot \mathcal{N}(x; \mu_2, \Sigma_2) \quad (1.75)$$

for some mean vector $\mu \in \mathbb{R}^n$ and covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$. That is, show that the product of two Gaussian PDFs is proportional to the PDF of a Gaussian.

1.8. Independence of Gaussians.

Show that two jointly Gaussian random vectors, X and Y , are independent if and only if X and Y are uncorrelated.

1.9. Marginal / conditional distribution of a Gaussian.

Prove Theorem 1.24. That is, show that

1. every marginal of a Gaussian is Gaussian; and

2. conditioning on a subset of variables of a joint Gaussian is Gaussian

by finding their corresponding PDFs.

Hint: You may use that for matrices Σ and Λ such that $\Sigma^{-1} = \Lambda$,

- if Σ and Λ are symmetric,

$$\begin{aligned} & \begin{bmatrix} \mathbf{x}_A \\ \mathbf{x}_B \end{bmatrix}^\top \begin{bmatrix} \Lambda_{AA} & \Lambda_{AB} \\ \Lambda_{BA} & \Lambda_{BB} \end{bmatrix} \begin{bmatrix} \mathbf{x}_A \\ \mathbf{x}_B \end{bmatrix} \\ &= \mathbf{x}_A^\top \Lambda_{AA} \mathbf{x}_A + \mathbf{x}_A^\top \Lambda_{AB} \mathbf{x}_B + \mathbf{x}_B^\top \Lambda_{BA} \mathbf{x}_A + \mathbf{x}_B^\top \Lambda_{BB} \mathbf{x}_B \\ &= \mathbf{x}_A^\top (\Lambda_{AA} - \Lambda_{AB} \Lambda_{BB}^{-1} \Lambda_{BA}) \mathbf{x}_A + \\ & \quad (\mathbf{x}_B + \Lambda_{BB}^{-1} \Lambda_{BA} \mathbf{x}_A)^\top \Lambda_{BB} (\mathbf{x}_B + \Lambda_{BB}^{-1} \Lambda_{BA} \mathbf{x}_A), \end{aligned}$$

- $\Lambda_{BB}^{-1} = \Sigma_{BB} - \Sigma_{BA} \Sigma_{AA}^{-1} \Sigma_{AB}$,
- $\Lambda_{BB}^{-1} \Lambda_{BA} = -\Sigma_{BA} \Sigma_{AA}^{-1}$.

The final two equations follow from the general characterization of the inverse of a block matrix (Petersen et al., 2008, section 9.1.3).

1.10. Closedness properties of Gaussians.

Recall the notion of a *moment-generating function* (MGF) of a random vector \mathbf{X} in \mathbb{R}^n which is defined as

$$\varphi_{\mathbf{X}}(\mathbf{t}) \doteq \mathbb{E} \exp(\mathbf{t}^\top \mathbf{X}), \quad \text{for all } \mathbf{t} \in \mathbb{R}^n. \quad (1.76)$$

An MGF uniquely characterizes a distribution. The MGF of the multivariate Gaussian $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is

$$\varphi_{\mathbf{X}}(\mathbf{t}) = \exp\left(\mathbf{t}^\top \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t}\right). \quad (1.77)$$

This generalizes the MGF of the univariate Gaussian from Equation (A.40).

Prove the following facts.

1. *Closedness under affine transformations:* Given an n -dimensional Gaussian $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, and $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$,

$$\mathbf{AX} + \mathbf{b} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top). \quad (1.78)$$

2. *Additivity:* Given two independent Gaussian random vectors $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and $\mathbf{X}' \sim \mathcal{N}(\boldsymbol{\mu}', \boldsymbol{\Sigma}')$ in \mathbb{R}^n ,

$$\mathbf{X} + \mathbf{X}' \sim \mathcal{N}(\boldsymbol{\mu} + \boldsymbol{\mu}', \boldsymbol{\Sigma} + \boldsymbol{\Sigma}'). \quad (1.79)$$

These properties are unique to Gaussians and a reason for why they are widely used for learning and inference.

1.11. Expectation and variance of Gaussians.

Derive that $\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu}$ and $\text{Var}[\mathbf{X}] = \boldsymbol{\Sigma}$ when $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

Hint: First derive the expectation and variance of a univariate standard normal random variable.

1.12. Non-affine transformations of Gaussians.

Answer the following questions with **yes** or **no**.

1. Does there exist any non-affine transformation of a Gaussian random vector which is Gaussian? If yes, give an example.
2. Let X, Y, Z be independent standard normal random variables. Is $\frac{X+YZ}{\sqrt{1+Z^2}}$ Gaussian?

1.13. Decision theory.

Derive the optimal decisions under the squared loss and the asymmetric loss from Example 1.29.

PART I

Probabilistic Machine Learning

Preface to Part I

As humans, we constantly learn about the world around us. We learn to interact with our physical surroundings. We deepen our understanding of the world by establishing relationships between actors, objects, and events. And we learn about ourselves by observing how we interact with the world and with ourselves. We then continuously use this knowledge to make inferences and predictions, be it about the weather, the movement of a ball, or the behavior of a friend.

With limited computational resources, limited genetic information, and limited life experience, we are not able to learn everything about the world to complete certainty. We saw in Chapter 1 that probability theory is the mathematical framework for reasoning with uncertainty in the same way that logic is the mathematical framework for reasoning with certainty. We will discuss two kinds of uncertainty: “aleatoric” uncertainty which cannot be reduced under computational constraints, and “epistemic” uncertainty which can be reduced by observing more data.

An important aspect of learning is that we do not just learn once, but continually. Bayes’ rule allows us to update our beliefs and reduce our uncertainty as we observe new data — a process that is called *probabilistic inference*. By taking the former posterior as the new prior, probabilistic inference can be performed continuously and repeated indefinitely as we observe more and more data.

Our sensory information is often noisy and imperfect, which is another source of uncertainty. The same is true for machines, even if they can sometimes sense aspects of the world more accurately than humans. We discuss how one can infer latent structure of the world from sensed data, such as the state of a dynamical system like a car, in a process that is called *filtering*.

In this first part of the manuscript, we examine how we can build machines that are capable of (continual) learning and inference. First, we introduce probabilistic inference in the context of linear models which

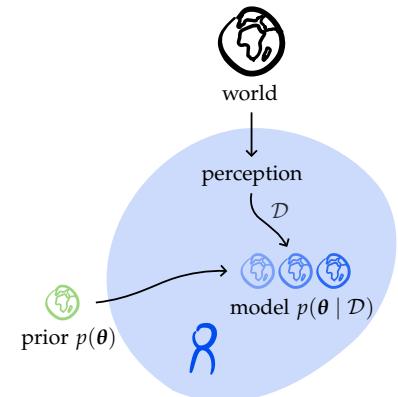


Figure 1.10: A schematic illustration of probabilistic inference in the context of the (supervised) learning of a model θ from perceived data D . The prior model $p(\theta)$ can equip the model with anything from substantial, to little, to no prior knowledge.

make predictions based on fixed (often hand-designed) features. We then discuss how probabilistic inference can be scaled to kernel methods and Gaussian processes which use a large (potentially infinite) number of features, and to deep neural networks which learn features dynamically from data. In these models, exact inference is typically intractable, and we discuss modern methods for approximate inference such as variational inference and Markov chain Monte Carlo. We highlight a tradeoff between curiosity (i.e., extrapolating beyond the given data) and conformity (i.e., fitting the given data), which surfaces as a fundamental principle of probabilistic inference in the regime where the data and our computational resources are limited.

2

Linear Regression

As a first example of probabilistic inference, we will study linear models for regression¹ which assume that the output $y \in \mathbb{R}$ is a linear function of the input $x \in \mathbb{R}^d$:

$$y \approx \mathbf{w}^\top \mathbf{x} + w_0$$

where $\mathbf{w} \in \mathbb{R}^d$ are the weights and $w_0 \in \mathbb{R}$ is the intercept. Observe that if we define the extended inputs $\mathbf{x}' \doteq (\mathbf{x}, 1)$ and $\mathbf{w}' \doteq (\mathbf{w}, w_0)$, then $\mathbf{w}'^\top \mathbf{x}' = \mathbf{w}^\top \mathbf{x} + w_0$, implying that without loss of generality it suffices to study linear functions without the intercept term w_0 . We will therefore consider the following function class of linear models

$$f(\mathbf{x}; \mathbf{w}) \doteq \mathbf{w}^\top \mathbf{x}.$$

We will consider the supervised learning task of learning weights \mathbf{w} from labeled training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$. We define the *design matrix*,

$$\mathbf{X} \doteq \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} \in \mathbb{R}^{n \times d}, \quad (2.1)$$

as the collection of inputs and the vector $\mathbf{y} \doteq [y_1 \cdots y_n]^\top \in \mathbb{R}^n$ as the collection of labels. For each noisy observation (\mathbf{x}_i, y_i) , we define the value of the approximation of our model, $f_i \doteq \mathbf{w}^\top \mathbf{x}_i$. Our model at the inputs \mathbf{X} is described by the vector $\mathbf{f} \doteq [f_1 \cdots f_n]^\top$ which can be expressed succinctly as $\mathbf{f} = \mathbf{X}\mathbf{w}$.

The most common way of estimating \mathbf{w} from data is the *least squares estimator*,

$$\hat{\mathbf{w}}_{\text{ls}} \doteq \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2, \quad (2.2)$$

¹ As we have discussed in Section 1.3, regression models can also be used for classification. The canonical example of a linear model for classification is logistic regression, which we will discuss in Section 5.1.1.

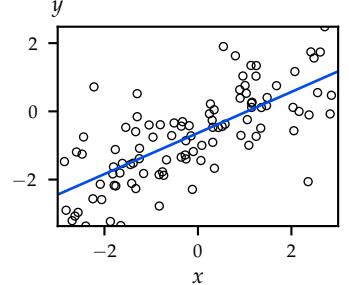


Figure 2.1: Example of linear regression with the least squares estimator (shown in blue).

minimizing the squared difference between the labels and predictions of the model. A slightly different estimator is used for *ridge regression*,

$$\hat{w}_{\text{ridge}} \doteq \arg \min_{w \in \mathbb{R}^d} \|y - Xw\|_2^2 + \lambda \|w\|_2^2 \quad (2.3)$$

where $\lambda > 0$. The squared L_2 regularization term $\lambda \|w\|_2^2$ penalizes large w and thus reduces the “complexity” of the resulting model.² It can be shown that the unique solutions to least squares and ridge regression are given by

$$\hat{w}_{\text{ls}} = (X^\top X)^{-1} X^\top y \quad \text{and} \quad (2.4)$$

$$\hat{w}_{\text{ridge}} = (X^\top X + \lambda I)^{-1} X^\top y, \quad (2.5)$$

respectively if the Hessian of the loss is positive definite (i.e., the loss is strictly convex) $\textcircled{?}$ which is the case as long as the columns of X are not linearly dependent. Least squares regression can be seen as finding the orthogonal projection of y onto the column space of X , as is illustrated in Figure 2.2 $\textcircled{?}$.

2.0.1 Maximum Likelihood Estimation

Since our function class comprises linear functions of the form $w^\top x$, the observation model from Equation (1.56) simplifies to

$$y_i = w^* \top x_i + \varepsilon_i \quad (2.6)$$

for some weight vector w , where for the purpose of this chapter we will additionally assume that $\varepsilon_i \sim \mathcal{N}(0, \sigma_n^2)$ is homoscedastic Gaussian noise.³ This observation model is equivalently characterized by the Gaussian likelihood,

$$y_i | x_i, w \sim \mathcal{N}(w^\top x_i, \sigma_n^2). \quad (2.7)$$

Based on this likelihood we can compute the MLE (1.57) of the weights:

$$\hat{w}_{\text{MLE}} = \arg \max_{w \in \mathbb{R}^d} \sum_{i=1}^n \log p(y_i | x_i, w) = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (y_i - w^\top x_i)^2.$$

Note that therefore $\hat{w}_{\text{MLE}} = \hat{w}_{\text{ls}}$.

In practice, the noise variance σ_n^2 is typically unknown and also has to be determined, for example, through maximum likelihood estimation. It is a straightforward exercise to check that the MLE of σ_n^2 given fixed weights w is $\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (y_i - w^\top x_i)^2$ $\textcircled{?}$.

2.1 Weight-space View

The most immediate and natural probabilistic interpretation of linear regression is to quantify uncertainty about the weights w . Recall

² Ridge regression is more robust to multicollinearity than standard linear regression. *Multicollinearity* occurs when multiple independent inputs are highly correlated. In this case, their individual effects on the predicted variable cannot be estimated well. Classical linear regression is highly volatile to small input changes. The regularization of ridge regression reduces this volatility by introducing a bias on the weights towards 0.

Problem 2.1 (1)

Problem 2.1 (2)

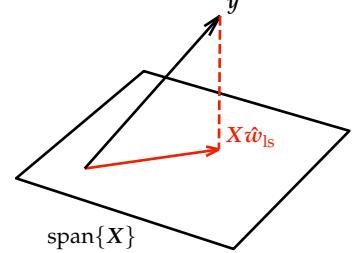


Figure 2.2: Least squares regression finds the orthogonal projection of y onto $\text{span}\{X\}$ (here illustrated as the plane).

³ ε_i is called *additive white Gaussian noise*.

using Equation (1.55)

plugging in the Gaussian likelihood and simplifying

Problem 2.2

that probabilistic inference requires specification of a generative model comprised of prior and likelihood. Throughout this chapter, we will use the Gaussian prior,

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I}), \quad (2.8)$$

and the Gaussian likelihood from Equation (2.7). We will discuss possible (probabilistic) strategies for choosing hyperparameters such as the prior variance σ_p^2 and the noise variance σ_n^2 in Section 4.4.

Remark 2.1: Why a Gaussian prior?

The choice of using a Gaussian prior may seem somewhat arbitrary at first sight, except perhaps for the nice analytical properties of Gaussians that we have seen in Section 1.2.3 and which will prove useful. The maximum entropy principle (cf. Section 1.2.1) provides a more fundamental justification for Gaussian priors since turns out that \mathcal{N} has the maximum entropy among all distributions on \mathbb{R}^d with known mean and variance $\textcircled{?}$.

Next, let us derive the posterior distribution over the weights.

$$\begin{aligned} & \log p(\mathbf{w} | \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) \\ &= \log p(\mathbf{w}) + \log p(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}, \mathbf{w}) + \text{const} \\ &= \log p(\mathbf{w}) + \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \mathbf{w}) + \text{const} \\ &= -\frac{1}{2} \left[\sigma_p^{-2} \|\mathbf{w}\|_2^2 + \sigma_n^{-2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 \right] + \text{const} \\ &= -\frac{1}{2} \left[\sigma_p^{-2} \|\mathbf{w}\|_2^2 + \sigma_n^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \right] + \text{const} \\ &= -\frac{1}{2} \left[\sigma_p^{-2} \mathbf{w}^\top \mathbf{w} + \sigma_n^{-2} (\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y}) \right] + \text{const} \\ &= -\frac{1}{2} \left[\mathbf{w}^\top (\sigma_n^{-2} \mathbf{X}^\top \mathbf{X} + \sigma_p^{-2} \mathbf{I}) \mathbf{w} - 2\sigma_n^{-2} \mathbf{y}^\top \mathbf{X} \mathbf{w} \right] + \text{const}. \end{aligned} \quad (2.9)$$

Observe that the log-posterior is a quadratic form in \mathbf{w} , so the posterior distribution must be Gaussian:

$$\mathbf{w} | \mathbf{x}_{1:n}, \mathbf{y}_{1:n} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (2.10a)$$

see Equation (A.12)

where we can read off the mean and variance to be

$$\boldsymbol{\mu} \doteq \sigma_n^{-2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y}, \quad (2.10b)$$

$$\boldsymbol{\Sigma} \doteq (\sigma_n^{-2} \mathbf{X}^\top \mathbf{X} + \sigma_p^{-2} \mathbf{I})^{-1}. \quad (2.10c)$$

This also shows that Gaussians with known variance and linear likelihood are self-conjugate, a property that we had hinted at in Section 1.2.2. It can be shown more generally that Gaussians with known

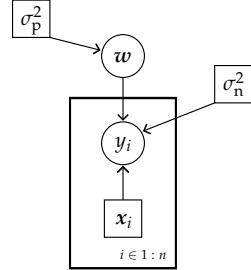


Figure 2.3: Directed graphical model of Bayesian linear regression in plate notation.

Problem 5.6

variance are self-conjugate to any Gaussian likelihood (Murphy, 2007). For other generative models, the posterior can typically not be expressed in closed-form — this is a very special property of Gaussians!

2.1.1 Maximum a Posteriori Estimation

Computing the MAP estimate for the weights,

$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \max_w \log p(y_{1:n} | \mathbf{x}_{1:n}, \mathbf{w}) + \log p(\mathbf{w}) \\ &= \arg \min_w \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{\sigma_n^2}{\sigma_p^2} \|\mathbf{w}\|_2^2,\end{aligned}\quad (2.11)$$

we observe that this is identical to ridge regression with weight decay $\lambda \doteq \sigma_n^2 / \sigma_p^2$: $\hat{\mathbf{w}}_{\text{MAP}} = \hat{\mathbf{w}}_{\text{ridge}}$. Equation (2.11) is simply the MLE loss with an additional L_2 -regularization (originating from the prior) that encourages keeping weights small. Recall that the MAP estimate corresponds to the mode of the posterior distribution, which in the case of a Gaussian is simply its mean μ . As to be expected, μ coincides with the analytical solution to ridge regression from Equation (2.5).

Example 2.2: Lasso as the MAP estimate with a Laplace prior

One problem with ridge regression is that the contribution of nearly-zero weights to the L_2 -regularization term is negligible. Thus, L_2 -regularization is typically not sufficient to perform variable selection (that is, set some weights to zero entirely), which is often desirable for interpretability of the model.

A commonly used alternative to ridge regression is the *least absolute shrinkage and selection operator* (or *lasso*), which regularizes with the L_1 -norm:

$$\hat{\mathbf{w}}_{\text{lasso}} \doteq \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1. \quad (2.12)$$

It turns out that lasso can also be viewed as probabilistic inference, using a Laplace prior $\mathbf{w} \sim \text{Laplace}(\mathbf{0}, h)$ with length scale h instead of a Gaussian prior.

Computing the MAP estimate for the weights yields,

$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \max_w \log p(y_{1:n} | \mathbf{x}_{1:n}, \mathbf{w}) + \log p(\mathbf{w}) \\ &= \arg \min_w \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \frac{\sigma_n^2}{h} \|\mathbf{w}\|_1\end{aligned}\quad (2.13)$$

which coincides with the lasso with weight decay $\lambda \doteq \sigma_n^2 / h$.

using that the likelihood and prior are Gaussian

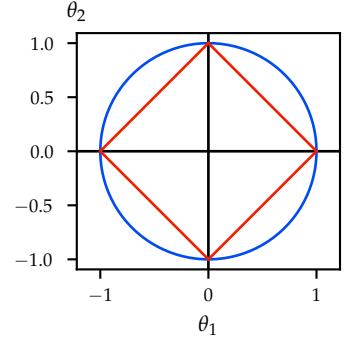


Figure 2.4: Level sets of L_2 - (blue) and L_1 -regularization (red), corresponding to Gaussian and Laplace priors, respectively. It can be seen that L_1 -regularization is more effective in encouraging sparse solutions (that is, solutions where many components are set to exactly 0).

using that the likelihood is Gaussian and the prior is Laplacian

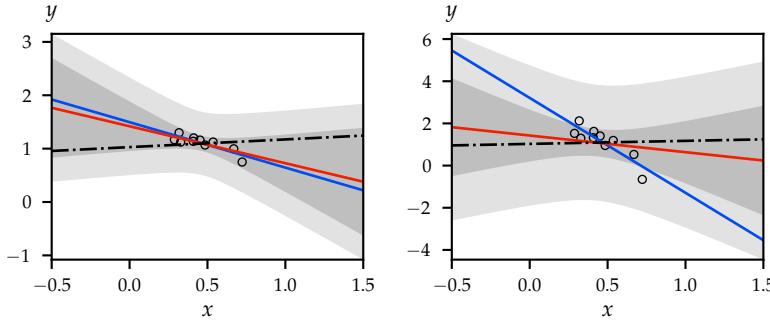
To make predictions at a test point x^* , we define the (model-)predicted point $f^* \doteq \hat{w}_{\text{MAP}}^\top x^*$ and obtain the label prediction

$$y^* | x^*, x_{1:n}, y_{1:n} \sim \mathcal{N}(f^*, \sigma_n^2). \quad (2.14)$$

Here we observe that using point estimates such as the MAP estimate does not quantify uncertainty in the weights. The MAP estimate simply collapses all mass of the posterior around its mode. This can be harmful when we are highly unsure about the best model, e.g., because we have observed insufficient data.

2.1.2 Probabilistic Inference

Rather than selecting a single weight vector \hat{w} to make predictions, we can use the full posterior distribution. This is known as *Bayesian linear regression* (BLR) and illustrated with an example in Figure 2.5.



To make predictions at a test point x^* , we let $f^* \doteq w^\top x^*$ which has the distribution

$$f^* | x^*, x_{1:n}, y_{1:n} \sim \mathcal{N}(\mu^\top x^*, x^{*\top} \Sigma x^*). \quad (2.15)$$

Note that this does not take into account the noise in the labels σ_n^2 . For the label prediction y^* , we obtain

$$y^* | x^*, x_{1:n}, y_{1:n} \sim \mathcal{N}(\mu^\top x^*, x^{*\top} \Sigma x^* + \sigma_n^2). \quad (2.16)$$

2.1.3 Recursive Probabilistic Inference

We have already discussed the recursive properties of probabilistic inference in Section 1.3.6. For Bayesian linear regression with a Gaussian prior and likelihood, this principle can be used to derive an efficient online algorithm since also the posterior is a Gaussian,

$$p^{(t)}(w) = \mathcal{N}(w; \mu^{(t)}, \Sigma^{(t)}), \quad (2.17)$$

which can be stored efficiently using only $O(d^2)$ parameters. This leads to an efficient online algorithm for Bayesian linear regression

Figure 2.5: Comparison of linear regression (MLE), ridge regression (MAP estimate), and Bayesian linear regression when the data is generated according to

$$y | w, x \sim \mathcal{N}(w^\top x, \sigma_n^2).$$

The **true mean** is shown in black, the MLE in blue, and the MAP estimate in red. The dark gray area denotes the epistemic uncertainty of Bayesian linear regression and the light gray area the additional homoscedastic noise. On the left, $\sigma_n = 0.15$. On the right, $\sigma_n = 0.7$.

using the closedness of Gaussians under linear transformations (1.78)

using additivity of Gaussians (1.79)

with time-independent(!) memory complexity $O(d)$ and round complexity $O(d^2)$ [?](#). The interpretation of Bayesian linear regression as an online algorithm also highlights similarities to other sequential models such as Kalman filters, which we discuss in Chapter 3. In Example 3.5, we will learn that online Bayesian linear regression is, in fact, an example of a Kalman filter.

Problem 2.5

2.2 Aleatoric and Epistemic Uncertainty

The predictive posterior distribution from Equation (2.16) highlights a decomposition of uncertainty wherein $x^{*\top} \Sigma x^*$ corresponds to the uncertainty about our model due to the lack of data (commonly referred to as the *epistemic uncertainty*) and σ_n^2 corresponds to the uncertainty about the labels that cannot be explained by the inputs and any model from the model class (commonly referred to as the *aleatoric uncertainty*, “irreducible noise”, or simply “(label) noise”) [?](#).

Problem 2.6

A natural probabilistic approach is to represent epistemic uncertainty with a probability distribution over models. Intuitively, the variance of this distribution measures our uncertainty about the model and its mode corresponds to our current best (point) estimate. The distribution over weights of a linear model is one example, and we will continue to explore this approach for other models in the following chapters.

It is a practical modeling choice how much inaccuracy to attribute to epistemic or aleatoric uncertainty. Generally, when a poor model is used to explain a process, more inaccuracy has to be attributed to irreducible noise. For example, when a linear model is used to “explain” a nonlinear process, most uncertainty is aleatoric as the model cannot explain the data well. As we use more expressive models, a larger portion of the uncertainty can be explained by the data.

Epistemic and aleatoric uncertainty can be formally defined in terms of the law of total variance (1.41),

$$\text{Var}[y^* | x^*] = \underbrace{\mathbb{E}_\theta[\text{Var}_{y^*}[y^* | x^*, \theta]]}_{\text{aleatoric uncertainty}} + \underbrace{\text{Var}_\theta[\mathbb{E}_{y^*}[y^* | x^*, \theta]]}_{\text{epistemic uncertainty}}. \quad (2.18)$$

Here, the mean variability of predictions y^* averaged across all models θ is the estimated *aleatoric uncertainty*. In contrast, the variability of the mean prediction y^* under each model θ is the estimated *epistemic uncertainty*. This decomposition of uncertainty will appear frequently throughout this manuscript.

2.3 Non-linear Regression

We can use linear regression not only to learn linear functions. The trick is to apply a nonlinear transformation $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^e$ to the features x_i , where d is the dimension of the input space and e is the dimension of the designed *feature space*. We denote the design matrix comprised of transformed features by $\Phi \in \mathbb{R}^{n \times e}$. Note that if the feature transformation ϕ is the identity function then $\Phi = X$.

Example 2.3: Polynomial regression

Let $\phi(x) \doteq [x^2, x, 1]$ and $w \doteq [a, b, c]$. Then the function that our model learns is given as

$$f = ax^2 + bx + c.$$

Thus, our model can exactly represent all polynomials up to degree 2.

However, to learn polynomials of degree m in d input dimensions, we need to apply the nonlinear transformation

$$\begin{aligned} \phi(x) = & [1, x_1, \dots, x_d, x_1^2, \dots, x_d^2, x_1 \cdot x_2, \dots, x_{d-1} \cdot x_d, \\ & \dots, \\ & x_{d-m+1} \cdots x_d]. \end{aligned}$$

Note that the feature dimension e is $\sum_{i=0}^m \binom{d+i-1}{i} = \Theta(d^m)$.⁴ Thus, the dimension of the feature space grows exponentially in the degree of polynomials and input dimensions. Even for relatively small m and d , this becomes completely unmanageable.

The example of polynomials highlights that it may be inefficient to keep track of the weights $w \in \mathbb{R}^e$ when e is large, and that it may be useful to instead consider a reparameterization which is of dimension n rather than of the feature dimension.

2.4 Function-space View

Let us now look at Bayesian linear regression through a different lens. Previously, we have been interpreting it as a distribution over the weights w of a linear function $f = \Phi w$. The key idea is that for a finite set of inputs (ensuring that the design matrix is well-defined), we can equivalently consider a distribution directly over the estimated function values f . We call this the *function-space view* of Bayesian linear regression.

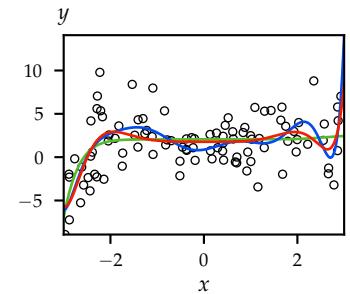


Figure 2.6: Applying linear regression with a feature space of polynomials of degree 10. The **least squares estimate** is shown in blue, **ridge regression** in red, and **lasso** in green.

⁴ Observe that the vector contains $\binom{d+i-1}{i}$ monomials of degree i as this is the number of ways to choose i times from d items with replacement and without consideration of order. To see this, consider the following encoding: We take a sequence of $d + i - 1$ spots. Selecting any subset of i spots, we interpret the remaining $d - 1$ spots as "barriers" separating each of the d items. The selected spots correspond to the number of times each item has been selected. For example, if 2 items are to be selected out of a total of 4 items with replacement, one possible configuration is "○ || ○ |" where ○ denotes a selected spot and | denotes a barrier. This configuration encodes that the first and third item have each been chosen once. The number of possible configurations — each encoding a unique outcome — is therefore $\binom{d+i-1}{i}$.

Instead of considering a prior over the weights $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$ as we have done previously, we now impose a prior directly on the values of our model at the observations. Using that Gaussians are closed under linear maps (1.78), we obtain the equivalent prior

$$\mathbf{f} | \mathbf{X} \sim \mathcal{N}(\Phi \mathbb{E}[\mathbf{w}], \Phi \text{Var}[\mathbf{w}] \Phi^\top) = \mathcal{N}(\mathbf{0}, \underbrace{\sigma_p^2 \Phi \Phi^\top}_K) \quad (2.19)$$

where $K \in \mathbb{R}^{n \times n}$ is the so-called *kernel matrix*. Observe that the entries of the kernel matrix can be expressed as $K(i, j) = \sigma_p^2 \cdot \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$.

You may say that nothing has changed, and you would be right — that is precisely the point. Note, however, that the shape of the kernel matrix is $n \times n$ rather than the $e \times e$ covariance matrix over weights, which becomes unmanageable when e is large. The kernel matrix K has entries only for the finite set of observed inputs. However, in principle, we could have observed any input, and this motivates the definition of the *kernel function*

$$k(\mathbf{x}, \mathbf{x}') \doteq \sigma_p^2 \cdot \phi(\mathbf{x})^\top \phi(\mathbf{x}') \quad (2.20)$$

for arbitrary inputs \mathbf{x} and \mathbf{x}' . A kernel matrix is simply a finite “view” of the kernel function,

$$K = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix} \quad (2.21)$$

Observe that by definition of the kernel matrix in Equation (2.19), the kernel matrix is a covariance matrix and the kernel function measures the covariance of the function values $f(\mathbf{x})$ and $f(\mathbf{x}')$ given inputs \mathbf{x} and \mathbf{x}' :

$$k(\mathbf{x}, \mathbf{x}') = \text{Cov}[f(\mathbf{x}), f(\mathbf{x}')]. \quad (2.22)$$

Moreover, note that we have reformulated⁵ the learning algorithm such that the feature space is now *implicit* in the choice of kernel, and the kernel is defined by inner products of (nonlinearly transformed) inputs. In other words, the choice of kernel implicitly determines the class of functions that f is sampled from (without expressing the functions explicitly in closed-form), which encodes our prior beliefs. This is known as the *kernel trick*.

2.4.1 Learning and Predictions

We have already kernelized the Bayesian linear regression prior. The posterior distribution $\mathbf{f} | \mathbf{X}, \mathbf{y}$ is again Gaussian due to the closedness

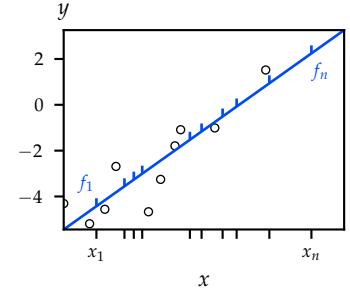


Figure 2.7: An illustration of the function-space view. The model is described by the points (\mathbf{x}_i, f_i) .

⁵ we often say “kernelized”

properties of Gaussians, analogously to our derivation of the prior kernel matrix in Equation (2.19).

It remains to show that we can also rely on the kernel trick for predictions. Given the test point x^* , we define

$$\tilde{\Phi} \doteq \begin{bmatrix} \Phi \\ \phi(x^*)^\top \end{bmatrix}, \quad \tilde{y} \doteq \begin{bmatrix} y \\ y^* \end{bmatrix}, \quad \tilde{f} \doteq \begin{bmatrix} f \\ f^* \end{bmatrix}.$$

We immediately obtain $\tilde{f} = \tilde{\Phi}w$. Analogously to our analysis of predictions from the weight-space view, we add the label noise to obtain the estimate $\tilde{y} = \tilde{f} + \tilde{\varepsilon}$ where $\tilde{\varepsilon} \doteq [\varepsilon_1 \dots \varepsilon_n \varepsilon^*]^\top \sim \mathcal{N}(\mathbf{0}, \sigma_n^2 I)$ is the independent label noise. Applying the same reasoning as we did for the prior, we obtain

$$\tilde{f} | X, x^* \sim \mathcal{N}(\mathbf{0}, \tilde{K}) \quad (2.23)$$

where $\tilde{K} \doteq \sigma_p^2 \tilde{\Phi} \tilde{\Phi}^\top$. Adding the label noise yields

$$\tilde{y} | X, x^* \sim \mathcal{N}(\mathbf{0}, \tilde{K} + \sigma_n^2 I). \quad (2.24)$$

Finally, we can conclude from the closedness of Gaussian random vectors under conditional distributions (1.53) that the predictive posterior $y^* | x^*, X, y$ follows again a normal distribution. We will do a full derivation of the posterior and predictive posterior in Section 4.1.

2.4.2 Efficient Polynomial Regression

But how does the kernel trick address our concerns about efficiency raised in Section 2.3? After all, computing the kernel for a feature space of dimension e still requires computing sums of length e which is prohibitive when e is large. The kernel trick opens up a couple of new doors for us:

1. For certain feature transformations ϕ , we may be able to find an easier to compute expression equivalent to $\phi(x)^\top \phi(x')$.
2. If this is not possible, we could approximate the inner product by an easier to compute expression.
3. Or, alternatively, we may decide not to care very much about the exact feature transformation and simply experiment with kernels that induce *some* feature space (which may even be infinitely dimensional).

We will explore the third approach when we revisit kernels in Section 4.3. A polynomial feature transformation can be computed efficiently in closed-form.

Fact 2.4. *For the polynomial feature transformation ϕ up to degree m from Example 2.3, it can be shown that up to constant factors,*

$$\phi(x)^\top \phi(x') = (1 + x^\top x')^m. \quad (2.25)$$

For example, for input dimension 2, the kernel $(1 + \mathbf{x}^\top \mathbf{x}')^2$ corresponds to the feature vector $\phi(\mathbf{x}) = [1 \ \sqrt{2}x_1 \ \sqrt{2}x_2 \ \sqrt{2}x_1x_2 \ x_1^2 \ x_2^2]^\top$.

Discussion

We have explored a probabilistic perspective on linear models, and seen that classical approaches such as least squares and ridge regression can be interpreted as approximate probabilistic inference. We then saw that we can even perform *exact* probabilistic inference efficiently if we adopt a Gaussian prior and Gaussian noise assumption. These are already powerful tools, which are often applied also to nonlinear models if we treat the latent feature space — which was either human-designed or learned via deep learning — as fixed. In the next chapter, we will digress briefly from the storyline on “learning” to see how we can adopt a similar probabilistic perspective to track latent states over time. Then, in Chapter 4, we will see how we can use the function-space view and kernel trick to learn flexible nonlinear models with exact probabilistic inference, without ever explicitly representing the feature space.

Problems

2.1. Closed-form linear regression.

1. Derive the unique solutions to least squares and ridge regression from Equations (2.4) and (2.5).
2. For an $n \times m$ matrix A and vector $\mathbf{x} \in \mathbb{R}^m$, we call $\Pi_A \mathbf{x}$ the orthogonal projection of \mathbf{x} onto $\text{span}\{A\} = \{Ax' \mid x' \in \mathbb{R}^m\}$. In particular, an orthogonal projection satisfies $\mathbf{x} - \Pi_A \mathbf{x} \perp Ax'$ for all $x' \in \mathbb{R}^m$. Show that $\hat{\mathbf{w}}_{\text{ls}}$ from Equation (2.4) is such that $X\hat{\mathbf{w}}_{\text{ls}}$ is the unique closest point to y on $\text{span}\{X\}$, i.e., it satisfies $X\hat{\mathbf{w}}_{\text{ls}} = \Pi_X y$.

2.2. MLE of noise variance.

Show that the MLE of σ_n^2 given fixed weights \mathbf{w} is

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2. \quad (2.26)$$

2.3. Variance of least squares around training data.

Show that the variance of a prediction at the point $[1 \ x^\star]^\top$ is smallest when x^\star is the mean of the training data. More formally, show that if inputs are of the form $\mathbf{x}_i = [1 \ x_i]^\top$ where $x_i \in \mathbb{R}$ and $\hat{\mathbf{w}}_{\text{ls}}$ is the least squares estimate, then $\text{Var}[y^\star \mid [1 \ x^\star]^\top, \hat{\mathbf{w}}_{\text{ls}}]$ is minimized for $x^\star = \frac{1}{n} \sum_{i=1}^n x_i$.

2.4. Bayesian linear regression.

Suppose you are given the following observations

$$\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 2 & 1 \\ 2 & 2 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 2.4 \\ 4.3 \\ 3.1 \\ 4.9 \end{bmatrix}$$

and assume the data follows a linear model with homoscedastic noise $\mathcal{N}(0, \sigma_n^2)$ where $\sigma_n^2 = 0.1$.

1. Find the maximum likelihood estimate $\hat{\mathbf{w}}_{\text{MLE}}$ given the data.
2. Now assume that we have a prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \sigma_p^2 \mathbf{I})$ with $\sigma_p^2 = 0.05$. Find the MAP estimate $\hat{\mathbf{w}}_{\text{MAP}}$ given the data and the prior.
3. Use the posterior $p(\mathbf{w} | \mathbf{X}, \mathbf{y})$ to get a posterior prediction for the label y^* at $\mathbf{x}^* = [3 \ 3]^\top$. Report the mean and the variance of this prediction.
4. How would you have to change the prior $p(\mathbf{w})$ such that

$$\hat{\mathbf{w}}_{\text{MAP}} \rightarrow \hat{\mathbf{w}}_{\text{MLE}}?$$

2.5. Online Bayesian linear regression.

1. Can you design an algorithm that updates the posterior (as opposed to recalculating it from scratch using Equation (2.10)) in a smarter way? The requirement is that the memory should not grow as $O(t)$.
2. If d is large, computing the inverse every round is very expensive. Can you use the recursive structure you found in the previous question to bring down the computational complexity of every round to $O(d^2)$?

The resulting efficient online algorithm is known as *online Bayesian linear regression*.

2.6. Aleatoric and epistemic uncertainty of BLR.

Prove for Bayesian linear regression that $\mathbf{x}^{*\top} \Sigma \mathbf{x}^*$ is the epistemic uncertainty and σ_n^2 the aleatoric uncertainty in y^* under the decomposition of Equation (2.18).

2.7. Hyperpriors.

We consider a dataset $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ of size n , where $\mathbf{x}_i \in \mathbb{R}^d$ denotes the feature vector and $y_i \in \mathbb{R}$ denotes the label of the i -th data point. Let ε_i be i.i.d. samples from the Gaussian distribution $\mathcal{N}(0, \lambda^{-1})$ for a given $\lambda > 0$. We collect the labels in a vector $\mathbf{y} \in \mathbb{R}^n$, the features in a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, and the noise in a vector $\boldsymbol{\varepsilon} \in \mathbb{R}^n$. The labels are generated according to $\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\varepsilon}$.

To perform Bayesian Linear Regression, we consider the prior distribution over the parameter vector w to be $\mathcal{N}(\mu, \lambda^{-1}I_d)$, where I_d denotes the d -dimensional identity matrix and $\mu \in \mathbb{R}^d$ is a hyperparameter.

1. Given this Bayesian data model, what is the conditional covariance matrix $\Sigma_y \doteq \text{Var}[y | X, \mu, \lambda]$?
2. Calculate the maximum likelihood estimate of the hyperparameter μ .
3. Since we are unsure about the hyperparameter μ , we decide to model our uncertainty about μ by placing the “hyperprior” $\mu \sim \mathcal{N}(\mathbf{0}, I_d)$. Is the posterior distribution $p(\mu | X, y, \lambda)$ a Gaussian distribution? If yes, what are its mean vector and covariance matrix?
4. What is the posterior distribution $p(\lambda | X, y, \mu)$?

Hint: For any $a \in \mathbb{R}$, $A \in \mathbb{R}^{n \times n}$ it holds that $\det(aA) = a^n \det(A)$.

3

Filtering

Before we continue in Chapter 4 with the function-space view of regression, we want to look at a seemingly different but very related problem. We will study Bayesian learning and inference in the *state space model*, where we want to keep track of the state of an agent over time based on noisy observations. In this model, we have a sequence of (hidden) states $(\mathbf{X}_t)_{t \in \mathbb{N}_0}$ where \mathbf{X}_t is in \mathbb{R}^d and a sequence of observations $(\mathbf{Y}_t)_{t \in \mathbb{N}_0}$ where \mathbf{Y}_t is in \mathbb{R}^m .

The process of keeping track of the state using noisy observations is also known as *Bayesian filtering* or *recursive Bayesian estimation*. Figure 3.1 illustrates this process, where an agent perceives the current state of the world and then updates its beliefs about the state based on this observation.

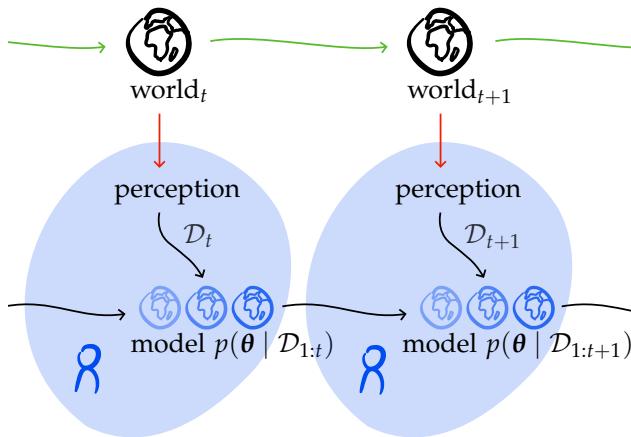


Figure 3.1: Schematic view of Bayesian filtering: An agent perceives the current state of the world and updates its belief accordingly.

We will discuss Bayesian filtering more broadly in the next section. A Kalman filter is an important special case of a Bayes' filter, which uses a Gaussian prior over the states and conditional linear Gaussians to describe the evolution of states and observations. Analogously to the

previous chapter, we will see that inference in this model is tractable due to the closedness properties of Gaussians.

Definition 3.1 (Kalman filter). A *Kalman filter* is specified by a Gaussian prior over the states,

$$\mathbf{X}_0 \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (3.1)$$

and a conditional linear Gaussian *motion model* and *sensor model*,

$$\mathbf{X}_{t+1} \doteq \mathbf{F}\mathbf{X}_t + \boldsymbol{\varepsilon}_t \quad \mathbf{F} \in \mathbb{R}^{d \times d}, \boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_x), \quad (3.2)$$

$$\mathbf{Y}_t \doteq \mathbf{H}\mathbf{X}_t + \boldsymbol{\eta}_t \quad \mathbf{H} \in \mathbb{R}^{m \times d}, \boldsymbol{\eta}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_y), \quad (3.3)$$

respectively. The motion model is sometimes also called *transition model* or *dynamics model*. Crucially, Kalman filters assume that \mathbf{F} and \mathbf{H} are known. In general, \mathbf{F} and \mathbf{H} may depend on t . Also, $\boldsymbol{\varepsilon}$ and $\boldsymbol{\eta}$ may have a non-zero mean, commonly called a “drift”.

Because Kalman filters use conditional linear Gaussians, which we have already seen in Equation (1.55), their joint distribution (over all variables) is also Gaussian. This means that predicting the future states of a Kalman filter is simply inference with multivariate Gaussians. In Bayesian filtering, however, we do not only want to make predictions occasionally. In Bayesian filtering, we want to *keep track* of states, that is, predict the current state of an agent online.¹ To do this efficiently, we need to update our *belief* about the state of the agent recursively, similarly to our recursive Bayesian updates in Bayesian linear regression (see Section 2.1.3).

From the directed graphical model of a Kalman filter shown in Figure 3.2, we can immediately gather the following conditional independence relations,²

$$\mathbf{X}_{t+1} \perp \mathbf{X}_{1:t-1}, \mathbf{Y}_{1:t-1} \mid \mathbf{X}_t, \quad (3.4)$$

$$\mathbf{Y}_t \perp \mathbf{X}_{1:t-1} \mid \mathbf{X}_t \quad (3.5)$$

$$\mathbf{Y}_t \perp \mathbf{Y}_{1:t-1} \mid \mathbf{X}_{t-1}. \quad (3.6)$$

The first conditional independence property is also known as the *Markov property*, which we will return to later in our discussion of Markov chains and Markov decision processes. This characterization of the Kalman filter, yields the following factorization of the joint distribution:

$$\begin{aligned} p(\mathbf{x}_{1:t}, \mathbf{y}_{1:t}) &= \prod_{i=1}^t p(\mathbf{x}_i \mid \mathbf{x}_{1:i-1}) p(\mathbf{y}_i \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:i-1}) \\ &= p(\mathbf{x}_1) p(\mathbf{y}_1 \mid \mathbf{x}_1) \prod_{i=2}^t p(\mathbf{x}_i \mid \mathbf{x}_{i-1}) p(\mathbf{y}_i \mid \mathbf{x}_i). \end{aligned} \quad (3.7)$$

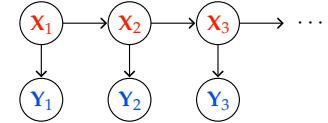


Figure 3.2: Directed graphical model of a Kalman filter with hidden states \mathbf{X}_t and observables \mathbf{Y}_t .

¹ Here, *online* is common terminology to say that we want to perform inference at time t without exposure to times $t+1, t+2, \dots$, so in “real-time”.

² Alternatively, they follow from the definition of the motion and sensor models as linear updates.

using the product rule (1.11)

using the conditional independence properties from (3.4), (3.5), and (3.6)

3.1 Conditioning and Prediction

We can describe Bayesian filtering by the following recursive scheme with the two phases, *conditioning* (also called “update”) and *prediction*:

Algorithm 3.2: Bayesian filtering

- 1 start with a prior over initial states $p(x_0)$
 - 2 **for** $t = 1$ **to** ∞ **do**
 - 3 | assume we have $p(x_t | y_{1:t-1})$
 - 4 | **conditioning:** compute $p(x_t | y_{1:t})$ using the new observation y_t
 - 5 | **prediction:** compute $p(x_{t+1} | y_{1:t})$
-

Let us consider the conditioning step first:

$$\begin{aligned} p(x_t | y_{1:t}) &= \frac{1}{Z} p(x_t | y_{1:t-1}) p(y_t | x_t, y_{1:t-1}) \\ &= \frac{1}{Z} p(x_t | y_{1:t-1}) p(y_t | x_t). \end{aligned} \quad (3.8)$$

For the prediction step, we obtain,

$$\begin{aligned} p(x_{t+1} | y_{1:t}) &= \int p(x_{t+1}, x_t | y_{1:t}) dx_t \\ &= \int p(x_{t+1} | x_t, y_{1:t}) p(x_t | y_{1:t}) dx_t \\ &= \int p(x_{t+1} | x_t) p(x_t | y_{1:t}) dx_t. \end{aligned} \quad (3.9)$$

In general, these distributions can be very complicated, but for Gaussians (i.e., Kalman filters) they can be expressed in closed-form.

Remark 3.3: Bayesian smoothing

Bayesian smoothing is a closely related task to Bayesian filtering. While Bayesian filtering methods estimate the current state based only on observations obtained before and at the current time step, Bayesian smoothing computes the distribution of $X_k | y_{1:t}$ where $t > k$. That is Bayesian smoothing estimates X_k based on data until and beyond time k . Note that if $k = t$, then Bayesian smoothing coincides with Bayesian filtering.

Analogously to Equation (3.8),

$$p(x_k | y_{1:t}) \propto p(x_k | y_{1:k}) p(y_{k+1:t} | x_k). \quad (3.10)$$

If we assume a Gaussian prior and conditional Gaussian transition and dynamics models (this is called *Kalman smoothing*), then

using Bayes' rule (1.45)

using the conditional independence structure (3.6)

using the sum rule (1.7)

using the product rule (1.11)

using the conditional independence structure (3.4)

by the closedness properties of Gaussians, $\mathbf{X}_k \mid \mathbf{y}_{1:t}$ is a Gaussian. Indeed, all terms of Equation (3.10) are Gaussian PDFs and as seen in Equation (1.75), the product of two Gaussian PDFs is again proportional to a Gaussian PDF.

The first term, $\mathbf{X}_k \mid \mathbf{y}_{1:k}$, is the marginal posterior of the hidden states of the Kalman filter which can be obtained with Bayesian filtering.

By conditioning on \mathbf{X}_{k+1} , we have for the second term,

$$\begin{aligned} p(\mathbf{y}_{k+1:t} \mid \mathbf{x}_k) &= \int p(\mathbf{y}_{k+1:t} \mid \mathbf{x}_k, \mathbf{x}_{k+1}) p(\mathbf{x}_{k+1} \mid \mathbf{x}_k) d\mathbf{x}_{k+1} \\ &= \int p(\mathbf{y}_{k+1:t} \mid \mathbf{x}_{k+1}) p(\mathbf{x}_{k+1} \mid \mathbf{x}_k) d\mathbf{x}_{k+1} \\ &= \int p(\mathbf{y}_{k+1} \mid \mathbf{x}_{k+1}) p(\mathbf{y}_{k+2:t} \mid \mathbf{x}_{k+1}) p(\mathbf{x}_{k+1} \mid \mathbf{x}_k) d\mathbf{x}_{k+1} \end{aligned} \quad (3.11)$$

using the sum rule (1.7) and product rule (1.11)

using the conditional independence structure (3.5)

using the conditional independence structure (3.6)

Let us have a look at the terms in the product:

- $p(\mathbf{y}_{k+1} \mid \mathbf{x}_{k+1})$ is obtained from the sensor model,
- $p(\mathbf{x}_{k+1} \mid \mathbf{x}_k)$ is obtained from the transition model, and
- $p(\mathbf{y}_{k+2:t} \mid \mathbf{x}_{k+1})$ can be computed recursively backwards in time.

This recursion results in linear equations resembling a Kalman filter running backwards in time.

Thus, in the setting of Kalman smoothing, both factors of Equation (3.10) can be computed efficiently: one using a (forward) Kalman filter; the other using a “backward” Kalman filter. More concretely, in time $O(t)$, we can compute the two factors for all $k \in [t]$. This approach is known as *two-filter smoothing* or the *forward-backward algorithm*.

3.2 Kalman Filters

Let us return to the setting of Kalman filters where priors and likelihoods are Gaussian. Here, we will see that the update and prediction steps can be computed in closed form.

3.2.1 Conditioning

The conditioning operation in Kalman filters is also called the Kalman update. Before introducing the general Kalman update, let us consider a simpler example:

Example 3.4: Random walk in 1d

We use the simple motion and sensor models,³

$$X_{t+1} \mid x_t \sim \mathcal{N}(x_t, \sigma_x^2), \quad (3.12a)$$

$$Y_t \mid x_t \sim \mathcal{N}(x_t, \sigma_y^2). \quad (3.12b)$$

Let $X_t \mid y_{1:t} \sim \mathcal{N}(\mu_t, \sigma_t^2)$ be our belief at time t . It can be shown that Bayesian filtering yields the belief $X_{t+1} \mid y_{1:t+1} \sim \mathcal{N}(\mu_{t+1}, \sigma_{t+1}^2)$ at time $t+1$ where ②

$$\mu_{t+1} \doteq \frac{\sigma_y^2 \mu_t + (\sigma_t^2 + \sigma_x^2) y_{t+1}}{\sigma_t^2 + \sigma_x^2 + \sigma_y^2}, \quad \sigma_{t+1}^2 \doteq \frac{(\sigma_t^2 + \sigma_x^2) \sigma_y^2}{\sigma_t^2 + \sigma_x^2 + \sigma_y^2}. \quad (3.13)$$

Although looking intimidating at first, this update has a very natural interpretation. Let us define the following quantity,

$$\lambda \doteq \frac{\sigma_t^2 + \sigma_x^2}{\sigma_t^2 + \sigma_x^2 + \sigma_y^2} = 1 - \frac{\sigma_y^2}{\sigma_t^2 + \sigma_x^2 + \sigma_y^2} \in [0, 1]. \quad (3.14)$$

Using λ , we can write the updated mean as a convex combination of the previous mean and the observation,

$$\mu_{t+1} = (1 - \lambda) \mu_t + \lambda y_{t+1} \quad (3.15)$$

$$= \mu_t + \lambda (y_{t+1} - \mu_t). \quad (3.16)$$

Intuitively, λ is a form of “gain” that influences how much of the new information should be incorporated into the updated mean. For this reason, λ is also called *Kalman gain*.

The updated variance can similarly be rewritten,

$$\sigma_{t+1}^2 = \lambda \sigma_y^2 = (1 - \lambda)(\sigma_t^2 + \sigma_x^2). \quad (3.17)$$

In particular, observe that if $\mu_t = y_{t+1}$ (i.e., we observe our prediction), we have $\mu_{t+1} = \mu_t$ as there is no new information. Similarly, for $\sigma_y^2 \rightarrow \infty$ (i.e., we do not trust our observations), we have

$$\lambda \rightarrow 0, \quad \mu_{t+1} = \mu_t, \quad \sigma_{t+1}^2 = \sigma_t^2 + \sigma_x^2.$$

In contrast, for $\sigma_y^2 \rightarrow 0$, we have

$$\lambda \rightarrow 1, \quad \mu_{t+1} = y_{t+1}, \quad \sigma_{t+1}^2 = 0.$$

³ This corresponds to $F = H = I$ and a drift of 0.

Problem 3.1

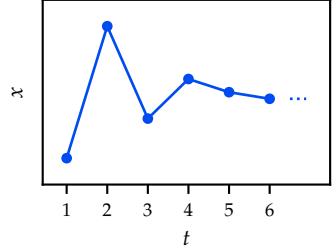


Figure 3.3: Hidden states during a random walk in one dimension.

The general formulas for the *Kalman update* follow the same logic as in the above example of a one-dimensional random walk. Given the

prior belief $\mathbf{X}_t \mid \mathbf{y}_{1:t} \sim \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$, we have

$$\mathbf{X}_{t+1} \mid \mathbf{y}_{1:t+1} \sim \mathcal{N}(\boldsymbol{\mu}_{t+1}, \boldsymbol{\Sigma}_{t+1}) \quad \text{where} \quad (3.18a)$$

$$\boldsymbol{\mu}_{t+1} \doteq \mathbf{F}\boldsymbol{\mu}_t + \mathbf{K}_{t+1}(\mathbf{y}_{t+1} - \mathbf{H}\mathbf{F}\boldsymbol{\mu}_t), \quad (3.18b)$$

$$\boldsymbol{\Sigma}_{t+1} \doteq (\mathbf{I} - \mathbf{K}_{t+1}\mathbf{H})(\mathbf{F}\boldsymbol{\Sigma}_t\mathbf{F}^\top + \boldsymbol{\Sigma}_x). \quad (3.18c)$$

Hereby, \mathbf{K}_{t+1} is the *Kalman gain*,

$$\mathbf{K}_{t+1} \doteq (\mathbf{F}\boldsymbol{\Sigma}_t\mathbf{F}^\top + \boldsymbol{\Sigma}_x)\mathbf{H}^\top(\mathbf{H}(\mathbf{F}\boldsymbol{\Sigma}_t\mathbf{F}^\top + \boldsymbol{\Sigma}_x)\mathbf{H}^\top + \boldsymbol{\Sigma}_y)^{-1} \in \mathbb{R}^{d \times m}. \quad (3.18d)$$

Note that $\boldsymbol{\Sigma}_t$ and \mathbf{K}_t can be computed offline as they are independent of the observation \mathbf{y}_{t+1} . $\mathbf{F}\boldsymbol{\mu}_t$ represents the expected state at time $t+1$, and hence, $\mathbf{H}\mathbf{F}\boldsymbol{\mu}_t$ corresponds to the expected observation. Therefore, the term $\mathbf{y}_{t+1} - \mathbf{H}\mathbf{F}\boldsymbol{\mu}_t$ measures the error in the predicted observation and the Kalman gain \mathbf{K}_{t+1} appears as a measure of relevance of the new observation compared to the prediction.

Example 3.5: Bayesian linear regression as a Kalman filter

Even though they arise from a rather different setting, it turns out that Kalman filters are a generalization of Bayesian linear regression! To see this, recall the online Bayesian linear regression algorithm from Section 2.1.3. Observe that by keeping attempting to estimate the (hidden) weights \mathbf{w}^* from sequential noisy observations \mathbf{y}_t , this algorithm performs Bayesian filtering! Moreover, we have used a Gaussian prior and likelihood. This is precisely the setting of a Kalman filter!

Concretely, we are estimating the constant (i.e., $\mathbf{F} = \mathbf{I}$, $\boldsymbol{\varepsilon} = \mathbf{0}$) hidden state $\mathbf{x}_t = \mathbf{w}^{(t)}$ with prior $\mathbf{w}^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$.

Our sensor model is time-dependent, since in each iteration we observe a different input \mathbf{x}_t . Furthermore, we only observe a scalar-valued label y_t .⁴ Formally, our sensor model is characterized by $\mathbf{h}_t = \mathbf{x}_t^\top$ and noise $\eta_t = \varepsilon_t$ with $\varepsilon_t \sim \mathcal{N}(0, \sigma_n^2)$.

You will show in ② that the Kalman update (3.18) is the online equivalent to computing the posterior of the weights in Bayesian linear regression.

⁴ That is, $m = 1$ in our general Kalman filter formulation from above.

Problem 3.2

3.2.2 Predicting

Using now that the marginal posterior of \mathbf{X}_t is a Gaussian due to the closedness properties of Gaussians, we have

$$\mathbf{X}_{t+1} \mid \mathbf{y}_{1:t} \sim \mathcal{N}(\hat{\boldsymbol{\mu}}_{t+1}, \hat{\boldsymbol{\Sigma}}_{t+1}), \quad (3.19)$$

and it suffices to compute the prediction mean $\hat{\mu}_{t+1}$ and covariance matrix $\hat{\Sigma}_{t+1}$.

For the mean,

$$\begin{aligned}\hat{\mu}_{t+1} &= \mathbb{E}[x_{t+1} | y_{1:t}] \\ &= \mathbb{E}[F\mathbf{x}_t + \boldsymbol{\varepsilon}_t | y_{1:t}] \\ &= F\mathbb{E}[\mathbf{x}_t | y_{1:t}] \\ &= F\boldsymbol{\mu}_t.\end{aligned}\tag{3.20}$$

using the motion model (3.2)

using linearity of expectation (1.20) and
 $\mathbb{E}[\boldsymbol{\varepsilon}_t] = \mathbf{0}$

using the mean of the Kalman update

For the covariance matrix,

$$\begin{aligned}\hat{\Sigma}_{t+1} &= \mathbb{E}\left[(x_{t+1} - \hat{\mu}_{t+1})(x_{t+1} - \hat{\mu}_{t+1})^\top \mid y_{1:t}\right] \\ &= F\mathbb{E}\left[(\mathbf{x}_t - \boldsymbol{\mu}_t)(\mathbf{x}_t - \boldsymbol{\mu}_t)^\top \mid y_{1:t}\right]F^\top + \mathbb{E}\left[\boldsymbol{\varepsilon}_t \boldsymbol{\varepsilon}_t^\top\right] \\ &= F\Sigma_t F^\top + \Sigma_x.\end{aligned}\tag{3.21}$$

using the definition of the covariance matrix (1.36)

using (3.20), the motion model (3.2) and that $\boldsymbol{\varepsilon}_t$ is independent of the observations

Optional Readings

Kalman filters and related models are often called *temporal models*.

For a broader look at such models, read chapter 15 of “Artificial intelligence: a modern approach” (Russell and Norvig, 2002).

Discussion

In this chapter, we have introduced Kalman filters as a special case of probabilistic filtering where probabilistic inference can be performed in closed form. Similarly to Bayesian linear regression, probabilistic inference is tractable due to assuming Gaussian priors and likelihoods. Indeed, learning linear models and Kalman filters are very closely related as seen in Example 3.5, and we will further explore this relationship in Problem 4.3. We will refer back to filtering in the second part of this manuscript when we discuss sequential decision-making with partial observability of the state space. Next, we return to the storyline on “learning” using exact probabilistic inference.

Problems

3.1. Kalman update.

Derive the predictive distribution $X_{t+1} | y_{1:t+1}$ (3.13) of the Kalman filter described in the above example using your knowledge about multivariate Gaussians from Section 1.2.3.

Hint: First compute the predictive distribution $X_{t+1} | y_{1:t}$.

3.2. Bayesian linear regression as a Kalman filter.

Recall the specific Kalman filter from Example 3.5. With this model the Kalman update (3.18) simplifies to

$$k_t = \frac{\Sigma_{t-1} x_t}{x_t^\top \Sigma_{t-1} x_t + \sigma_n^2}, \quad (3.22a)$$

$$\mu_t = \mu_{t-1} + k_t (y_t - x_t^\top \mu_{t-1}), \quad (3.22b)$$

$$\Sigma_t = \Sigma_{t-1} - k_t x_t^\top \Sigma_{t-1}, \quad (3.22c)$$

with $\mu_0 = \mathbf{0}$ and $\Sigma_0 = \sigma_p^2 I$. Note that the Kalman gain k_t is a vector in \mathbb{R}^d . We assume $\sigma_n^2 = \sigma_p^2 = 1$ for simplicity.

Prove by induction that the (μ_t, Σ_t) produced by the Kalman update are equivalent to (μ, Σ) from the posterior of Bayesian linear regression (2.10) given $x_{1:t}, y_{1:t}$. You may use that $\Sigma_t^{-1} k_t = x_t$.

Hint: In the inductive step, first prove the equivalence of Σ_t and then expand $\Sigma_t^{-1} \mu_t$ to prove the equivalence of μ_t .

3.3. Parameter estimation using Kalman filters.

Suppose that we want to estimate the value of an unknown constant π using uncorrelated measurements

$$y_t = \pi + \eta_t, \quad \eta_t \sim \mathcal{N}(0, \sigma_y^2).$$

1. How can this problem be formulated as a Kalman filter? Compute closed form expressions for the Kalman gain and the variance of the estimation error σ_t^2 in terms of t , σ_y^2 , and σ_0^2 .
2. What is the Kalman filter when $t \rightarrow \infty$?
3. Suppose that one has no prior assumptions on π , meaning that $\mu_0 = 0$ and $\sigma_0^2 \rightarrow \infty$. Which well-known estimator does the Kalman filter reduce to in this case?

4

Gaussian Processes

Let us remember our first attempt from Chapter 2 at scaling up Bayesian linear regression to nonlinear functions. We saw that we can model nonlinear functions by transforming the input space to a suitable higher-dimensional feature space, but found that this approach scales poorly if we require a large number of features. We then found something remarkable: by simply changing our perspective from a weight-space view to a function-space view, we could implement Bayesian linear regression without ever needing to compute the features explicitly. Under the function-space view, the key object describing the class of functions we can model is not the features $\phi(\mathbf{x})$, but instead the kernel function which only implicitly defines a feature space. Our key observation in this chapter is that we can therefore stop reasoning about feature spaces, and instead directly work with kernel functions that describe “reasonable” classes of functions.

We are still concerned with the problem of estimating the value of a function $f : \mathcal{X} \rightarrow \mathbb{R}$ at arbitrary points $\mathbf{x}^* \in \mathcal{X}$ given training data $\{\mathbf{x}_i, y_i\}_{i=1}^n$, where the labels are assumed to be corrupted by homoscedastic Gaussian noise with variance σ_n^2 ,

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma_n^2).$$

As in Chapter 2 on Bayesian linear regression, we denote by \mathbf{X} the design matrix (collection of training inputs) and by \mathbf{y} the vector of training labels. We will represent the unknown function value at a point $\mathbf{x} \in \mathcal{X}$ by the random variable $f_{\mathbf{x}} \doteq f(\mathbf{x})$. The collection of these random variables is then called a Gaussian process if any finite subset of them is jointly Gaussian:

Definition 4.1 (Gaussian process, GP). A *Gaussian process* is an infinite set of random variables such that any finite number of them are jointly Gaussian and such that they are consistent under marginalization.¹

¹ That is, if you take a joint distribution for n variables and marginalize out one of them, you should recover the joint distribution for the remaining $n - 1$ variables.

The fact that with a Gaussian process, any finite subset of the random variables is jointly Gaussian is the key property allowing us to perform exact probabilistic inference. Intuitively, a Gaussian process can be interpreted as a normal distribution over functions — and is therefore often called an “infinite-dimensional Gaussian”.

A Gaussian process is characterized by a *mean function* $\mu : \mathcal{X} \rightarrow \mathbb{R}$ and a *covariance function* (or *kernel function*) $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ such that for any set of points $A \doteq \{x_1, \dots, x_m\} \subseteq \mathcal{X}$, we have

$$f_A \doteq [f_{x_1} \ \dots \ f_{x_m}]^\top \sim \mathcal{N}(\boldsymbol{\mu}_A, \mathbf{K}_{AA}) \quad (4.1)$$

where

$$\boldsymbol{\mu}_A \doteq \begin{bmatrix} \mu(x_1) \\ \vdots \\ \mu(x_m) \end{bmatrix}, \quad \mathbf{K}_{AA} \doteq \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_m) \\ \vdots & \ddots & \vdots \\ k(x_m, x_1) & \dots & k(x_m, x_m) \end{bmatrix}. \quad (4.2)$$

We write $f \sim \mathcal{GP}(\mu, k)$. In particular, given a mean function, covariance function, and using the homoscedastic noise assumption,

$$y^* \mid x^* \sim \mathcal{N}(\mu(x^*), k(x^*, x^*) + \sigma_n^2). \quad (4.3)$$

Commonly, for notational simplicity, the mean function is taken to be zero. Note that for a fixed mean this is not a restriction, as we can simply apply the zero-mean Gaussian process to the difference between the mean and the observations.²

4.1 Learning and Inference

First, let us look at learning and inference in the context of Gaussian processes. With slight abuse of our previous notation, let us denote the set of observed points by $A \doteq \{x_1, \dots, x_n\}$. Given a prior $f \sim \mathcal{GP}(\mu, k)$ and the noisy observations $y_i = f(x_i) + \varepsilon_i$ with $\varepsilon_i \sim \mathcal{N}(0, \sigma_n^2)$, we can then write the joint distribution of the observations $y_{1:n}$ and the noise-free prediction f^* at a test point x^* as

$$\begin{bmatrix} \mathbf{y} \\ f^* \end{bmatrix} \mid x^*, \mathbf{X} \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\mathbf{K}}), \quad \text{where} \quad (4.4)$$

$$\tilde{\boldsymbol{\mu}} \doteq \begin{bmatrix} \boldsymbol{\mu}_A \\ \mu(x^*) \end{bmatrix}, \quad \tilde{\mathbf{K}} \doteq \begin{bmatrix} \mathbf{K}_{AA} + \sigma_n^2 \mathbf{I} & \mathbf{k}_{x^*, A} \\ \mathbf{k}_{x^*, A}^\top & k(x^*, x^*) \end{bmatrix}, \quad \mathbf{k}_{x^*, A} \doteq \begin{bmatrix} k(x, x_1) \\ \vdots \\ k(x, x_n) \end{bmatrix}. \quad (4.5)$$

Deriving the conditional distribution using (1.53), we obtain that the Gaussian process posterior is given by

$$f \mid x_{1:n}, y_{1:n} \sim \mathcal{GP}(\boldsymbol{\mu}', \mathbf{k}'), \quad \text{where} \quad (4.6)$$

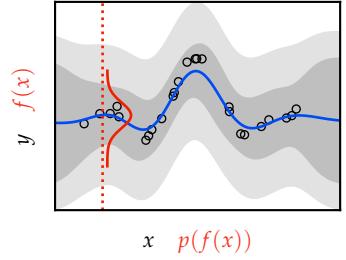


Figure 4.1: A Gaussian process can be interpreted as an infinite-dimensional Gaussian over functions. At any location x in the domain, this yields a distribution over values $f(x)$ shown in red. The blue line corresponds to the MAP estimate (i.e., mean function of the Gaussian process), the dark gray region corresponds to the epistemic uncertainty and the light gray region denotes the additional aleatoric uncertainty.

² For alternative ways of representing a mean function, refer to section 2.7 of “Gaussian processes for machine learning” (Williams and Rasmussen, 2006).

$$\mu'(\mathbf{x}) \doteq \mu(\mathbf{x}) + \mathbf{k}_{\mathbf{x}, A}^\top (\mathbf{K}_{AA} + \sigma_n^2 \mathbf{I})^{-1} (\mathbf{y}_A - \boldsymbol{\mu}_A), \quad (4.7)$$

$$k'(\mathbf{x}, \mathbf{x}') \doteq k(\mathbf{x}, \mathbf{x}') - \mathbf{k}_{\mathbf{x}, A}^\top (\mathbf{K}_{AA} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}_{\mathbf{x}', A}. \quad (4.8)$$

Observe that analogously to Bayesian linear regression, the posterior covariance can only decrease when conditioning on additional data, and is independent of the observations y_i .

We already studied inference in the function-space view of Bayesian linear regression, but did not make the predictive posterior explicit. Using Equation (4.6), the predictive posterior at \mathbf{x}^* is simply

$$f^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n} \sim \mathcal{N}(\mu'(\mathbf{x}^*), k'(\mathbf{x}^*, \mathbf{x}^*)). \quad (4.9)$$

4.2 Sampling

Often, we are not interested in the full predictive posterior distribution, but merely want to obtain samples of our Gaussian process model. We will briefly examine two approaches.

1. For the first approach, consider a discretized subset of points

$$\mathbf{f} \doteq [f_1, \dots, f_n]$$

that we want to sample.³ Note that $\mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$. We have already seen in Equation (1.54) that

$$\mathbf{f} = \mathbf{K}^{1/2} \boldsymbol{\varepsilon} + \boldsymbol{\mu} \quad (4.10)$$

where $\mathbf{K}^{1/2}$ is the square root of \mathbf{K} and $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is standard Gaussian noise.⁴ However, computing the square root of \mathbf{K} takes $O(n^3)$ time.

2. For the second approach, recall the product rule (1.11),

$$p(f_1, \dots, f_n) = \prod_{i=1}^n p(f_i | f_{1:i-1}).$$

That is the joint distribution factorizes neatly into a product where each factor only depends on the “outcomes” of preceding factors. We can therefore obtain samples one-by-one, each time conditioning on one more observation:

$$\begin{aligned} f_1 &\sim p(f_1) \\ f_2 &\sim p(f_2 | f_1) \\ f_3 &\sim p(f_3 | f_1, f_2) \\ &\vdots \end{aligned} \quad (4.11)$$

This general approach is known as *forward sampling*. Due to the matrix inverse in the formula of the GP posterior (4.6), this approach also takes $O(n^3)$ time.

We will discuss more efficient approximate sampling methods in Section 4.5.

³ For example, if we want to render the function, the length of this vector could be guided by the screen resolution.

⁴ We discuss square roots of matrices in Appendix A.2.

4.3 Kernel Functions

We have seen that kernel functions are the key object describing the class of functions a Gaussian process can model. Depending on the kernel function, the “shape” of functions that are realized from a Gaussian process varies greatly. Let us recap briefly from Section 2.4 what a kernel function is:

Definition 4.2 (Kernel function). A *kernel function* $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ satisfies

- $k(x, x') = k(x', x)$ for any $x, x' \in \mathcal{X}$ (symmetry), and
- K_{AA} is positive semi-definite for any $A \subseteq \mathcal{X}$.

The two defining conditions ensure that for any $A \subseteq \mathcal{X}$, K_{AA} is a valid covariance matrix. We say that a kernel function is *positive definite* if K_{AA} is positive definite for any $A \subseteq \mathcal{X}$.

Intuitively, the kernel function evaluated at locations x and x' describes how $f(x)$ and $f(x')$ are related, which we can express formally as

$$k(x, x') = \text{Cov}[f(x), f(x')]. \quad (4.12)$$

If x and x' are “close”, then $f(x)$ and $f(x')$ are usually taken to be positively correlated, encoding a “smooth” function.

In the following, we will discuss some of the most common kernel functions, how they can be combined to create “new” kernels, and how we can characterize the class of functions they can model.

4.3.1 Common Kernels

First, we look into some of the most commonly used kernels. Often an additional factor σ^2 (*output scale*) is added, which we assume here to be 1 for simplicity.

1. The *linear kernel* is defined as

$$k(x, x'; \phi) \doteq \phi(x)^\top \phi(x') \quad (4.13)$$

where ϕ is a nonlinear transformation as introduced in Section 2.3 or the identity.

Remark 4.3: GPs with linear kernel and BLR

A Gaussian process with a linear kernel is equivalent to Bayesian linear regression. This follows directly from the function-space view of Bayesian linear regression (see Section 2.4) and comparing the derived kernel function (2.20) with the definition of the linear kernel (4.13).

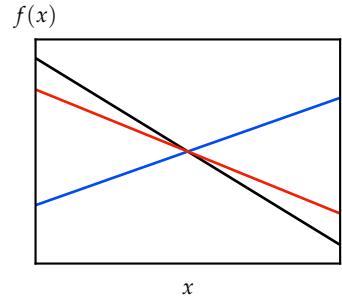


Figure 4.2: Functions sampled according to a Gaussian process with a linear kernel and $\phi = \text{id}$.

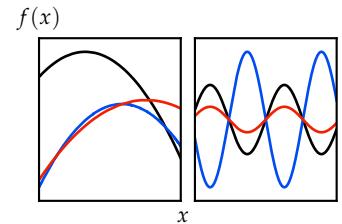
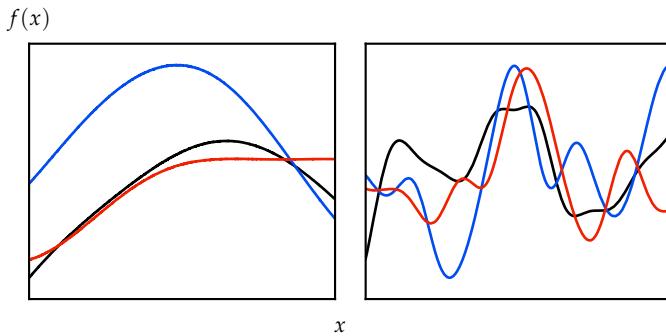


Figure 4.3: Functions sampled according to a Gaussian process with a linear kernel and $\phi(x) = [1, x, x^2]$ (left) and $\phi(x) = \sin(x)$ (right).

2. The *Gaussian kernel* (also known as *squared exponential kernel* or *radial basis function (RBF) kernel*) is defined as

$$k(\mathbf{x}, \mathbf{x}'; h) \doteq \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2h^2}\right) \quad (4.14)$$

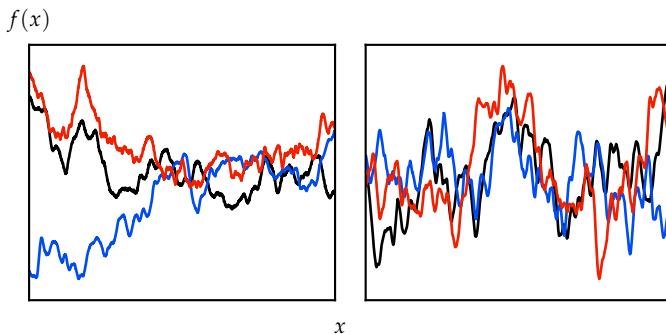
where h is its *length scale*. The larger the length scale h , the smoother the resulting functions.⁵ Furthermore, it turns out that the feature space (think back to Section 2.4!) corresponding to the Gaussian kernel is “infinitely dimensional”, as you will show in (?). So the Gaussian kernel already encodes a function class that we were not able to model under the weight-space view of Bayesian linear regression.



3. The *Laplace kernel* (also known as *exponential kernel*) is defined as

$$k(\mathbf{x}, \mathbf{x}'; h) \doteq \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2}{h}\right). \quad (4.15)$$

As can be seen in Figure 4.7, samples from a GP with Laplace kernel are non-smooth as opposed to the samples from a GP with Gaussian kernel.



4. The *Matérn kernel* trades the smoothness of the Gaussian and the Laplace kernels. As such, it is frequently used in practice to model

⁵ As the length scale is increased, the exponent of the exponential increases, resulting in a higher dependency between locations.

Problem 4.1

Figure 4.4: Functions sampled according to a Gaussian process with a Gaussian kernel and length scales $h = 5$ (left) and $h = 1$ (right).

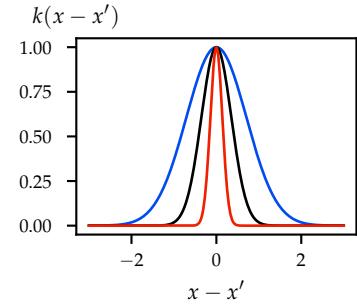


Figure 4.5: Gaussian kernel with length scales $h = 1$, $h = 0.5$, and $h = 0.2$.

Figure 4.7: Functions sampled according to a Gaussian process with a Laplace kernel and length scales $h = 10000$ (left) and $h = 10$ (right).

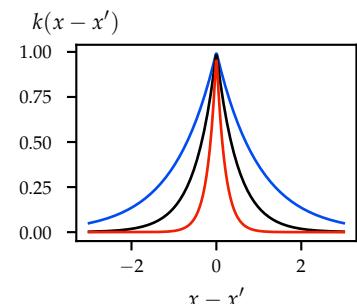


Figure 4.6: Laplace kernel with length scales $h = 1$, $h = 0.5$, and $h = 0.2$.

“real world” functions that are relatively smooth. It is defined as

$$k(\mathbf{x}, \mathbf{x}'; \nu, h) \doteq \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|\mathbf{x} - \mathbf{x}'\|_2}{h} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu} \|\mathbf{x} - \mathbf{x}'\|_2}{h} \right) \quad (4.16)$$

where Γ is the Gamma function, K_ν the modified Bessel function of the second kind, and h a length scale parameter. For $\nu = 1/2$, the Matérn kernel is equivalent to the Laplace kernel. For $\nu \rightarrow \infty$, the Matérn kernel is equivalent to the Gaussian kernel. The resulting functions are $\lceil \nu \rceil - 1$ times mean square differentiable.⁶ In particular, GPs with a Gaussian kernel are infinitely many times mean square differentiable whereas GPs with a Laplace kernel are mean square continuous but not mean square differentiable.

4.3.2 Composing Kernels

Given two kernels $k_1 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and $k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, they can be composed to obtain a new kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ in the following ways:

- $k(\mathbf{x}, \mathbf{x}') \doteq k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$,
- $k(\mathbf{x}, \mathbf{x}') \doteq k_1(\mathbf{x}, \mathbf{x}') \cdot k_2(\mathbf{x}, \mathbf{x}')$,
- $k(\mathbf{x}, \mathbf{x}') \doteq c \cdot k_1(\mathbf{x}, \mathbf{x}')$ for any $c > 0$,
- $k(\mathbf{x}, \mathbf{x}') \doteq f(k_1(\mathbf{x}, \mathbf{x}'))$ for any polynomial f with positive coefficients or $f = \exp$.

For example, the additive structure of a function $f(\mathbf{x}) \doteq f_1(\mathbf{x}) + f_2(\mathbf{x})$ can be easily encoded in GP models. Suppose that $f_1 \sim \mathcal{GP}(\mu_1, k_1)$ and $f_2 \sim \mathcal{GP}(\mu_2, k_2)$, then the distribution of the sum of those two functions $f = f_1 + f_2 \sim \mathcal{GP}(\mu_1 + \mu_2, k_1 + k_2)$ is another GP.⁷

Whereas the addition of two kernels k_1 and k_2 can be thought of as an *OR* operation (i.e., the kernel has high value if either k_1 or k_2 have high value), the multiplication of k_1 and k_2 can be thought of as an *AND* operation (i.e., the kernel has high value if both k_1 and k_2 have high value). For example, the product of two linear kernels results in functions which are quadratic.

As mentioned previously, the constant c of a scaled kernel function $k'(\mathbf{x}, \mathbf{x}') \doteq c \cdot k(\mathbf{x}, \mathbf{x}')$ is generally called the *output scale* of a kernel, and it scales the variance $\text{Var}[f(\mathbf{x})] = c \cdot k(\mathbf{x}, \mathbf{x})$ of the predictions $f(\mathbf{x})$ from $\mathcal{GP}(\mu, k')$.

⁶ Refer to Remark A.12 for the definitions of mean square continuity and differentiability.

⁷ We use $f \doteq f_1 + f_2$ to denote the function $f(\cdot) = f_1(\cdot) + f_2(\cdot)$.

Optional Readings

For a broader introduction to how kernels can be used and combined to model certain classes of functions, read

- chapter 2 of “Automatic model construction with Gaussian processes” (Duvenaud, 2014) also known as the “kernel cookbook”,
- chapter 4 of “Gaussian processes for machine learning” (Williams and Rasmussen, 2006).

4.3.3 Stationarity and Isotropy

Kernel functions are commonly classified according to two properties:

Definition 4.4 (Stationarity and isotropy). A kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is called

- *stationary* (or *shift-invariant*) if there exists a function \tilde{k} such that $\tilde{k}(\mathbf{x} - \mathbf{x}') = k(\mathbf{x}, \mathbf{x}')$, and
- *isotropic* if there exists a function \tilde{k} such that $\tilde{k}(\|\mathbf{x} - \mathbf{x}'\|) = k(\mathbf{x}, \mathbf{x}')$ with $\|\cdot\|$ any norm.

Note that stationarity is a necessary condition for isotropy. In other words, isotropy implies stationarity.

Example 4.5: Stationarity and isotropy of kernels

	stationary	isotropic
linear kernel	no	no
Gaussian kernel	yes	yes
$k(\mathbf{x}, \mathbf{x}') \doteq \exp(-\ \mathbf{x} - \mathbf{x}'\ _M^2)$ where M is positive semi-definite	yes	no

$\|\cdot\|_M$ denotes the Mahalanobis norm induced by matrix M

For $\mathbf{x}' = \mathbf{x}$, stationarity implies that the kernel must only depend on $\mathbf{0}$. In other words, a stationary kernel must depend on relative locations only. This is clearly not the case for the linear kernel, which depends on the absolute locations of \mathbf{x} and \mathbf{x}' . Therefore, the linear kernel cannot be isotropic either.

For the Gaussian kernel, isotropy follows immediately from its definition.

The last kernel is clearly stationary by definition, but not isotropic for general matrices M . Note that for $M = I$ it is indeed isotropic.

Stationarity encodes the idea that relative location matters more than absolute location: the process “looks the same” no matter where we shift it in the input space. This is often appropriate when we believe the same statistical behavior holds across the entire domain (e.g., no region is special). Isotropy goes one step further by requiring that

the kernel depends only on the distance between points, so that all directions in the space are treated equally. In other words, there is no preferred orientation or axis. This is especially useful in settings where we expect uniform behavior in every direction (as with the Gaussian kernel). Such kernels are simpler to specify and interpret since we only need a single “scale” (like a length scale) rather than multiple parameters or directions.

4.3.4 Reproducing Kernel Hilbert Spaces

We can characterize the precise class of functions that can be modeled by a Gaussian process with a given kernel function. This corresponding function space is called a *reproducing kernel Hilbert space* (RKHS), and we will discuss it briefly in this section.

Recall that Gaussian processes keep track of a posterior distribution $f \mid \mathbf{x}_{1:n}, y_{1:n}$ over functions. We will in fact show later that the corresponding MAP estimate \hat{f} corresponds to the solution to a regularized optimization problem in the RKHS space of functions. This duality is similar to the duality between the MAP estimate of Bayesian linear regression and ridge regression we observed in Chapter 2. So what is the reproducing kernel Hilbert space of a kernel function k ?

Definition 4.6 (Reproducing kernel Hilbert space, RKHS). Given a kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, its corresponding *reproducing kernel Hilbert space* is the space of functions f defined as

$$\mathcal{H}_k(\mathcal{X}) \doteq \left\{ f(\cdot) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \cdot) : n \in \mathbb{N}, \mathbf{x}_i \in \mathcal{X}, \alpha_i \in \mathbb{R} \right\}. \quad (4.17)$$

The inner product of the RKHS is defined as

$$\langle f, g \rangle_k \doteq \sum_{i=1}^n \sum_{j=1}^{n'} \alpha_i \alpha'_j k(\mathbf{x}_i, \mathbf{x}'_j), \quad (4.18)$$

where $g(\cdot) = \sum_{j=1}^{n'} \alpha'_j k(\mathbf{x}'_j, \cdot)$, and induces the norm $\|f\|_k = \sqrt{\langle f, f \rangle_k}$. You can think of the norm as measuring the “smoothness” or “complexity” of f . $\textcircled{?}$

Problem 4.4 (2)

It is straightforward to check that for all $\mathbf{x} \in \mathcal{X}$, $k(\mathbf{x}, \cdot) \in \mathcal{H}_k(\mathcal{X})$. Moreover, the RKHS inner product $\langle \cdot, \cdot \rangle_k$ satisfies for all $\mathbf{x} \in \mathcal{X}$ and $f \in \mathcal{H}_k(\mathcal{X})$ that $f(\mathbf{x}) = \langle f(\cdot), k(\mathbf{x}, \cdot) \rangle_k$ which is also known as the *reproducing property* $\textcircled{?}$. That is, evaluations of RKHS functions f are inner products in $\mathcal{H}_k(\mathcal{X})$ parameterized by the “feature map” $k(\mathbf{x}, \cdot)$.

Problem 4.4 (1)

The *representer theorem* (Schölkopf et al., 2001) characterizes the solution to regularized optimization problems in RKHSs:

Theorem 4.7 (Representer theorem). **?** Let k be a kernel and let $\lambda > 0$. For $f \in \mathcal{H}_k(\mathcal{X})$ and training data $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$, let $\mathcal{L}(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) \in \mathbb{R} \cup \{\infty\}$ denote any loss function which depends on f only through its evaluation at the training points. Then, any minimizer

$$\hat{f} \in \arg \min_{f \in \mathcal{H}_k(\mathcal{X})} \mathcal{L}(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) + \lambda \|f\|_k^2 \quad (4.19)$$

admits a representation of the form

$$\hat{f}(\mathbf{x}) = \hat{\alpha}^\top \mathbf{k}_{\mathbf{x}, \{\mathbf{x}_i\}_{i=1}^n} = \sum_{i=1}^n \hat{\alpha}_i k(\mathbf{x}, \mathbf{x}_i) \quad \text{for some } \hat{\alpha} \in \mathbb{R}^n. \quad (4.20)$$

This statement is remarkable: the solutions to general regularized optimization problems over the generally infinite-dimensional space of functions $\mathcal{H}_k(\mathcal{X})$ can be represented as a linear combination of the kernel functions evaluated at the training points. The representer theorem can be used to show that the MAP estimate of a Gaussian process corresponds to the solution of a regularized linear regression problem in the RKHS of the kernel function, namely, **?**

Problem 4.5

Problem 4.6

$$\hat{f} \doteq \arg \min_{f \in \mathcal{H}_k(\mathcal{X})} -\log p(y_{1:n} | \mathbf{x}_{1:n}, f) + \frac{1}{2} \|f\|_k^2. \quad (4.21)$$

Here, the first term corresponds to the likelihood, measuring the “quality of fit”. The regularization term limits the “complexity” of \hat{f} . Regularization is necessary to prevent overfitting since in an expressive RKHSs, there may be many functions that interpolate the training data perfectly. This shows the close link between Gaussian process regression and Bayesian linear regression, with the kernel function k generalizing the inner product of feature maps to feature spaces of possibly “infinite dimensionality”. Because solutions can be represented as linear combinations of kernel evaluations at the training points, Gaussian processes remain computationally tractable even though they can model functions over “infinite-dimensional” feature spaces.

4.4 Model Selection

We have not yet discussed how to pick the hyperparameters θ (e.g., parameters of kernels). A common technique in supervised learning is to select hyperparameters θ , such that the resulting function estimate \hat{f}_θ leads to the most accurate predictions on hold-out validation data. After reviewing this approach, we contrast it with a probabilistic approach to model selection, which avoids using point estimates of \hat{f}_θ and rather utilizes the full posterior.

4.4.1 Optimizing Validation Set Performance

A common approach to model selection is to split our data \mathcal{D} into separate training set $\mathcal{D}^{\text{train}} \doteq \{(x_i^{\text{train}}, y_i^{\text{train}})\}_{i=1}^n$ and validation sets $\mathcal{D}^{\text{val}} \doteq \{(x_i^{\text{val}}, y_i^{\text{val}})\}_{i=1}^m$. We then optimize the model for a parameter candidate θ_j using the training set. This is usually done by picking a point estimate (like the MAP estimate),

$$\hat{f}_j \doteq \arg \max_f p(f \mid x_{1:n}^{\text{train}}, y_{1:n}^{\text{train}}). \quad (4.22)$$

Then, we score θ_j according to the performance of \hat{f}_j on the validation set,

$$\hat{\theta} \doteq \arg \max_{\theta_j} p(y_{1:m}^{\text{val}} \mid x_{1:m}^{\text{val}}, \hat{f}_j). \quad (4.23)$$

This ensures that \hat{f}_j does not depend on \mathcal{D}^{val} .

Remark 4.8: Approximating population risk

Why is it useful to separate the data into a training and a validation set? Recall from Appendix A.3.5 that minimizing the empirical risk without separating training and validation data may lead to overfitting as both the loss and \hat{f}_j depend on the same data \mathcal{D} . In contrast, using independent training and validation sets, \hat{f}_j does not depend on \mathcal{D}^{val} , and we have that

$$\frac{1}{m} \sum_{i=1}^m \ell(y_i^{\text{val}} \mid x_i^{\text{val}}, \hat{f}_j) \approx \mathbb{E}_{(x,y) \sim \mathcal{P}} [\ell(y \mid x, \hat{f}_j)], \quad (4.24)$$

using Monte Carlo sampling.⁸ In words, for reasonably large m , minimizing the empirical risk as we do in Equation (4.23) approximates minimizing the population risk.

While this approach often is quite effective at preventing overfitting as compared to using the same data for training and picking $\hat{\theta}$, it still collapses the uncertainty in f into a point estimate. Can we do better?

4.4.2 Maximizing the Marginal Likelihood

We have already seen for Bayesian linear regression, that picking a point estimate loses a lot of information. Instead of optimizing the effects of θ for a specific point estimate \hat{f} of the model f , *maximizing the marginal likelihood* optimizes the effects of θ across all realizations of f . In this approach, we obtain our hyperparameter estimate via

$$\hat{\theta}_{\text{MLE}} \doteq \arg \max_{\theta} p(y_{1:n} \mid x_{1:n}, \theta) \quad (4.25)$$

⁸ We generally assume $\mathcal{D} \stackrel{\text{iid}}{\sim} \mathcal{P}$, in particular, we assume that the individual samples of the data are i.i.d.. Recall that in this setting, Hoeffding's inequality (A.41) can be used to gauge how large m should be.

using the definition of marginal likelihood in Bayes' rule (1.45)

$$\begin{aligned}
&= \arg \max_{\theta} \int p(y_{1:n}, f | x_{1:n}, \theta) df \\
&= \arg \max_{\theta} \int p(y_{1:n} | x_{1:n}, f, \theta) p(f | \theta) df. \quad (4.26)
\end{aligned}$$

Remarkably, this approach typically avoids overfitting even though we do not use a separate training and validation set. The following table provides an intuitive argument for why maximizing the marginal likelihood is a good strategy.

	likelihood	prior
“underfit” model (too simple θ)	small for “almost all” f	large
“overfit” model (too complex θ)	large for “few” f small for “most” f	small
“just right”	moderate for “many” f	moderate

For an “underfit” model, the likelihood is mostly small as the data cannot be well described, while the prior is large as there are “fewer” functions to choose from. For an “overfit” model, the likelihood is large for “some” functions (which would be picked if we were only minimizing the training error and not doing cross validation) but small for “most” functions. The prior is small, as the probability mass has to be distributed among “more” functions. Thus, in both cases, one term in the product will be small. Hence, maximizing the marginal likelihood naturally encourages trading between a large likelihood and a large prior.

In the context of Gaussian process regression, recall from Equation (4.3) that

$$y_{1:n} | x_{1:n}, \theta \sim \mathcal{N}(\mathbf{0}, K_{f,\theta} + \sigma_n^2 I) \quad (4.27)$$

where $K_{f,\theta}$ denotes the kernel matrix at the inputs $x_{1:n}$ depending on the kernel function parameterized by θ . We write $K_{y,\theta} \doteq K_{f,\theta} + \sigma_n^2 I$. Continuing from Equation (4.25), we obtain

$$\begin{aligned}
\hat{\theta}_{\text{MLE}} &= \arg \max_{\theta} \mathcal{N}(y; \mathbf{0}, K_{y,\theta}) \\
&= \arg \min_{\theta} \frac{1}{2} y^\top K_{y,\theta}^{-1} y + \frac{1}{2} \log \det(K_{y,\theta}) + \frac{n}{2} \log 2\pi \quad (4.28)
\end{aligned}$$

$$\begin{aligned}
&= \arg \min_{\theta} \frac{1}{2} y^\top K_{y,\theta}^{-1} y + \frac{1}{2} \log \det(K_{y,\theta}) \quad (4.29)
\end{aligned}$$

The first term of the optimization objective describes the “goodness of fit” (i.e., the “alignment” of y with $K_{y,\theta}$). The second term characterizes the “volume” of the model class. Thus, this optimization naturally trades the aforementioned objectives.

by conditioning on f using the sum rule (1.7)

using the product rule (1.11)

Table 4.1: The table gives an intuitive explanation of effects of parameter choices θ on the marginal likelihood. Note that words in quotation marks refer to intuitive quantities, as we have infinitely many realizations of f .

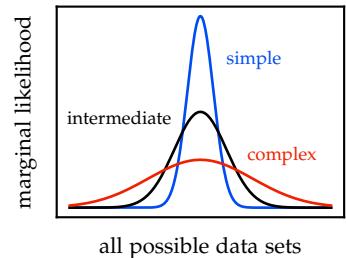


Figure 4.8: A schematic illustration of the marginal likelihood of a simple, intermediate, and complex model across all possible data sets.

taking the negative logarithm

the last term is independent of θ

Marginal likelihood maximization is an empirical Bayes method. Often it is simply referred to as *empirical Bayes*. It also has the nice property that the gradient of its objective (the MLL loss) can be expressed in closed-form (??),

$$\frac{\partial}{\partial \theta_j} \log p(y_{1:n} | \mathbf{x}_{1:n}, \boldsymbol{\theta}) = \frac{1}{2} \text{tr} \left((\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - \mathbf{K}_{y,\boldsymbol{\theta}}^{-1}) \frac{\partial \mathbf{K}_{y,\boldsymbol{\theta}}}{\partial \theta_j} \right) \quad (4.30)$$

where $\boldsymbol{\alpha} \doteq \mathbf{K}_{y,\boldsymbol{\theta}}^{-1} \mathbf{y}$ and $\text{tr}(M)$ is the trace of a matrix M . This optimization problem is, in general, non-convex. Figure 4.10 gives an example of two local optima according to empirical Bayes.

Taking a step back, observe that taking a probabilistic perspective on model selection naturally led us to consider all realizations of our model f instead of using point estimates. However, we are still using point estimates for our model parameters $\boldsymbol{\theta}$. Continuing on our probabilistic adventure, we could place a prior $p(\boldsymbol{\theta})$ on them too.⁹ We could use it to obtain the MAP estimate (still a point estimate!) which adds an additional regularization term

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} \doteq \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | \mathbf{x}_{1:n}, y_{1:n}) \quad (4.31)$$

$$= \arg \min_{\boldsymbol{\theta}} -\log p(\boldsymbol{\theta}) - \log p(y_{1:n} | \mathbf{x}_{1:n}, \boldsymbol{\theta}). \quad (4.32)$$

An alternative approach is to consider the full posterior distribution over parameters $\boldsymbol{\theta}$. The resulting predictive distribution is, however, intractable,

$$p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}) = \int \int p(y^* | \mathbf{x}^*, f) \cdot p(f | \mathbf{x}_{1:n}, y_{1:n}, \boldsymbol{\theta}) \cdot p(\boldsymbol{\theta}) df d\boldsymbol{\theta}. \quad (4.33)$$

Recall that as the mode of Gaussians coincides with their mean, the MAP estimate corresponds to the mean of the predictive posterior.

As a final note, observe that in principle, there is nothing stopping us from descending deeper in the probabilistic hierarchy. The prior on the model parameters $\boldsymbol{\theta}$ is likely to have parameters too. Ultimately, we need to break out of this hierarchy of dependencies and choose a prior.

4.5 Approximations

To learn a Gaussian process, we need to invert $n \times n$ matrices, hence the computational cost is $O(n^3)$. Compare this to Bayesian linear regression which allows us to learn a regression model in $O(nd^2)$ time (even online) where d is the feature dimension. It is therefore natural to look for ways of approximating a Gaussian process.

Problem 4.7

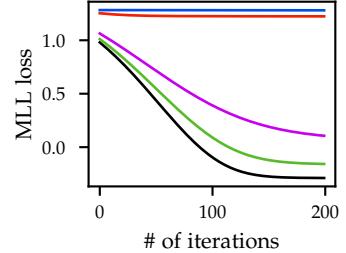


Figure 4.9: An example of model selection by maximizing the log likelihood (without hyperpriors) using a **linear**, **quadratic**, **Laplace**, **Matérn** ($v = 3/2$), and **Gaussian** kernel, respectively. They are used to learn the function

$$x \mapsto \frac{\sin(x)}{x} + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 0.01)$$

using SGD with learning rate 0.1.

⁹ Such a prior is called *hyperprior*.

using Bayes' rule (1.45) and then taking the negative logarithm

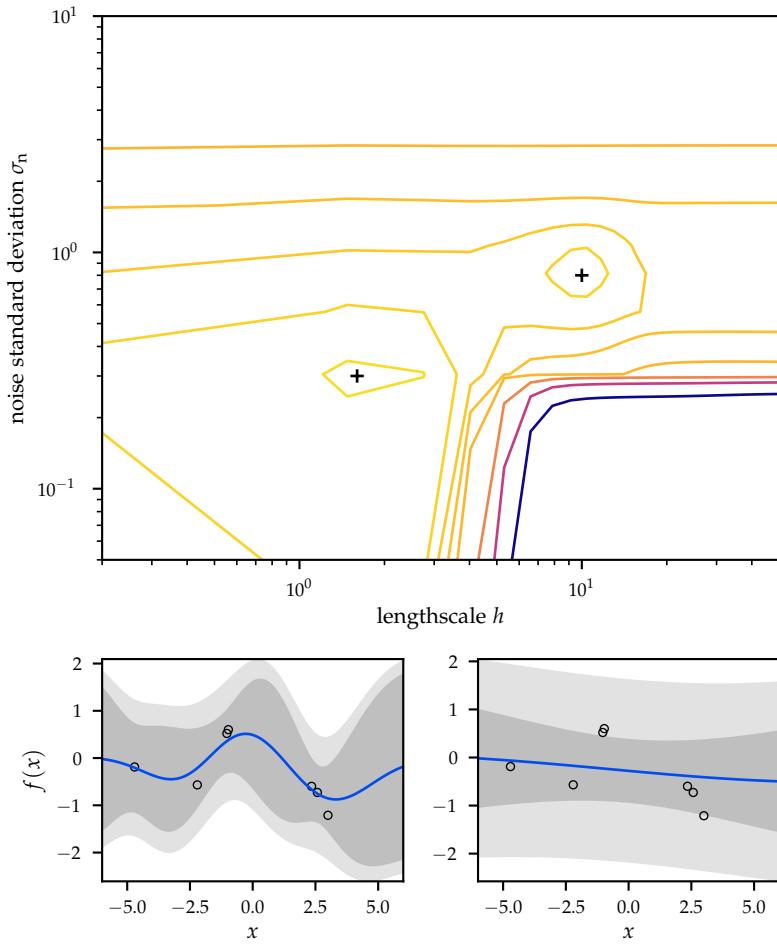


Figure 4.10: The top plot shows contour lines of an empirical Bayes with two local optima. The bottom two plots show the Gaussian processes corresponding to the two optimal models. The left model with smaller lengthscale is chosen within a more flexible class of models, while the right model explains more observations through noise. Adapted from figure 5.5 of “Gaussian processes for machine learning” (Williams and Rasmussen, 2006).

4.5.1 Local Methods

Recall that during forward sampling, we had to condition on a larger and larger number of previous samples. When sampling at a location x , a very simple approximation is to only condition on those samples x' that are “close” (where $|k(x, x')| \geq \tau$ for some $\tau > 0$). Essentially, this method “cuts off the tails” of the kernel function k . However, τ has to be chosen carefully as if τ is chosen too large, samples become essentially independent.

This is one example of a *sparse approximation* of a Gaussian process. We will discuss more advanced sparse approximations known as “inducing point methods” in Section 4.5.3.

4.5.2 Kernel Function Approximation

Another method is to approximate the kernel function directly. The idea is to construct a “low-dimensional” feature map $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ that approximates the kernel,

$$k(x, x') \approx \phi(x)^\top \phi(x'). \quad (4.34)$$

Then, we can apply Bayesian linear regression, resulting in a time complexity of $O(nm^2 + m^3)$.

One example of this approach are *random Fourier features*, which we will discuss in the following.

Remark 4.9: Fourier transform

First, let us remind ourselves of Fourier transformations. The Fourier transform is a method of decomposing frequencies into their individual components.

Recall *Euler's formula* which states that for any $x \in \mathbb{R}$,

$$e^{ix} = \cos x + i \sin x \quad (4.35)$$

where i is the imaginary unit of complex numbers. The formula is illustrated in Figure 4.11. Note that $e^{-i2\pi x}$ corresponds to rotating clockwise around the unit circle in \mathbb{R}^2 — completing a rotation whenever $x \in \mathbb{R}$ reaches the next natural number.

We can scale x by a frequency ξ : $e^{-i2\pi\xi x}$. If $x \in \mathbb{R}^d$, we can also scale each component j of x by a different frequency $\xi(j)$. Multiplying a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ with the rotation around the unit circle with given frequencies ξ , yields a quantity that describes the

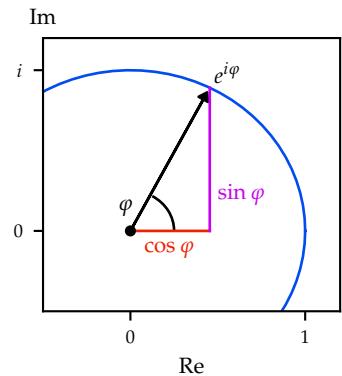


Figure 4.11: Illustration of Euler's formula. It can be seen that $e^{i\varphi}$ corresponds to a (counter-clockwise) rotation on the unit circle as φ varies from 0 to 2π .

amplitude of the frequencies ξ ,

$$\hat{f}(\xi) \doteq \int_{\mathbb{R}^d} f(x) e^{-i2\pi\xi^\top x} dx. \quad (4.36)$$

\hat{f} is called the *Fourier transform* of f . f is called the *inverse Fourier transform* of \hat{f} , and can be computed using

$$f(x) = \int_{\mathbb{R}^d} \hat{f}(\xi) e^{i2\pi\xi^\top x} d\xi. \quad (4.37)$$

It is common to write $\omega \doteq 2\pi\xi$. See Figure 4.12 for an example.

Refer to “But what is the Fourier Transform? A visual introduction” (Sanderson, 2018) for a visual introduction.

Because a stationary kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ can be interpreted as a function in one variable, it has an associated Fourier transform which we denote by $p(\omega)$. That is,

$$k(x - x') = \int_{\mathbb{R}^d} p(\omega) e^{i\omega^\top (x - x')} d\omega. \quad (4.38)$$

Fact 4.10 (Bochner’s theorem). *A continuous stationary kernel on \mathbb{R}^d is positive definite if and only if its Fourier transform $p(\omega)$ is non-negative.*

Bochner’s theorem implies that when a continuous and stationary kernel is positive definite and scaled appropriately, its Fourier transform $p(\omega)$ is a proper probability distribution. In this case, $p(\omega)$ is called the *spectral density* of the kernel k .

Remark 4.11: Eigenvalue spectrum of stationary kernels

When a kernel k is stationary (i.e., a univariate function of $x - x'$), its eigenfunctions (with respect to the usual Lebesgue measure) turn out to be the complex exponentials $\exp(i\omega^\top (x - x'))$. In simpler terms, you can think of these exponentials as “building blocks” at different frequencies ω . The spectral density $p(\omega)$ associated with the kernel tells you how strongly each frequency contributes, i.e., how large the corresponding eigenvalue is.

A key insight of this analysis is that the rate at which these magnitudes $p(\omega)$ decay with increasing frequency ω reveals the smoothness of the processes governed by the kernel. If a kernel allocates more “power” to high frequencies (meaning the spectral density decays slowly), the resulting processes will appear “rougher”. Conversely, if high-frequency components are suppressed, the process will appear “smoother”.

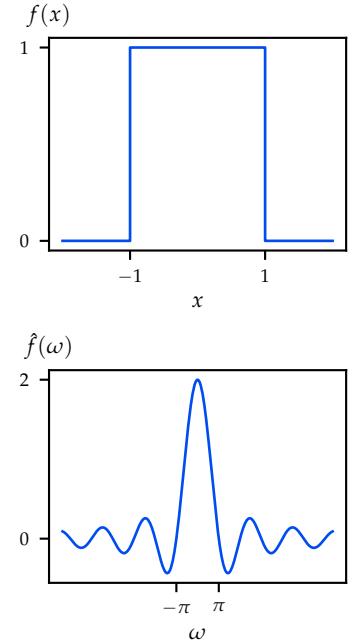


Figure 4.12: The Fourier transform of a rectangular pulse,

$$f(x) \doteq \begin{cases} 1 & x \in [-1, 1] \\ 0 & \text{otherwise,} \end{cases}$$

is given by

$$\begin{aligned} \hat{f}(\omega) &= \int_{-1}^1 e^{-i\omega x} dx = \frac{1}{i\omega} (e^{i\omega} - e^{-i\omega}) \\ &= \frac{2\sin(\omega)}{\omega}. \end{aligned}$$

For an in-depth introduction to the eigenfunction analysis of kernels, refer to section 4.3 of “Gaussian processes for machine learning” (Williams and Rasmussen, 2006).

Example 4.12: Spectral density of the Gaussian kernel

The Gaussian kernel with length scale h has the spectral density

$$\begin{aligned} p(\omega) &= \int_{\mathbb{R}^d} k(x - x'; h) e^{-i\omega^\top(x - x')} d(x - x') \\ &= \int_{\mathbb{R}^d} \exp\left(-\frac{\|x\|_2^2}{2h^2} - i\omega^\top x\right) dx \\ &= (2h^2\pi)^{d/2} \exp\left(-h^2\frac{\|\omega\|_2^2}{2}\right). \end{aligned} \quad (4.39)$$

using the definition of the Fourier transform (4.36)

using the definition of the Gaussian kernel (4.14)

The key idea is now to interpret the kernel as an expectation,

$$\begin{aligned} k(x - x') &= \int_{\mathbb{R}^d} p(\omega) e^{i\omega^\top(x - x')} d\omega \\ &= \mathbb{E}_{\omega \sim p} [e^{i\omega^\top(x - x')}] \\ &= \mathbb{E}_{\omega \sim p} [\cos(\omega^\top x - \omega^\top x') + i \sin(\omega^\top x - \omega^\top x')]. \end{aligned}$$

from Equation (4.38)

by the definition of expectation (1.19)

using Euler’s formula (4.35)

Observe that as both k and p are real, convergence of the integral implies $\mathbb{E}_{\omega \sim p} [\sin(\omega^\top x - \omega^\top x')] = 0$. Hence,

$$\begin{aligned} &= \mathbb{E}_{\omega \sim p} [\cos(\omega^\top x - \omega^\top x')] \\ &= \mathbb{E}_{\omega \sim p} \mathbb{E}_{b \sim \text{Unif}([0, 2\pi])} [\cos((\omega^\top x + b) - (\omega^\top x' + b))] \\ &= \mathbb{E}_{\omega \sim p} \mathbb{E}_{b \sim \text{Unif}([0, 2\pi])} [\cos(\omega^\top x + b) \cos(\omega^\top x' + b) \\ &\quad + \sin(\omega^\top x + b) \sin(\omega^\top x' + b)] \\ &= \mathbb{E}_{\omega \sim p} \mathbb{E}_{b \sim \text{Unif}([0, 2\pi])} [2 \cos(\omega^\top x + b) \cos(\omega^\top x' + b)] \\ &= \mathbb{E}_{\omega \sim p, b \sim \text{Unif}([0, 2\pi])} [z_{\omega, b}(x) \cdot z_{\omega, b}(x')] \end{aligned} \quad (4.40)$$

expanding with $b - b$

using the angle subtraction identity,
 $\cos(\alpha - \beta) = \cos \alpha \cos \beta + \sin \alpha \sin \beta$

using

$$\begin{aligned} &\mathbb{E}_b [\cos(\alpha + b) \cos(\beta + b)] \\ &= \mathbb{E}_b [\sin(\alpha + b) \sin(\beta + b)] \end{aligned}$$

for $b \sim \text{Unif}([0, 2\pi])$

using Monte Carlo sampling to estimate the expectation, see Example A.6

where $z_{\omega, b}(x) \doteq \sqrt{2} \cos(\omega^\top x + b)$,

$$\approx \frac{1}{m} \sum_{i=1}^m z_{\omega^{(i)}, b^{(i)}}(x) \cdot z_{\omega^{(i)}, b^{(i)}}(x') \quad (4.41)$$

for independent samples $\omega^{(i)} \stackrel{\text{iid}}{\sim} p$ and $b^{(i)} \stackrel{\text{iid}}{\sim} \text{Unif}([0, 2\pi])$,

$$= z(x)^\top z(x') \quad (4.42)$$

where the (randomized) feature map of random Fourier features is

$$z(x) \doteq \frac{1}{\sqrt{m}} [z_{\omega^{(1)}, b^{(1)}}(x), \dots, z_{\omega^{(m)}, b^{(m)}}(x)]^\top. \quad (4.43)$$

Intuitively, each component of the feature map $z(x)$ projects x onto a random direction ω drawn from the (inverse) Fourier transform $p(\omega)$ of $k(x - x')$, and wraps this line onto the unit circle in \mathbb{R}^2 . After transforming two points x and x' in this way, their inner product is an unbiased estimator of $k(x - x')$. The mapping $z_{\omega,b}(x) = \sqrt{2} \cos(\omega^\top x + b)$ additionally rotates the circle by a random amount b and projects the points onto the interval $[0, 1]$.

Rahimi et al. (2007) show that Bayesian linear regression with the feature map z approximates Gaussian processes with a stationary kernel:

Theorem 4.13 (Uniform convergence of Fourier features). *Suppose \mathcal{M} is a compact subset of \mathbb{R}^d with diameter $\text{diam}(\mathcal{M})$. Then for a stationary kernel k , the random Fourier features z , and any $\epsilon > 0$ it holds that*

$$\begin{aligned} \mathbb{P}\left(\sup_{x,x' \in \mathcal{M}} |z(x)^\top z(x') - k(x - x')| \geq \epsilon\right) \\ \leq 2^8 \left(\frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon}\right)^2 \exp\left(-\frac{m\epsilon^2}{8(d+2)}\right) \end{aligned} \quad (4.44)$$

where $\sigma_p^2 \doteq \mathbb{E}_{\omega \sim p}[\omega^\top \omega]$ is the second moment of p , m is the dimension of $z(x)$, and d is the dimension of x . $\textcircled{?}$

Note that the error probability decays exponentially fast in the dimension of the Fourier feature space.

4.5.3 Data Sampling

Another natural approach is to only consider a (random) subset of the training data during learning. The naive approach is to subsample uniformly at random. Not very surprisingly, we can do much better.

One subsampling method is the *inducing points method* (Quinonero-Candela and Rasmussen, 2005). The idea is to summarize the data around so-called inducing points.¹⁰ For now, let us consider an arbitrary set of inducing points,

$$U \doteq \{\bar{x}_1, \dots, \bar{x}_k\}.$$

Then, the original Gaussian process can be recovered using marginalization,

$$p(f^*, f) = \int_{\mathbb{R}^k} p(f^*, f, u) du = \int_{\mathbb{R}^k} p(f^*, f | u) p(u) du, \quad (4.45)$$

where $f \doteq [f(x_1) \dots f(x_n)]^\top$ and $f^* \doteq f(x^*)$ at some evaluation point $x^* \in \mathcal{X}$. We use $u \doteq [f(\bar{x}_1) \dots f(\bar{x}_k)]^\top \in \mathbb{R}^k$ to denote the predictions of the model at the inducing points U . Due to the marginalization property of Gaussian processes (4.1), we have that $u \sim \mathcal{N}(\mathbf{0}, K_{UU})$.

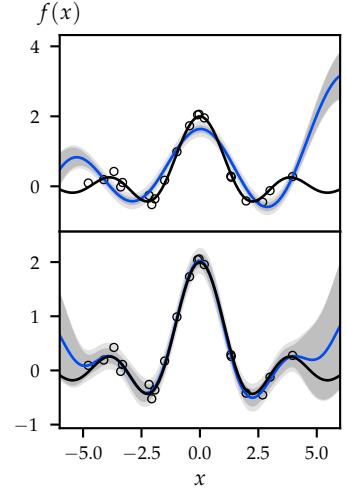


Figure 4.13: Example of random Fourier features with where the number of features m is 5 (top) and 10 (bottom), respectively. The noise-free true function is shown in black and the mean of the Gaussian process is shown in blue.

Problem 4.8

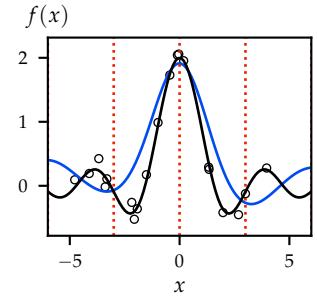


Figure 4.14: Inducing points u are shown as vertical dotted red lines. The noise-free true function is shown in black and the mean of the Gaussian process is shown in blue. Observe that the true function is approximated “well” around the inducing points.

¹⁰The inducing points can be treated as hyperparameters.

using the sum rule (1.7) and product rule (1.11)

The key idea is to approximate the joint prior, assuming that f^* and f are conditionally independent given \mathbf{u} ,

$$p(f^*, f) \approx \int_{\mathbb{R}^k} p(f^* | \mathbf{u}) p(f | \mathbf{u}) p(\mathbf{u}) d\mathbf{u}. \quad (4.46)$$

Here, $p(f | \mathbf{u})$ and $p(f^* | \mathbf{u})$ are commonly called the *training conditional* and the *testing conditional*, respectively. Still denoting the observations by $A = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ and defining $\star \doteq \{\mathbf{x}^*\}$, we know, using the closed-form expression for conditional Gaussians (1.53),

$$p(f | \mathbf{u}) \sim \mathcal{N}(f; \mathbf{K}_{AU} \mathbf{K}_{UU}^{-1} \mathbf{u}, \mathbf{K}_{AA} - \mathbf{Q}_{AA}), \quad (4.47a)$$

$$p(f^* | \mathbf{u}) \sim \mathcal{N}(f^*; \mathbf{K}_{\star U} \mathbf{K}_{UU}^{-1} \mathbf{u}, \mathbf{K}_{\star\star} - \mathbf{Q}_{\star\star}) \quad (4.47b)$$

where $\mathbf{Q}_{ab} \doteq \mathbf{K}_{aU} \mathbf{K}_{UU}^{-1} \mathbf{K}_{Ub}$. Intuitively, \mathbf{K}_{AA} represents the prior covariance and \mathbf{Q}_{AA} represents the covariance “explained” by the inducing points.¹¹

Computing the full covariance matrix is expensive. In the following, we mention two approximations to the covariance of the training conditional (and testing conditional).

Example 4.14: Subset of regressors

The *subset of regressors* (SoR) approximation is defined as

$$q_{\text{SoR}}(f | \mathbf{u}) \doteq \mathcal{N}(f; \mathbf{K}_{AU} \mathbf{K}_{UU}^{-1} \mathbf{u}, \mathbf{0}), \quad (4.48a)$$

$$q_{\text{SoR}}(f^* | \mathbf{u}) \doteq \mathcal{N}(f^*; \mathbf{K}_{\star U} \mathbf{K}_{UU}^{-1} \mathbf{u}, \mathbf{0}). \quad (4.48b)$$

Comparing to Equation (4.47), SoR simply forgets about all variance and covariance.

Example 4.15: Fully independent training conditional

The *fully independent training conditional* (FITC) approximation is defined as

$$q_{\text{FITC}}(f | \mathbf{u}) \doteq \mathcal{N}(f; \mathbf{K}_{AU} \mathbf{K}_{UU}^{-1} \mathbf{u}, \text{diag}\{\mathbf{K}_{AA} - \mathbf{Q}_{AA}\}), \quad (4.49a)$$

$$q_{\text{FITC}}(f^* | \mathbf{u}) \doteq \mathcal{N}(f^*; \mathbf{K}_{\star U} \mathbf{K}_{UU}^{-1} \mathbf{u}, \text{diag}\{\mathbf{K}_{\star\star} - \mathbf{Q}_{\star\star}\}). \quad (4.49b)$$

In contrast to SoR, FITC keeps track of the variances but forgets about the covariance.

The computational cost for inducing point methods SoR and FITC is dominated by the cost of inverting \mathbf{K}_{UU} . Thus, the time complexity is cubic in the number of inducing points, but only linear in the number of data points.

¹¹ For more details, refer to section 2 of “A unifying view of sparse approximate Gaussian process regression” (Quinonero-Candela and Rasmussen, 2005).

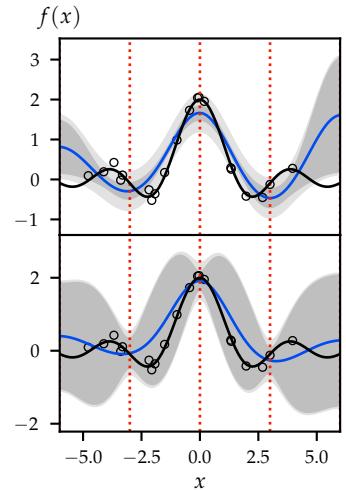


Figure 4.15: Comparison of SoR (top) and FITC (bottom). The inducing points \mathbf{u} are shown as vertical dotted red lines. The noise-free true function is shown in black and the mean of the Gaussian process is shown in blue.

Discussion

This chapter introduced Gaussian processes which leverage the function-space view on linear regression to perform exact probabilistic inference with flexible, nonlinear models. A Gaussian process can be seen as a non-parametric model since it can represent an infinite-dimensional parameter space. Instead, as we saw with the representer theorem, such non-parametric (i.e., “function-space”) models are directly represented as functions of the data points. While this can make these models more flexible than a simple linear parametric model in input space, it also makes them computationally expensive as the number of data points grows. To this end, we discussed several ways of approximating Gaussian processes.

Nevertheless, for today’s internet-scale datasets, modern machine learning typically relies on large parametric models that learn features from data. These models can effectively amortize the cost of inference during training by encoding information into a fixed set of parameters. In the following chapters, we will start to explore approaches to approximate probabilistic inference that can be applied to such models.

Problems

4.1. Feature space of Gaussian kernel.

1. Show that the univariate Gaussian kernel with length scale $h = 1$ implicitly defines a feature space with basis vectors

$$\boldsymbol{\phi}(x) = \begin{bmatrix} \phi_0(x) \\ \phi_1(x) \\ \vdots \end{bmatrix} \quad \text{with} \quad \phi_j(x) = \frac{1}{\sqrt{j!}} e^{-\frac{x^2}{2}} x^j.$$

Hint: Use the Taylor series expansion of the exponential function, $e^x = \sum_{j=0}^{\infty} \frac{x^j}{j!}$.

2. Note that the vector $\boldsymbol{\phi}(x)$ is ∞ -dimensional. Thus, taking the function-space view allows us to perform regression in an infinite-dimensional feature space. What is the effective dimension when regressing n univariate data points with a Gaussian kernel?

4.2. Kernels on the circle.

Consider a dataset $\{(x_i, y_i)\}_{i=1}^n$ with labels $y_i \in \mathbb{R}$ and inputs x_i which lie on the unit circle $S \subset \mathbb{R}^2$. In particular, any element of S can be identified with points in \mathbb{R}^2 of form $(\cos(\theta), \sin(\theta))$ or with the respective angles $\theta \in [0, 2\pi]$.

You now want to use GP regression to learn an unknown mapping from S to \mathbb{R} using this dataset. Thus, you need a valid kernel $k :$

$S \times S \rightarrow \mathbb{R}$. First, we look at kernels k which can be understood as analogous to the Gaussian kernel.

1. You think of the “extrinsic” kernel $k_e : S \times S \rightarrow \mathbb{R}$ defined by

$$k_e(\theta, \theta') \doteq \exp\left(-\frac{\|x(\theta) - x(\theta')\|_2^2}{2\kappa^2}\right),$$

where $x(\theta) \doteq (\cos(\theta), \sin(\theta))$. Is k_e positive semi-definite for all values of $\kappa > 0$?

2. Then, you think of an “intrinsic” kernel $k_i : S \times S \rightarrow \mathbb{R}$ defined by

$$k_i(\theta, \theta') \doteq \exp\left(-\frac{d(\theta, \theta')^2}{2\kappa^2}\right)$$

where $d(\theta, \theta') \doteq \min(|\theta - \theta'|, |\theta - \theta' - 2\pi|, |\theta - \theta' + 2\pi|)$ is the standard arc length distance on the circle S .

You would now like to test whether this kernel is positive semi-definite. We pick $\kappa = 2$ and compute the kernel matrix K for the points corresponding to the angles $\{0, \pi/2, \pi, 3\pi/2\}$. This kernel matrix K has eigenvectors $(1, 1, 1, 1)$ and $(-1, 1, -1, 1)$.

Now compute the eigenvalue corresponding to the eigenvector $(-1, 1, -1, 1)$.

3. Is k_i positive semi-definite for $\kappa = 2$?
4. A mathematician friend of yours suggests to you yet another kernel for points on the circle S , called the *heat kernel*. The kernel itself has a complicated expression but can be accurately approximated by

$$k_h(\theta, \theta') \doteq \frac{1}{C_\kappa} \left(1 + \sum_{l=1}^{L-1} e^{-\frac{\kappa^2}{2} l^2} 2 \cos(l(\theta - \theta')) \right),$$

where $L \in \mathbb{N}$ controls the quality of approximation and $C_\kappa > 0$ is a normalizing constant that depends only on κ .

Is k_h positive semi-definite for all values of $\kappa > 0$ and $L \in \mathbb{N}$?

Hint: Recall that $\cos(a - b) = \cos(a)\cos(b) + \sin(a)\sin(b)$.

4.3. A Kalman filter as a Gaussian process.

Next we will show that the Kalman filter from Example 3.4 can be seen as a Gaussian process. To this end, we define

$$f : \mathbb{N}_0 \rightarrow \mathbb{R}, \quad t \mapsto X_t. \tag{4.50}$$

Assuming that $X_0 \sim \mathcal{N}(0, \sigma_0^2)$ and $X_{t+1} \doteq X_t + \varepsilon_t$ with independent noise $\varepsilon_t \sim \mathcal{N}(0, \sigma_x^2)$, show that

$$f \sim \mathcal{GP}(0, k_{KF}) \quad \text{where} \tag{4.51}$$

$$k_{KF}(t, t') \doteq \sigma_0^2 + \sigma_x^2 \min\{t, t'\}. \tag{4.52}$$

This particular kernel $k(t, t') \doteq \min\{t, t'\}$ but over the continuous-time domain defines the *Wiener process* (also known as Brownian motion).

4.4. Reproducing property and RKHS norm.

1. Derive the reproducing property.

Hint: Use $k(\mathbf{x}, \mathbf{x}') = \langle k(\mathbf{x}, \cdot), k(\mathbf{x}', \cdot) \rangle_k$.

2. Show that the RKHS norm $\|\cdot\|_k$ is a measure of smoothness by proving that for any $f \in \mathcal{H}_k(\mathcal{X})$ and $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ it holds that

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq \|f\|_k \|k(\mathbf{x}, \cdot) - k(\mathbf{y}, \cdot)\|_k.$$

4.5. Representer theorem.

With this, we can now derive the representer theorem (4.20).

Hint: Recall

1. the reproducing property $f(\mathbf{x}) = \langle f, k(\mathbf{x}, \cdot) \rangle_k$ with $k(\mathbf{x}, \cdot) \in \mathcal{H}_k(\mathcal{X})$ which holds for all $f \in \mathcal{H}_k(\mathcal{X})$ and $\mathbf{x} \in \mathcal{H}_k(\mathcal{X})$, and
2. that the norm after projection is smaller or equal the norm before projection.

Then decompose f into parallel and orthogonal components with respect to $\text{span}\{k(\mathbf{x}_1, \cdot), \dots, k(\mathbf{x}_n, \cdot)\}$.

4.6. MAP estimate of Gaussian processes.

Let us denote by $A = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ the set of training points. We will now show that the MAP estimate of GP regression corresponds to the solution of the regularized linear regression problem in the RKHS stated in Equation (4.21):

$$\hat{f} \doteq \arg \min_{f \in \mathcal{H}_k(\mathcal{X})} -\log p(y_{1:n} | \mathbf{x}_{1:n}, f) + \frac{1}{2} \|f\|_k^2.$$

In the following, we abbreviate $\mathbf{K} = \mathbf{K}_{AA}$. We will also assume that the GP has a zero mean function.

1. Show that Equation (4.21) is equivalent to

$$\hat{\boldsymbol{\alpha}} \doteq \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{K}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_{\mathbf{K}}^2 \quad (4.53)$$

for some $\lambda > 0$ which is also known as *kernel ridge regression*. Determine λ .

2. Show that Equation (4.53) with the λ determined in (1) is equivalent to the MAP estimate of GP regression.

Hint: Recall from Equation (4.6) that the MAP estimate at a point \mathbf{x}^ is $\mathbb{E}[f^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y}] = \mathbf{k}_{\mathbf{x}^*, A}^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$.*

4.7. Gradient of the marginal likelihood.

In this exercise, we derive Equation (4.30).

Recall that we were considering a dataset (\mathbf{X}, \mathbf{y}) of noise-perturbed evaluations $y_i = f(\mathbf{x}_i) + \varepsilon_i$ where $\varepsilon_i \sim \mathcal{N}(0, \sigma_n^2)$ and f is an unknown

function. We make the hypothesis $f \sim \mathcal{GP}(0, k_\theta)$ with a zero mean function and the covariance function k_θ . We are interested in finding the hyperparameters θ that maximize the marginal likelihood $p(\mathbf{y} | \mathbf{X}, \theta)$.

1. Derive Equation (4.30).

Hint: You can use the following identities:

- (a) *for any invertible matrix M ,*

$$\frac{\partial}{\partial \theta_j} M^{-1} = -M^{-1} \frac{\partial M}{\partial \theta_j} M^{-1} \quad \text{and} \quad (4.54)$$

- (b) *for any symmetric positive definite matrix M ,*

$$\frac{\partial}{\partial \theta_j} \log \det(M) = \text{tr}\left(M^{-1} \frac{\partial M}{\partial \theta_j}\right). \quad (4.55)$$

2. Assume now that the covariance function for the noisy targets (i.e., including the noise contribution) can be expressed as

$$k_{y,\theta}(\mathbf{x}, \mathbf{x}') = \theta_0 \tilde{k}(\mathbf{x}, \mathbf{x}')$$

where \tilde{k} is a valid kernel independent of θ_0 .¹²

Show that $\frac{\partial}{\partial \theta_0} \log p(\mathbf{y} | \mathbf{X}, \theta) = 0$ admits a closed-form solution for θ_0 which we denote by θ_0^* .

3. How should the optimal parameter θ_0^* be scaled if we scale the labels \mathbf{y} by a scalar s ?

4.8. Uniform convergence of Fourier features.

In this exercise, we will prove Theorem 4.13.

Let $s(\mathbf{x}, \mathbf{x}') \doteq \mathbf{z}(\mathbf{x})^\top \mathbf{z}(\mathbf{x}')$ and $f(\mathbf{x}, \mathbf{x}') \doteq s(\mathbf{x}, \mathbf{x}') - k(\mathbf{x}, \mathbf{x}')$. Observe that both functions are shift invariant, and we will therefore denote them as univariate functions with argument $\Delta \equiv \mathbf{x} - \mathbf{x}' \in \mathcal{M}_\Delta$. Notice that our goal is to bound the probability of the event $\sup_{\Delta \in \mathcal{M}_\Delta} |f(\Delta)| \geq \epsilon$.

1. Show that for all $\Delta \in \mathcal{M}_\Delta$, $\mathbb{P}(|f(\Delta)| \geq \epsilon) \leq 2 \exp\left(-\frac{m\epsilon^2}{4}\right)$.

What we have derived in (1) is known as a *pointwise convergence* guarantee. However, we are interested in bounding the *uniform convergence* over the compact set \mathcal{M}_Δ .

Our approach will be to “cover” the compact set \mathcal{M}_Δ using T balls of radius r whose centers we denote by $\{\Delta_i\}_{i=1}^T$. It can be shown that this is possible for some $T \leq (4 \operatorname{diam}(\mathcal{M})/r)^d$. It can furthermore be shown that

$$\forall i. |f(\Delta_i)| < \frac{\epsilon}{2} \text{ and } \|\nabla f(\Delta^*)\|_2 < \frac{\epsilon}{2r} \implies \sup_{\Delta \in \mathcal{M}_\Delta} |f(\Delta)| < \epsilon$$

where $\Delta^* = \arg \max_{\Delta \in \mathcal{M}_\Delta} \|\nabla f(\Delta)\|_2$.

¹² That is, $K_{y,\theta}(i, j) = k_{y,\theta}(\mathbf{x}_i, \mathbf{x}_j)$.

2. Prove $\mathbb{P}(\|\nabla f(\Delta^*)\|_2 \geq \frac{\epsilon}{2r}) \leq \left(\frac{2r\sigma_p}{\epsilon}\right)^2$.
Hint: Recall that the random Fourier feature approximation is unbiased, i.e., $\mathbb{E}[s(\Delta)] = k(\Delta)$.
3. Prove $\mathbb{P}\left(\bigcup_{i=1}^T |f(\Delta_i)| \geq \frac{\epsilon}{2}\right) \leq 2T \exp\left(-\frac{m\epsilon^2}{16}\right)$.
4. Combine the results from (2) and (3) to prove Theorem 4.13.
Hint: You may use that
 - (a) $\alpha r^{-d} + \beta r^2 = 2\beta^{\frac{d}{d+2}} \alpha^{\frac{2}{d+2}}$ for $r = (\alpha/\beta)^{\frac{1}{d+2}}$ and
 - (b) $\frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon} \geq 1$.
5. Show that for the Gaussian kernel (4.14), $\sigma_p^2 = \frac{d}{h^2}$.
Hint: First show $\sigma_p^2 = -\text{tr}(\mathbf{H}_\Delta k(\mathbf{0}))$.

4.9. Subset of regressors.

1. Using an SoR approximation, prove the following:

$$q_{\text{SoR}}(f, f^*) = \mathcal{N}\left(\begin{bmatrix} f \\ f^* \end{bmatrix}; \mathbf{0}, \begin{bmatrix} Q_{AA} & Q_{A^*} \\ Q_{A^*} & Q_{**} \end{bmatrix}\right) \quad (4.56)$$

$$q_{\text{SoR}}(f^* \mid \mathbf{y}) = \mathcal{N}(f^*; Q_{*A} \tilde{Q}_{AA}^{-1} \mathbf{y}, Q_{**} - Q_{*A} \tilde{Q}_{AA}^{-1} Q_{A^*}) \quad (4.57)$$

where $\tilde{Q}_{ab} \doteq Q_{ab} + \sigma_n^2$.

2. Derive that the resulting model is a degenerate Gaussian process with covariance function

$$k_{\text{SoR}}(\mathbf{x}, \mathbf{x}') \doteq \mathbf{k}_{x,U}^\top \mathbf{K}_{UU}^{-1} \mathbf{k}_{x',U}. \quad (4.58)$$

5

Variational Inference

We have seen how to perform (efficient) probabilistic inference with Gaussians, exploiting their closed-form formulas for marginal and conditional distributions. But what if we work with other distributions?

In this and the following chapter, we will discuss two methods of approximate inference. We begin by discussing variational (probabilistic) inference, which aims to find a good approximation of the posterior distribution from which it is easy to sample. In Chapter 6, we discuss Markov chain Monte Carlo methods, which approximate the sampling from the posterior distribution directly.

The fundamental idea behind variational inference is to approximate the true posterior distribution using a “simpler” posterior that is as close as possible to the true posterior:

$$p(\boldsymbol{\theta} \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \frac{1}{Z} p(\boldsymbol{\theta}, \mathbf{y}_{1:n} \mid \mathbf{x}_{1:n}) \approx q(\boldsymbol{\theta} \mid \boldsymbol{\lambda}) \doteq q_{\boldsymbol{\lambda}}(\boldsymbol{\theta}) \quad (5.1)$$

where $\boldsymbol{\lambda}$ represents the parameters of the *variational posterior* $q_{\boldsymbol{\lambda}}$, also called *variational parameters*. In doing so, variational inference reduces probabilistic inference — where the fundamental difficulty lies in solving high-dimensional integrals — to an optimization problem. Optimizing (stochastic) objectives is a well-understood problem with efficient algorithms that perform well in practice.¹

5.1 Laplace Approximation

Before introducing a general framework of variational inference, we discuss a simpler method of approximate inference known as *Laplace’s method*. This method was proposed as a method of approximating integrals as early as 1774 by Pierre-Simon Laplace. The idea is to use a Gaussian approximation (that is, a second-order Taylor approxima-

¹ We provide an overview of first-order methods such as stochastic gradient descent in Appendix A.4.

tion) of the posterior distribution around its mode. Let

$$\psi(\boldsymbol{\theta}) \doteq \log p(\boldsymbol{\theta} \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) \quad (5.2)$$

denote the log-posterior. Then, using a second-order Taylor approximation (A.53) around the mode $\hat{\boldsymbol{\theta}}$ of ψ (i.e., the MAP estimate), we obtain the approximation $\hat{\psi}$ which is accurate for $\boldsymbol{\theta} \approx \hat{\boldsymbol{\theta}}$:

$$\begin{aligned} \psi(\boldsymbol{\theta}) \approx \hat{\psi}(\boldsymbol{\theta}) &\doteq \psi(\hat{\boldsymbol{\theta}}) + (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \nabla \psi(\hat{\boldsymbol{\theta}}) + \frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H}_\psi(\hat{\boldsymbol{\theta}}) (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \\ &= \psi(\hat{\boldsymbol{\theta}}) + \frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H}_\psi(\hat{\boldsymbol{\theta}}) (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}). \end{aligned} \quad (5.3) \quad \text{using } \nabla \psi(\hat{\boldsymbol{\theta}}) = 0$$

Compare this expression to the log-PDF of a Gaussian:

$$\log \mathcal{N}(\boldsymbol{\theta}; \hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda}^{-1}) = -\frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \boldsymbol{\Lambda} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) + \text{const.} \quad (5.4)$$

Since $\psi(\hat{\boldsymbol{\theta}})$ is constant with respect to $\boldsymbol{\theta}$,

$$\hat{\psi}(\boldsymbol{\theta}) = \log \mathcal{N}(\boldsymbol{\theta}; \hat{\boldsymbol{\theta}}, -\mathbf{H}_\psi(\hat{\boldsymbol{\theta}})^{-1}) + \text{const.} \quad (5.5)$$

The *Laplace approximation* q of p is

$$q(\boldsymbol{\theta}) \doteq \mathcal{N}(\boldsymbol{\theta}; \hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda}^{-1}) \propto \exp(\hat{\psi}(\boldsymbol{\theta})) \quad \text{where} \quad (5.6a)$$

$$\boldsymbol{\Lambda} \doteq -\mathbf{H}_\psi(\hat{\boldsymbol{\theta}}) = -\mathbf{H}_\theta \log p(\boldsymbol{\theta} \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}. \quad (5.6b)$$

Recall that for this approximation to be well-defined, the covariance matrix $\boldsymbol{\Lambda}^{-1}$ (or equivalently the precision matrix $\boldsymbol{\Lambda}$) needs to be symmetric and positive semi-definite. Let us verify that this is indeed the case for sufficiently smooth ψ .² In this case, the Hessian $\boldsymbol{\Lambda}$ is symmetric since the order of differentiation does not matter. Moreover, by the second-order optimality condition, $\mathbf{H}_\psi(\hat{\boldsymbol{\theta}})$ is negative semi-definite since $\hat{\boldsymbol{\theta}}$ is a maximum of ψ , which implies that $\boldsymbol{\Lambda}$ is positive semi-definite.

² ψ being twice continuously differentiable around $\hat{\boldsymbol{\theta}}$ is sufficient.

Example 5.1: Laplace approximation of a Gaussian

Consider approximating the Gaussian density $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ using a Laplace approximation.

We know that the mode of p is $\boldsymbol{\mu}$, which we can verify by computing the gradient,

$$\nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}) = -\frac{1}{2} (2\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta} - 2\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}) \stackrel{!}{=} 0 \iff \boldsymbol{\theta} = \boldsymbol{\mu}. \quad (5.7)$$

For the Hessian of $\log p(\boldsymbol{\theta})$, we get

$$\mathbf{H}_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}) = (\mathbf{D}_{\boldsymbol{\theta}}(\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} - \boldsymbol{\Sigma}^{-1}\boldsymbol{\theta}))^\top = -(\boldsymbol{\Sigma}^{-1})^\top = -\boldsymbol{\Sigma}^{-1}. \quad (5.8)$$

using $(\mathbf{A}^{-1})^\top = (\mathbf{A}^\top)^{-1}$ and symmetry of $\boldsymbol{\Sigma}$

We see that the Laplace approximation of a Gaussian $p(\theta)$ is exact, which should not come as a surprise since the second-order Taylor approximation of $\log p(\theta)$ is exact for Gaussians.

The Laplace approximation matches the shape of the true posterior around its mode but may not represent it accurately elsewhere — often leading to extremely overconfident predictions. An example is given in Figure 5.1. Nevertheless, the Laplace approximation has some desirable properties such as being relatively easy to apply in a post-hoc manner, that is, after having already computed the MAP estimate. It preserves the MAP point estimate as its mean and just “adds” a little uncertainty around it. However, the fact that it can be arbitrarily different from the true posterior makes it unsuitable for approximate probabilistic inference.

5.1.1 Example: Bayesian Logistic Regression

As an example, we will look at Laplace approximation in the context of Bayesian logistic regression. Logistic regression learns a classifier that decides for a given input whether it belongs to one of two classes. A sigmoid function, typically the *logistic function*,

$$\sigma(z) \doteq \frac{1}{1 + \exp(-z)} \in (0, 1), \quad z = \mathbf{w}^\top \mathbf{x}, \quad (5.9)$$

is used to obtain the class probabilities. *Bayesian logistic regression* corresponds to Bayesian linear regression with a Bernoulli likelihood,

$$y | \mathbf{x}, \mathbf{w} \sim \text{Bern}(\sigma(\mathbf{w}^\top \mathbf{x})), \quad (5.10)$$

where $y \in \{-1, 1\}$ is the binary class label.³ Observe that given a data point (\mathbf{x}, y) , the probability of a correct classification is

$$p(y | \mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(\mathbf{w}^\top \mathbf{x}) & \text{if } y = 1 \\ 1 - \sigma(\mathbf{w}^\top \mathbf{x}) & \text{if } y = -1 \end{cases} = \sigma(y \mathbf{w}^\top \mathbf{x}), \quad (5.11)$$

as the logistic function σ is symmetric around 0. Also, recall that Bayesian linear regression used the prior

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \sigma_p^2 \mathbf{I}) \propto \exp\left(-\frac{1}{2\sigma_p^2} \|\mathbf{w}\|_2^2\right).$$

Let us first find the posterior mode, that is, the MAP estimate of the weights:

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} p(\mathbf{w} | \mathbf{x}_{1:n}, y_{1:n})$$

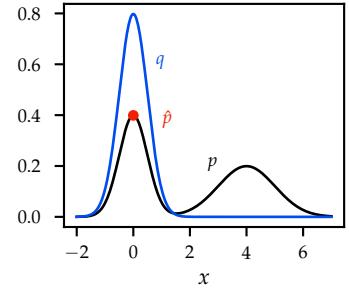


Figure 5.1: The Laplace approximation q greedily selects the mode of the true posterior distribution p and matches the curvature around the mode \hat{p} . As shown here, the Laplace approximation can be extremely overconfident when p is not approximately Gaussian.

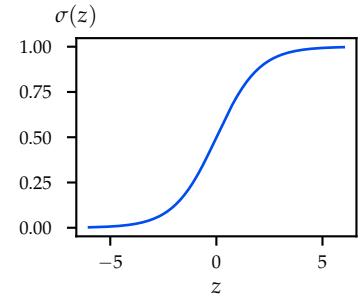


Figure 5.2: The logistic function squashes the linear function $\mathbf{w}^\top \mathbf{x}$ onto the interval $(0, 1)$.

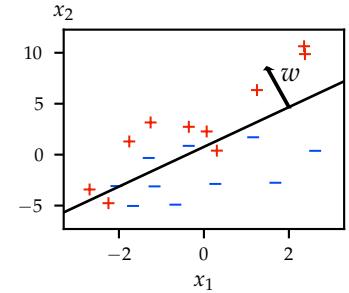


Figure 5.3: Logistic regression classifies data into two classes with a linear decision boundary.

³ The same approach extends to Gaussian processes where it is known as *Gaussian process classification*, see Problem 5.2 and Hensman et al. (2015).

$$\begin{aligned}
&= \arg \max_{\mathbf{w}} p(\mathbf{w}) p(y_{1:n} | \mathbf{x}_{1:n}, \mathbf{w}) \\
&= \arg \max_{\mathbf{w}} \log p(\mathbf{w}) + \log p(y_{1:n} | \mathbf{x}_{1:n}, \mathbf{w}) \\
&= \arg \max_{\mathbf{w}} -\frac{1}{2\sigma_p^2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \log \sigma(y_i \mathbf{w}^\top \mathbf{x}_i) \\
&= \arg \min_{\mathbf{w}} \frac{1}{2\sigma_p^2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)). \quad (5.12)
\end{aligned}$$

Note that for $\lambda = 1/2\sigma_p^2$, the above optimization is equivalent to standard (regularized) logistic regression where

$$\ell_{\text{log}}(\mathbf{w}^\top \mathbf{x}; y) \doteq \log(1 + \exp(-y \mathbf{w}^\top \mathbf{x})) \quad (5.13)$$

is called *logistic loss*. The gradient of the logistic loss is given by [\(?\)](#)

$$\nabla_{\mathbf{w}} \ell_{\text{log}}(\mathbf{w}^\top \mathbf{x}; y) = -y \mathbf{x} \cdot \sigma(-y \mathbf{w}^\top \mathbf{x}). \quad (5.14)$$

Recall that due to the symmetry of σ around 0, $\sigma(-y \mathbf{w}^\top \mathbf{x})$ is the probability that \mathbf{x} was *not* classified as y . Intuitively, if the model is “surprised” by the label, the gradient is large.

We can therefore use SGD with the (regularized) gradient step and with batch size 1,

$$\mathbf{w} \leftarrow \mathbf{w}(1 - 2\lambda\eta_t) + \eta_t y \mathbf{x} \sigma(-y \mathbf{w}^\top \mathbf{x}), \quad (5.15)$$

for the data point (\mathbf{x}, y) picked uniformly at random from the training data. Here, $2\lambda\eta_t$ is due to the gradient of the regularization term, in effect, performing weight decay.

Example 5.2: Laplace approx. of Bayesian logistic regression

We have already found the mode of the posterior distribution, $\hat{\mathbf{w}}$.

Let us denote by

$$\pi_i \doteq \mathbb{P}(y_i = 1 | \mathbf{x}_i, \hat{\mathbf{w}}) = \sigma(\hat{\mathbf{w}}^\top \mathbf{x}_i) \quad (5.16)$$

the probability of x_i belonging to the positive class under the model given by the MAP estimate of the weights. For the precision matrix, we then have

$$\begin{aligned}
\Lambda &= -H_{\mathbf{w}} \log p(\mathbf{w} | \mathbf{x}_{1:n}, y_{1:n})|_{\mathbf{w}=\hat{\mathbf{w}}} \\
&= -H_{\mathbf{w}} \log p(y_{1:n} | \mathbf{x}_{1:n}, \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} - H_{\mathbf{w}} \log p(\mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} \\
&= \sum_{i=1}^n H_{\mathbf{w}} \ell_{\text{log}}(\mathbf{w}^\top \mathbf{x}_i; y_i)|_{\mathbf{w}=\hat{\mathbf{w}}} + \sigma_p^{-2} \mathbf{I} \\
&= \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \pi_i (1 - \pi_i) + \sigma_p^{-2} \mathbf{I}
\end{aligned}$$

using Bayes' rule (1.45)

taking the logarithm

using independence of the observations and Equation (5.11)

using the definition of σ (5.9)

Problem 5.1 (1)

using the definition of the logistic loss (5.13)

using the Hessian of the logistic loss (5.73) which you derive in [Problem 5.1 \(2\)](#)

$$= \mathbf{X}^\top \text{diag}_{i \in [n]} \{\pi_i(1 - \pi_i)\} \mathbf{X} + \sigma_p^{-2} \mathbf{I}. \quad (5.17)$$

Observe that $\pi_i(1 - \pi_i) \approx 0$ if $\pi_i \approx 1$ or $\pi_i \approx 0$. That is, if a training example is “well-explained” by $\hat{\mathbf{w}}$, then its contribution to the precision matrix is small. In contrast, we have $\pi_i(1 - \pi_i) = 0.25$ for $\pi_i = 0.5$. Importantly, Λ does not depend on the normalization constant of the posterior distribution which is hard to compute.

In summary, we have that $\mathcal{N}(\hat{\mathbf{w}}, \Lambda^{-1})$ is the Laplace approximation of $p(\mathbf{w} | \mathbf{x}_{1:n}, y_{1:n})$.

5.2 Predictions with a Variational Posterior

How can we make predictions using our variational approximation? We simply approximate the (intractable) true posterior with our variational posterior:

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}) &= \int p(y^* | \mathbf{x}^*, \theta) p(\theta | \mathbf{x}_{1:n}, y_{1:n}) d\theta \\ &\approx \int p(y^* | \mathbf{x}^*, \theta) q_\lambda(\theta) d\theta. \end{aligned} \quad (5.18)$$

using the sum rule (1.7)

A straightforward approach is to observe that Equation (5.18) can be viewed as an expectation over the variational posterior q_λ and approximated via Monte Carlo sampling:

$$= \mathbb{E}_{\theta \sim q_\lambda} [p(y^* | \mathbf{x}^*, \theta)] \quad (5.19)$$

$$\approx \frac{1}{m} \sum_{j=1}^m p(y^* | \mathbf{x}^*, \theta_j) \quad (5.20)$$

where $\theta_j \stackrel{\text{iid}}{\sim} q_\lambda$.

Example 5.3: Predictions in Bayesian logistic regression

In the case of Bayesian logistic regression with a Gaussian approximation of the posterior, we can obtain more accurate predictions.

Observe that the final prediction y^* is conditionally independent of the model parameters \mathbf{w} given the “latent value” $f^* = \mathbf{w}^\top \mathbf{x}^*$:

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}) &\approx \int p(y^* | \mathbf{x}^*, \mathbf{w}) q_\lambda(\mathbf{w}) d\mathbf{w} \\ &= \int \int p(y^* | f^*) p(f^* | \mathbf{x}^*, \mathbf{w}) q_\lambda(\mathbf{w}) d\mathbf{w} df^* \\ &= \int p(y^* | f^*) \int p(f^* | \mathbf{x}^*, \mathbf{w}) q_\lambda(\mathbf{w}) d\mathbf{w} df^*. \end{aligned} \quad (5.21)$$

once more, using the sum rule (1.7)

rearranging terms

The outer integral can be readily approximated since it is only one-dimensional! The challenging part is the inner integral, which is a high-dimensional integral over the model weights w . Since the posterior over weights $q_\lambda(w) = \mathcal{N}(w; \hat{w}, \Lambda^{-1})$ is a Gaussian, we have due to the closedness properties of Gaussians (1.78) that

$$\int p(f^* | \mathbf{x}^*, w) q_\lambda(w) dw = \mathcal{N}(\hat{w}^\top \mathbf{x}^*, \mathbf{x}^{*\top} \Lambda^{-1} \mathbf{x}^*). \quad (5.22)$$

Crucially, this is a one-dimensional Gaussian in function-space as opposed to the d -dimensional Gaussian q_λ in weight-space!

As we have seen in Equation (5.11), for Bayesian logistic regression, the prediction y^* depends deterministically on the predicted latent value f^* : $p(y^* | f^*) = \sigma(y^* f^*)$. Combining Equations (5.21) and (5.22), we obtain

$$p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}) \approx \int \sigma(y^* f^*) \cdot \mathcal{N}(f^*; \hat{w}^\top \mathbf{x}^*, \mathbf{x}^{*\top} \Lambda^{-1} \mathbf{x}^*) df^*. \quad (5.23)$$

We have replaced the high-dimensional integral over the model parameters θ by the one-dimensional integral over the prediction of our variational posterior f^* . While this integral is generally still intractable, it can be approximated efficiently using numerical quadrature methods such as the Gauss-Legendre quadrature or alternatively with Monte Carlo sampling.

5.3 Blueprint of Variational Inference

General probabilistic inference poses the challenge of approximating the posterior distribution with limited memory and computation, resource constraints also present in humans and other intelligent systems. These resource constraints require information to be compressed, and as we will see, such a compression poses a fundamental tradeoff between model accuracy (on the observed data) and model complexity (to avoid overfitting).

Laplace approximation approximates the true (intractable) posterior with a simpler one, by greedily matching mode and curvature around it. Can we find “less greedy” approaches? We can view variational probabilistic inference more generally as a family of approaches aiming to approximate the true posterior distribution by one that is closest (according to some criterion) among a “simpler” class of distributions. To this end, we need to fix a class of distributions and define suitable criteria, which we can then optimize numerically. The key ben-

efit is that we can reduce the (generally intractable) problem of high-dimensional integration to the (often much more tractable) problem of optimization.

Definition 5.4 (Variational family). Let \mathcal{P} be the class of all probability distributions. A *variational family* $\mathcal{Q} \subseteq \mathcal{P}$ is a class of distributions such that each distribution $q \in \mathcal{Q}$ is characterized by unique variational parameters $\lambda \in \Lambda$.

Example 5.5: Family of independent Gaussians

A straightforward example for a variational family is the family of independent Gaussians,

$$\mathcal{Q} \doteq \left\{ q(\theta) = \mathcal{N}(\theta; \mu, \text{diag}_{i \in [d]} \{\sigma_i^2\}) \right\}, \quad (5.24)$$

which is parameterized by $\lambda \doteq [\mu_{1:d}, \sigma_{1:d}^2]$. Such a multivariate distribution where all variables are independent is called a *mean-field distribution*. Importantly, this family of distributions is characterized by only $2d$ parameters!

Note that Figure 5.4 is a generalization of the canonical distinction between estimation error and approximation error from Figure 1.7, only that here, we operate in the space of distributions over functions as opposed the space of functions. A common notion of distance between two distributions q and p is the Kullback-Leibler divergence $\text{KL}(q \| p)$ which we will define in the next section. Using this notion of distance, we need to solve the following optimization problem:

$$q^* \doteq \arg \min_{q \in \mathcal{Q}} \text{KL}(q \| p) = \arg \min_{\lambda \in \Lambda} \text{KL}(q_\lambda \| p). \quad (5.25)$$

In Section 5.4, we introduce information theory and the Kullback-Leibler divergence. Then, in Section 5.5, we discuss how the optimization problem of Equation (5.25) can be solved efficiently.

5.4 Information Theoretic Aspects of Uncertainty

One of our main objectives throughout this manuscript is to capture the “uncertainty” about events A in an appropriate probability space. One very natural measure of uncertainty is their probability, $\mathbb{P}(A)$. In this section, we will introduce an alternative measure of uncertainty, namely the so-called “surprise” about the event A .

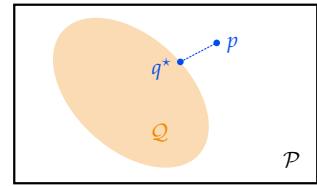


Figure 5.4: An illustration of variational inference in the space of distributions \mathcal{P} . The variational distribution $q^* \in \mathcal{Q}$ is the optimal approximation of the true posterior p .

5.4.1 Surprise

The *surprise* about an event with probability u is defined as

$$S[u] \doteq -\log u. \quad (5.26)$$

Observe that the surprise is a function from $\mathbb{R}_{\geq 0}$ to \mathbb{R} , where we let $S[0] \equiv \infty$. Moreover, for a discrete random variable X , we have that $p(x) \leq 1$, and hence, $S[p(x)] \geq 0$. But why is it reasonable to measure surprise by $-\log u$?

Remarkably, it can be shown that the following natural axiomatic characterization leads to exactly this definition of surprise.

Theorem 5.6 (Axiomatic characterization of surprise). *The axioms*

1. $S[u] > S[v] \implies u < v$ (*anti-monotonicity*)
 2. S is continuous,
 3. $S[uv] = S[u] + S[v]$ for independent events,
- characterize S up to a positive constant factor.

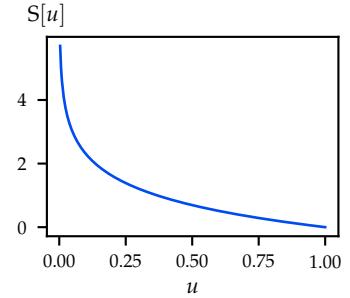


Figure 5.5: Surprise $S[u]$ associated with an event of probability u .

*we are more surprised by unlikely events
no jumps in surprise for infinitesimal
changes of probability
the surprise of independent events is
additive*

Proof. Observe that the third condition looks similar to the product rule of logarithms: $\log(uv) = \log v + \log u$. We can formalize this intuition by remembering Cauchy's functional equation, $f(x+y) = f(x) + f(y)$, which has the unique family of solutions $\{f : x \mapsto cx : c \in \mathbb{R}\}$ if f is required to be continuous. Such a solution is called an "additive function". Consider the function $g(x) \doteq f(e^x)$. Then, g is additive if and only if

$$f(e^x e^y) = f(e^{x+y}) = g(x+y) = g(x) + g(y) = f(e^x) + f(e^y).$$

This is precisely the third axiom of surprise for $f = S$ and $e^x = u$! Hence, the second and third axioms of surprise imply that g must be additive and that $g(x) = S[e^x] = cx$ for any $c \in \mathbb{R}$. If we replace e^x by u , we obtain $S[u] = c \log u$. The first axiom of surprise implies that $c < 0$, and thus, $S[u] = -c' \log u$ for any $c' > 0$. \square

Importantly, surprise offers a different perspective on uncertainty as opposed to probability: the uncertainty about an event can either be interpreted in terms of its probability or in terms of its surprise, and the two "spaces of uncertainty" are related by a log-transform. This relationship is illustrated in Figure 5.6. Information theory is the study of uncertainty in terms of surprise.

Throughout this manuscript we will see many examples where modeling uncertainty in terms of surprise (i.e., the information-theoretic interpretation of uncertainty) is useful. One example where we have

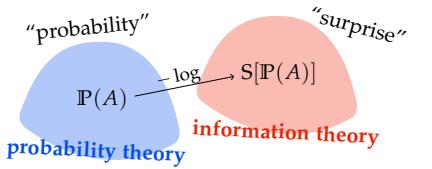


Figure 5.6: Illustration of the probability space and the corresponding "surprise space".

already encountered the “surprise space” was in the context of likelihood maximization (cf. Section 1.3.1) where we used that the log-transform linearizes products of probabilities. We will see later in Chapter 6 that in many cases the surprise $S[p(x)]$ can also be interpreted as a “cost” or “energy” associated with the state x .

5.4.2 Entropy

The *entropy* of a distribution p is the average surprise about samples from p . In this way, entropy is a notion of uncertainty associated with the distribution p : if the entropy of p is large, we are more uncertain about $x \sim p$ than if the entropy of p were low. Formally,

$$H[p] \doteq \mathbb{E}_{x \sim p}[S[p(x)]] = \mathbb{E}_{x \sim p}[-\log p(x)]. \quad (5.27)$$

When $\mathbf{X} \sim p$ is a random vector distributed according to p , we write $H[\mathbf{X}] \doteq H[p]$. Observe that by definition, if p is discrete then $H[p] \geq 0$ as $p(x) \leq 1 (\forall x)$.⁴ For discrete distributions it is common to use the logarithm with base 2 rather than the natural logarithm:⁵

$$H[p] = - \sum_x p(x) \log_2 p(x) \quad (\text{if } p \text{ is discrete}), \quad (5.28a)$$

$$H[p] = - \int p(x) \log p(x) dx \quad (\text{if } p \text{ is continuous}). \quad (5.28b)$$

Let us briefly recall Jensen’s inequality, which is a useful tool when working with expectations of convex functions such as entropy:⁶

Fact 5.7 (Jensen’s Inequality). *Given a random variable X and a convex function $g : \mathbb{R} \rightarrow \mathbb{R}$, we have*

$$g(\mathbb{E}[X]) \leq \mathbb{E}[g(X)]. \quad (5.29)$$

Example 5.8: Examples of entropy

- Fair Coin $H[\text{Bern}(0.5)] = -2(0.5 \log_2 0.5) = 1$.
- Unfair Coin

$$H[\text{Bern}(0.1)] = -0.1 \log_2 0.1 - 0.9 \log_2 0.9 \approx 0.469.$$

- Uniform Distribution

$$H[\text{Unif}(\{1, \dots, n\})] = - \sum_{i=1}^n \frac{1}{n} \log_2 \frac{1}{n} = \log_2 n.$$

The uniform distribution has the maximum entropy among all discrete distributions supported on $\{1, \dots, n\}$ *?*. Note that a fair coin corresponds to a uniform distribution with $n = 2$. Also

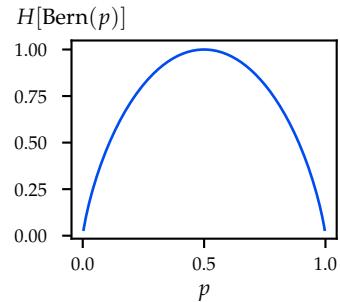


Figure 5.7: Entropy of a Bernoulli experiment with success probability p .

⁴ The entropy of a continuous distribution can be negative. For example,

$$\begin{aligned} H[\text{Unif}([a, b])] &= - \int \frac{1}{b-a} \log \frac{1}{b-a} dx \\ &= \log(b-a) \end{aligned}$$

which is negative if $b - a < 1$.

⁵ Recall that $\log_2 x = \frac{\log x}{\log 2}$, that is, logarithms to a different base only differ by a constant factor.

⁶ The surprise $S[u]$ is convex in u .

Problem 5.3 (1)

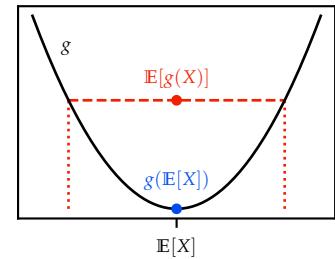


Figure 5.8: An illustration of Jensen’s inequality. Due to the convexity of g , we have that g evaluated at $\mathbb{E}[X]$ will always be below the average of evaluations of g .

Problem 5.3 (2)

observe that $\log_2 n$ corresponds to the number of bits required to encode the outcome of the experiment.

In general, the entropy $H[p]$ of a discrete distribution p can be interpreted as the average number of bits required to encode a sample $x \sim p$, or in other words, the average “information” carried by a sample x .

Example 5.9: Entropy of a Gaussian

Let us derive the entropy of a univariate Gaussian. Recall the PDF,

$$\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{Z} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

where $Z = \sqrt{2\pi\sigma^2}$. Using the definition of entropy (5.28b), we obtain,

$$\begin{aligned} H[\mathcal{N}(\mu, \sigma^2)] &= - \int \frac{1}{Z} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \\ &\quad \cdot \log\left(\frac{1}{Z} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)\right) dx \\ &= \log Z \underbrace{\int \frac{1}{Z} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) dx}_1 \\ &\quad + \int \frac{1}{Z} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \frac{(x - \mu)^2}{2\sigma^2} dx \\ &= \log Z + \frac{1}{2\sigma^2} \mathbb{E}[(x - \mu)^2] \\ &= \log(\sigma\sqrt{2\pi}) + \frac{1}{2} \\ &= \log(\sigma\sqrt{2\pi e}). \end{aligned} \tag{5.30}$$

using LOTUS (1.22)

using $\mathbb{E}[(x - \mu)^2] = \text{Var}[x] = \sigma^2$ (1.34)

using $\log \sqrt{e} = 1/2$

In general, the entropy of a Gaussian is

$$H[\mathcal{N}(\mu, \Sigma)] = \frac{1}{2} \log \det(2\pi e \Sigma) = \frac{1}{2} \log((2\pi e)^d \det(\Sigma)). \tag{5.31}$$

Note that the entropy is a function of the determinant of the covariance matrix Σ . In general, there are various ways of “scalarizing” the notion of uncertainty for a multivariate distribution. The determinant of Σ measures the volume of the credible sets around the mean μ , and is also called the *generalized variance*. Next to entropy and generalized variance (which are closely related for Gaussians), a common scalarization is the trace of Σ , which is also called the *total variance*.

5.4.3 Cross-Entropy

How can we use entropy to measure our average surprise when assuming the data follows some distribution q but in reality the data follows a different distribution p ?

Definition 5.10 (Cross-entropy). The *cross-entropy* of a distribution q relative to the distribution p is

$$H[p\|q] \doteq \mathbb{E}_{x\sim p}[S[q(x)]] = \mathbb{E}_{x\sim p}[-\log q(x)]. \quad (5.32)$$

Cross-entropy can also be expressed in terms of the KL-divergence (cf. Section 5.4.4) $KL(p\|q)$ which measures how “different” the distribution q is from a reference distribution p ,

$$H[p\|q] = H[p] + KL(p\|q) \geq H[p]. \quad (5.33)$$

$KL(p\|q) \geq 0$ is shown in Problem 5.5

Quite intuitively, the average surprise in samples from p with respect to the distribution q is given by the inherent uncertainty in p and the additional surprise that is due to us assuming the wrong data distribution q . The “closer” q is to the true data distribution p , the smaller is the additional average surprise.

5.4.4 Kullback-Leibler Divergence

As mentioned, the Kullback-Leibler divergence is a (non-metric) measure of distance between distributions. It is defined as follows.

Definition 5.11 (Kullback-Leibler divergence, KL-divergence). Given two distributions p and q , the *Kullback-Leibler divergence* (or *relative entropy*) of q with respect to p ,

$$KL(p\|q) \doteq H[p\|q] - H[p] \quad (5.34)$$

$$= \mathbb{E}_{\theta\sim p}[S[q(\theta)] - S[p(\theta)]] \quad (5.35)$$

$$= \mathbb{E}_{\theta\sim p}\left[\log \frac{p(\theta)}{q(\theta)}\right], \quad (5.36)$$

measures how different q is from a reference distribution p .

In words, $KL(p\|q)$ measures the *additional* expected surprise when observing samples from p that is due to assuming the (wrong) distribution q and which is not inherent in the distribution p already.⁷

The KL-divergence has the following properties:

- $KL(p\|q) \geq 0$ for any distributions p and q $\textcircled{?}$,
- $KL(p\|q) = 0$ if and only if $p = q$ almost surely $\textcircled{?}$, and
- there exist distributions p and q such that $KL(p\|q) \neq KL(q\|p)$.

⁷ The KL-divergence only captures the additional expected surprise since the surprise inherent in p (as measured by $H[p]$) is subtracted.

Problem 5.5 (1)

Problem 5.5 (2)

The KL-divergence can simply be understood as a shifted version of cross-entropy, which is zero if we consider the divergence between two identical distributions.

We will briefly look at another interpretation for how KL-divergence measures “distance” between distributions. Suppose we are presented with a sequence $\theta_1, \dots, \theta_n$ of independent samples from either a distribution p or a distribution q , both of which are known. Which of p or q was used to generate the data is, however, unknown to us, and we would like to find out. A natural approach is to choose the distribution whose data likelihood is larger. That is, we choose p if $p(\theta_{1:n}) > q(\theta_{1:n})$ and vice versa. Assuming that the samples are independent and rewriting the inequality slightly, we choose p if

$$\prod_{i=1}^n \frac{p(\theta_i)}{q(\theta_i)} > 1, \quad \text{or equivalently if } \sum_{i=1}^n \log \frac{p(\theta_i)}{q(\theta_i)} > 0. \quad (5.37)$$

taking the logarithm

Assume without loss of generality that $\theta_i \sim p$. Then, using the law of large numbers (A.36),

$$\frac{1}{n} \sum_{i=1}^n \log \frac{p(\theta_i)}{q(\theta_i)} \xrightarrow{\text{a.s.}} \mathbb{E}_{\theta \sim p} \left[\log \frac{p(\theta)}{q(\theta)} \right] = \text{KL}(p \| q) \quad (5.38)$$

as $n \rightarrow \infty$. Plugging this into our decision criterion from Equation (5.37), we find that

$$\mathbb{E} \left[\sum_{i=1}^n \log \frac{p(\theta_i)}{q(\theta_i)} \right] = n \text{KL}(p \| q). \quad (5.39)$$

In this way, $\text{KL}(p \| q)$ measures the observed “distance” between p and q . Recall that assuming $p \neq q$ we have that $\text{KL}(p \| q) > 0$ with probability 1, and therefore we correctly choose p with probability 1 as $n \rightarrow \infty$. Moreover, Hoeffding’s inequality (A.41) can be used to determine “how quickly” samples converge to this limit, that is, how quickly we can distinguish between p and q .

Example 5.12: KL-divergence of Bernoulli random variables

Suppose we are given two Bernoulli distributions $\text{Bern}(p)$ and $\text{Bern}(q)$. Then, their KL-divergence is

$$\begin{aligned} \text{KL}(\text{Bern}(p) \| \text{Bern}(q)) &= \sum_{x \in \{0,1\}} \text{Bern}(x; p) \log \frac{\text{Bern}(x; p)}{\text{Bern}(x; q)} \\ &= p \log \frac{p}{q} + (1-p) \log \frac{(1-p)}{(1-q)}. \end{aligned} \quad (5.40)$$

Observe that $\text{KL}(\text{Bern}(p) \| \text{Bern}(q)) = 0$ if and only if $p = q$.

Example 5.13: KL-divergence of Gaussians

Suppose we are given two Gaussian distributions $p \doteq \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$ and $q \doteq \mathcal{N}(\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q)$ with dimension d . The KL-divergence of p and q is given by (?)

$$\begin{aligned} \text{KL}(p\|q) &= \frac{1}{2} \left(\text{tr}(\boldsymbol{\Sigma}_q^{-1} \boldsymbol{\Sigma}_p) + (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)^\top \boldsymbol{\Sigma}_q^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q) \right. \\ &\quad \left. - d + \log \frac{\det(\boldsymbol{\Sigma}_q)}{\det(\boldsymbol{\Sigma}_p)} \right). \end{aligned} \quad (5.41)$$

For independent Gaussians with unit variance, $\boldsymbol{\Sigma}_p = \boldsymbol{\Sigma}_q = \mathbf{I}$, the expression simplifies to the squared Euclidean distance,

$$\text{KL}(p\|q) = \frac{1}{2} \|\boldsymbol{\mu}_q - \boldsymbol{\mu}_p\|_2^2. \quad (5.42)$$

If we approximate independent Gaussians with variances σ_i^2 ,

$$p \doteq \mathcal{N}(\boldsymbol{\mu}, \text{diag}\{\sigma_1^2, \dots, \sigma_d^2\}),$$

by a standard normal distribution, $q \doteq \mathcal{N}(\mathbf{0}, \mathbf{I})$, the expression simplifies to

$$\text{KL}(p\|q) = \frac{1}{2} \sum_{i=1}^d (\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2). \quad (5.43)$$

Here, the term μ_i^2 penalizes a large mean of p , the term σ_i^2 penalizes a large variance of p , and the term $-\log \sigma_i^2$ penalizes a small variance of p . As expected, $\text{KL}(p\|q)$ is proportional to the amount of information we lose by approximating p with the simpler distribution q .

5.4.5 Forward and Reverse KL-divergence

$\text{KL}(p\|q)$ is also called the *forward* (or *inclusive*) KL-divergence. In contrast, $\text{KL}(q\|p)$ is called the *reverse* (or *exclusive*) KL-divergence. Figure 5.9 shows the approximations of a general Gaussian obtained when \mathcal{Q} is the family of diagonal (independent) Gaussians. Thereby,

$$q_1^* \doteq \arg \min_{q \in \mathcal{Q}} \text{KL}(p\|q) \quad \text{and} \quad q_2^* \doteq \arg \min_{q \in \mathcal{Q}} \text{KL}(q\|p).$$

q_1^* is the result when using the forward KL-divergence and q_2^* is the result when using reverse KL-divergence. It can be seen that the reverse KL-divergence tends to greedily select the mode and underestimating the variance which, in this case, leads to an overconfident predic-

Problem 5.8

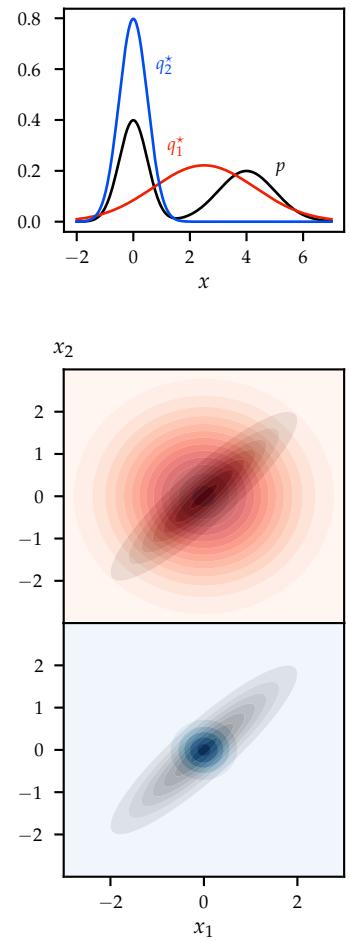


Figure 5.9: Comparison of the **forward** KL-divergence q_1^* and the **reverse** KL-divergence q_2^* when used to approximate the **true posterior** p . The first plot shows the PDFs in a one-dimensional feature space where p is a mixture of two univariate Gaussians. The second plot shows contour lines of the PDFs in a two-dimensional feature space where the non-diagonal Gaussian p is approximated by diagonal Gaussians q_1^* and q_2^* . It can be seen that q_1^* selects the variance and q_2^* selects the mode of p . The approximation q_1^* is more conservative than the (overconfident) approximation q_2^* .

tion. The forward KL-divergence, in contrast, is more conservative and yields what one could consider the “desired” approximation.

Recall that in the blueprint of variational inference (5.25) we used the reverse KL-divergence. This is for computational reasons. Observe that to approximate the KL-divergence $\text{KL}(p\|q)$ using Monte Carlo sampling, we would need to obtain samples from p yet p is the intractable posterior distribution which we were trying to approximate in the first place. Crucially, observe that if the true posterior $p(\cdot | \mathbf{x}_{1:n}, y_{1:n})$ is in the variational family \mathcal{Q} , then

$$\arg \min_{q \in \mathcal{Q}} \text{KL}(q \| p(\cdot | \mathbf{x}_{1:n}, y_{1:n})) \stackrel{a.s.}{=} p(\cdot | \mathbf{x}_{1:n}, y_{1:n}), \quad (5.44)$$

as $\min_{q \in \mathcal{Q}} \text{KL}(q \| p(\cdot | \mathbf{x}_{1:n}, y_{1:n})) \stackrel{a.s.}{=} 0$

so minimizing reverse-KL still recovers the true posterior almost surely.

Remark 5.14: Greediness of reverse-KL

As in the previous example, consider the independent Gaussians

$$p \doteq \mathcal{N}(\boldsymbol{\mu}, \text{diag}_{i \in [d]} \{\sigma_i^2\}),$$

which we seek to approximate by a standard normal distribution $q \doteq \mathcal{N}(\mathbf{0}, \mathbf{I})$. Using (5.41), we obtain for the reverse KL-divergence,

$$\begin{aligned} \text{KL}(q \| p) &= \frac{1}{2} \left(\text{tr} \left(\text{diag}\{\sigma_i^{-2}\} \right) + \boldsymbol{\mu}^\top \text{diag}\{\sigma_i^{-2}\} \boldsymbol{\mu} - d \right. \\ &\quad \left. + \log \det \left(\text{diag}\{\sigma_i^2\} \right) \right) \\ &= \frac{1}{2} \sum_{i=1}^d \left(\sigma_i^{-2} + \frac{\mu_i^2}{\sigma_i^2} - 1 + \log \sigma_i^2 \right). \end{aligned} \quad (5.45)$$

Here, σ_i^{-2} penalizes small variance, μ_i^2/σ_i^2 penalizes a large mean, and $\log \sigma_i^2$ penalizes large variance. Compare this to the expression for the forward KL-divergence $\text{KL}(p\|q)$ that we have seen in Equation (5.43). In particular, observe that reverse-KL penalizes large variance less strongly than forward-KL.

Note, however, that reverse-KL is not greedy in the same sense as Laplace approximation, as it does still take the variance into account and does not *purely* match the mode of p .

5.4.6 Interlude: Minimizing Forward KL-Divergence

Before completing the blueprint of variational inference in Section 5.5 by showing how *reverse-KL* can be efficiently minimized, we will digress briefly and relate minimizing *forward-KL* to two other well-known

inference algorithms. This discussion will deepen our understanding of the KL-divergence and its role in probabilistic inference, but feel free to skip ahead to Section 5.5 if you are eager to complete the blueprint.

Minimizing forward-KL as maximum likelihood estimation: First, we observe that minimizing the forward KL-divergence is equivalent to maximum likelihood estimation on an infinitely large sample size. The classical application of this result is in the setting where $p(\mathbf{x})$ is a generative model, and we aim to estimate its density with the parameterized model q_λ .

Lemma 5.15 (Forward KL-divergence as MLE). *Given some generative model $p(\mathbf{x})$ and a likelihood $q_\lambda(\mathbf{x}) = q(\mathbf{x} \mid \lambda)$ (that we use to approximate the true data distribution), we have*

$$\arg \min_{\lambda \in \Lambda} \text{KL}(p \parallel q_\lambda) \stackrel{\text{a.s.}}{=} \arg \max_{\lambda \in \Lambda} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \log q(\mathbf{x}_i \mid \lambda), \quad (5.46)$$

where $\mathbf{x}_i \stackrel{\text{iid}}{\sim} p$ are independent samples from the true data distribution.

Proof.

$$\begin{aligned} \text{KL}(p \parallel q_\lambda) &= H[p \parallel q_\lambda] - H[p] \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim p} [-\log q(\mathbf{x} \mid \lambda)] + \text{const} \\ &\stackrel{\text{a.s.}}{=} -\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \log q(\mathbf{x}_i \mid \lambda) + \text{const} \end{aligned}$$

where $\mathbf{x}_i \stackrel{\text{iid}}{\sim} p$ are independent samples. \square

using the definition of KL-divergence (5.34)

dropping $H[p]$ and using the definition of cross-entropy (5.32)

using Monte Carlo sampling, i.e., the law of large numbers (A.36)

This tells us that *any* maximum likelihood estimate q_λ minimizes the forward KL-divergence to the empirical data distribution. Note that here, we aim to learn model parameters λ for estimating the probability of \mathbf{x} , whereas in the setting of variational probabilistic inference, we want to learn parameters λ of a distribution over θ and θ parameterizes a *distribution over \mathbf{x}* . This interpretation is therefore not immediately useful for probabilistic inference (i.e., in the setting where p is a posterior distribution over model parameters θ) as a maximum likelihood estimate requires i.i.d. samples from p which we cannot easily obtain in this case.⁸

Example 5.16: Minimizing cross-entropy

Minimizing the KL-divergence between p and q_λ is equivalent to minimizing cross-entropy since $\text{KL}(p \parallel q_\lambda) = H[p \parallel q_\lambda] - H[p]$ and $H[p]$ is constant with respect to p .

⁸ It is possible to obtain “approximate” samples using Markov chain Monte Carlo (MCMC) methods which we discuss in Chapter 6.

Lets consider an example in a binary classification problem with the label $y \in \{0, 1\}$ and predicted class probability $\hat{y} \in [0, 1]$ for some fixed input. It is natural to use cross-entropy as a measure of dissimilarity between y and \hat{y} ,

$$\begin{aligned}\ell_{\text{bce}}(\hat{y}; y) &\doteq H[\text{Bern}(y) \parallel \text{Bern}(\hat{y})] \\ &= - \sum_{x \in \{0, 1\}} \text{Bern}(x; y) \log \text{Bern}(x; \hat{y}) \quad (5.47) \\ &= -y \log \hat{y} - (1 - y) \log(1 - \hat{y}).\end{aligned}$$

This loss function is also known as the *binary cross-entropy loss* and we will discuss it in more detail in Section 7.1.3 in the context of neural networks.

Minimizing forward-KL as moment matching: Now to a second interpretation of minimizing forward-KL. *Moment matching* (also known as the *method of moments*) is a technique for approximating an unknown distribution p with a parameterized distribution q_λ where λ is chosen such that q_λ matches the (estimated) moments of p . For example, given the estimates a and B of the first and second moment of p ,⁹ and if q_λ is a Gaussian with parameters $\lambda = \{\mu, \Sigma\}$, then moment matching chooses λ as the solution to

$$\begin{aligned}\mathbb{E}_p[\theta] &\approx a \stackrel{!}{=} \mu = \mathbb{E}_{q_\lambda}[\theta] \\ \mathbb{E}_p[\theta\theta^\top] &\approx B \stackrel{!}{=} \Sigma + \mu\mu^\top = \mathbb{E}_{q_\lambda}[\theta\theta^\top].\end{aligned}$$

In general the number of moments to be matched (i.e., the number of equations) is adjusted such that it is equal to the number of parameters to be estimated. We will see now that the “matching” of moments is also ensured when q_λ is obtained by minimizing the forward KL-divergence within the family of Gaussians.

The Gaussian PDF can be expressed as

$$\mathcal{N}(\theta; \mu, \Sigma) = \frac{1}{Z(\lambda)} \exp(\lambda^\top s(\theta)) \quad \text{where} \quad (5.48)$$

$$\lambda \doteq \begin{bmatrix} \Sigma^{-1}\mu \\ \text{vec}[\Sigma^{-1}] \end{bmatrix} \quad (5.49)$$

$$s(\theta) \doteq \begin{bmatrix} \theta \\ \text{vec}[-\frac{1}{2}\theta\theta^\top] \end{bmatrix} \quad (5.50)$$

and $Z(\lambda) \doteq \int \exp(\lambda^\top s(\theta)) d\theta$, and we will confirm this in just a moment.¹⁰ The family of distributions with densities of the form (5.48) — with an additional scaling constant $h(\theta)$ which is often 1 — is called

⁹ These estimates are computed using the samples from p . For example, using a sample mean and a sample variance to compute the estimates of the first and second moment.

using the definition of variance (1.35)

¹⁰ Given a matrix $A \in \mathbb{R}^{n \times m}$, we use

$$\text{vec}[A] \in \mathbb{R}^{n \cdot m}$$

to denote the row-by-row concatenation of A yielding a vector of length $n \cdot m$.

the *exponential family* of distributions. Here, $s(\theta)$ are the *sufficient statistics*, λ are called the *natural parameters*, and $Z(\lambda)$ is the normalizing constant. In this context, $Z(\lambda)$ is often called the *partition function*.

To see that the Gaussian is indeed part of the exponential family as promised in Equations (5.49) and (5.50), consider

$$\begin{aligned}\mathcal{N}(\theta; \mu, \Sigma) &\propto \exp\left(-\frac{1}{2}(\theta - \mu)^\top \Sigma^{-1}(\theta - \mu)\right) \\ &\propto \exp\left(\text{tr}\left(-\frac{1}{2}\theta^\top \Sigma^{-1}\theta\right) + \theta^\top \Sigma^{-1}\mu\right) \\ &= \exp\left(\text{tr}\left(-\frac{1}{2}\theta\theta^\top \Sigma^{-1}\right) + \theta^\top \Sigma^{-1}\mu\right) \\ &= \exp\left(\text{vec}\left[-\frac{1}{2}\theta\theta^\top\right]^\top \text{vec}[\Sigma^{-1}] + \theta^\top \Sigma^{-1}\mu\right).\end{aligned}$$

This allows us to express the forward KL-divergence as

$$\begin{aligned}\text{KL}(p\|q_\lambda) &= \int p(\theta) \log \frac{p(\theta)}{q_\lambda(\theta)} d\theta \\ &= - \int p(\theta) \cdot \lambda^\top s(\theta) d\theta + \log Z(\lambda) + \text{const.}\end{aligned}$$

expanding the inner product and using that $\text{tr}(x) = x$ for all $x \in \mathbb{R}$

using that the trace is invariant under cyclic permutations

using $\text{tr}(AB) = \text{vec}[A]^\top \text{vec}[B]$ for any $A, B \in \mathbb{R}^{n \times n}$

using that $\int p(\theta) \log p(\theta) d\theta$ is constant

Differentiating with respect to the natural parameters λ gives

$$\begin{aligned}\nabla_\lambda \text{KL}(p\|q_\lambda) &= - \int p(\theta) s(\theta) d\theta + \frac{1}{Z(\lambda)} \int s(\theta) \exp(\lambda^\top s(\theta)) d\theta \\ &= -\mathbb{E}_{\theta \sim p}[s(\theta)] + \mathbb{E}_{\theta \sim q_\lambda}[s(\theta)].\end{aligned}$$

Hence, for any minimizer of $\text{KL}(p\|q_\lambda)$, we have that the sufficient statistics under p and q_λ match:

$$\mathbb{E}_p[s(\theta)] = \mathbb{E}_{q_\lambda}[s(\theta)]. \quad (5.51)$$

Therefore, in the Gaussian case,

$$\mathbb{E}_p[\theta] = \mathbb{E}_{q_\lambda}[\theta] \quad \text{and} \quad \mathbb{E}_p\left[-\frac{1}{2}\theta\theta^\top\right] = \mathbb{E}_{q_\lambda}\left[-\frac{1}{2}\theta\theta^\top\right],$$

implying that

$$\mathbb{E}_p[\theta] = \mu \quad \text{and} \quad \text{Var}_p[\theta] = \mathbb{E}_p\left[\theta\theta^\top\right] - \mathbb{E}_p[\theta] \cdot \mathbb{E}_p[\theta]^\top = \Sigma \quad (5.52)$$

using Equation (1.35)

where μ and Σ are the mean and variance of the approximation q_λ , respectively. That is, a Gaussian q_λ minimizing $\text{KL}(p\|q_\lambda)$ has the same first and second moment as p . Combining this insight with our observation from Equation (5.46) that minimizing forward-KL is equivalent to maximum likelihood estimation, we see that if we use MLE to fit a Gaussian to given data, this Gaussian will eventually match the first and second moments of the data distribution.

5.5 Evidence Lower Bound

Let us return to the blueprint of variational inference from Section 5.3. To complete this blueprint, it remains to show that the reverse KL-divergence can be minimized efficiently. We have

$$\begin{aligned}
 \text{KL}(q\|p(\cdot | \mathbf{x}_{1:n}, y_{1:n})) &= \mathbb{E}_{\theta \sim q} \left[\log \frac{q(\theta)}{p(\theta | \mathbf{x}_{1:n}, y_{1:n})} \right] \\
 &= \mathbb{E}_{\theta \sim q} \left[\log \frac{p(y_{1:n} | \mathbf{x}_{1:n})q(\theta)}{p(y_{1:n}, \theta | \mathbf{x}_{1:n})} \right] \\
 &= \log p(y_{1:n} | \mathbf{x}_{1:n}) \\
 &\quad - \underbrace{\mathbb{E}_{\theta \sim q} [\log p(y_{1:n}, \theta | \mathbf{x}_{1:n})] - H[q]}_{-L(q, p; \mathcal{D}_n)}
 \end{aligned}$$

using the definition of the KL-divergence (5.34)

using the definition of conditional probability (1.8)

using linearity of expectation (1.20)

where $L(q, p; \mathcal{D}_n)$ is called the *evidence lower bound* (ELBO) given the data $\mathcal{D}_n = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$. This gives the relationship

$$L(q, p; \mathcal{D}_n) = \underbrace{\log p(y_{1:n} | \mathbf{x}_{1:n})}_{\text{const}} - \text{KL}(q\|p(\cdot | \mathbf{x}_{1:n}, y_{1:n})). \quad (5.53)$$

Thus, maximizing the ELBO coincides with minimizing reverse-KL.

Maximizing the ELBO,

$$L(q, p; \mathcal{D}_n) \doteq \mathbb{E}_{\theta \sim q} [\log p(y_{1:n}, \theta | \mathbf{x}_{1:n})] + H[q], \quad (5.54)$$

selects q that has large joint likelihood $p(y_{1:n}, \theta | \mathbf{x}_{1:n})$ and large entropy $H[q]$. The ELBO can also be expressed in various other forms:

$$L(q, p; \mathcal{D}_n) = \mathbb{E}_{\theta \sim q} [\log p(y_{1:n}, \theta | \mathbf{x}_{1:n}) - \log q(\theta)] \quad (5.55a)$$

$$= \mathbb{E}_{\theta \sim q} [\log p(y_{1:n} | \mathbf{x}_{1:n}, \theta) + \log p(\theta) - \log q(\theta)] \quad (5.55b)$$

$$= \underbrace{\mathbb{E}_{\theta \sim q} [\log p(y_{1:n} | \mathbf{x}_{1:n}, \theta)]}_{\text{log-likelihood}} - \underbrace{\text{KL}(q\|p(\cdot))}_{\text{proximity to prior}}. \quad (5.55c)$$

using the definition of entropy (5.27)

using the product rule (1.11)

using the definition of KL-divergence (5.34)

where we denote by $p(\cdot)$ the prior distribution. Equation (5.55c) highlights the connection to probabilistic inference, namely that maximizing the ELBO selects a variational distribution q that is close to the prior distribution $p(\cdot)$ while also maximizing the *average* likelihood of the data $p(y_{1:n} | \mathbf{x}_{1:n}, \theta)$ for $\theta \sim q$. This is in contrast to maximum a posteriori estimation, which picks a *single* model θ that maximizes the likelihood and proximity to the prior. As an example, let us look at the case where the prior is noninformative, i.e., $p(\cdot) \propto 1$. In this case, the ELBO simplifies to $\mathbb{E}_{\theta \sim q} [\log p(y_{1:n} | \mathbf{x}_{1:n}, \theta)] + H[q] + \text{const}$. That is, maximizing the ELBO maximizes the average likelihood of the data under the variational distribution q while regularizing q to have high

entropy. Why is it reasonable to maximize the entropy of q ? Consider two distributions q_1 and q_2 under which the data is “equally” likely and which are “equally” close to the prior. Maximizing the entropy selects the distribution that exhibits the most uncertainty which is in accordance with the maximum entropy principle.¹¹

Recalling that KL-divergence is non-negative, it follows from Equation (5.53) that the evidence lower bound is a (uniform¹²) lower bound to the *evidence* $\log p(y_{1:n} | \mathbf{x}_{1:n})$:

$$\log p(y_{1:n} | \mathbf{x}_{1:n}) \geq L(q, p; \mathcal{D}_n). \quad (5.56)$$

This indicates that maximizing the evidence lower bound is an adequate method of model selection which can be used instead of maximizing the evidence (marginal likelihood) directly (as was discussed in Section 4.4.2). Note that this inequality lower bounds the logarithm of an integral by an expectation of a logarithm over some variational distribution q . Hence, the ELBO is a family of lower bounds — one for each variational distribution. Such inequalities are called variational inequalities.

Example 5.17: Gaussian VI vs Laplace approximation

Consider maximizing the ELBO within the variational family of Gaussians $\mathcal{Q} \doteq \{q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}$. How does this relate to Laplace approximation which also fits a Gaussian approximation to the posterior

$$p(\boldsymbol{\theta} | \mathcal{D}) \propto p(y_{1:n}, \boldsymbol{\theta} | \mathbf{x}_{1:n})?$$

It turns out that both approximations are closely related. Indeed, it can be shown that while the Laplace approximation is fitted *locally* at the MAP estimate $\hat{\boldsymbol{\theta}}$, satisfying

$$\begin{aligned} \mathbf{0} &= \nabla_{\boldsymbol{\theta}} \log p(y_{1:n}, \boldsymbol{\theta} | \mathbf{x}_{1:n}) \\ \boldsymbol{\Sigma}^{-1} &= -H_{\boldsymbol{\theta}} \log p(y_{1:n}, \boldsymbol{\theta} | \mathbf{x}_{1:n})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}, \end{aligned} \quad (5.57)$$

Gaussian variational inference satisfies the conditions of the Laplace approximation *on average* with respect to the approximation q (??):

$$\begin{aligned} \mathbf{0} &= \mathbb{E}_{\boldsymbol{\theta} \sim q} [\nabla_{\boldsymbol{\theta}} \log p(y_{1:n}, \boldsymbol{\theta} | \mathbf{x}_{1:n})] \\ \boldsymbol{\Sigma}^{-1} &= -\mathbb{E}_{\boldsymbol{\theta} \sim q} [H_{\boldsymbol{\theta}} \log p(y_{1:n}, \boldsymbol{\theta} | \mathbf{x}_{1:n})]. \end{aligned} \quad (5.58)$$

For this reason, the Gaussian variational approximation does not suffer from the same overconfidence as the Laplace approximation.¹³

¹¹ In Section 1.2.1, we discussed the maximum entropy principle and the related principle of indifference at length. In simple terms, the maximum entropy principle states that without information, we should choose the distribution that is maximally uncertain.

¹² That is, the bound holds for any variational distribution q (with full support).

Problem 5.10

¹³ see Figure 5.1

Example 5.18: ELBO for Bayesian logistic regression

Recall that Bayesian logistic regression uses the prior distribution $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.¹⁴

Suppose we use the variational family \mathcal{Q} of all diagonal Gaussians from Example 5.5. We have already seen in Equation (5.43) that for a prior $p \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and a variational distribution

$$q_\lambda \sim \mathcal{N}(\boldsymbol{\mu}, \text{diag}_{i \in [d]} \{\sigma_i^2\}),$$

we have

$$\text{KL}(q \| p(\cdot)) = \frac{1}{2} \sum_{i=1}^d (\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2).$$

It remains to find the expected likelihood under models from our approximate posterior:

$$\begin{aligned} \mathbb{E}_{\mathbf{w} \sim q_\lambda} [\log p(y_{1:n} | \mathbf{x}_{1:n}, \mathbf{w})] &= \mathbb{E}_{\mathbf{w} \sim q_\lambda} \left[\sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \mathbf{w}) \right] \\ &= \mathbb{E}_{\mathbf{w} \sim q_\lambda} \left[- \sum_{i=1}^n \ell_{\log}(\mathbf{w}; \mathbf{x}_i, y_i) \right]. \end{aligned} \tag{5.59}$$

¹⁴ We omit the scaling factor σ_p^2 here for simplicity.

using independence of the data

substituting the logistic loss (5.13)

5.5.1 Gradient of Evidence Lower Bound

We have yet to discuss how the optimization problem of maximizing the ELBO can be solved efficiently. A suitable tool is stochastic gradient descent (SGD), however, SGD requires unbiased gradient estimates of the loss $\ell(\lambda; \mathcal{D}_n) \doteq -L(q_\lambda, p; \mathcal{D}_n)$. That is, we need to obtain gradient estimates of

$$\nabla_\lambda L(q_\lambda, p; \mathcal{D}_n) = \nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [\log p(y_{1:n} | \mathbf{x}_{1:n}, \theta)] - \nabla_\lambda \text{KL}(q_\lambda \| p(\cdot)). \tag{5.60}$$

using the definition of the ELBO (5.55c)

Typically, the KL-divergence (and its gradient) can be computed exactly for commonly used variational families. For example, we have already seen a closed-form expression of the KL-divergence for Gaussians in Equation (5.41).

Obtaining the gradient of the expected log-likelihood is more difficult. This is because the expectation integrates over the measure q_λ , which depends on the variational parameters λ . Thus, we cannot move the gradient operator inside the expectation as commonly done (cf. Appendix A.1.5). There are two main techniques which are used to

rewrite the gradient in such a way that Monte Carlo sampling becomes possible.

One approach is to use *score gradients* via the “score function trick”:

$$\begin{aligned} \nabla_{\lambda} \mathbb{E}_{\theta \sim q_{\lambda}} [\log p(y_{1:n} | x_{1:n}, \theta)] \\ = \mathbb{E}_{\theta \sim q_{\lambda}} [\log p(y_{1:n} | x_{1:n}, \theta) \underbrace{\nabla_{\lambda} \log q_{\lambda}(\theta)}_{\text{score function}}], \end{aligned} \quad (5.61)$$

which we introduce in Section 12.3.2 in the context of reinforcement learning. More common in the context of variational inference is the so-called “reparameterization trick”.

Theorem 5.19 (Reparameterization trick). *Given a random variable $\varepsilon \sim \phi$ (which is independent of λ) and given a differentiable and invertible function $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$. We let $\theta \doteq g(\varepsilon; \lambda)$. Then,*

$$q_{\lambda}(\theta) = \phi(\varepsilon) \cdot |\det(D_{\varepsilon}g(\varepsilon; \lambda))|^{-1}, \quad (5.62)$$

$$\mathbb{E}_{\theta \sim q_{\lambda}} [f(\theta)] = \mathbb{E}_{\varepsilon \sim \phi} [f(g(\varepsilon; \lambda))] \quad (5.63)$$

for a “nice” function $f : \mathbb{R}^d \rightarrow \mathbb{R}^e$.

Proof. By the change of variables formula (1.43) and using $\varepsilon = g^{-1}(\theta; \lambda)$,

$$\begin{aligned} q_{\lambda}(\theta) &= \phi(\varepsilon) \cdot \left| \det(D_{\theta}g^{-1}(\theta; \lambda)) \right| \\ &= \phi(\varepsilon) \cdot \left| \det((D_{\varepsilon}g(\varepsilon; \lambda))^{-1}) \right| \\ &= \phi(\varepsilon) \cdot |\det(D_{\varepsilon}g(\varepsilon; \lambda))|^{-1}. \end{aligned}$$

by the inverse function theorem,
 $Dg^{-1}(y) = Dg(x)^{-1}$
using $\det(A^{-1}) = \det(A)^{-1}$

Equation (5.63) is a direct consequence of the law of the unconscious statistician (1.22). \square

In other words, the reparameterization trick allows a change of “densities” by finding a function $g(\cdot; \lambda)$ and a reference density ϕ such that $q_{\lambda} = g(\cdot; \lambda)_\sharp \phi$ is the pushforward of ϕ under perturbation g . Applying the reparameterization trick, we can swap the order of gradient and expectation,

$$\nabla_{\lambda} \mathbb{E}_{\theta \sim q_{\lambda}} [f(\theta)] = \mathbb{E}_{\varepsilon \sim \phi} [\nabla_{\lambda} f(g(\varepsilon; \lambda))]. \quad (5.64)$$

using Equation (A.5)

We call a distribution q_{λ} *reparameterizable* if it admits reparameterization, i.e., if we can find g and a suitable reference density ϕ which is independent of λ .

Example 5.20: Reparameterization trick for Gaussians

Suppose we use a Gaussian variational approximation,

$$q_\lambda(\theta) \doteq \mathcal{N}(\theta; \mu, \Sigma),$$

where we assume Σ to have full rank (i.e., be invertible). We have seen in Equation (1.78) that a Gaussian random vector $\varepsilon \sim \mathcal{N}(\mathbf{0}, I)$ following a standard normal distribution can be transformed to follow the Gaussian distribution q_λ by using the linear transformation,

$$\theta = g(\varepsilon; \lambda) \doteq \Sigma^{1/2}\varepsilon + \mu. \quad (5.65)$$

In particular, we have

$$\varepsilon = g^{-1}(\theta; \lambda) = \Sigma^{-1/2}(\theta - \mu) \quad \text{and} \quad (5.66)$$

$$\phi(\varepsilon) = q_\lambda(\theta) \cdot |\det(\Sigma^{1/2})|. \quad (5.67)$$

In the following, we write $C \doteq \Sigma^{1/2}$. Let us now derive the gradient estimate for the evidence lower bound assuming the Gaussian variational approximation from Example 5.20. This approach extends to any reparameterizable distribution.

$$\begin{aligned} & \nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [\log p(y_{1:n} | x_{1:n}, \theta)] \\ &= \nabla_{C, \mu} \mathbb{E}_{\varepsilon \sim \mathcal{N}(\mathbf{0}, I)} \left[\log p(y_{1:n} | x_{1:n}, \theta) \Big|_{\theta=C\varepsilon+\mu} \right] \end{aligned} \quad (5.68)$$

$$\begin{aligned} &= n \cdot \nabla_{C, \mu} \mathbb{E}_{\varepsilon \sim \mathcal{N}(\mathbf{0}, I)} \left[\frac{1}{n} \sum_{i=1}^n \log p(y_i | x_i, \theta) \Big|_{\theta=C\varepsilon+\mu} \right] \\ &= n \cdot \nabla_{C, \mu} \mathbb{E}_{\varepsilon \sim \mathcal{N}(\mathbf{0}, I)} \mathbb{E}_{i \sim \text{Unif}([n])} \left[\log p(y_i | x_i, \theta) \Big|_{\theta=C\varepsilon+\mu} \right] \\ &= n \cdot \mathbb{E}_{\varepsilon \sim \mathcal{N}(\mathbf{0}, I)} \mathbb{E}_{i \sim \text{Unif}([n])} \left[\nabla_{C, \mu} \log p(y_i | x_i, \theta) \Big|_{\theta=C\varepsilon+\mu} \right] \end{aligned} \quad (5.69)$$

$$\approx n \cdot \frac{1}{m} \sum_{j=1}^m \nabla_{C, \mu} \log p(y_{i_j} | x_{i_j}, \theta) \Big|_{\theta=C\varepsilon_j+\mu} \quad (5.70)$$

by solving Equation (5.65) for ε

using the reparameterization trick (i.e., the change of variables formula) (5.62)

using the reparameterization trick (5.63)

using independence of the data and extending with n/n

interpreting the sum as an expectation

using Equation (A.5)

using Monte Carlo sampling

where $\varepsilon_j \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, I)$ and $i_j \stackrel{\text{iid}}{\sim} \text{Unif}([n])$. This yields an unbiased gradient estimate, which we can use with stochastic gradient descent to maximize the evidence lower bound. We have successfully recast the difficult problems of learning and inference as an optimization problem!

The procedure of approximating the true posterior using a variational posterior by maximizing the evidence lower bound using stochastic optimization is also called *black box stochastic variational inference* (Ranganath et al., 2014; Titsias and Lázaro-Gredilla, 2014; Duvenaud and

Adams, 2015). The only requirement is that we can obtain unbiased gradient estimates from the evidence lower bound (and the likelihood). We have just discussed one of many approaches to obtain such gradient estimates (Mohamed et al., 2020). If we use the variational family of diagonal Gaussians, we only require twice as many parameters as other inference techniques like MAP estimation. The performance can be improved by using natural gradients and variance reduction techniques for the gradient estimates such as control variates.

5.5.2 Minimizing Surprise via Exploration and Exploitation

Now that we have established a way to optimize the ELBO, let us dwell a bit more on its interpretation. Observe that the evidence can also be interpreted as the negative surprise about the observations under the prior distribution $p(\theta)$, and by negating Equation (5.56), we obtain the variational upper bound

$$\begin{aligned} S[p(y_{1:n} | \mathbf{x}_{1:n})] &\leq \underbrace{\mathbb{E}_{\theta \sim q}[S[p(y_{1:n}, \theta | \mathbf{x}_{1:n})]] - H[q]}_{\text{called energy}} \\ &= -L(q, p; \mathcal{D}_n). \end{aligned} \quad (5.71)$$

Here, $-L(q, p; \mathcal{D}_n)$ is commonly called the (*variational*) free energy with respect to q . Free energy can also be characterized as

$$-L(q, p; \mathcal{D}_n) = \underbrace{\mathbb{E}_{\theta \sim q}[S[p(y_{1:n} | \mathbf{x}_{1:n}, \theta)]]}_{\text{average surprise}} + \underbrace{\text{KL}(q \| p(\cdot))}_{\text{proximity to prior}}, \quad (5.72)$$

analogously to Equation (5.55c)

and therefore its minimization is minimizing the average surprise about the data under the variational distribution q while maximizing proximity to the prior $p(\cdot)$. Systems that minimize the surprise in their observations are widely studied in many areas of science.¹⁵

To minimize surprise, the free energy makes apparent a natural trade-off between two extremes: On the one hand, models $q(\theta)$ that “overfit” to observations (e.g., by using a point estimate of θ), and hence, result in a large surprise when new observations deviate from this specific belief. On the other hand, models $q(\theta)$ that “underfit” to observations (e.g., by expecting outcomes with equal probability), and hence, any observation results in a non-negligible surprise. Either of these extremes is undesirable.

¹⁵ Refer to the *free energy principle* which was originally introduced by the neuroscientist Karl Friston (Friston, 2010).

As we have alluded to previously when introducing the ELBO, the two terms constituting free energy map neatly onto this tradeoff. Therein, the *entropy* (which is maximized) encourages q to have uncertainty, in other words, to “explore” beyond the finite data. In contrast, the *energy* (which is minimized) encourages q to fit the observed data closely, in

other words, to “exploit” the finite data. This tradeoff is ubiquitous in approximations to probabilistic inference that deal with limited computational resources and limited time, and we will encounter it many times more. We will point out these connections as we go along.

Since this tradeoff is so fundamental and appears in many branches of science under different names, it is difficult to give it an appropriate unifying name. The essence of this tradeoff can be captured as a *principle of curiosity and conformity*, which suggests that reasoning under uncertainty requires curiosity to entertain and pursue alternative explanations of the data *and* conformity to make consistent predictions.

Discussion

We have explored variational inference, where we approximate the intractable posterior distribution of probabilistic inference with a simpler distribution. We operationalized this idea by turning the inference problem, which requires computing high-dimensional integrals, into a tractable optimization problem. Gaussians are frequently used as variational distributions due to their versatility and compact representation.

Nevertheless, recall from Figure 5.4 that while the estimation error in variational inference can be small, choosing a variational family that is too simple can lead to a large approximation error. We have seen that for posteriors that are multimodal or have heavy tails, Gaussians may not provide a good approximation. In the next chapter, we will explore alternative techniques for approximate inference that can handle more complex posteriors.

Problems

5.1. Logistic loss.

1. Derive the gradient of ℓ_{\log} as given in Equation (5.14).
2. Show that

$$H_w \ell_{\log}(\mathbf{w}^\top \mathbf{x}; y) = \mathbf{x} \mathbf{x}^\top \cdot \sigma(\mathbf{w}^\top \mathbf{x}) \cdot (1 - \sigma(\mathbf{w}^\top \mathbf{x})). \quad (5.73)$$

Hint: Begin by deriving the first derivative of the logistic function, and use the chain rule of multivariate calculus,

$$\underbrace{D_x(f \circ g)}_{\mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}} = \underbrace{(D_{g(x)} f \circ g) \cdot D_x g}_{\mathbb{R}^n \rightarrow \mathbb{R}^{m \times k}, \mathbb{R}^{k \times n}} \quad (5.74)$$

where $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$ and $f : \mathbb{R}^k \rightarrow \mathbb{R}^m$.

3. Is the logistic loss ℓ_{\log} convex in \mathbf{w} ?

5.2. Gaussian process classification.

In this exercise, we will study the use of Gaussian processes for classification tasks, commonly called *Gaussian process classification* (GPC). Linear logistic regression is extended to GPC by replacing the Gaussian prior over weights with a GP prior on f ,

$$f \sim \mathcal{GP}(0, k), \quad y | x, f \sim \text{Bern}(\sigma(f(x))) \quad (5.75)$$

where $\sigma : \mathbb{R} \rightarrow (0, 1)$ is some logistic-type function. Note that Bayesian logistic regression is the special case where k is the linear kernel and σ is the logistic function. This is analogous to the relationship of Bayesian linear regression and Gaussian process regression.

In the GP regression setting of Chapter 4, y_i was assumed to be a noisy observation of $f(x_i)$. In the classification setting, we now have that $y_i \in \{-1, +1\}$ is a binary class label and $f(x_i) \in \mathbb{R}$ is a latent value. We study the setting where $\sigma(z) = \Phi(z; 0, \sigma_n^2)$ is the CDF of a univariate Gaussian with mean 0 and variance σ_n^2 , also called a *probit likelihood*.

To make probabilistic predictions for a query point x^* , we first compute the distribution of the latent variable f^* ,

$$p(f^* | x_{1:n}, y_{1:n}, x^*) = \int p(f^* | x_{1:n}, x^*, f) p(f | x_{1:n}, y_{1:n}) df \quad (5.76)$$

where $p(f | x_{1:n}, y_{1:n})$ is the posterior over the latent variables.

1. Assuming that we can efficiently compute $p(f^* | x_{1:n}, y_{1:n}, x^*)$ (approximately), describe how we can find the predictive posterior $p(y^* = +1 | x_{1:n}, y_{1:n}, x^*)$.
2. The posterior over the latent variables is not a Gaussian as we used a non-Gaussian likelihood, and hence, the integral of the latent predictive posterior (5.76) is analytically intractable. A common technique is to approximate the latent posterior $p(f | x_{1:n}, y_{1:n})$ with a Gaussian using a Laplace approximation $q = \mathcal{N}(\hat{f}, \Lambda^{-1})$. It is generally not possible to obtain an analytical representation of the mode of the Laplace approximation \hat{f} . Instead, \hat{f} is commonly found using a second-order optimization scheme such as Newton's method.
 - (a) Find the precision matrix Λ of the Laplace approximation.

Hint: Observe that for a label $y_i \in \{-1, +1\}$, the probability of a correct classification given the latent value f_i is $p(y_i | f_i) = \sigma(y_i f_i)$, where we use the symmetry of the probit likelihood around 0.
- (b) Assume that $k(x, x') = x^\top x'$ is the linear kernel ($\sigma_p = 1$) and that σ is the logistic function (5.9). Show for this setting that the matrix Λ derived in (a) is equivalent to the precision matrix

using sum rule (1.7) and product rule (1.11), and $f^* \perp y_{1:n} | f$

Λ' of the Laplace approximation of Bayesian logistic regression (5.17).¹⁶ You may assume that $\hat{f}_i = \hat{\mathbf{w}}^\top \mathbf{x}_i$.

Hint: First derive under which condition Λ and Λ' are “equivalent”.

- (c) Observe that the (approximate) latent predictive posterior

$$q(f^* | \mathbf{x}_{1:n}, y_{1:n}, \mathbf{x}^*) \doteq \int p(f^* | \mathbf{x}_{1:n}, \mathbf{x}^*, f) q(f | \mathbf{x}_{1:n}, y_{1:n}) df$$

which uses the Laplace approximation of the latent posterior is Gaussian.¹⁷ Determine its mean and variance.

Hint: Condition on the latent variables f using the laws of total expectation and variance.

- (d) Compare the prediction $p(f^* | \mathbf{x}_{1:n}, y_{1:n}, \mathbf{x}^*)$ you obtained in (1) (but now using the Laplace-approximated latent predictive posterior) to the prediction $\sigma(\mathbb{E}_{f^* \sim q}[f^*])$. Are they identical? If not, describe how they are different.
3. The use of the probit likelihood may seem arbitrary. Consider the following model which may be more natural,

$$f \sim \mathcal{GP}(0, k), \quad y = \mathbb{1}\{\underbrace{f(\mathbf{x}) + \varepsilon}_{\text{GP regression}} \geq 0\}, \quad \varepsilon \sim \mathcal{N}(0, \sigma_n^2). \quad (5.77)$$

Show that the model from Equation (5.75) using a noise-free latent process with probit likelihood $\Phi(z; 0, \sigma_n^2)$ is equivalent (in expectation over ε) to the model from Equation (5.77).

5.3. Jensen's inequality.

1. Prove the finite form of Jensen's inequality.

Theorem 5.21 (Jensen's inequality, finite form). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. Suppose that $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^n$ and $\theta_1, \dots, \theta_k \geq 0$ with $\theta_1 + \dots + \theta_k = 1$. Then,*

$$f(\theta_1 \mathbf{x}_1 + \dots + \theta_k \mathbf{x}_k) \leq \theta_1 f(\mathbf{x}_1) + \dots + \theta_k f(\mathbf{x}_k). \quad (5.78)$$

Observe that if X is a random variable with finite support, the above two versions of Jensen's inequality are equivalent.

2. Show that for any discrete distribution p supported on a finite domain of size n , $H[p] \leq \log_2 n$. This implies that the uniform distribution has maximum entropy.

5.4. Binary cross-entropy loss.

Show that the logistic loss (5.13) is equivalent to the binary cross-entropy loss with $\hat{y} = \sigma(\hat{f})$. That is,

$$\ell_{\text{log}}(\hat{f}; y) = \ell_{\text{bce}}(\hat{y}; y). \quad (5.79)$$

¹⁶ This should not be surprising since — as already mentioned — Gaussian process classification is a generalization of Bayesian logistic regression.

¹⁷ Using the Laplace-approximated latent posterior, $[f^* | f]$ are jointly Gaussian. Thus, it directly follows from Theorem 1.24 that the marginal distribution over f^* is also a Gaussian.

5.5. Gibbs' inequality.

1. Prove $\text{KL}(p\|q) \geq 0$ which is also known as *Gibbs' inequality*.
2. Let p and q be discrete distributions with finite identical support
A. Show that $\text{KL}(p\|q) = 0$ if and only if $p \equiv q$.

Hint: Use that if a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is strictly convex and $x_1, \dots, x_k \in \mathbb{R}^n$, $\theta_1, \dots, \theta_k \geq 0$, $\theta_1 + \dots + \theta_k = 1$, we have that

$$f(\theta_1 x_1 + \dots + \theta_k x_k) = \theta_1 f(x_1) + \dots + \theta_k f(x_k) \quad (5.80)$$

iff $x_1 = \dots = x_k$. This is a slight adaptation of Jensen's inequality in finite-form, which you proved in Problem 5.3.

5.6. Maximum entropy principle.

In this exercise we will prove that the normal distribution is the distribution with maximal entropy among all (univariate) distributions supported on \mathbb{R} with fixed mean μ and variance σ^2 . Let $g(x) \doteq \mathcal{N}(x; \mu, \sigma^2)$, and $f(x)$ be any distribution on \mathbb{R} with mean μ and variance σ^2 .

1. Prove that $\text{KL}(f\|g) = H[g] - H[f]$.

Hint: Equivalently, show that $H[f\|g] = H[g]$. That is, the expected surprise evaluated based on the Gaussian g is invariant to the true distribution f .

2. Conclude that $H[g] \geq H[f]$.

5.7. Probabilistic inference as a consequence of the maximum entropy principle.

Consider the family of generative models of the random vectors \mathbf{X} in \mathcal{X} and \mathbf{Y} in \mathcal{Y} :

$$\Delta^{\mathcal{X} \times \mathcal{Y}} = \left\{ q : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0} \mid \int_{\mathcal{X} \times \mathcal{Y}} q(x, y) dx dy = 1 \right\}.$$

Suppose that we observe \mathbf{Y} to be \mathbf{y}' , and are looking for a (new) generative model that is consistent with this information, that is,

$$q(\mathbf{y}) = \int_{\mathcal{X}} q(\mathbf{x}, \mathbf{y}) d\mathbf{x} = \delta_{\mathbf{y}'}(\mathbf{y}) \quad \text{using the sum rule (1.7)}$$

where $\delta_{\mathbf{y}'}$ denotes the point density at \mathbf{y}' . The product rule (1.11) implies that $q(\mathbf{x}, \mathbf{y}) = \delta_{\mathbf{y}'}(\mathbf{y}) \cdot q(\mathbf{x} \mid \mathbf{y})$, but any choice of $q(\mathbf{x} \mid \mathbf{y})$ is possible.

We will derive that given any fixed generative model $p_{\mathbf{X}, \mathbf{Y}}$, the “posterior” distribution $q_{\mathbf{X}}(\cdot) = p_{\mathbf{X}|\mathbf{Y}}(\cdot \mid \mathbf{y}')$ minimizes the relative entropy $\text{KL}(q_{\mathbf{X}, \mathbf{Y}} \| p_{\mathbf{X}, \mathbf{Y}})$ subject to the constraint $\mathbf{Y} = \mathbf{y}'$. In other words, among all distributions $q_{\mathbf{X}, \mathbf{Y}}$ that are consistent with the observation $\mathbf{Y} = \mathbf{y}'$, the posterior distribution $q_{\mathbf{X}}(\cdot) = p_{\mathbf{X}|\mathbf{Y}}(\cdot \mid \mathbf{y}')$ is the one with “maximum entropy”.

1. Show that the optimization problem

$$\begin{aligned} & \arg \min_{q \in \Delta^{\mathcal{X} \times \mathcal{Y}}} \text{KL}(q_{\mathbf{X}, \mathbf{Y}} \| p_{\mathbf{X}, \mathbf{Y}}) \\ & \text{subject to } q(\mathbf{y}) = \delta_{\mathbf{y}'}(\mathbf{y}) \quad \forall \mathbf{y} \in \mathcal{Y} \end{aligned}$$

is solved by $q(\mathbf{x}, \mathbf{y}) = \delta_{\mathbf{y}'}(\mathbf{y}) \cdot p(\mathbf{x} | \mathbf{y})$.

Hint: Solve the dual problem.

2. Conclude that $q(\mathbf{x}) = p(\mathbf{x} | \mathbf{y}')$.

5.8. KL-divergence of Gaussians.

Derive Equation (5.41).

Hint: If $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ in d dimensions, then we have that for any $\mathbf{m} \in \mathbb{R}^d$ and $\mathbf{A} \in \mathbb{R}^{d \times d}$,

$$\mathbb{E}\left[(\mathbf{X} - \mathbf{m})^\top \mathbf{A}(\mathbf{X} - \mathbf{m})\right] = (\boldsymbol{\mu} - \mathbf{m})^\top \mathbf{A}(\boldsymbol{\mu} - \mathbf{m}) + \text{tr}(\mathbf{A}\boldsymbol{\Sigma}) \quad (5.81)$$

5.9. Forward vs reverse KL.

1. Consider a factored approximation $q(x, y) = q(x)q(y)$ to a joint distribution $p(x, y)$. Show that to minimize the forward $\text{KL}(p \| q)$ we should set $q(x) = p(x)$ and $q(y) = p(y)$, i.e., the optimal approximation is a product of marginals.
2. Consider the following joint distribution p , where the rows represent y and the columns x :

	1	2	3	4
1	1/8	1/8	0	0
2	1/8	1/8	0	0
3	0	0	1/4	0
4	0	0	0	1/4

Show that the reverse $\text{KL}(q \| p)$ for this p has three distinct minima.

Identify those minima and evaluate $\text{KL}(q \| p)$ at each of them.

3. What is the value of $\text{KL}(q \| p)$ if we use the approximation $q(x, y) = p(x)p(y)$?

5.10. Gaussian VI vs Laplace approximation.

In this exercise, we compare the Laplace approximation from Section 5.1 to variational inference with the variational family of Gaussians,

$$\mathcal{Q} \doteq \{q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}.$$

1. Let p be any distribution on \mathbb{R} , and let $q^* = \arg \min_{q \in \mathcal{Q}} \text{KL}(p \| q)$. Show that q^* differs from the Laplace approximation of p .

Minimizing forward-KL is typically intractable, and we have seen that it is therefore common to minimize the reverse-KL instead:

$$\tilde{q} = \arg \min_{q \in \mathcal{Q}} \text{KL}(q \| p(\cdot | \mathcal{D}_n)).$$

2. Show that $\tilde{q} = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ satisfies Equation (5.58):

$$\begin{aligned}\mathbf{0} &= \mathbb{E}_{\boldsymbol{\theta} \sim \tilde{q}} [\nabla_{\boldsymbol{\theta}} \log p(y_{1:n}, \boldsymbol{\theta} | \mathbf{x}_{1:n})] \\ \boldsymbol{\Sigma}^{-1} &= -\mathbb{E}_{\boldsymbol{\theta} \sim \tilde{q}} [\mathbf{H}_{\boldsymbol{\theta}} \log p(y_{1:n}, \boldsymbol{\theta} | \mathbf{x}_{1:n})].\end{aligned}$$

Hint 1: For any positive definite and symmetric matrix \mathbf{A} , it holds that

$$\nabla_{\mathbf{A}} \log \det(\mathbf{A}) = \mathbf{A}^{-1}.$$

Hint 2: For any function f and Gaussian $p = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$,

$$\begin{aligned}\nabla_{\boldsymbol{\mu}} \mathbb{E}_{\mathbf{x} \sim p}[f(\mathbf{x})] &= \mathbb{E}_{\mathbf{x} \sim p}[\nabla_{\mathbf{x}} f(\mathbf{x})] \\ \nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{\mathbf{x} \sim p}[f(\mathbf{x})] &= \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p}[\mathbf{H}_{\mathbf{x}} f(\mathbf{x})].\end{aligned}\tag{5.82}$$

Recall the conditions satisfied by the Laplace approximation of the posterior $p(\boldsymbol{\theta} | \mathcal{D}) \propto p(y_{1:n}, \boldsymbol{\theta} | \mathbf{x}_{1:n})$ as detailed in Equation (5.57). The Laplace approximation is fitted *locally* at the MAP estimate $\hat{\boldsymbol{\theta}}$. Comparing Equation (5.57) to Equation (5.58), we see that Gaussian variational inference satisfies the conditions of the Laplace approximation *on average*. For more details, refer to Opper and Archambeau (2009).

5.11. Gradient of reverse-KL.

Suppose $p \doteq \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$ and a tractable distribution described by

$$q_{\boldsymbol{\lambda}} \doteq \mathcal{N}(\boldsymbol{\mu}, \text{diag}\{\sigma_1^2, \dots, \sigma_d^2\})$$

where $\boldsymbol{\mu} \doteq [\mu_1 \dots \mu_d]$ and $\boldsymbol{\lambda} \doteq [\mu_1 \dots \mu_d \ \sigma_1 \dots \sigma_d]$. Show that the gradient of $\text{KL}(q_{\boldsymbol{\lambda}} \| p(\cdot))$ with respect to $\boldsymbol{\lambda}$ is given by

$$\nabla_{\boldsymbol{\mu}} \text{KL}(q_{\boldsymbol{\lambda}} \| p(\cdot)) = \sigma_p^{-2} \boldsymbol{\mu}, \quad \text{and} \tag{5.83a}$$

$$\nabla_{[\sigma_1 \dots \sigma_d]} \text{KL}(q_{\boldsymbol{\lambda}} \| p(\cdot)) = \left[\frac{\sigma_1}{\sigma_p^2} - \frac{1}{\sigma_1} \quad \dots \quad \frac{\sigma_d}{\sigma_p^2} - \frac{1}{\sigma_d} \right]. \tag{5.83b}$$

5.12. Reparameterizable distributions.

1. Let $X \sim \text{Unif}([a, b])$ for any $a \leq b$. That is,

$$p_X(x) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{otherwise.} \end{cases} \tag{5.84}$$

Show that X can be reparameterized in terms of $\text{Unif}([0, 1])$. Hint:

You may use that for any $Y \sim \text{Unif}([a, b])$ and $c \in \mathbb{R}$,

- $Y + c \sim \text{Unif}([a + c, b + c])$ and
- $cY \sim \text{Unif}([c \cdot a, c \cdot b])$.

2. Let $Z \sim \mathcal{N}(\mu, \sigma^2)$ and $X \doteq e^Z$. That is, X is logarithmically normal distributed with parameters μ and σ^2 . Show that X can be reparameterized in terms of $\mathcal{N}(0, 1)$.
3. Show that Cauchy($0, 1$) can be reparameterized in terms of Unif($[0, 1]$). Finally, let us apply the reparameterization trick to compute the gradient of an expectation.
4. Let $\text{ReLU}(z) \doteq \max\{0, z\}$ and $w > 0$. Show that

$$\frac{d}{d\mu} \mathbb{E}_{x \sim \mathcal{N}(\mu, 1)} \text{ReLU}(wx) = w\Phi(\mu)$$

where Φ denotes the CDF of the standard normal distribution.

6

Markov Chain Monte Carlo Methods

Variational inference approximates the entire posterior distribution. However, note that the key challenge in probabilistic inference is not learning the posterior distribution, but using the posterior distribution for predictions,

$$p(y^* \mid \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}) = \int p(y^* \mid \mathbf{x}^*, \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathbf{x}_{1:n}, y_{1:n}) d\boldsymbol{\theta}. \quad (6.1)$$

This integral can be interpreted as an expectation over the posterior distribution,

$$= \mathbb{E}_{\boldsymbol{\theta} \sim p(\cdot \mid \mathbf{x}_{1:n}, y_{1:n})} [p(y^* \mid \mathbf{x}^*, \boldsymbol{\theta})]. \quad (6.2)$$

Observe that the likelihood $f(\boldsymbol{\theta}) \doteq p(y^* \mid \mathbf{x}^*, \boldsymbol{\theta})$ is easy to evaluate. The difficulty lies in sampling from the posterior distribution. Assuming we can obtain independent samples from the posterior distribution, we can use Monte Carlo sampling to obtain an unbiased estimate of the expectation,

$$\approx \frac{1}{m} \sum_{i=1}^m f(\boldsymbol{\theta}^{(i)}) \quad (6.3)$$

for independent $\boldsymbol{\theta}^{(i)} \stackrel{\text{iid}}{\sim} p(\cdot \mid \mathbf{x}_{1:n}, y_{1:n})$. The law of large numbers (A.36) and Hoeffding's inequality (A.41) imply that this estimator is consistent and sharply concentrated.¹

¹ For more details, see Appendix A.3.3.

Obtaining samples of the posterior distribution is therefore sufficient to perform approximate inference. Recall that the difficulty of computing the posterior p exactly, was in finding the normalizing constant Z ,

$$p(x) = \frac{1}{Z} q(x). \quad (6.4)$$

The joint likelihood q is typically easy to obtain. Note that $q(x)$ is proportional to the probability density associated with x , but q does not

integrate to 1. Such functions are also called a *finite measure*. Without normalizing q , we cannot directly sample from it.

Remark 6.1: The difficulty of sampling — even with a PDF

Even a decent approximation of Z does not yield a general efficient sampling method. For example, one very common approach to sampling is *inverse transform sampling* (cf. Appendix A.1.3) which requires an (approximate) quantile function. Computing the quantile function given an arbitrary PDF requires solving integrals over the domain of the PDF which is what we were trying to avoid in the first place.

The key idea of Markov chain Monte Carlo methods is to construct a Markov chain, which is efficient to simulate and has the stationary distribution p .

6.1 Markov Chains

To start, let us revisit the fundamental theory behind Markov chains.

Definition 6.2 (Markov chain). A (*finite and discrete-time*) *Markov chain* over the state space

$$S \doteq \{0, \dots, n - 1\} \quad (6.5)$$

is a stochastic process² $(X_t)_{t \in \mathbb{N}_0}$ valued in S such that the *Markov property* is satisfied:

$$X_{t+1} \perp X_{0:t-1} \mid X_t. \quad (6.6)$$

Intuitively, the Markov property states that future behavior is independent of past states given the present state.

Remark 6.3: Generalizations of Markov chains

One can also define continuous-state Markov chains (for example, where states are vectors in \mathbb{R}^d) and the results which we state for (finite) Markov chains will generally carry over. For a survey, refer to “General state space Markov chains and MCMC algorithms” (Roberts and Rosenthal, 2004).

Moreover, one can also consider continuous-time Markov chains. One example of such a continuous-space and continuous-time Markov chain is the *Wiener process* (cf. Remark 6.23).

We restrict our attention to *time-homogeneous* Markov chains,³ which

² A *stochastic process* is a sequence of random variables.

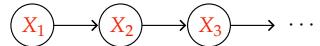


Figure 6.1: Directed graphical model of a Markov chain. The random variable X_{t+1} is conditionally independent of the random variables $X_{0:t-1}$ given X_t .

³ That is, the transition probabilities do not change over time.

can be characterized by a *transition function*,

$$p(x' | x) \doteq \mathbb{P}(X_{t+1} = x' | X_t = x). \quad (6.7)$$

As the state space is finite, we can describe the transition function by the *transition matrix*,

$$\mathbf{P} \doteq \begin{bmatrix} p(x_1 | x_1) & \cdots & p(x_n | x_1) \\ \vdots & \ddots & \vdots \\ p(x_1 | x_n) & \cdots & p(x_n | x_n) \end{bmatrix} \in \mathbb{R}^{n \times n}. \quad (6.8)$$

Note that each row of \mathbf{P} must always sum to 1. Such matrices are also called *stochastic*.

The *transition graph* of a Markov chain is a directed graph consisting of vertices S and weighted edges represented by the adjacency matrix \mathbf{P} .

The current state of the Markov chain at time t is denoted by the probability distribution q_t over states S , that is, $X_t \sim q_t$. In the finite setting, q_t is a PMF, which is often written explicitly as the row vector $q_t \in \mathbb{R}^{1 \times |S|}$. The initial state (or prior) of the Markov chain is given as $X_0 \sim q_0$. One iteration of the Markov chain can then be expressed as follows: $\textcircled{?}$

Problem 6.1

$$q_{t+1} = q_t \mathbf{P}. \quad (6.9)$$

It is implied directly that we can write the state of the Markov chain at time $t + k$ as

$$q_{t+k} = q_t \mathbf{P}^k. \quad (6.10)$$

The entry $\mathbf{P}^k(x, x')$ corresponds to the probability of transitioning from state $x \in S$ to state $x' \in S$ in exactly k steps $\textcircled{?}$. We denote this entry by $p^{(k)}(x' | x)$.

Problem 6.2

In the analysis of Markov chains, there are two main concepts of interest: stationarity and convergence. We begin by introducing stationarity.

6.1.1 Stationarity

Definition 6.4 (Stationary distribution). A distribution π is *stationary* with respect to the transition function p iff

$$\pi(x) = \sum_{x' \in S} p(x | x') \pi(x') \quad (6.11)$$

holds for all $x \in S$. It follows from Equation (6.9) that equivalently, π is stationary w.r.t. a transition matrix \mathbf{P} iff

$$\pi = \pi \mathbf{P}. \quad (6.12)$$

After entering a stationary distribution π , a Markov chain will always remain in the stationary distribution. In particular, suppose that X_t is distributed according to π , then for all $k \geq 0$, $X_{t+k} \sim \pi$.

Remark 6.5: When does a stationary distribution exist?

In general, there are Markov chains with infinitely many stationary distributions or no stationary distribution at all. You can find some examples in Figure 6.2.

It can be shown that there exists a unique stationary distribution π if the Markov chain is *irreducible*, that is, if every state is reachable from every other state with a positive probability when the Markov chain is run for enough steps. Formally,

$$\forall x, x' \in S. \exists k \in \mathbb{N}. p^{(k)}(x' | x) > 0. \quad (6.13)$$

Equivalently, a Markov chain is irreducible iff its transition graph is strongly connected.

6.1.2 Convergence

Let us now consider Markov chains with a unique stationary distribution.⁴ A natural next question is whether this Markov chain converges to its stationary distribution. We say that a Markov chain converges to its stationary distribution iff we have

$$\lim_{t \rightarrow \infty} q_t = \pi, \quad (6.14)$$

irrespectively of the initial distribution q_0 .

⁴Observe that the stationary distribution of an irreducible Markov chain must have full support, that is, assign positive probability to every state.

Remark 6.6: When does a Markov chain converge?

Even if a Markov chain has a unique stationary distribution, it does not have to converge to it. Consider example (3) in Figure 6.2. Clearly, $\pi = (\frac{1}{2}, \frac{1}{2})$ is the unique stationary distribution. However, observe that if we start with a suitable initial distribution such as $q_0 = (1, 0)$, at no point in time will the probability of all states be positive, and in particular, the chain will not converge to π . Instead, the chain behaves periodically, i.e., its state distributions are $q_{2t} = (0, 1)$ and $q_{2t+1} = (1, 0)$ for all $t \in \mathbb{N}$. It turns out that if we exclude such “periodic” Markov chains, then the remaining (irreducible) Markov chains will always converge to their stationary distribution.

Formally, a Markov chain is *aperiodic* if for all states $x \in S$,

$$\exists k_0 \in \mathbb{N}. \forall k \geq k_0. p^{(k)}(x | x) > 0. \quad (6.15)$$

In words, a Markov chain is aperiodic iff for every state x , the transition graph has a closed path from x to x with length k for all $k \in \mathbb{N}$ greater than some $k_0 \in \mathbb{N}$.

This additional property leads to the concept of ergodicity.

Definition 6.7 (Ergodicity). A Markov chain is *ergodic* iff there exists a $t \in \mathbb{N}_0$ such that for any $x, x' \in S$ we have

$$p^{(t)}(x' | x) > 0, \quad (6.16)$$

whereby $p^{(t)}(x' | x)$ is the probability to reach x' from x in exactly t steps. Equivalent conditions are

1. that there exists some $t \in \mathbb{N}_0$ such that all entries of P^t are strictly positive; and
2. that it is irreducible and aperiodic.

Example 6.8: Making a Markov chain ergodic

A commonly used strategy to ensure that a Markov chain is ergodic is to add “self-loops” to every vertex in the transition graph. That is, to ensure that at any point in time, the Markov chain remains with positive probability in its current state.

Take a (not necessarily ergodic) but irreducible Markov chain with transition matrix P . We define the new Markov chain

$$P' = \frac{1}{2}P + \frac{1}{2}I. \quad (6.17)$$

It is a simple exercise to confirm that P' is stochastic, and hence a valid transition matrix. Also, it follows directly that P' is irreducible (as P is irreducible) and aperiodic as every vertex has a closed path of length 1 to itself, and therefore the chain is ergodic.

Take now π to be a stationary distribution of P . We have that π is also a stationary distribution of P' as

$$\pi P' = \frac{1}{2}\pi P + \frac{1}{2}\pi I = \frac{1}{2}\pi + \frac{1}{2}\pi = \pi. \quad (6.18)$$

Fact 6.9 (Fundamental theorem of ergodic Markov chains, theorem 4.9 of Levin and Peres (2017)). *An ergodic Markov chain has a unique stationary distribution π (with full support) and*

$$\lim_{t \rightarrow \infty} q_t = \pi \quad (6.19)$$

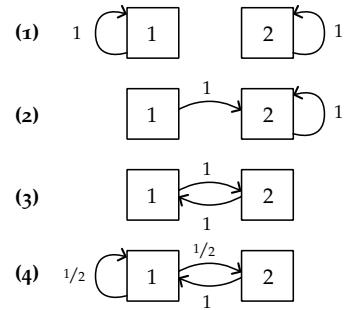


Figure 6.2: Transition graphs of Markov chains: (1) is not ergodic as its transition diagram is not strongly connected; (2) is not ergodic for the same reason; (3) is irreducible but periodic and therefore not ergodic; (4) is ergodic with stationary distribution $\pi(1) = 2/3, \pi(2) = 1/3$.

using (6.12)

irrespectively of the initial distribution q_0 .

This naturally suggests constructing an ergodic Markov chain such that its stationary distribution coincides with the posterior distribution. If we then sample “sufficiently long”, X_t is drawn from a distribution that is “very close” to the posterior distribution.

Remark 6.10: How quickly does a Markov chain converge?

The convergence speed of Markov chains is a rich field of research. “Sufficiently long” and “very close” are commonly made precise by the notions of *rapidly mixing* Markov chains and *total variation distance*.

Definition 6.11 (Total variation distance). The total variation distance between two probability distributions μ and ν on \mathcal{A} is defined by

$$\|\mu - \nu\|_{\text{TV}} \doteq 2 \sup_{A \subseteq \mathcal{A}} |\mu(A) - \nu(A)|. \quad (6.20)$$

It defines the distance between μ and ν to be the maximum difference between the probabilities that μ and ν assign to the same event.

As opposed to the KL-divergence (5.34), the total variation distance is a metric. In particular, it is symmetric and satisfies the triangle inequality. It can be shown that

$$\|\mu - \nu\|_{\text{TV}} \leq \sqrt{2\text{KL}(\mu\|\nu)} \quad (6.21)$$

which is known as *Pinsker’s inequality*. Moreover, if μ and ν are discrete distributions over the set S , it can be shown that

$$\|\mu - \nu\|_{\text{TV}} = \sum_{i \in S} |\mu(i) - \nu(i)|. \quad (6.22)$$

Definition 6.12 (Mixing time). For a Markov chain with stationary distribution π , its *mixing time* with respect to the total variation distance for any $\epsilon > 0$ is

$$\tau_{\text{TV}}(\epsilon) \doteq \min\{t \mid \forall q_0 : \|q_t - \pi\|_{\text{TV}} \leq \epsilon\}. \quad (6.23)$$

Thus, the mixing time measures the time required by a Markov chain for the distance to stationarity to be small. A Markov chain is typically said to be *rapidly mixing* if for any $\epsilon > 0$,

$$\tau_{\text{TV}}(\epsilon) \in O(\text{poly}(n, \log(1/\epsilon))). \quad (6.24)$$

That is, a rapidly mixing Markov chain on n states needs to be simulated for at most $\text{poly}(n)$ steps to obtain a “good” sample from its stationary distribution π .

You can find a thorough introduction to mixing times in chapter 4 of “Markov chains and mixing times” (Levin and Peres, 2017). Later chapters introduce methods for showing that a Markov chain is rapidly mixing.

6.1.3 Detailed Balance Equation

How can we confirm that the stationary distribution of a Markov chain coincides with the posterior distribution? The detailed balance equation yields a very simple method.

Definition 6.13 (Detailed balance equation / reversibility). A Markov chain satisfies the *detailed balance equation* with respect to a distribution π iff

$$\pi(x)p(x' | x) = \pi(x')p(x | x') \quad (6.25)$$

holds for any $x, x' \in S$. A Markov chain that satisfies the detailed balance equation with respect to π is called *reversible* with respect to π .

Lemma 6.14. *Given a finite Markov chain, if the Markov chain is reversible with respect to π then π is a stationary distribution.*⁵

Proof. Let $\pi \doteq q_t$. We have,

$$\begin{aligned} q_{t+1}(x) &= \sum_{x' \in S} p(x | x')q_t(x') \\ &= \sum_{x' \in S} p(x | x')\pi(x') \\ &= \sum_{x' \in S} p(x' | x)\pi(x) \\ &= \pi(x) \sum_{x' \in S} p(x' | x) \\ &= \pi(x). \end{aligned}$$

⁵ Note that reversibility of π is only a sufficient condition for stationarity of π , it is not necessary! In particular, there are irreversible ergodic Markov chains.

using the Markov property (6.6)

using the detailed balance equation (6.25)

using that $\sum_{x' \in S} p(x' | x) = 1$

That is, if we can show that the detailed balance equation (6.25) holds for some distribution q , then we know that q is the stationary distribution of the Markov chain.

Next, reconsider our posterior distribution $p(x) = \frac{1}{Z}q(x)$ from Equation (6.4). If we substitute the posterior for π in the detailed balance equation, we obtain

$$\frac{1}{Z}q(x)p(x' | x) = \frac{1}{Z}q(x')p(x | x'), \quad (6.26)$$

or equivalently,

$$q(x)p(x' | x) = q(x')p(x | x'). \quad (6.27)$$

In words, we do not need to know the true posterior p to check that the stationary distribution of our Markov chain coincides with p , it suffices to know the finite measure q !

6.1.4 Ergodic Theorem

If we now suppose that we can construct a Markov chain whose stationary distribution coincides with the posterior distribution — we will see later that this is possible — it is not apparent that this allows us to estimate expectations over the posterior distribution. Note that although constructing such a Markov chain allows us to obtain samples from the posterior distribution, they are *not* independent. In fact, due to the structure of a Markov chain, by design, they are strongly dependent. Thus, the law of large numbers and Hoeffding's inequality do not apply. By itself, it is not even clear that an estimator relying on samples from a single Markov chain will be unbiased.

Theoretically, we could simulate many Markov chains separately and obtain one sample from each of them. This, however, is extremely inefficient. It turns out that there is a way to generalize the (strong) law of large numbers to Markov chains.

Theorem 6.15 (Ergodic theorem, appendix C of Levin and Peres (2017)). *Given an ergodic Markov chain $(X_t)_{t \in \mathbb{N}_0}$ over a finite state space S with stationary distribution π and a function $f : S \rightarrow \mathbb{R}$,*

$$\frac{1}{n} \sum_{i=1}^n f(x_i) \xrightarrow{a.s.} \sum_{x \in S} \pi(x) f(x) = \mathbb{E}_{x \sim \pi}[f(x)] \quad (6.28)$$

as $n \rightarrow \infty$ where $x_i \sim X_i | x_{i-1}$.

This result is the fundamental reason for why Markov chain Monte Carlo methods are possible. There are analogous results for continuous domains.

Note, however, that the ergodic theorem only tells us that simulating a single Markov chain yields an unbiased estimator. It does not tell us anything about the rate of convergence and variance of such an estimator. The convergence rate depends on the mixing time of the Markov chain, which is difficult to establish in general.

In practice, one observes that Markov chain Monte Carlo methods have a so-called “burn-in” time during which the distribution of the Markov

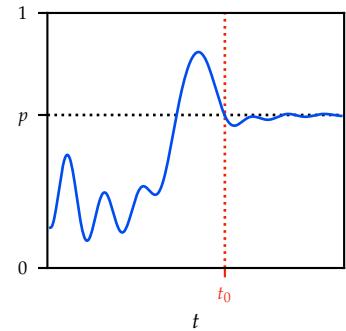


Figure 6.3: Illustration of the “burn-in” time t_0 of a Markov chain approximating the posterior $p(y^* = 1 | X, y)$ of Bayesian logistic regression. The true posterior p is shown in gray. The distribution of the Markov chain at time t is shown in red.

chain does not yet approximate the posterior distribution well. Typically, the first t_0 samples are therefore discarded,

$$\mathbb{E}[f(X)] \approx \frac{1}{T-t_0} \sum_{t=t_0+1}^T f(X_t). \quad (6.29)$$

It is not clear in general how T and t_0 should be chosen such that the estimator is unbiased, rather they have to be tuned.

Another widely used heuristic is to first find the mode of the posterior distribution and then start the Markov chain at that point. This tends to increase the rate of convergence drastically, as the Markov chain does not have to “walk to the location in the state space where most probability mass will be located”.

6.2 Elementary Sampling Methods

We will now examine methods for constructing and sampling from a Markov chain with the goal of approximating samples from the posterior distribution p . Note that in this setting the state space of the Markov chain is \mathbb{R}^n and a single state at time t is described by the random vector $\mathbf{X} \doteq [X_1, \dots, X_n]$.

6.2.1 Metropolis-Hastings Algorithm

Suppose we are given a *proposal distribution* $r(x' | x)$ which, given we are in state x , proposes a new state x' . Metropolis and Hastings showed that using the *acceptance distribution* $\text{Bern}(\alpha(x' | x))$ where

$$\alpha(x' | x) \doteq \min\left\{1, \frac{p(x')r(x | x')}{p(x)r(x' | x)}\right\} \quad (6.30)$$

$$= \min\left\{1, \frac{q(x')r(x | x')}{q(x)r(x' | x)}\right\} \quad (6.31)$$

to decide whether to follow the proposal yields a Markov chain with stationary distribution $p(x) = \frac{1}{Z}q(x)$.

similarly to the detailed balance equation, the normalizing constant Z cancels

Algorithm 6.16: Metropolis-Hastings algorithm

```

1 initialize  $x \in \mathbb{R}^n$ 
2 for  $t = 1$  to  $T$  do
3   sample  $x' \sim r(x' | x)$ 
4   sample  $u \sim \text{Unif}(0, 1)$ 
5   if  $u \leq \alpha(x' | x)$  then update  $x \leftarrow x'$ 
6   else update  $x \leftarrow x$ 
```

Intuitively, the acceptance distribution corrects for the bias in the proposal distribution. That is, if the proposal distribution r is likely to propose states with low probability under p , the acceptance distribution will reject these proposals frequently. The following theorem formalizes this intuition.

Theorem 6.17 (Metropolis-Hastings theorem). *Given an arbitrary proposal distribution r whose support includes the support of q , the stationary distribution of the Markov chain simulated by the Metropolis-Hastings algorithm is $p(x) = \frac{1}{Z}q(x)$.*

Proof. First, let us define the transition probabilities of the Markov chain. The probability of transitioning from a state x to a state x' is given by $r(x' | x)\alpha(x' | x)$ if $x \neq x'$ and the probability of proposing to remain in state x , $r(x | x)$, plus the probability of denying the proposal, otherwise.

$$p(x' | x) = \begin{cases} r(x' | x)\alpha(x' | x) & \text{if } x \neq x' \\ r(x | x) + \sum_{x'' \neq x} r(x'' | x)(1 - \alpha(x'' | x)) & \text{otherwise.} \end{cases} \quad (6.32)$$

We will show that the stationary distribution is p by showing that p satisfies the detailed balance equation (6.25). Let us fix arbitrary states x and x' . First, observe that if $x = x'$, then the detailed balance equation is trivially satisfied. Without loss of generality we assume

$$\alpha(x | x') = 1, \quad \alpha(x' | x) = \frac{q(x')r(x | x')}{q(x)r(x' | x)}.$$

For $x \neq x'$, we then have,

$$\begin{aligned} p(x) \cdot p(x' | x) &= \frac{1}{Z}q(x)p(x' | x) \\ &= \frac{1}{Z}q(x)r(x' | x)\alpha(x' | x) \\ &= \frac{1}{Z}q(x)r(x' | x) \frac{q(x')r(x | x')}{q(x)r(x' | x)} \\ &= \frac{1}{Z}q(x')r(x | x') \\ &= \frac{1}{Z}q(x')r(x | x')\alpha(x | x') \\ &= \frac{1}{Z}q(x')p(x | x') \\ &= p(x') \cdot p(x | x'). \end{aligned} \quad \square$$

using the definition of the distribution p

using the transition probabilities of the Markov chain

using the definition of the acceptance distribution α

using the definition of the acceptance distribution α

using the transition probabilities of the Markov chain

using the definition of the distribution p

Note that by the fundamental theorem of ergodic Markov chains (6.19), for convergence to the stationary distribution, it is sufficient for the

Markov chain to be ergodic. Ergodicity follows immediately when the transition probabilities $p(\cdot | \mathbf{x})$ have full support. For example, if the proposal distribution $r(\cdot | \mathbf{x})$ has full support, the full support of $p(\cdot | \mathbf{x})$ follows immediately from Equation (6.32). The rate of convergence of Metropolis-Hastings depends strongly on the choice of the proposal distribution, and we will explore different choices of proposal distribution in the following.

6.2.2 Gibbs Sampling

A popular example of a Metropolis-Hastings algorithm is Gibbs sampling as presented in Algorithm 6.18.

Algorithm 6.18: Gibbs sampling

```

1 initialize  $\mathbf{x} = [x_1, \dots, x_n] \in \mathbb{R}^n$ 
2 for  $t = 1$  to  $T$  do
3   pick a variable  $i$  uniformly at random from  $\{1, \dots, n\}$ 
4   set  $\mathbf{x}_{-i} \doteq [x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n]$ 
5   update  $x_i$  by sampling according to the posterior distribution
      $p(x_i | \mathbf{x}_{-i})$ 
```

Intuitively, by re-sampling single coordinates according to the posterior distribution given the other coordinates, Gibbs sampling finds states that are successively “more” likely. Selecting the index i uniformly at random ensures that the underlying Markov chain is ergodic provided the conditional distributions $p(\cdot | \mathbf{x}_{-i})$ have full support.

Theorem 6.19 (Gibbs sampling as Metropolis-Hastings). *Gibbs sampling is a Metropolis-Hastings algorithm. For any fixed $i \in [n]$, it has proposal distribution*

$$r_i(\mathbf{x}' | \mathbf{x}) \doteq \begin{cases} p(x'_i | \mathbf{x}'_{-i}) & \mathbf{x}' \text{ differs from } \mathbf{x} \text{ only in entry } i \\ 0 & \text{otherwise} \end{cases} \quad (6.33)$$

and acceptance distribution $\alpha_i(\mathbf{x}' | \mathbf{x}) \doteq 1$.

Proof. We show that $\alpha_i(\mathbf{x}' | \mathbf{x}) = 1$ follows from the definition of an acceptance distribution in Metropolis-Hastings (6.31) and the choice of proposal distribution (6.33).

By (6.31),

$$\alpha_i(\mathbf{x}' | \mathbf{x}) = \min \left\{ 1, \frac{p(\mathbf{x}') r_i(\mathbf{x} | \mathbf{x}')}{p(\mathbf{x}) r_i(\mathbf{x}' | \mathbf{x})} \right\}$$

Note that $p(\mathbf{x}) = p(x_i, \mathbf{x}_{-i}) = p(x_i | \mathbf{x}_{-i})p(\mathbf{x}_{-i})$ using the product rule (1.11). Therefore,

$$\begin{aligned}
&= \min \left\{ 1, \frac{p(x'_i | \mathbf{x}'_{-i})p(\mathbf{x}'_{-i})r_i(\mathbf{x} | \mathbf{x}')}{p(x_i | \mathbf{x}_{-i})p(\mathbf{x}_{-i})r_i(\mathbf{x}' | \mathbf{x})} \right\} \\
&= \min \left\{ 1, \frac{p(x'_i | \mathbf{x}'_{-i})p(\mathbf{x}'_{-i})p(x_i | \mathbf{x}_{-i})}{p(x_i | \mathbf{x}_{-i})p(\mathbf{x}_{-i})p(x'_i | \mathbf{x}'_{-i})} \right\} \\
&= \min \left\{ 1, \frac{p(\mathbf{x}'_{-i})}{p(\mathbf{x}_{-i})} \right\} \\
&= 1.
\end{aligned}$$

using the proposal distribution (6.33)

□ using that $p(\mathbf{x}'_{-i}) = p(\mathbf{x}_{-i})$

If the index i is chosen uniformly at random as in Algorithm 6.18, then the proposal distribution is $r(\mathbf{x}' | \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n r_i(\mathbf{x}' | \mathbf{x})$, which analogously to Theorem 6.19 has the associated acceptance distribution

$$\alpha(\mathbf{x}' | \mathbf{x}) = \min \left\{ 1, \frac{p(\mathbf{x}')r(\mathbf{x} | \mathbf{x}')}{p(\mathbf{x})r(\mathbf{x}' | \mathbf{x})} \right\} = \min \left\{ 1, \frac{\sum_{i=1}^n p(\mathbf{x}')r_i(\mathbf{x} | \mathbf{x}')}{\sum_{i=1}^n p(\mathbf{x})r_i(\mathbf{x}' | \mathbf{x})} \right\} = 1.$$

Corollary 6.20 (Convergence of Gibbs sampling). *As Gibbs sampling is a specific example of an MH-algorithm, the stationary distribution of the simulated Markov chain is $p(\mathbf{x})$.*

Note that for the proposals of Gibbs sampling, we have

$$p(x_i | \mathbf{x}_{-i}) = \frac{p(x_i, \mathbf{x}_{-i})}{\sum_{x_i} p(x_i, \mathbf{x}_{-i})} = \frac{q(x_i, \mathbf{x}_{-i})}{\sum_{x_i} q(x_i, \mathbf{x}_{-i})}. \quad (6.34)$$

using the definition of condition probability (1.8) and the sum rule (1.7)

Under many models, this probability can be efficiently evaluated due to the conditioning on the remaining coordinates \mathbf{x}_{-i} . If X_i has finite support, the normalizer can be computed exactly.

6.3 Sampling using Gradients

In this section, we discuss more advanced sampling methods. The main idea that we will study is the interpretation of sampling as an optimization problem. We will build towards an optimization view of sampling step-by-step, and first introduce what is commonly called the “energy” of a distribution.

6.3.1 Gibbs Distributions

Gibbs distributions are a special class of distributions that are widely used in machine learning, and which are characterized by an energy.

Definition 6.21 (Gibbs distribution). Formally, a *Gibbs distribution* (also called a *Boltzmann distribution*) is a continuous distribution p whose PDF is of the form

$$p(\mathbf{x}) = \frac{1}{Z} \exp(-f(\mathbf{x})). \quad (6.35)$$

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ is also called an *energy function*. When the energy function f is convex, its Gibbs distribution is called *log-concave*.

A useful property is that Gibbs distributions always have full support.⁶ It is often easier to reason about “energies” rather than probabilities as they are neither restricted to be non-negative nor do they have to integrate to 1. Note that the Gibbs distribution belongs to the exponential family (5.48) with sufficient statistic $-f(\mathbf{x})$.

⁶ This can easily be seen as $\exp(\cdot) > 0$.

Observe that the posterior distribution can always be interpreted as a Gibbs distribution as long as prior and likelihood have full support,

$$\begin{aligned} p(\boldsymbol{\theta} | \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) &= \frac{1}{Z} p(\boldsymbol{\theta}) p(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}, \boldsymbol{\theta}) \\ &= \frac{1}{Z} \exp(-[-\log p(\boldsymbol{\theta}) - \log p(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}, \boldsymbol{\theta})]). \end{aligned} \quad (6.36)$$

using Bayes' rule (1.45)

Thus, defining the energy function

$$f(\boldsymbol{\theta}) \doteq -\log p(\boldsymbol{\theta}) - \log p(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}, \boldsymbol{\theta}) \quad (6.37)$$

$$= -\log p(\boldsymbol{\theta}) - \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}), \quad (6.38)$$

yields

$$p(\boldsymbol{\theta} | \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \frac{1}{Z} \exp(-f(\boldsymbol{\theta})). \quad (6.39)$$

Note that f coincides with the loss function used for MAP estimation (1.62). For a noninformative prior, the regularization term vanishes and the energy reduces to the negative log-likelihood $\ell_{\text{nll}}(\boldsymbol{\theta}; \mathcal{D})$ (i.e., the loss function of maximum likelihood estimation (1.57)).

Using that the posterior is a Gibbs distribution, we can rewrite the acceptance distribution of Metropolis-Hastings,

$$\alpha(\mathbf{x}' | \mathbf{x}) = \min \left\{ 1, \frac{r(\mathbf{x} | \mathbf{x}')}{r(\mathbf{x}' | \mathbf{x})} \exp(f(\mathbf{x}) - f(\mathbf{x}')) \right\}. \quad (6.40)$$

this is obtained by substituting the PDF of a Gibbs distribution for the posterior

Example 6.22: Metropolis-Hastings with Gaussian proposals

Let us consider Metropolis-Hastings with the Gaussian proposal distribution,

$$r(\mathbf{x}' | \mathbf{x}) \doteq \mathcal{N}(\mathbf{x}'; \mathbf{x}, \tau \mathbf{I}). \quad (6.41)$$

Due to the symmetry of Gaussians, we have

$$\frac{r(\mathbf{x} | \mathbf{x}')}{r(\mathbf{x}' | \mathbf{x})} = \frac{\mathcal{N}(\mathbf{x}; \mathbf{x}', \tau \mathbf{I})}{\mathcal{N}(\mathbf{x}'; \mathbf{x}, \tau \mathbf{I})} = 1.$$

Hence, the acceptance distribution is defined by

$$\alpha(x' | x) = \min\{1, \exp(f(x) - f(x'))\}. \quad (6.42)$$

Intuitively, when a state with lower energy is proposed, that is $f(x') \leq f(x)$, then the proposal will always be accepted. In contrast, if the energy of the proposed state is higher, the acceptance probability decreases exponentially in the difference of energies $f(x) - f(x')$. Thus, Metropolis-Hastings minimizes the energy function, which corresponds to minimizing the negative log-likelihood and negative log-prior. The variance in the proposals τ helps in getting around local optima, but the search direction is uniformly random (i.e., “uninformed”).

6.3.2 From Energy to Surprise (and back)

Energy-based models are a well-known class of models in machine learning where an energy function f is learned from data. These energy functions do not need to originate from a probability distribution, yet they induce a probability distribution via their Gibbs distribution $p(x) \propto \exp(-f(x))$. As we will see in Problem 6.7, this Gibbs distribution is the associated maximum entropy distribution. Observe that the surprise about x under this distribution is given by

$$S[p(x)] = f(x) + \log Z. \quad (6.43)$$

That is, up to a constant shift, the energy of x coincides with the surprise about x . Energies are therefore sufficient for comparing the “likelihood” of points, and they do not require normalization.⁷

What kind of energies could we use? In Section 6.3.1, we discussed the use of the negative log-posterior or negative log-likelihood as energies. In general, *any* loss function $\ell(x)$ can be thought of as an energy function with an associated maximum entropy distribution $p(x) \propto \exp(-\ell(x))$.

6.3.3 Langevin Dynamics

Until now, we have looked at Metropolis-Hastings algorithms with proposal distributions that do not explicitly take into account the curvature of the energy function around the current state. Langevin dynamics adapts the Gaussian proposals of the Metropolis-Hastings algorithm we have seen in Example 6.22 to search the state space in an “informed” direction. The simple idea is to bias the sampling towards states with lower energy, thereby making it more likely that a proposal is accepted.

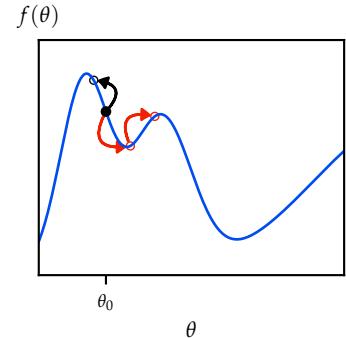


Figure 6.4: Metropolis-Hastings and Langevin dynamics minimize the energy function $f(\theta)$ shown in blue. Suppose we start at the black dot θ_0 , then the black and red arrows denote possible subsequent samples. Metropolis-Hastings uses an “uninformed” search direction, whereas Langevin dynamics uses the gradient of $f(\theta)$ to make “more promising” proposals. The random proposals help get past local optima.

⁷ Intuitively, an energy can be used to compare the “likelihood” of two points x and x' whereas the probability x makes a statement about the “likelihood” of x relative to all other points.

A natural idea is to shift the proposal distribution perpendicularly to the gradient of the energy function. This yields the following proposal distribution,

$$r(\mathbf{x}' \mid \mathbf{x}) = \mathcal{N}(\mathbf{x}'; \mathbf{x} - \eta_t \nabla f(\mathbf{x}), 2\eta_t \mathbf{I}). \quad (6.44)$$

The resulting variant of Metropolis-Hastings is known as the *Metropolis adjusted Langevin algorithm* (MALA) or *Langevin Monte Carlo* (LMC). It can be shown that, as $\eta_t \rightarrow 0$, we have for the acceptance probability $\alpha(\mathbf{x}' \mid \mathbf{x}) \rightarrow 1$ using that the acceptance probability is 1 if $\mathbf{x}' = \mathbf{x}$. Hence, the Metropolis-Hastings acceptance step can be omitted once the rejection probability becomes negligible. The algorithm which always accepts the proposal of Equation (6.44) is known as the *unadjusted Langevin algorithm* (ULA).

Observe that if the stationary distribution is the posterior distribution

$$p(\boldsymbol{\theta} \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \frac{1}{Z} \exp \underbrace{\left(\log p(\boldsymbol{\theta}) + \sum_{i=1}^n \log p(y_i \mid \mathbf{x}_i, \boldsymbol{\theta}) \right)}_{-f(\boldsymbol{\theta})}$$

with energy f as we discussed in Section 6.3.1, then the proposal $\boldsymbol{\theta}'$ of MALA/LMC can be equivalently formulated as

$$\begin{aligned} \boldsymbol{\theta}' &\leftarrow \boldsymbol{\theta} - \eta_t \nabla f(\boldsymbol{\theta}) + \boldsymbol{\varepsilon} \\ &= \boldsymbol{\theta} + \eta_t \left(\nabla \log p(\boldsymbol{\theta}) + \sum_{i=1}^n \nabla \log p(y_i \mid \mathbf{x}_i, \boldsymbol{\theta}) \right) + \boldsymbol{\varepsilon} \end{aligned} \quad (6.45)$$

where $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, 2\eta_t \mathbf{I})$.

Remark 6.23: ULA is a discretized diffusion

The unadjusted Langevin algorithm can be seen as a discretization of Langevin dynamics, which is a continuous-time stochastic process with a drift and with random stationary and independent Gaussian increments. The randomness is modeled by a Wiener process.

Definition 6.24 (Wiener process). The *Wiener process* (also known as *Brownian motion*) is a sequence of random vectors $\{\mathbf{W}_t\}_{t \geq 0}$ such that

1. $\mathbf{W}_0 = \mathbf{0}$,
2. with probability 1, \mathbf{W}_t is continuous in t ,
3. the process has independent increments,⁸ and
4. $\mathbf{W}_{t+u} - \mathbf{W}_t \sim \mathcal{N}(\mathbf{0}, u\mathbf{I})$.

Consider the continuous-time stochastic process $\boldsymbol{\theta}$ defined by the

⁸ That is, the “future” increments $\mathbf{W}_{t+u} - \mathbf{W}_t$ for $u \geq 0$ are independent of past values \mathbf{W}_s for $s < t$.

stochastic differential equation (SDE)

$$d\theta_t = \underbrace{-\nabla f(\theta_t) dt}_{\text{drift}} + \underbrace{\sqrt{2} d\mathbf{W}_t}_{\text{noise}} \quad (6.46)$$

Such a stochastic process is also called a *diffusion (process)* and Equation (6.46) specifically is called *Langevin dynamics*. Here, the first term is called the “drift” of the process, and the second term is called its “noise”. Note that if the noise term is zero then Equation (6.46) is simply an ordinary differential equation (ODE).

A diffusion can be discretized using the Euler-Maruyama method (also called “forward Euler”) to obtain a discrete approximation $\theta_k \approx \theta(\tau_k)$ where τ_k denotes the k -th time step. Choosing the time steps such that $\Delta t_k \doteq \tau_{k+1} - \tau_k = \eta_k$ yields the approximation

$$\theta_{k+1} = \theta_k - \nabla f(\theta_k) \Delta t_k + \sqrt{2} \Delta \mathbf{W}_k \quad (6.47)$$

where $\Delta \mathbf{W}_k \doteq \mathbf{W}_{\tau_{k+1}} - \mathbf{W}_{\tau_k} \sim \mathcal{N}(\mathbf{0}, \Delta t_k \mathbf{I})$. Observe that this coincides with the update rule of Langevin dynamics from Equation (6.45).

The appearance of drift and noise is closely related to the *principle of curiosity and conformity*, we encountered in the previous chapter on variational inference. The noise term (also called the *diffusion*) leads to exploration of the state space (i.e., curiosity about alternative explanations of the data), whereas the drift term (also called the *distillation*) leads to minimizing the energy or loss (i.e., conformity to the data). Interestingly, this same principle appears in both variational inference and MCMC, two very different approaches to approximate probabilistic inference. In the remainder of this manuscript we will find this to be a reoccurring theme.

For log-concave distributions, the mixing time of the Markov chain underlying Langevin dynamics can be shown to be polynomial in the dimension n (Vempala and Wibisono, 2019). You will prove this result for strongly log-concave distributions in (?) and see that the analysis is analogous to the canonical convergence analysis of classical optimization schemes.

Problem 6.9

6.3.4 Stochastic Gradient Langevin Dynamics

Note that computing the gradient of the energy function, which corresponds to computing exact gradients of the log-prior and log-likelihood, in every step can be expensive. The proposal step of MALA/LMC can be made more efficient by approximating the gradient with an unbi-

ased gradient estimate, leading to *stochastic gradient Langevin dynamics* (SGLD) shown in Algorithm 6.25 (Welling and Teh, 2011). Observe that SGLD (6.48) differs from MALA/LMC (6.45) only by using a sample-based approximation of the gradient.

Algorithm 6.25: Stochastic gradient Langevin dynamics, SGLD

```

1 initialize  $\theta$ 
2 for  $t = 1$  to  $T$  do
3   sample  $i_1, \dots, i_m \sim \text{Unif}(\{1, \dots, n\})$  independently
4   sample  $\varepsilon \sim \mathcal{N}(\mathbf{0}, 2\eta_t \mathbf{I})$ 
5    $\theta \leftarrow \theta + \eta_t \left( \nabla \log p(\theta) + \frac{n}{m} \sum_{j=1}^m \nabla \log p(y_{i_j} | x_{i_j}, \theta) \right) + \varepsilon$  // (6.48)

```

Intuitively, in the initial phase of the algorithm, the stochastic gradient term dominates, and therefore, SGLD corresponds to a variant of stochastic gradient ascent. In the later phase, the update rule is dominated by the injected noise ε , and will effectively be Langevin dynamics. SGLD transitions smoothly between the two phases.

Under additional assumptions, SGLD is guaranteed to converge to the posterior distribution for decreasing learning rates $\eta_t = \Theta(t^{-1/3})$ (Raginsky et al., 2017; Xu et al., 2018). SGLD does not use the acceptance step from Metropolis-Hastings as asymptotically, SGLD corresponds to Langevin dynamics and the Metropolis-Hastings rejection probability goes to zero for a decreasing learning rate.

6.3.5 Hamiltonian Monte Carlo

As MALA and SGLD can be seen as a sampling-based analogue of GD and SGD, a similar analogue for (stochastic) gradient descent with momentum is the (stochastic gradient) *Hamiltonian Monte Carlo* (HMC) algorithm, which we discuss in the following (Duane et al., 1987; Chen et al., 2014).

We have seen that if we want to sample from a distribution

$$p(x) \propto \exp(-f(x))$$

with energy function f , we can construct a Markov chain whose distribution converges to p . We have also seen that for this approach to work, the chain must move through all areas of significant probability with reasonable speed.

If one is faced with a distribution p which is multimodal (i.e., that has several “peaks”), one has to ensure that the chain will explore all modes, and can therefore “jump between different areas of the space”.

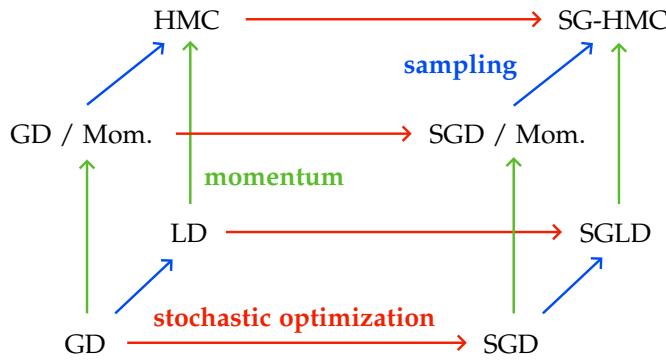


Figure 6.5: A commutative diagram of sampling and optimization algorithms. Langevin dynamics (LD) is the non-stochastic variant of SGLD.

So in general, *local* updates are doomed to fail. Methods such as Metropolis-Hastings with Gaussian proposals, or even Langevin Monte Carlo might face this issue, as they do not jump to distant areas of the state space with significant acceptance probability. It will therefore take a long time to move from one peak to another.

The HMC algorithm is an instance of Metropolis-Hastings which uses momentum to propose distant points that conserve energy, with high acceptance probability. The general idea of HMC is to *lift* samples x to a higher-order space by considering an auxiliary variable y with the same dimension as x . We also lift the distribution p to a distribution on the (x, y) -space by defining a distribution $p(y | x)$ and setting $p(x, y) \doteq p(y | x)p(x)$. It is common to pick $p(y | x)$ to be a Gaussian with zero mean and variance mI . Hence,

$$p(x, y) \propto \exp\left(-\frac{1}{2m} \|y\|_2^2 - f(x)\right). \quad (6.49)$$

Physicists might recognize the above as the canonical distribution of a Newtonian system if one takes x as the position and y as the momentum. $H(x, y) \doteq \frac{1}{2m} \|y\|_2^2 + f(x)$ is called the *Hamiltonian*. HMC then takes a step in this higher-order space according to the Hamiltonian dynamics,⁹

$$\frac{dx}{dt} = \nabla_y H, \quad \frac{dy}{dt} = -\nabla_x H, \quad (6.50)$$

reaching some new point (x', y') and *projecting* back to the state space by selecting x' as the new sample. This is illustrated in Figure 6.6. In the next iteration, we resample the momentum $y' \sim p(\cdot | x')$ and repeat the procedure.

In an implementation of this algorithm, one has to solve Equation (6.50) numerically rather than exactly. Typically, this is done using the *Leapfrog*

⁹ That is, HMC follows the trajectory of these dynamics for some time.

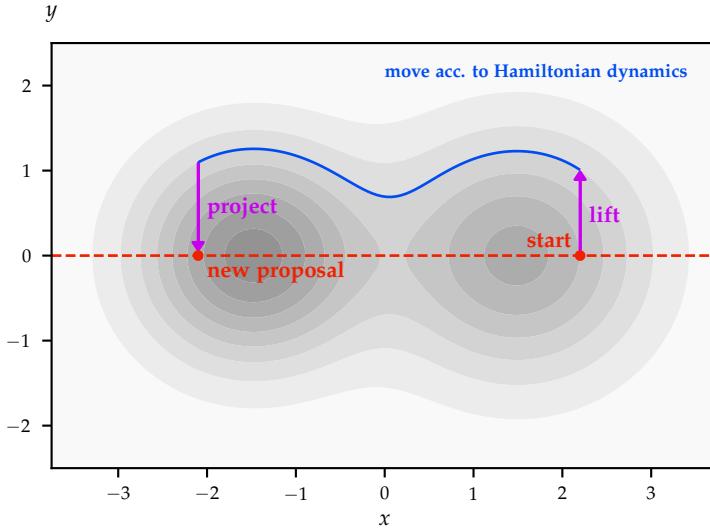


Figure 6.6: Illustration of Hamiltonian Monte Carlo. Shown is the contour plot of a distribution p , which is a mixture of two Gaussians, in the (x, y) -space.

First, the initial point in the state space is lifted to the (x, y) -space. Then, we move according to Hamiltonian dynamics and finally project back onto the state space.

method, which for a step size τ computes

$$\mathbf{y}(t + \tau/2) = \mathbf{y}(t) - \frac{\tau}{2} \nabla_{\mathbf{x}} f(\mathbf{x}(t)) \quad (6.51a)$$

$$\mathbf{x}(t + \tau) = \mathbf{x}(t) + \frac{\tau}{m} \mathbf{y}(t + \tau/2) \quad (6.51b)$$

$$\mathbf{y}(t + \tau) = \mathbf{y}(t + \tau/2) - \frac{\tau}{2} \nabla_{\mathbf{x}} f(\mathbf{x}(t + \tau)). \quad (6.51c)$$

Then, one repeats this procedure L times to arrive at a point $(\mathbf{x}', \mathbf{y}')$. To correct for the resulting discretization error, the proposal is either accepted or rejected in a final Metropolis-Hastings acceptance step. If the proposal distribution is symmetric (which we will confirm in a moment), the acceptance probability is

$$\alpha((\mathbf{x}', \mathbf{y}') \mid (\mathbf{x}, \mathbf{y})) \doteq \min\{1, \exp(H(\mathbf{x}', \mathbf{y}') - H(\mathbf{x}, \mathbf{y}))\}. \quad (6.52)$$

It follows that $p(\mathbf{x}, \mathbf{y})$ is the stationary distribution of the Markov chain underlying HMC. Due to the independence of \mathbf{x} and \mathbf{y} , this also implies that the projection to \mathbf{x} yields a Markov chain with stationary distribution $p(\mathbf{x})$.

So why is the proposal distribution symmetric? This follows from the time-reversibility of Hamiltonian dynamics. It is straightforward to check that the dynamics from Equation (6.50) are identical if we replace t with $-t$ and \mathbf{y} with $-\mathbf{y}$. Intuitively, unlike the position \mathbf{x} , the momentum \mathbf{y} is reversed when time is reversed as it depends on the velocity which is the time-derivative of the position.¹⁰ In simpler terms, time-reversibility states that if we observe the evolution of a system (e.g., two billiard balls colliding), we cannot distinguish whether

¹⁰The momentum of the i -th coordinate is $y_i = m_i v_i$ where m_i is the mass and $v_i = \frac{dx_i}{dt}$ is the velocity.

we are observing the system evolve forward or backward in time. The Leapfrog method maintains the time-reversibility of the dynamics.

Symmetry of the proposal distribution is ensured by proposing the point $(x', -y')$.¹¹ Intuitively, this is simply to ensure that the system is run backward in time as often as it is run forward in time. Recall that the momentum is resampled before each iteration (i.e., the proposed momentum is “discarded”) and observe that $p(x', -y') = p(x', y')$,¹² so we can safely disregard the direction of time when computing the acceptance probability in Equation (6.52).

Discussion

To summarize, Markov chain Monte Carlo methods use a Markov chain to approximately sample from an intractable distribution. Note that unlike for variational inference, the convergence of many methods can be guaranteed. Moreover, for log-concave distributions (e.g., with Bayesian logistic regression), the underlying Markov chain converges quickly to the stationary distribution. Methods such as Langevin dynamics and Hamiltonian Monte Carlo aim to accelerate mixing by proposing points with a higher acceptance probability than Metropolis-Hastings with “undirected” Gaussian proposals. Nevertheless, in general, the convergence (mixing time) may be slow, meaning that, in practice, accuracy and efficiency have to be traded.

Optional Readings

- Ma, Chen, Jin, Flammarion, and Jordan (2019).
Sampling can be faster than optimization.
- Teh, Thiery, and Vollmer (2016).
Consistency and fluctuations for stochastic gradient Langevin dynamics.
- Chen, Fox, and Guestrin (2014).
Stochastic gradient hamiltonian monte carlo.

Problems

6.1. Markov chain update.

Prove Equation (6.9), i.e., that one iteration of the Markov chain can be expressed as $q_{t+1} = q_t P$.

6.2. k -step transitions.

Prove that the entry $P^k(x, x')$ corresponds to the probability of transitioning from state $x \in S$ to state $x' \in S$ in exactly k steps.

¹¹ More formally, the proposal distribution is the Dirac delta at $(x', -y')$.

¹² We use here that $p(y | x)$ was chosen to be symmetric around zero.

6.3. Finding stationary distributions.

A news station classifies each day as “good”, “fair”, or “poor” based on its daily ratings which fluctuate with what is occurring in the news. Moreover, the following table shows the probabilistic relationship between the type of the current day and the probability of the type of the next day conditioned on the type of the current day.

		next day		
		good	fair	poor
current day	good	0.60	0.30	0.10
	fair	0.50	0.25	0.25
	poor	0.20	0.40	0.40

In the long run, what percentage of news days will be classified as “good”?

6.4. Example of Metropolis-Hastings.

Consider the state space $\{0, 1\}^n$ of binary strings having length n . Let the proposal distribution be $r(x' | x) = 1/n$ if x' differs from x in exactly one bit and $r(x' | x) = 0$ otherwise. Suppose we desire a stationary distribution p for which $p(x)$ is proportional to the number of ones that occur in the bit string x . For example, in the long run, a random walk should visit a string having five 1s five times as often as it visits a string having only a single 1. Provide a general formula for the acceptance probability $\alpha(x' | x)$ that would be used if we were to obtain the desired stationary distribution used the Metropolis-Hastings algorithm.

6.5. Practical examples of Gibbs sampling.

In this exercise, we look at some examples where Gibbs sampling is useful.

1. Consider the distribution

$$p(x, y) \doteq \binom{n}{x} y^{x+\alpha-1} (1-y)^{n-x+\beta-1}, \quad x \in [n], y \in [0, 1].$$

Convince yourself that it is hard to sample directly from p and prove that it is an easy task if one uses Gibbs sampling. That is, show that the conditional distributions $p(x | y)$ and $p(y | x)$ are easy to sample from.

Hint: Take a look at the Beta distribution (1.50).

2. Consider the following generative model $p(\mu, \lambda, x_{1:n})$ given by the likelihood $x_{1:n} | \mu, \lambda \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \lambda^{-1})$ and the independent priors

$$\mu \sim \mathcal{N}(\mu_0, \lambda_0^{-1}) \quad \text{and} \quad \lambda \sim \text{Gamma}(\alpha, \beta).$$

Gamma distribution The PDF of the *gamma distribution* $\text{Gamma}(\alpha, \beta)$ is defined as

$$\text{Gamma}(x; \alpha, \beta) \propto x^{\alpha-1} e^{-\beta x}, \quad x \in \mathbb{R}_{>0}.$$

A random variable $X \sim \text{Gamma}(\alpha, \beta)$ measures the waiting time until $\alpha > 0$ events occur in a Poisson process with rate $\beta > 0$. In particular, when $\alpha = 1$ then the gamma distribution coincides with the exponential distribution with rate β .

We would like to sample from the posterior $p(\mu, \lambda \mid x_{1:n})$. Show that

$$\mu \mid \lambda, x_{1:n} \sim \mathcal{N}(m_\lambda, l_\lambda^{-1}) \quad \text{and} \quad \lambda \mid \mu, x_{1:n} \sim \text{Gamma}(a_\mu, b_\mu),$$

and derive $m_\lambda, l_\lambda, a_\mu, b_\mu$. Such a prior is called a *semi-conjugate prior* to the likelihood, as the prior on μ is conjugate for any fixed value of λ and vice-versa.

3. Let us assume that $x_{1:n} \mid \alpha, c \stackrel{\text{iid}}{\sim} \text{Pareto}(\alpha, c)$ and assume the improper prior $p(\alpha, c) \propto \mathbb{1}\{\alpha, c > 0\}$ which corresponds to a noninformative prior. Derive the posterior $p(\alpha, c \mid x_{1:n})$. Then, also derive the conditional distributions $p(\alpha \mid c, x_{1:n})$ and $p(c \mid \alpha, x_{1:n})$, and observe that they correspond to known distributions / are easy to sample from.

6.6. Energy function of Bayesian logistic regression.

Recall from Equation (5.12) that the energy function of Bayesian logistic regression is

$$f(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)), \quad (6.53)$$

which coincided with the standard optimization objective of (regularized) logistic regression.

Show that the posterior distribution of Bayesian logistic regression is log-concave.

6.7. Maximum entropy property of Gibbs distribution.

1. Let X be a random variable supported on the finite set $\mathcal{T} \subset \mathbb{R}$.¹³ Show that the Gibbs distribution with energy function $\frac{1}{T}f(x)$ for some temperature scalar $T \in \mathbb{R}$ is the distribution with maximum entropy of all distributions supported on \mathcal{T} that satisfy the constraint $\mathbb{E}[f(X)] < \infty$.

Hint: Solve the dual problem (analogously to Problem 5.7).

2. What happens for $T \rightarrow \{0, \infty\}$?

¹³ The same result can be shown to hold for arbitrary compact subsets.

6.8. Energy reduction of Gibbs sampling.

Let $p(\mathbf{x})$ be a probability density over \mathbb{R}^d , which we want to sample from. Assume that p is a Gibbs distribution with energy function $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

In this exercise, we will study a single round of Gibbs sampling with initial state \mathbf{x} and final state \mathbf{x}' where

$$x'_j = \begin{cases} x'_i & \text{if } j = i \\ x_j & \text{otherwise} \end{cases}$$

for some fixed index i and $x'_i \sim p(\cdot | \mathbf{x}_{-i})$.

Show that

$$\mathbb{E}_{x'_i \sim p(\cdot | \mathbf{x}_{-i})} [f(\mathbf{x}')] \leq f(\mathbf{x}) - S[p(x_i | \mathbf{x}_{-i})] + H[p(\cdot | \mathbf{x}_{-i})]. \quad (6.54)$$

That is, the energy is expected to decrease if the surprise of x_i given \mathbf{x}_{-i} is larger than the expected surprise of the new x'_i given \mathbf{x}_{-i} , i.e., $S[p(x_i | \mathbf{x}_{-i})] \geq H[p(\cdot | \mathbf{x}_{-i})]$.

Hint: Recall the framing of Gibbs sampling as a variant of Metropolis-Hastings and relate this to the acceptance distribution of Metropolis-Hastings when p is a Gibbs distribution.

6.9. Mixing time of Langevin dynamics.

In this exercise, we will show that for certain Gibbs distributions, $p(\theta) \propto \exp(-f(\theta))$, Langevin dynamics is rapidly mixing. To do this, we will observe that Langevin dynamics can be seen as a continuous-time optimization algorithm in the space of distributions.

First, we consider a simpler and more widely-known optimization algorithm, namely the *gradient flow*

$$d\mathbf{x}_t = -\nabla f(\mathbf{x}_t) dt. \quad (6.55)$$

Note that gradient descent is simply the discrete-time approximation of gradient flow just as ULA is the discrete-time approximation of Langevin dynamics. In the analysis of ODEs such as the gradient flow, so-called *Lyapunov functions* are commonly used to prove convergence of \mathbf{x}_t to a fixed point (also called an *equilibrium*).

Let us assume that f is α -strongly convex for some $\alpha > 0$, that is,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\alpha}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n. \quad (6.56)$$

In words, f is lower bounded by a quadratic function with curvature α . Moreover, assume w.l.o.g. that f minimized at $f(\mathbf{0}) = 0$.¹⁴

1. Show that f satisfies the *Polyak-Łojasiewicz (PL) inequality*, i.e.,

$$f(\mathbf{x}) \leq \frac{1}{2\alpha} \|\nabla f(\mathbf{x})\|_2^2 \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (6.57)$$

2. Prove $\frac{d}{dt} f(\mathbf{x}_t) \leq -2\alpha f(\mathbf{x}_t)$.

Thus, $\mathbf{0}$ is the fixed point of Equation (6.55) and the Lyapunov function f is monotonically decreasing along the trajectory of \mathbf{x}_t . We recall *Grönwall's inequality* which states that for any real-valued continuous functions $g(t)$ and $\beta(t)$ on the interval $[0, T] \subset \mathbb{R}$ such that $\frac{d}{dt} g(t) \leq \beta(t)g(t)$ for all $t \in [0, T]$ we have

$$g(t) \leq g(0) \exp\left(\int_0^t \beta(s) ds\right) \quad \forall t \in [0, T]. \quad (6.58)$$

¹⁴ This can always be achieved by shifting the coordinate system and subtracting a constant from f .

3. Conclude that $f(\mathbf{x}_t) \leq e^{-2\alpha t} f(\mathbf{x}_0)$.

Now that we have proven the convergence of gradient flow using f as Lyapunov function, we will follow the same template to prove the convergence of Langevin dynamics to the distribution $p(\theta) \propto \exp(-f(\theta))$. We will use that the evolution of $\{\theta_t\}_{t \geq 0}$ following the Langevin dynamics (6.46) is equivalently characterized by their densities $\{q_t\}_{t \geq 0}$ following the *Fokker-Planck equation*

$$\frac{\partial q_t}{\partial t} = \nabla \cdot (q_t \nabla f) + \Delta q_t. \quad (6.59)$$

Here, $\nabla \cdot$ and Δ are the divergence and Laplacian operators, respectively.¹⁵ Intuitively, the first term of the Fokker-Planck equation corresponds to the drift and its second term corresponds to the diffusion (i.e., the Gaussian noise).

Remark 6.26: Intuition on vector calculus

Recall that the divergence $\nabla \cdot \mathbf{F}$ of a vector field \mathbf{F} measures the change of volume under the flow of \mathbf{F} . That is, if in the small neighborhood of a point x , \mathbf{F} points towards x , then the divergence at x is negative as the volume shrinks. If \mathbf{F} points away from x , then the divergence at x is positive as the volume increases.

The Laplacian $\Delta \varphi = \nabla \cdot (\nabla \varphi)$ of a scalar field φ can be understood intuitively as measuring “heat dissipation”. That is, if $\varphi(x)$ is smaller than the average value of φ in a small neighborhood of x , then the Laplacian at x is positive.

Regarding the Fokker-Planck equation (6.59), the second term Δq_t can therefore be understood as locally dissipating the probability mass of q_t (which is due to the diffusion term in the SDE). On the other hand, the term $\nabla \cdot (q_t \nabla f)$ can be understood as a Laplacian of f “weighted” by q_t . Intuitively, the vector field ∇f moves flow in the direction of high energy, and hence, its divergence is larger in regions of lower energy and smaller in regions of higher energy. This term therefore corresponds to a drift from regions of high energy to regions of low energy.

4. Show that $\Delta q_t = \nabla \cdot (q_t \nabla \log q_t)$, implying that the Fokker-Planck equation simplifies to

$$\frac{\partial q_t}{\partial t} = \nabla \cdot \left(q_t \nabla \log \frac{q_t}{p} \right). \quad (6.60)$$

Hint: The Laplacian of a scalar field φ is $\Delta \varphi \doteq \nabla \cdot (\nabla \varphi)$.

Observe that the Fokker-Planck equation already implies that p is indeed a stationary distribution, as if $q_t = p$ then $\frac{\partial q_t}{\partial t} = 0$. Moreover, note

¹⁵ For ease of notation, we omit the explicit dependence of q_t , p , and f on θ .

the similarity of the integrand of $\text{KL}(q_t\|p)$, $q_t \log \frac{q_t}{p}$, to Equation (6.60).

We will therefore use the KL-divergence as Lyapunov function.

5. Prove $\frac{d}{dt} \text{KL}(q_t\|p) = -J(q_t\|p)$. Here,

$$J(q_t\|p) \doteq \mathbb{E}_{\theta \sim q_t} \left[\left\| \nabla \log \frac{q_t(\theta)}{p(\theta)} \right\|_2^2 \right] \quad (6.61)$$

denotes the *relative Fisher information* of p with respect to q_t . Hint:
For any distribution q on \mathbb{R}^n ,

$$\int_{\mathbb{R}^n} (\nabla \cdot q \mathbf{F}) \varphi \, dx = - \int_{\mathbb{R}^n} q \, \nabla \varphi \cdot \mathbf{F} \, dx \quad (6.62)$$

follows for any vector field \mathbf{F} and scalar field φ from the divergence theorem and the product rule of the divergence operator.

Thus, the relative Fisher information can be seen as the negated time-derivative of the KL-divergence, and as $J(q_t\|p) \geq 0$ it follows that the KL-divergence is decreasing along the trajectory.

The *log-Sobolev inequality* (LSI) is satisfied by a distribution p with a constant $\alpha > 0$ if for all q :

$$\text{KL}(q\|p) \leq \frac{1}{2\alpha} J(q\|p). \quad (6.63)$$

It is a classical result that if f is α -strongly convex then p satisfies the LSI with constant α (Bakry and Émery, 2006).

6. Show that if f is α -strongly convex for some $\alpha > 0$ (we say that p is “strongly log-concave”), then $\text{KL}(q_t\|p) \leq e^{-2\alpha t} \text{KL}(q_0\|p)$.
7. Conclude that under the same assumption on f , Langevin dynamics is rapidly mixing, i.e., $\tau_{\text{TV}}(\epsilon) \in O(\text{poly}(n, \log(1/\epsilon)))$.

To summarize, we have seen that Langevin dynamics is an optimization scheme in the space of distributions, and that its convergence can be analyzed analogously to classical optimization schemes. Notably, in this exercise we have studied continuous-time Langevin dynamics. Convergence guarantees for discrete-time approximations can be derived using the same techniques. If this interests you, refer to “Rapid convergence of the unadjusted Langevin algorithm: Isoperimetry suffices” (Vempala and Wibisono, 2019).

6.10. Hamiltonian Monte Carlo.

1. Prove that if the dynamics are solved exactly (as opposed to numerically using the Leapfrog method), then the acceptance probability of the MH-step is always 1.
2. Prove that the Langevin Monte Carlo algorithm from (6.44) can be seen as a special case of HMC if only one Leapfrog step is used ($L = 1$) and $m = 1$.

7

Deep Learning

We began our journey through probabilistic machine learning with linear models. In most practical applications, however, it is seen that models perform better when labels may nonlinearly depend on the inputs, and for this reason linear models are often used in conjunction with “hand-designed” nonlinear features. Designing these features is costly, time-consuming, and requires expert knowledge. Moreover, the design of such features is inherently limited by human comprehension and ingenuity.

7.1 Artificial Neural Networks

One widely used family of nonlinear functions are *artificial “deep” neural networks*,¹

$$f : \mathbb{R}^d \rightarrow \mathbb{R}^k, \quad f(\mathbf{x}; \boldsymbol{\theta}) \doteq \varphi(W_L \varphi(W_{L-1}(\cdots \varphi(W_1 \mathbf{x})))) \quad (7.1)$$

where $\boldsymbol{\theta} \doteq [W_1, \dots, W_L]$ is a vector of *weights* (written as matrices $W_l \in \mathbb{R}^{n_l \times n_{l-1}}$)² and $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is a component-wise nonlinear function. Thus, a deep neural network can be seen as nested (“deep”) linear functions composed with nonlinearities. This simple kind of neural network is also called a *multilayer perceptron*.

In this chapter, we will be focusing mostly on performing probabilistic inference with a *given* neural network architecture. To this end, understanding the basic architecture of a multilayer perceptron will be sufficient for us. For a more thorough introduction to the field of deep learning and various architectures, you may refer to the books of Goodfellow et al. (2016) and Prince (2023).

A neural network can be visualized by a *computation graph*. An example for such a computation graph is given in Figure 7.1. The columns of the computation graph are commonly called *layers*, whereby the

¹ In the following, we will refrain from using the characterizations “artificial” and “deep” for better readability.

² where $n_0 = d$ and $n_L = k$

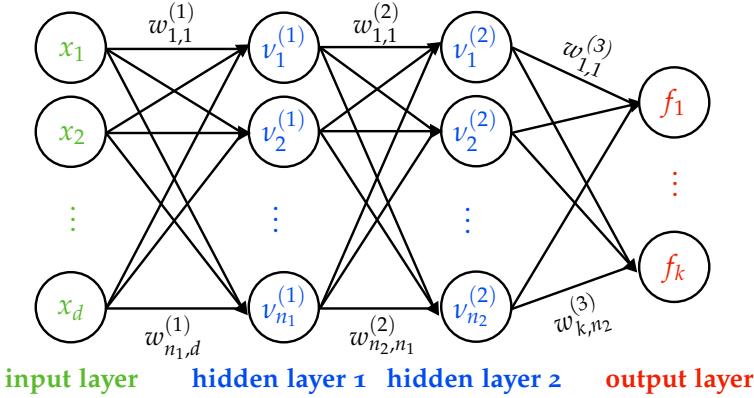


Figure 7.1: Computation graph of a neural network with two hidden layers.

left-most column is the *input layer*, the right-most column is the *output layer*, and the remaining columns are the *hidden layers*. The inputs are (as we have previously) referred to as $x \doteq [x_1, \dots, x_d]$. The outputs (i.e., vertices of the output layer) are often referred to as *logits* and named $f \doteq [f_1, \dots, f_k]$. The *activations* of an individual (hidden) layer l of the neural network are described by

$$\nu^{(l)} \doteq \varphi(W_l \nu^{(l-1)}) \quad (7.2)$$

where $\nu^{(0)} = x$. The activation of the i -th node is $\nu_i^{(l)} = \nu^{(l)}(i)$.

7.1.1 Activation Functions

The nonlinearity φ is called an *activation function*. The following two activation functions are particularly common:

1. The *hyperbolic tangent* (Tanh) is defined as

$$\text{Tanh}(z) \doteq \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} \in (-1, 1). \quad (7.3)$$

Tanh is a scaled and shifted variant of the sigmoid function (5.9) which we have previously seen in the context of logistic regression as $\text{Tanh}(z) = 2\sigma(2z) - 1$.

2. The *rectified linear unit* (ReLU) is defined as

$$\text{ReLU}(z) \doteq \max\{z, 0\} \in [0, \infty). \quad (7.4)$$

In particular, the ReLU activation function leads to “sparser” gradients as it selects the halfspace of inputs with positive sign. Moreover, the gradients of ReLU do not “vanish” as $z \rightarrow \pm\infty$ which can lead to faster training.

It is important that the activation function is nonlinear because otherwise, any composition of layers would still represent a linear function. Non-linear activation functions allow the network to represent arbitrary functions. This is known as the *universal approximation theorem*,

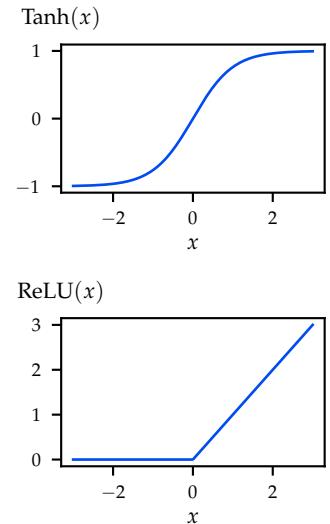


Figure 7.2: The Tanh and ReLU activation functions, respectively.

and it states that any artificial neural network with just a single hidden layer (with arbitrary width) and non-polynomial activation function φ can approximate any continuous function to an arbitrary accuracy.

7.1.2 Classification

Although we mainly focus on regression, neural networks can equally well be used for classification. If we want to classify inputs into c separate classes, we can simply construct a neural network with c outputs, $f = [f_1, \dots, f_c]$, and normalize them into a probability distribution. Often, the *softmax function* is used for normalization,

$$\sigma_i(f) \doteq \frac{\exp(f_i)}{\sum_{j=1}^c \exp(f_j)} \quad (7.5)$$

where $\sigma_i(f)$ corresponds to the probability mass of class i . The softmax is a generalization of the logistic function (5.9) to more than two classes $\textcircled{2}$. Note that the softmax corresponds to a Gibbs distribution with energies $-f$.

7.1.3 Maximum Likelihood Estimation

We will study neural networks under the lens of supervised learning (cf. Section 1.3) where we are provided some independently-sampled (noisy) data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ generated according to an unknown process $y \sim p(\cdot | x, \theta^*)$, which we wish to approximate.

Upon initialization, the network does generally not approximate this process well, so a key element of deep learning is “learning” a parameterization θ that is a good approximation. To this end, one typically considers a loss function $\ell(\theta; y)$ which quantifies the “error” of the network outputs $f(x; \theta)$. In the classical setting of regression, i.e., $y \in \mathbb{R}$ and $k = 1$, ℓ is often taken to be the (*empirical*) *mean squared error*,

$$\ell_{\text{mse}}(\theta; \mathcal{D}) \doteq \frac{1}{n} \sum_{i=1}^n (f(x_i; \theta) - y_i)^2. \quad (7.6)$$

As we have already seen in Section 2.0.1 in the context of linear regression, minimizing mean squared error corresponds to maximum likelihood estimation under a Gaussian likelihood.

In the setting of classification where $y \in \{0, 1\}^c$ is a one-hot encoding of class membership,³ it is instead common to interpret the outputs of a neural network as probabilities akin to our discussion in Section 7.1.2. We denote by $q_\theta(\cdot | x)$ the resulting probability distribution over classes with PMF $[\sigma_1(f(x; \theta)), \dots, \sigma_c(f(x; \theta))]$, and aim to find θ such that $q_\theta(y | x) \approx p(y | x)$. In this context, it is common to

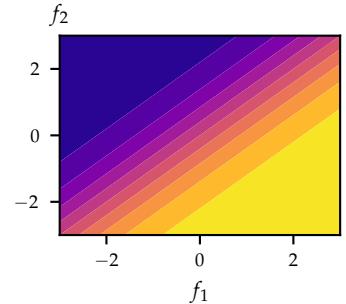


Figure 7.3: Softmax $\sigma_1(f_1, f_2)$ for a binary classification problem. Blue denotes a small probability and yellow denotes a large probability of belonging to class 1, respectively.

Problem 7.1

³ That is, exactly one component of y is 1 and all others are 0, indicating to which class the given example belongs.

minimize the cross-entropy,

$$\begin{aligned} H[p\|q_{\theta}] &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p} [-\log q_{\theta}(\mathbf{y} \mid \mathbf{x})] \\ &\approx -\frac{1}{n} \sum_{i=1}^n \underbrace{\log q_{\theta}(\mathbf{y}_i \mid \mathbf{x}_i)}_{\doteq \ell_{\text{ce}}(\boldsymbol{\theta}; \mathcal{D})} \end{aligned} \quad (7.7)$$

which can be understood as minimizing the surprise about the training data under the model. ℓ_{ce} is called the *cross-entropy loss*. Disregarding the constant $1/n$, we can rewrite the cross-entropy loss as

$$\propto -\sum_{i=1}^n \log q_{\theta}(\mathbf{y}_i \mid \mathbf{x}_i) = \ell_{\text{nll}}(\boldsymbol{\theta}; \mathcal{D}) \quad (7.8)$$

Recall that $\ell_{\text{nll}}(\boldsymbol{\theta}; \mathcal{D})$ is the *negative log-likelihood* of the training data, and thus, empirically minimizing cross-entropy can equivalently be interpreted as maximum likelihood estimation.⁴ Furthermore, recall from Problem 5.1 that for a two-class classification problem the cross-entropy loss is equivalent to the logistic loss.

7.1.4 Backpropagation

A crucial property of neural networks is that they are differentiable. That is, we can compute gradients $\nabla_{\theta} \ell$ of f with respect to the parameterization of the model $\boldsymbol{\theta} = W_{1:L}$ and some loss function $\ell(f; \mathbf{y})$. Being able to obtain these gradients efficiently allows for “learning” a particular function from data using first-order optimization methods.

The algorithm for computing gradients of a neural network is called *backpropagation* and is essentially a repeated application of the chain rule. Note that using the chain rule for every path through the network is computationally infeasible, as this quickly leads to a combinatorial explosion as the number of hidden layers is increased. The key insight of backpropagation is that we can use the *feed-forward* structure of our neural network to memoize computations of the gradient, yielding a linear time algorithm. Obtaining gradients by backpropagation is often called *automatic differentiation (auto-diff)*. For more details, refer to Goodfellow et al. (2016).

Computing the exact gradient for each data point is still fairly expensive when the size of the neural network is large. Typically, stochastic gradient descent is used to obtain unbiased gradient estimates using batches of only m of the n data points, where $m \ll n$.

7.2 Bayesian Neural Networks

How can we perform probabilistic inference in neural networks? We adopt the same strategy which we already used for Bayesian linear

using the definition of cross-entropy
(5.32)

using Monte Carlo sampling

⁴ We have previously seen this equivalence of MLE and empirically minimizing KL-divergence in Section 5.4.6 (minimizing the cross-entropy $H[p\|q_{\theta}]$ is equivalent to minimizing forward-KL $KL(p\|q_{\theta})$). Note that this interpretation is not exclusive to the canonical cross-entropy loss from Equation (7.7), but holds for any MLE. For example, minimizing mean squared error corresponds to empirically minimizing the KL-divergence with a Gaussian likelihood.

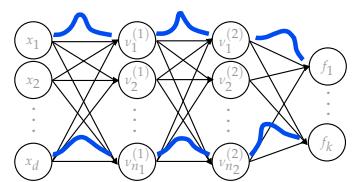


Figure 7.4: Bayesian neural networks model a distribution over the weights of a neural network.

regression, we impose a Gaussian prior on the weights,

$$\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I}). \quad (7.9)$$

Similarly, we can use a Gaussian likelihood to describe how well the data is described by the model f ,

$$y | \mathbf{x}, \boldsymbol{\theta} \sim \mathcal{N}(f(\mathbf{x}; \boldsymbol{\theta}), \sigma_n^2). \quad (7.10)$$

Thus, instead of considering weights as point estimates which are learned exactly, *Bayesian neural networks* learn a distribution over the weights of the network. In principle, other priors and likelihoods can be used, yet Gaussians are typically chosen due to their closedness properties, which we have seen in Section 1.2.3 and many times since.

7.2.1 Maximum a Posteriori Estimation

Before studying probabilistic inference, let us first consider MAP estimation in the context of neural networks.

The MAP estimate of the weights is obtained by

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \arg \max_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}) + \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}). \quad (7.11)$$

In Section 7.1.3, we have seen that the negative log-likelihood under a Gaussian likelihood (7.10) is the squared error between label and prediction,

$$\log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = -\frac{(y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2}{2\sigma_n^2} + \text{const.} \quad (7.12)$$

Obtaining the MAP estimate instead, simply corresponds to adding an L_2 -regularization term to the squared error loss,

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_p^2} \|\boldsymbol{\theta}\|_2^2 + \frac{1}{2\sigma_n^2} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2. \quad (7.13)$$

As we have already seen in Remark A.31 and the context of Bayesian linear regression (and ridge regression), using a Gaussian prior is equivalent to applying weight decay.⁵ Using gradient ascent, we obtain the following update rule,

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}(1 - \lambda \eta_t) + \eta_t \sum_{i=1}^n \nabla \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) \quad (7.14)$$

where $\lambda \doteq 1/\sigma_p^2$. The gradients of the likelihood can be obtained using automatic differentiation.

using that for an isotropic Gaussian prior, $\log p(\boldsymbol{\theta}) = -\frac{1}{2\sigma_p^2} \|\boldsymbol{\theta}\|_2^2 + \text{const}$

⁵ Recall that weight decay regularizes weights by shrinking them towards zero.

7.2.2 Heteroscedastic Noise

Equation (7.10) uses the scalar parameter σ_n^2 to model the aleatoric uncertainty (label noise), similarly to how we modeled the label noise in Bayesian linear regression and Gaussian processes. Such a noise model is called *homoscedastic noise* as the noise is assumed to be uniform across the domain. In many settings, however, the noise is inherently *heteroscedastic noise*. That is, the noise varies depending on the input and which “region” of the domain the input is from. This behavior is visualized in Figure 7.5.

There is a natural way of modeling heteroscedastic noise with Bayesian neural networks. We use a neural network with two outputs f_1 and f_2 and define

$$y | \mathbf{x}, \boldsymbol{\theta} \sim \mathcal{N}(\mu(\mathbf{x}; \boldsymbol{\theta}), \sigma^2(\mathbf{x}; \boldsymbol{\theta})) \quad \text{where} \quad (7.15a)$$

$$\mu(\mathbf{x}; \boldsymbol{\theta}) \doteq f_1(\mathbf{x}; \boldsymbol{\theta}), \quad (7.15b)$$

$$\sigma^2(\mathbf{x}; \boldsymbol{\theta}) \doteq \exp(f_2(\mathbf{x}; \boldsymbol{\theta})). \quad (7.15c)$$

Using this model, the likelihood term from Equation (7.11) is

$$\begin{aligned} \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) &= \log \mathcal{N}(y_i; \mu(\mathbf{x}_i; \boldsymbol{\theta}), \sigma^2(\mathbf{x}_i; \boldsymbol{\theta})) \\ &= \log \frac{1}{\sqrt{2\pi\sigma^2(\mathbf{x}_i; \boldsymbol{\theta})}} - \frac{(y_i - \mu(\mathbf{x}_i; \boldsymbol{\theta}))^2}{2\sigma^2(\mathbf{x}_i; \boldsymbol{\theta})} \\ &= \underbrace{\log \frac{1}{\sqrt{2\pi}}}_{\text{const}} - \frac{1}{2} \left[\log \sigma^2(\mathbf{x}_i; \boldsymbol{\theta}) + \frac{(y_i - \mu(\mathbf{x}_i; \boldsymbol{\theta}))^2}{\sigma^2(\mathbf{x}_i; \boldsymbol{\theta})} \right]. \end{aligned} \quad (7.16)$$

Hence, the model can either explain the label y_i by an accurate model $\mu(\mathbf{x}_i; \boldsymbol{\theta})$ or by a large variance $\sigma^2(\mathbf{x}_i; \boldsymbol{\theta})$, yet, it is penalized for choosing large variances. Intuitively, this allows to attenuate losses for some data points by attributing them to large variance (when no model reflecting all data points simultaneously can be found). This allows the model to “learn” its aleatoric uncertainty. However, recall that the MAP estimate still corresponds to a point estimate of the weights, so we forgo modeling the epistemic uncertainty.

7.3 Approximate Probabilistic Inference

Naturally, we want to understand the epistemic uncertainty of our model too. However, learning and inference in Bayesian neural networks are generally intractable (even when using a Gaussian prior and likelihood) when the noise is not assumed to be homoscedastic and known.⁶ Thus, we are led to the techniques of approximate inference, which we discussed in the previous two chapters.

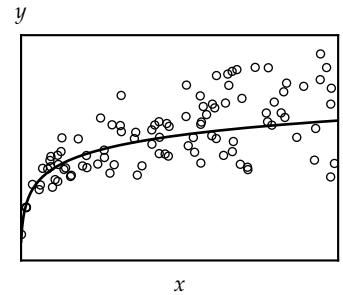


Figure 7.5: Illustration of data with variable (heteroscedastic) noise. The noise increases as the inputs increase in magnitude. The noise-free function is shown in black.

we exponentiate f_2 to ensure non-negativity of the variance

note that the normalizing constant depends on the noise model!

⁶ In this case, the conjugate prior to a Gaussian likelihood is not a Gaussian. See, e.g., https://en.wikipedia.org/wiki/Conjugate_prior.

7.3.1 Variational Inference

As we have discussed in Chapter 5, we can apply black box stochastic variational inference which — in the context of neural networks — is also known as *Bayes by Backprop*. As variational family, we use the family of independent Gaussians which we have already encountered in Example 5.5.⁷ Recall the fundamental objective of variational inference (5.25),

$$\begin{aligned} & \arg \min_{q \in \mathcal{Q}} \text{KL}(q \| p(\cdot | \mathbf{x}_{1:n}, y_{1:n})) \\ &= \arg \max_{q \in \mathcal{Q}} L(q, p; \mathcal{D}) \\ &= \arg \max_{q \in \mathcal{Q}} \mathbb{E}_{\theta \sim q} [\log p(y_{1:n} | \mathbf{x}_{1:n}, \theta)] - \text{KL}(q \| p(\cdot)). \end{aligned}$$

The KL-divergence $\text{KL}(q \| p(\cdot))$ can be expressed in closed-form for Gaussians.⁸ Recall that we can obtain unbiased gradient estimates of the expectation using the reparameterization trick (5.68),

$$\mathbb{E}_{\theta \sim q} [\log p(y_{1:n} | \mathbf{x}_{1:n}, \theta)] = \mathbb{E}_{\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\log p(y_{1:n} | \mathbf{x}_{1:n}, \theta) \Big|_{\theta = \Sigma^{1/2} \varepsilon + \mu} \right].$$

As Σ is the diagonal matrix $\text{diag}\{\sigma_1^2, \dots, \sigma_d^2\}$, $\Sigma^{1/2} = \text{diag}\{\sigma_1, \dots, \sigma_d\}$. The gradients of the likelihood can be obtained using backpropagation. We can now use the variational posterior q_λ to perform approximate inference,

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}) &= \int p(y^* | \mathbf{x}^*, \theta) p(\theta | \mathbf{x}_{1:n}, y_{1:n}) d\theta \\ &= \mathbb{E}_{\theta \sim p(\cdot | \mathbf{x}_{1:n}, y_{1:n})} [p(y^* | \mathbf{x}^*, \theta)] \\ &\approx \mathbb{E}_{\theta \sim q_\lambda} [p(y^* | \mathbf{x}^*, \theta)] \\ &\approx \frac{1}{m} \sum_{i=1}^m p(y^* | \mathbf{x}^*, \theta^{(i)}) \end{aligned} \tag{7.17}$$

for independent samples $\theta^{(i)} \stackrel{\text{iid}}{\sim} q_\lambda$,

$$= \frac{1}{m} \sum_{i=1}^m \mathcal{N}(y^*; \mu(\mathbf{x}^*; \theta^{(i)}), \sigma^2(\mathbf{x}^*; \theta^{(i)})). \tag{7.18}$$

Intuitively, variational inference in Bayesian neural networks can be interpreted as averaging the predictions of multiple neural networks drawn according to the variational posterior q_λ .

Using the Monte Carlo samples $\theta^{(i)}$, we can also estimate the mean of our predictions,

$$\mathbb{E}[y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}] \approx \frac{1}{m} \sum_{i=1}^m \mu(\mathbf{x}^*; \theta^{(i)}) \doteq \bar{\mu}(\mathbf{x}^*), \tag{7.19}$$

⁷ Independent Gaussians are useful because they can be encoded using only $2d$ parameters, which is crucial when the size of the neural network is large.

using Equation (5.53)

using Equation (5.55c)

⁸ see Equation (5.41)

using the sum rule (1.7) and product rule (1.11)

interpreting the integral as an expectation over the posterior

approximating the posterior with the variational posterior q_λ

using Monte Carlo sampling

using the neural network

and the variance of our predictions,

$$\begin{aligned}\text{Var}[y^* \mid \mathbf{x}^*, \mathbf{x}_{1:n}, \mathbf{y}_{1:n}] &= \mathbb{E}_\theta[\text{Var}_{y^*}[y^* \mid \mathbf{x}^*, \theta]] + \text{Var}_\theta[\mathbb{E}_{y^*}[y^* \mid \mathbf{x}^*, \theta]] \\ &= \mathbb{E}_\theta[\sigma^2(\mathbf{x}^*; \theta)] + \text{Var}_\theta[\mu(\mathbf{x}^*; \theta)].\end{aligned}$$

Recall from Equation (2.18) that the first term corresponds to the aleatoric uncertainty of the data and the second term corresponds to the epistemic uncertainty of the model. We can approximate them using the Monte Carlo samples $\theta^{(i)}$,

$$\begin{aligned}\text{Var}[y^* \mid \mathbf{x}^*, \mathbf{x}_{1:n}, \mathbf{y}_{1:n}] &\approx \frac{1}{m} \sum_{i=1}^m \sigma^2(\mathbf{x}^*; \theta^{(i)}) \\ &\quad + \frac{1}{m-1} \sum_{i=1}^m (\mu(\mathbf{x}^*; \theta^{(i)}) - \bar{\mu}(\mathbf{x}^*))^2\end{aligned}\tag{7.20}$$

using a sample mean (A.14) and sample variance (A.16).

7.3.2 Markov Chain Monte Carlo

As we have discussed in Chapter 6, an alternative method to approximating the full posterior distribution is to sample from it directly. By the ergodic theorem (6.28), we can use any of the discussed sampling strategies to obtain samples $\theta^{(t)}$ such that

$$p(y^* \mid \mathbf{x}^*, \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) \approx \frac{1}{T} \sum_{t=1}^T p(y^* \mid \mathbf{x}^*, \theta^{(t)}). \quad \text{see (6.29)}$$

Here, we omit the offset t_0 which is commonly used to avoid the “burn-in” period for simplicity. Algorithms such as SGLD or SG-HMC are often used as they rely only on stochastic gradients of the loss function which can be computed efficiently using automatic differentiation.

Typically, for large networks, we cannot afford to store all T samples of models. Thus, we need to summarize the iterates.⁹ One approach is to keep track of m snapshots of weights $[\theta^{(1)}, \dots, \theta^{(m)}]$ according to some schedule and use those for inference (e.g., by averaging the predictions of the corresponding neural networks). This approach of sampling a subset of some data is generally called *subsampling*.

Another approach is to summarize (that is, approximate) the weights using sufficient statistics (e.g., a Gaussian).¹⁰ In other words, we learn the Gaussian approximation,

$$\theta \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad \text{where} \tag{7.21a}$$

$$\boldsymbol{\mu} \doteq \frac{1}{T} \sum_{i=1}^T \theta^{(i)}, \tag{7.21b}$$

using the law of total variance (1.41)

using the likelihood (7.15a)

see (6.29)

⁹ That is, combine the individual samples of weights $\theta^{(i)}$.

¹⁰ A statistic is *sufficient* for a family of probability distributions if the samples from which it is calculated yield no more information than the statistic with respect to the learned parameters. We provide a formal definition in Section 8.2.

using a sample mean (A.14)

$$\Sigma \doteq \frac{1}{T-1} \sum_{i=1}^T (\theta^{(i)} - \mu)(\theta^{(i)} - \mu)^\top, \quad (7.21c)$$

using sample means and sample (co)variances. This can be implemented efficiently using running averages of the first and second moments,

$$\mu \leftarrow \frac{1}{T+1}(T\mu + \theta) \quad \text{and} \quad A \leftarrow \frac{1}{T+1}(TA + \theta\theta^\top) \quad (7.22)$$

upon observing the fresh sample θ . Σ can easily be calculated from these moments,

$$\Sigma = \frac{T}{T-1}(A - \mu\mu^\top). \quad (7.23)$$

To predict, we can sample weights θ from the learned Gaussian.

Remark 7.1: Stochastic weight averaging

It turns out that this approach works well even without injecting additional Gaussian noise during training, e.g., when using SGD rather than SGLD. Simply taking the mean of the iterates of SGD is called *stochastic weight averaging* (SWA). Describing the iterates of SGD by Gaussian sufficient statistics (analogously to Equation (7.21)), is known as the *stochastic weight averaging-Gaussian* (SWAG) method (Izmailov et al., 2018).

using a sample variance (A.16)

using the characterization of sample variance in terms of estimators of the first and second moment (A.17)

7.3.3 Dropout and Dropconnect

We will now discuss two approximate inference techniques that are tailored to neural networks. The first is “dropout”/“dropconnect” regularization. Traditionally, *dropout regularization* (Hinton et al., 2012; Srivastava et al., 2014) randomly omits vertices of the computation graph during training, see Figure 7.7. In contrast, *dropconnect regularization* (Wan et al., 2013) randomly omits edges of the computation graph. The key idea that we will present here is to interpret this type of regularization as performing variational inference.

In our exposition, we will focus on dropconnect, but the same approach also applies to dropout (Gal and Ghahramani, 2016). Suppose that we omit an edge of the computation graph (i.e., set its weight to zero) with probability p . Then our variational posterior is given by

$$q(\theta | \lambda) \doteq \prod_{j=1}^d q_j(\theta_j | \lambda_j) \quad (7.24)$$

where d is the number of weights of the neural network and

$$q_j(\theta_j | \lambda_j) \doteq p\delta_0(\theta_j) + (1-p)\delta_{\lambda_j}(\theta_j). \quad (7.25)$$

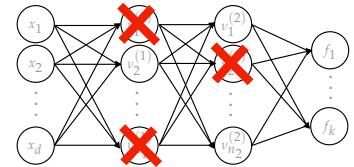


Figure 7.6: Illustration of dropout regularization. Some vertices of the computation graph are randomly omitted. In contrast, dropconnect regularization randomly omits edges of the computation graph.

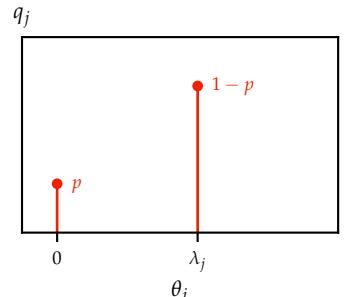


Figure 7.7: Interpretation of dropconnect regularization as variational inference. The only coordinates where the variational posterior q_j has positive probability are 0 and λ_j .

Here, δ_α is the Dirac delta function with point mass at α .¹¹ The variational parameters λ correspond to the “original” weights of the network. In words, the variational posterior expresses that the j -th weight has value 0 with probability p and value λ_j with probability $1 - p$. For fixed weights λ , sampling from the variational posterior q_λ corresponds to sampling a vector z with entries $z(i) \sim \text{Bern}(p)$, yielding $z \odot \lambda$ which is one of 2^d possible subnetworks.¹²

The weights λ can be learned by maximizing the ELBO, analogously to black-box variational inference. The KL-divergence term of the ELBO is not tractable for the variational family described by Equation (7.25), instead a common approach is to use a mixture of Gaussians:

$$q_j(\theta_j | \lambda_j) \doteq p\mathcal{N}(\theta_j; 0, 1) + (1 - p)\mathcal{N}(\theta_j; \lambda_j, 1). \quad (7.26)$$

In this case, it can be shown that $\text{KL}(q_\lambda \| p(\cdot)) \approx \frac{p}{2} \|\lambda\|_2^2$ for sufficiently large d (Gal and Ghahramani, 2015, proposition 1). Thus,

$$\begin{aligned} & \arg \max_{\lambda \in \Lambda} L(q_\lambda, p; \mathcal{D}) \\ &= \arg \max_{\lambda \in \Lambda} \mathbb{E}_{\theta \sim q_\lambda} [\log p(y_{1:n} | x_{1:n}, \theta)] - \text{KL}(q \| p(\cdot)) && \text{using Equation (5.55c)} \\ &\approx \arg \min_{\lambda \in \Lambda} -\frac{1}{m} \sum_{i=1}^m \log p(y_{1:n} | x_{1:n}, \theta) \Big|_{\theta=z^{(i)} \odot \lambda + \varepsilon^{(i)}} + \frac{p}{2} \|\lambda\|_2^2 && \text{using Monte Carlo sampling} \end{aligned} \quad (7.27)$$

where we reparameterize $\theta \sim q_\lambda$ by $\theta = z \odot \lambda + \varepsilon$ with $z(i) \sim \text{Bern}(p)$ and $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Equation (7.27) is the standard L_2 -regularized loss function of a neural network with weights λ and dropconnect, and it is straightforward to obtain unbiased gradient estimates by automatic differentiation.

Crucially, for the interpretation of dropconnect regularization as variational inference to be valid, we also need to perform dropconnect regularization during inference,

$$\begin{aligned} p(y^* | x^*, x_{1:n}, y_{1:n}) &\approx \mathbb{E}_{\theta \sim q_\lambda} [p(y^* | x^*, \theta)] \\ &\approx \frac{1}{m} \sum_{i=1}^m p(y^* | x^*, \theta^{(i)}) && \text{using Monte Carlo sampling} \end{aligned} \quad (7.28)$$

where $\theta^{(i)} \stackrel{\text{iid}}{\sim} q_\lambda$ are independent samples. This coincides with our earlier discussion of variational inference for Bayesian neural networks in Equation (7.17). In words, we average the predictions of m neural networks for each of which we randomly “drop out” weights.

¹¹ see Appendix A.1.4

¹² $A \odot B$ denotes the Hadamard (element-wise) product.

Remark 7.2: Maskensembles

A practical problem of dropout is that for any reasonable choice of dropout probability, the dropout masks $z^{(i)}$ will overlap significantly, which tends to make the predictions $p(y^* | \mathbf{x}^*, \theta^{(i)})$ highly correlated. This can lead to underestimation of epistemic uncertainty. *Maskensembles* (Durasov et al., 2021) mitigate this issue by choosing a fixed set of pre-defined dropout masks, which have controlled overlap, and alternating between them during training. In the extreme case of “infinitely many” masks, maskensembles are equivalent to dropout since each mask is only seen once.

7.3.4 Probabilistic Ensembles

We have seen that variational inference in the context of Bayesian neural networks can be interpreted as averaging the predictions of m neural networks drawn according to the variational posterior.

A natural adaptation of this idea is to immediately learn the weights of m neural networks. The idea is to randomly choose m training sets by sampling uniformly from the data with replacement. Then, using our analysis from Section 7.2.1, we obtain m MAP estimates of the weights $\theta^{(i)}$, yielding the approximation

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) &= \mathbb{E}_{\theta \sim p(\cdot | \mathbf{x}_{1:n}, \mathbf{y}_{1:n})} [p(y^* | \mathbf{x}^*, \theta)] \\ &\approx \frac{1}{m} \sum_{i=1}^m p(y^* | \mathbf{x}^*, \theta^{(i)}). \end{aligned} \quad (7.29)$$

using bootstrapping

Here, the randomly chosen m training sets induce “diversity” of the models. In practice, in the context of deep neural networks where the global minimizer of the loss can rarely be identified, it is common to use the full training data to train each of the m neural networks. Random initialization and random shuffling of the training data is typically enough to ensure some degree of diversity between the individual models (Lakshminarayanan et al., 2017). We can connect ensembles to the other approximate inference techniques we have discussed: First, ensembles can be seen as a specific kind of maskensemble, where the masks are non-overlapping,¹³ which mitigates the issue of correlated predictions from Remark 7.2. Second, ensembling can be combined with other approximate inference techniques such as variational inference, Laplace approximation, or SWAG to get a mixture of Gaussians as the posterior approximation.

Note that Equation (7.29) is not equivalent to Monte Carlo sampling, although it looks very similar. The key difference is that this approach does not sample from the true posterior distribution p , but instead

¹³ That is, the m models do not share any of their parameters. Ensembles and dropout lie on opposite ends of this spectrum.

from the empirical posterior distribution \hat{p} given the (re-sampled) MAP estimates. Intuitively, this can be understood as the difference between sampling from a distribution p directly (Monte Carlo sampling) versus sampling from an approximate (empirical) distribution \hat{p} (corresponding to the training data), which itself is constructed from samples of the true distribution p . This approach is known as *bootstrapping* or *bagging* (short for *bootstrap aggregating*) and plays a central role in model-free reinforcement learning. We will return to this concept in Section 11.4.1.

7.3.5 Diverse Probabilistic Ensembles

Probabilistic ensembles can be loosely interpreted as randomly initializing m “particles” $\{\theta^{(i)}\}_{i=1}^m$ and then pushing each particle towards regions of high posterior probability. A potential issue with this approach is that if the initialization of particles is not sufficiently diverse, the particles may “collapse” since every particle eventually converges to a local optimum of the loss function. A natural approach to mitigate this issue is to alter the objective of each particle from simply aiming to minimize the loss $-\log p(\theta | \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$, which we will abbreviate by $-\log p(\theta)$, to also “push” the particles away from each other. We will see next that this can be interpreted as a form of variational inference under a very flexible variational family.

In our discussion of *variational inference*, we have seen that minimizing reverse-KL is equivalent to maximizing the evidence lower bound, and used this to derive an optimization algorithm to compute approximate posteriors. We will discuss the alternative approach of directly computing the gradient of the KL-divergence. Consider the variational family of all densities that can be expressed as smooth transformations of points sampled from a reference density ϕ . That is, we consider $\mathcal{Q}_\phi \doteq \{\mathbf{T}_\sharp \phi \mid \mathbf{T} : \Theta \rightarrow \Theta \text{ is smooth}\}$ where $\mathbf{T}_\sharp \phi$ is the density of the random vector $\theta' = \mathbf{T}(\theta)$ with $\theta \sim \phi$.¹⁴ The density ϕ can be thought of as the initial distribution of the particles, and the smooth map \mathbf{T} as the dynamics that push the particles towards the target density p . It can be shown that for “almost any” reference density ϕ , this variational family \mathcal{Q}_ϕ is expressive enough to closely approximate “almost arbitrary” distributions.¹⁵ A natural approach is therefore to *learn* the appropriate smooth map \mathbf{T} between the reference density ϕ and the target density p .

¹⁴ Refer back to Section 1.1.11 for a recap on pushforwards.

¹⁵ For a more detailed discussion, refer to “Stein variational gradient descent: A general purpose Bayesian inference algorithm” (Liu and Wang, 2016).

Example 7.3: Gaussian variational inference

We have seen in Section 5.5 that if ϕ is standard normal and $\mathbf{T}(\mathbf{x}; \{\mu, \Sigma^{1/2}\}) = \mu + \Sigma^{1/2}\mathbf{x}$ an affine transformation then the ELBO

can be maximized efficiently using stochastic gradient descent. However, in this case, \mathcal{Q}_ϕ can only approximate Gaussian-like distributions since the expressivity of the map \mathbf{T} is limited under the fixed reference density ϕ .

An alternative approach to Gaussian variational inference is the following algorithm known as *Stein variational gradient descent* (SVGD), where we recursively apply carefully chosen smooth maps to the current variational approximation:

$$q_0 \xrightarrow{\mathbf{T}_0^*} q_1 \xrightarrow{\mathbf{T}_1^*} q_2 \xrightarrow{\mathbf{T}_2^*} \dots \quad \text{where } q_{t+1} \doteq \mathbf{T}_t^* q_t. \quad (7.30)$$

We consider maps $\mathbf{T} = \mathbf{id} + f$ where $\mathbf{id}(\theta) \doteq \theta$ denotes the identity map and $f(x)$ represents a (small) perturbation. Recall that at time t we seek to minimize $\text{KL}(\mathbf{T}_\sharp q_t \| p)$, so we choose the smooth map as

$$\mathbf{T}_t^* \doteq \mathbf{id} - \eta_t \nabla_f \text{KL}(\mathbf{T}_\sharp q_t \| p) \Big|_{f=0} \quad (7.31)$$

where η_t is a step size. In this way, the SVGD update (7.31) can be interpreted as a step of “functional” gradient descent.

To be able to compute the gradient of the KL-divergence, we need to make some structural assumptions on the perturbation f . SVGD assumes that $f = [f_1 \cdots f_d]^\top$ with $f_i \in \mathcal{H}_k(\Theta)$ is from some reproducing kernel Hilbert space $\mathcal{H}_k(\Theta)$ of a positive definite kernel k ; we say $f \in \mathcal{H}_k^d(\Theta)$. Within the RKHS, we can compute the gradient of the KL-divergence exactly. Liu and Wang (2016) show that in this case, the functional gradient of the KL-divergence can be expressed in closed-form as $-\nabla_f \text{KL}(\mathbf{T}_\sharp q \| p) \Big|_{f=0} = \boldsymbol{\varphi}_{q,p}^*$ where

$$\boldsymbol{\varphi}_{q,p}^*(\cdot) \doteq \mathbb{E}_{\theta \sim q}[k(\cdot, \theta) \nabla_\theta \log p(\theta) + \nabla_\theta k(\cdot, \theta)]. \quad (7.32)$$

SVGD then approximates q using the particles $\{\theta^{(i)}\}_{i=1}^m$ as follows:

Algorithm 7.4: Stein variational gradient descent, SVGD

```

1 initialize particles  $\{\theta^{(i)}\}_{i=1}^m$ 
2 repeat
3   for each particle  $i \in [m]$  do
4      $\theta^{(i)} \leftarrow \theta^{(i)} + \eta_t \hat{\boldsymbol{\varphi}}_{q,p}^*(\theta^{(i)})$  where
5      $\hat{\boldsymbol{\varphi}}_{q,p}^*(\theta) \doteq \frac{1}{m} \sum_{j=1}^m [k(\theta, \theta^{(j)}) \nabla_\theta \log p(\theta) + \nabla_\theta k(\theta, \theta^{(j)})]$ 
5 until converged

```

Often, a Gaussian kernel (4.14) with length scale h is used to model the perturbations, in which case the repulsion term is

$$\nabla_{\theta^{(j)}} k(\theta, \theta^{(j)}) = \frac{1}{h^2} (\theta - \theta^{(j)}) k(\theta, \theta^{(j)}) \quad (7.33)$$

and the negative functional gradient simplifies to

$$\hat{\varphi}_{q,p}^*(\theta) = \frac{1}{m} \sum_{j=1}^m k(\theta, \theta^{(j)}) \left[\underbrace{\nabla_\theta \log p(\theta)}_{\text{drift}} + \underbrace{h^{-2}(\theta - \theta^{(j)})}_{\text{repulsion}} \right]. \quad (7.34)$$

Note that SVGD has similarities to Langevin dynamics, which as seen in (?) can also be interpreted as following a gradient of the KL-divergence. Whereas Langevin dynamics perturbs particles according to a drift towards regions of high posterior probability and some random diffusion (cf. Equation (6.46)), the first term of $\hat{\varphi}_{q,p}^*(\theta)$ perturbs particles to drift towards regions of high posterior probability while the second term leads to a mutual “repulsion” of particles. Notably, the perturbations of Langevin dynamics are noisy, while SVGD perturbs particles deterministically and the randomness is exclusively in the initialization of particles. The repulsion term prevents particles from collapsing to a single mode of the posterior distribution, which is a possible failure mode of other particle-based posterior approximations such as ensembles.

Problem 6.9

Note that the above decomposition of $\hat{\varphi}_{q,p}^*(\theta)$ is once more an example of the *principle of curiosity and conformity* which we have seen to be a recurring theme in approaches to approximate inference. The repulsion term leads to exploration of the particles (i.e., “curiosity” about alternative explanations), while the drift term leads to minimization of the loss (i.e., “conformity” to the data).

Lu et al. (2019) show that under some assumptions, SVGD converges asymptotically to the target density p as $\eta_t \rightarrow 0$. SVGD’s name originates from *Stein’s method* which is a general-purpose approach for characterizing convergence in distribution.¹⁶

7.4 Calibration

A key challenge of Bayesian deep learning (and also other probabilistic methods) is the calibration of models. We say that a model is *well-calibrated* if its confidence coincides with its accuracy across many predictions. Consider a classification model that predicts that the label of a given input belongs to some class with probability 80%. If the model is well-calibrated, then the prediction is correct about 80% of the time. In other words, during calibration, we adjust the probability estimation of the model.

¹⁶ For an introduction to Stein’s method, read “Measuring sample quality with Stein’s method” (Gorham and Mackey, 2015).

We will first mention two methods of estimating the calibration of a model, namely the marginal likelihood and reliability diagrams. Then, in Section 7.4.3, we survey commonly used heuristics for empirically improving the calibration.

7.4.1 Evidence

A popular method (which we already encountered multiple times) is to use the evidence of a validation set $x_{1:m}^{\text{val}}$ of size m given the training set $x_{1:n}^{\text{train}}$ of size n for estimating the model calibration. Here, the evidence can be understood as describing how well the validation set is described by the model trained on the training set. We obtain,

$$\begin{aligned} & \log p(y_{1:m}^{\text{val}} | x_{1:m}^{\text{val}}, x_{1:n}^{\text{train}}, y_{1:n}^{\text{train}}) \\ &= \log \int p(y_{1:m}^{\text{val}} | x_{1:m}^{\text{val}}, \theta) p(\theta | x_{1:n}^{\text{train}}, y_{1:n}^{\text{train}}) d\theta \\ &\approx \log \int p(y_{1:m}^{\text{val}} | x_{1:m}^{\text{val}}, \theta) q_\lambda(\theta) d\theta \\ &= \log \int \prod_{i=1}^m p(y_i^{\text{val}} | x_i^{\text{val}}, \theta) q_\lambda(\theta) d\theta \end{aligned} \tag{7.35}$$

using the sum rule (1.7) and product rule (1.11)

approximating with the variational posterior

using the independence of the data

The resulting integrals are typically very small which leads to numerical instabilities. Therefore, it is common to maximize a lower bound to the evidence instead,

$$\begin{aligned} &= \log \mathbb{E}_{\theta \sim q_\lambda} \left[\prod_{i=1}^m p(y_i^{\text{val}} | x_i^{\text{val}}, \theta) \right] \\ &\geq \mathbb{E}_{\theta \sim q_\lambda} \left[\sum_{i=1}^m \log p(y_i^{\text{val}} | x_i^{\text{val}}, \theta) \right] \end{aligned} \tag{7.36}$$

interpreting the integral as an expectation over the variational posterior

using Jensen's inequality (5.29)

$$\approx \frac{1}{k} \sum_{j=1}^k \sum_{i=1}^m \log p(y_i^{\text{val}} | x_i^{\text{val}}, \theta^{(j)}) \tag{7.37}$$

using Monte Carlo sampling

for independent samples $\theta^{(j)} \stackrel{\text{iid}}{\sim} q_\lambda$.

7.4.2 Reliability Diagrams

Reliability diagrams take a frequentist perspective to estimate the calibration of a model. For simplicity, we assume a calibration problem with two classes, 1 and -1 (similarly to logistic regression).¹⁷

We group the predictions of a validation set into M interval bins of size $1/M$ according to the class probability predicted by the model, $\mathbb{P}(Y_i = 1 | x_i)$. We then compare within each bin, how often the model thought the inputs belonged to the class (confidence) with how often the inputs actually belonged to the class (frequency). Formally, we

¹⁷ Reliability diagrams generalize beyond this restricted example.

define B_m as the set of samples falling into bin m and let

$$\text{freq}(B_m) \doteq \frac{1}{|B_m|} \sum_{i \in B_m} \mathbb{1}\{Y_i = 1\} \quad (7.38)$$

be the proportion of samples in bin m that belong to class 1 and let

$$\text{conf}(B_m) \doteq \frac{1}{|B_m|} \sum_{i \in B_m} \mathbb{P}(Y_i = 1 \mid \mathbf{x}_i) \quad (7.39)$$

be the average confidence of samples belonging to class 1 within the bin m .

Thus, a model is well calibrated if $\text{freq}(B_m) \approx \text{conf}(B_m)$ for each bin $m \in [M]$. There are two common metrics of calibration that quantify how “close” a model is to being well calibrated.

1. The *expected calibration error* (ECE) is the average deviation of a model from perfect calibration,

$$\ell_{\text{ECE}} \doteq \sum_{m=1}^M \frac{|B_m|}{n} |\text{freq}(B_m) - \text{conf}(B_m)| \quad (7.40)$$

where n is the size of the validation set.

2. The *maximum calibration error* (MCE) is the maximum deviation of a model from perfect calibration among all bins,

$$\ell_{\text{MCE}} \doteq \max_{m \in [M]} |\text{freq}(B_m) - \text{conf}(B_m)|. \quad (7.41)$$

7.4.3 Heuristics for Improving Calibration

We now survey a few heuristics which can be used empirically to improve model calibration.

1. *Histogram binning* assigns a calibrated score $q_m \doteq \text{freq}(B_m)$ to each bin during validation. Then, during inference, we return the calibrated score q_m of the bin corresponding to the prediction of the model.
2. *Isotonic regression* extends histogram binning by using variable bin boundaries. We find a piecewise constant function $f \doteq [f_1, \dots, f_M]$ that minimizes the bin-wise squared loss,

$$\min_{M, f, a} \sum_{m=1}^M \sum_{i=1}^n \mathbb{1}\{a_m \leq \mathbb{P}(Y_i = 1 \mid \mathbf{x}_i) < a_{m+1}\} (f_m - y_i)^2 \quad (7.42a)$$

$$\text{subject to } 0 = a_1 \leq \dots \leq a_{M+1} = 1, \quad (7.42b)$$

$$f_1 \leq \dots \leq f_M \quad (7.42c)$$

where f are the calibrated scores and $a \doteq [a_1, \dots, a_{M+1}]$ are the bin boundaries. We then return the calibrated score f_m of the bin corresponding to the prediction of the model.

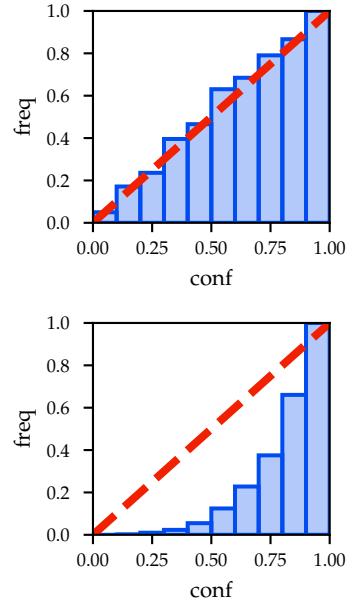


Figure 7.8: Examples of reliability diagrams with ten bins. A perfectly calibrated model approximates the diagonal dashed red line. The first reliability diagram shows a well-calibrated model. In contrast, the second reliability diagram shows an overconfident model.

3. *Platt scaling* adjusts the logits z_i of the output layer to

$$q_i \doteq \sigma(az_i + b) \quad (7.43)$$

and then learns parameters $a, b \in \mathbb{R}$ to maximize the likelihood.

4. *Temperature scaling* is a special and widely used instance of Platt scaling where $a \doteq 1/T$ and $b \doteq 0$ for some temperature scalar $T > 0$,

$$q_i \doteq \sigma\left(\frac{z_i}{T}\right). \quad (7.44)$$

Intuitively, for a larger temperature T , the probability is distributed more evenly among the classes (without changing the ranking), yielding a more uncertain prediction. In contrast, for a lower temperature T , the probability is concentrated more towards the top choices, yielding a less uncertain prediction. As seen in Problem 6.7, temperature scaling can be motivated as tuning the mean of the softmax distribution.

Optional Readings

- Guo, Pleiss, Sun, and Weinberger (2017).
On calibration of modern neural networks.
- Blundell, Cornebise, Kavukcuoglu, and Wierstra (2015).
Weight uncertainty in neural network.
- Kendall and Gal (2017).
What uncertainties do we need in Bayesian deep learning for computer vision?.

Discussion

This chapter concludes our discussion of (approximate) probabilistic inference. Across the last three chapters, we have seen numerous methods for approximating the posterior distributions of deep neural networks:

- Methods such as dropout and stochastic weight averaging are frequently used in practice. Other particle-based approaches such as ensembles and SVGD are used less frequently since they are computationally more expensive to train, but are some of the most effective methods in estimating uncertainty.
- Recently, Laplace approximations regained interest since they can be applied “post-hoc” after training simply by computing or approximating the Hessian of the loss function (Daxberger et al., 2021; Antorán et al., 2022). Still, Laplace approximations come with the limitations inherent to unimodal Gaussian approximations.

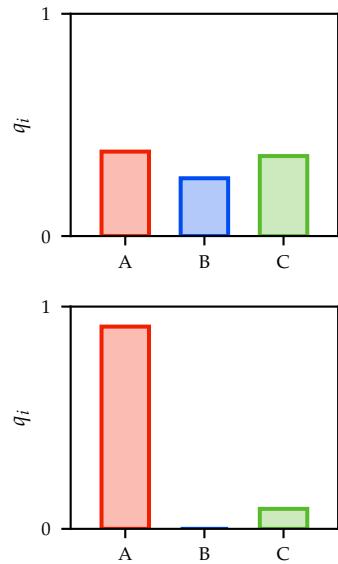


Figure 7.9: Illustration of temperature scaling for a classifier with three classes. On the top, we have a prediction with a high temperature, yielding a very uncertain prediction (in favor of class A). Below, we have a prediction with a low temperature, yielding a prediction that is strongly in favor of class A. Note that the ranking ($A \succ C \succ B$) is preserved.

- Other work, particularly in fine-tuning, has explored approximating the posterior distribution of deep neural networks by treating them as linear functions in a fixed learned feature space, in which case one can use the tools for exact probabilistic inference from Chapters 2 and 4 (e.g., Hübotter et al., 2025).

Despite large progress in approximate inference over the past decade, efficient and reliable uncertainty estimation of large models remains an important open challenge.

Problems

7.1. Softmax is a generalization of the logistic function.

Show that for a two-class classification problem (i.e., $c = 2$), the softmax function is equivalent to the logistic function (5.9) for the univariate model $f \doteq f_1 - f_0$. That is, $\sigma_1(f) = \sigma(f)$ and $\sigma_0(f) = 1 - \sigma(f)$.

Thus, the softmax function is a generalization of the logistic function to more than two classes.

PART II

Sequential Decision-Making

Preface to Part II

In the first part of the manuscript, we have learned about how we can build machines that are capable of updating their beliefs and reducing their epistemic uncertainty through probabilistic inference. We have also discussed ways of keeping track of the world through noisy sensory information by filtering. An important aspect of intelligence is to use this acquired knowledge for making decisions and taking actions that have a positive impact on the world.

Already today, we are surrounded by machines that make decisions and take actions; that is, exhibit some degree of agency. Be it a search engine producing a list of search results, a chatbot answering a question, or a driving-assistance system steering a car: these systems are all perceiving the world, making decisions, and then taking actions that in turn have an effect on the world. Figure 7.10 illustrates this *perception-action loop*.

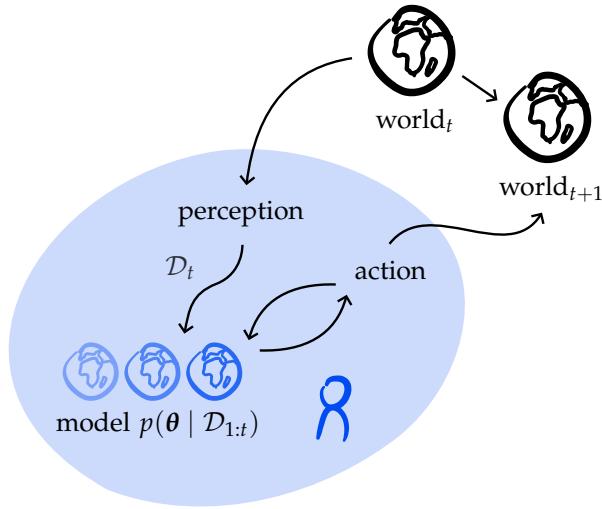


Figure 7.10: An illustration of the perception-action loop. This is a straightforward extension of our view of probabilistic inference from Figure 1.10 with the addition of an “action” component which is capable of “adaptively” interacting with the outside world and the internal world model.

In the second part of this manuscript, we will discuss the underpinning principles of building machines that are capable of making se-

quential decisions. We will see that decision-making itself can be cast as probabilistic inference, obeying the same mechanisms that we used in the first part to build learning systems.

We discuss various ways of addressing the question:

How to act, given that computational resources and time are limited?

One approach is to act with the aim to reduce epistemic uncertainty, which is the topic of active learning. Another approach is to act with the aim to maximize some reward signal, which is the topic of bandits, Bayesian optimization, and reinforcement learning.

This surfaces the exploration-exploitation dilemma where the agent has to prioritize either maximizing its immediate rewards or reducing its uncertainty about the world which might pay off in the future. We discuss that this dilemma is, in fact, in direct correspondence to the principle of curiosity and conformity which we discussed extensively throughout Part I.

Since time is limited, it is critical to be sample-efficient when learning the most important aspects of the world. At the same time, interactions with the world are often complex, and some interactions might even be harmful. We discuss how an agent can use its epistemic uncertainty to guide the exploration of its environment while mitigating risks and reasoning about safety.

Active Learning

By now, we have seen that probabilistic machine learning is very useful for estimating the uncertainty in our models (epistemic uncertainty) and in the data (aleatoric uncertainty). We have been focusing on the setting of supervised learning where we are given a set $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ of labeled data, yet we often encounter settings where we have only little data and acquiring new data is costly.

In this chapter — and in the following chapter on Bayesian optimization — we will discuss how one can use uncertainty to effectively collect more data. In other words, we want to figure out where in the domain we should sample to obtain the most useful information. Throughout most of this chapter, we focus on the most common way of quantifying “useful information”, namely the expected reduction in entropy which is also called the *mutual information*.

8.1 Conditional Entropy

We begin by introducing the notion of conditional entropy. Recall that the entropy $H[\mathbf{X}]$ of a random vector \mathbf{X} can be interpreted as our average surprise when observing realizations $\mathbf{x} \sim \mathbf{X}$. Thus, entropy can be considered as a quantification of the uncertainty about a random vector (or equivalently, its distribution).¹

A natural extension is to consider the entropy of \mathbf{X} given the occurrence of an event corresponding to another random variable (e.g., $\mathbf{Y} = \mathbf{y}$ for a random vector \mathbf{Y}),

$$H[\mathbf{X} \mid \mathbf{Y} = \mathbf{y}] \doteq \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{y})}[-\log p(\mathbf{x} \mid \mathbf{y})]. \quad (8.1)$$

Instead of averaging over the surprise of samples from the distribution $p(\mathbf{x})$ (like the entropy $H[\mathbf{X}]$), this quantity simply averages over the surprise of samples from the conditional distribution $p(\mathbf{x} \mid \mathbf{y})$.

¹ We discussed entropy extensively in Section 5.4.

Definition 8.1 (Conditional entropy). The *conditional entropy* of a random vector \mathbf{X} given the random vector \mathbf{Y} is defined as

$$H[\mathbf{X} \mid \mathbf{Y}] \doteq \mathbb{E}_{y \sim p(y)}[H[\mathbf{X} \mid \mathbf{Y} = y]] \quad (8.2)$$

$$= \mathbb{E}_{(x,y) \sim p(x,y)}[-\log p(x \mid y)]. \quad (8.3)$$

Intuitively, the conditional entropy of \mathbf{X} given \mathbf{Y} describes our average surprise about realizations of \mathbf{X} given a particular realization of \mathbf{Y} , averaged over all such possible realizations of \mathbf{Y} . In other words, conditional entropy corresponds to the expected remaining uncertainty in \mathbf{X} after we observe \mathbf{Y} . Note that, in general, $H[\mathbf{X} \mid \mathbf{Y}] \neq H[\mathbf{Y} \mid \mathbf{X}]$.

It is crucial to stress the difference between $H[\mathbf{X} \mid \mathbf{Y} = y]$ and the conditional entropy $H[\mathbf{X} \mid \mathbf{Y}]$. The former simply corresponds to a probabilistic update of our uncertainty in \mathbf{X} after we have observed the realization $y \sim \mathbf{Y}$. In contrast, conditional entropy *predicts* how much uncertainty will remain about \mathbf{X} (in expectation) after we *will* observe \mathbf{Y} .

Definition 8.2 (Joint entropy). One can also define the *joint entropy* of random vectors \mathbf{X} and \mathbf{Y} ,

$$H[\mathbf{X}, \mathbf{Y}] \doteq \mathbb{E}_{(x,y) \sim p(x,y)}[-\log p(x, y)], \quad (8.4)$$

as the combined uncertainty about \mathbf{X} and \mathbf{Y} . Observe that joint entropy is symmetric.

This gives the *chain rule for entropy*,

$$H[\mathbf{X}, \mathbf{Y}] = H[\mathbf{Y}] + H[\mathbf{X} \mid \mathbf{Y}] \quad (8.5)$$

$$= H[\mathbf{X}] + H[\mathbf{Y} \mid \mathbf{X}]. \quad (8.6)$$

using the product rule (1.11) and the definition of conditional entropy (8.2)
using symmetry of joint entropy

That is, the joint entropy of \mathbf{X} and \mathbf{Y} is given by the uncertainty about \mathbf{X} and the additional uncertainty about \mathbf{Y} given \mathbf{X} . Moreover, this also yields *Bayes' rule for entropy*,

$$H[\mathbf{X} \mid \mathbf{Y}] = H[\mathbf{Y} \mid \mathbf{X}] + H[\mathbf{X}] - H[\mathbf{Y}]. \quad (8.7)$$

using the chain rule for entropy (8.5)
twice

A very intuitive property of entropy is its monotonicity: when conditioning on additional observations the entropy can never increase,

$$H[\mathbf{X} \mid \mathbf{Y}] \leq H[\mathbf{X}]. \quad (8.8)$$

Colloquially, this property is also called the “*information never hurts*” principle. We will derive a proof in the following section.

8.2 Mutual Information

Recall that our fundamental objective is to reduce entropy, as this corresponds to reduced uncertainty in the variables, which we want to predict. Thus, we are interested in how much information we “gain” about the random vector \mathbf{X} by choosing to observe a random vector \mathbf{Y} . By our interpretation of conditional entropy from the previous section, this is described by the following quantity.

Definition 8.3 (Mutual information, MI). The *mutual information* of \mathbf{X} and \mathbf{Y} (also known as the *information gain*) is defined as

$$I(\mathbf{X}; \mathbf{Y}) \doteq H[\mathbf{X}] - H[\mathbf{X} | \mathbf{Y}] \quad (8.9)$$

$$= H[\mathbf{X}] + H[\mathbf{Y}] - H[\mathbf{X}, \mathbf{Y}]. \quad (8.10)$$

In words, we subtract the uncertainty left about \mathbf{X} after observing \mathbf{Y} from our initial uncertainty about \mathbf{X} . This measures the reduction in our uncertainty in \mathbf{X} (as measured by entropy) upon observing \mathbf{Y} . Unlike conditional entropy, it follows from the definition that mutual information is symmetric. That is,

$$I(\mathbf{X}; \mathbf{Y}) = I(\mathbf{Y}; \mathbf{X}). \quad (8.11)$$

Thus, the mutual information between \mathbf{X} and \mathbf{Y} can be understood as the approximation error (or information loss) when assuming that \mathbf{X} and \mathbf{Y} are independent.

In particular, using Gibbs’ inequality (cf. Problem 5.5), this relationship shows that $I(\mathbf{X}; \mathbf{Y}) \geq 0$ with equality when \mathbf{X} and \mathbf{Y} are independent, and also proves the *information never hurts* principle (8.8) as

$$0 \leq I(\mathbf{X}; \mathbf{Y}) = H[\mathbf{X}] - H[\mathbf{X} | \mathbf{Y}]. \quad (8.12)$$

Example 8.4: Mutual information of Gaussians

Given the Gaussian random vector $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and the noisy observation $\mathbf{Y} = \mathbf{X} + \boldsymbol{\varepsilon}$ where $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I})$, we want to find the information gain of \mathbf{X} when observing \mathbf{Y} . Using our definitions from this chapter, we obtain

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y}) &= I(\mathbf{Y}; \mathbf{X}) \\ &= H[\mathbf{Y}] - H[\mathbf{Y} | \mathbf{X}] \\ &= H[\mathbf{Y}] - H[\boldsymbol{\varepsilon}] \\ &= \frac{1}{2} \log((2\pi e)^d \det(\boldsymbol{\Sigma} + \sigma_n^2 \mathbf{I})) - \frac{1}{2} \log((2\pi e)^d \det(\sigma_n^2 \mathbf{I})) \end{aligned}$$

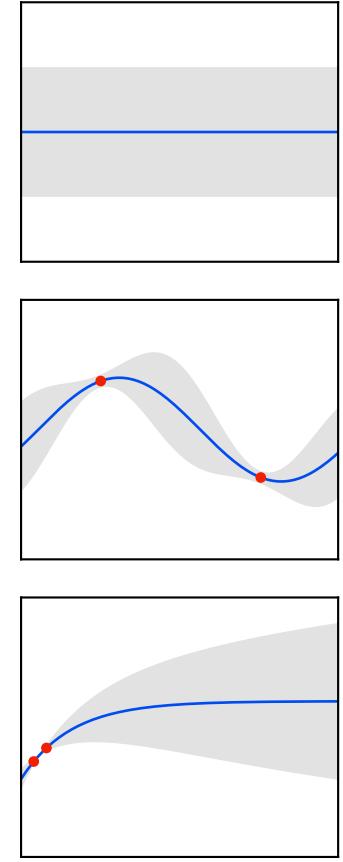


Figure 8.1: Information gain. The first graph shows the prior. The second graph shows a selection of samples with large information gain (large uncertainty reduction). The third graph shows a selection of samples with small information gain (small uncertainty reduction).

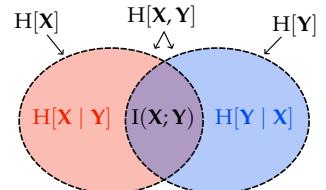


Figure 8.2: Relationship between mutual information and entropy, expressed as a Venn diagram.

using symmetry (8.11)

by mutual information (8.9)

given \mathbf{X} , the only randomness in \mathbf{Y} originates from $\boldsymbol{\varepsilon}$

using the entropy of Gaussians (5.31)

$$\begin{aligned}
&= \frac{1}{2} \log \frac{\det(\boldsymbol{\Sigma} + \sigma_n^2 \mathbf{I})}{\det(\sigma_n^2 \mathbf{I})} \\
&= \frac{1}{2} \log \det(\mathbf{I} + \sigma_n^{-2} \boldsymbol{\Sigma}). \tag{8.13}
\end{aligned}$$

Intuitively, the larger the noise σ_n^2 in relation to the covariance of \mathbf{X} , the smaller the information gain.

8.2.1 Synergy and Redundancy

It is sometimes useful to write down the mutual information of \mathbf{X} and \mathbf{Y} conditioned (in expectation) on a third random vector \mathbf{Z} . This leads us to the following definition.

Definition 8.5 (Conditional mutual information). The *conditional mutual information* of \mathbf{X} and \mathbf{Y} given \mathbf{Z} is defined as

$$I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) \doteq H[\mathbf{X} | \mathbf{Z}] - H[\mathbf{X} | \mathbf{Y}, \mathbf{Z}]. \tag{8.14}$$

$$= H[\mathbf{X}, \mathbf{Z}] + H[\mathbf{Y}, \mathbf{Z}] - H[\mathbf{Z}] - H[\mathbf{X}, \mathbf{Y}, \mathbf{Z}] \tag{8.15}$$

$$= I(\mathbf{X}; \mathbf{Y}, \mathbf{Z}) - I(\mathbf{X}; \mathbf{Z}). \tag{8.16}$$

using the relationship of joint and conditional entropy (8.5)

Thus, the conditional mutual information corresponds to the reduction of uncertainty in \mathbf{X} when observing \mathbf{Y} , given we already observed \mathbf{Z} . It also follows that conditional mutual information is symmetric:

$$I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) = I(\mathbf{Y}; \mathbf{X} | \mathbf{Z}). \tag{8.17}$$

We have seen in this chapter that entropy is monotonically decreasing as we condition on new information, and called this the “information never hurts” principle (8.8). However, the same does not hold for mutual information! That is, information about a random vector \mathbf{Z} may reduce the mutual information between random vectors \mathbf{X} and \mathbf{Y} (??).

Problem 8.2

Remark 8.6: Sufficient statistics and data processing inequality

A related concept is the data processing inequality (8.42) which you prove in (??) and which allows us to formalize a concept which we have seen multiple times already, namely the notion of a sufficient statistic. Consider the Markov chain $\lambda \rightarrow \mathbf{X} \rightarrow s(\mathbf{X})$, for example, λ may be parameters of the distribution of \mathbf{X} . By the data processing inequality (8.42), $I(\lambda; s(\mathbf{X})) \leq I(\lambda; \mathbf{X})$. If the data processing inequality is satisfied with equality then $s(\mathbf{X})$ is called a *sufficient statistic* of \mathbf{X} for the inference of λ .

Problem 8.2 (3)

To understand the behavior of mutual information under conditioning, it is helpful to consider the *interaction information*

$$I(\mathbf{X}; \mathbf{Y}; \mathbf{Z}) \doteq I(\mathbf{X}; \mathbf{Y}) - I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}). \quad (8.18)$$

If the interaction is positive then some information about \mathbf{X} that is provided by \mathbf{Y} is also provided by \mathbf{Z} (i.e., conditioning on \mathbf{Z} decreases MI between \mathbf{X} and \mathbf{Y}), and we say that there is *redundancy* between \mathbf{Y} and \mathbf{Z} (with respect to \mathbf{X}). Conversely, if the interaction is negative then learning about \mathbf{Z} increases what can be learned from \mathbf{Y} about \mathbf{X} , and we say that there is *synergy* between \mathbf{Y} and \mathbf{Z} . We will see later in this chapter that the absence of synergies can lead to efficient algorithms for maximizing mutual information.

8.2.2 Mutual Information as Utility Function

Following our introduction of mutual information, it is natural to answer the question “where should I collect data?” by saying “wherever mutual information is maximized”. More concretely, assume we are given a set \mathcal{X} of possible observations of f , where y_x denotes a single such observation at $x \in \mathcal{X}$,

$$y_x \doteq f_x + \varepsilon_x, \quad (8.19)$$

$f_x \doteq f(x)$, and ε_x is some zero-mean Gaussian noise. For a set of observations $S = \{x_1, \dots, x_n\}$, we can write $\mathbf{y}_S = f_S + \boldsymbol{\varepsilon}$ where

$$\mathbf{y}_S \doteq \begin{bmatrix} y_{x_1} \\ \vdots \\ y_{x_n} \end{bmatrix}, \quad f_S \doteq \begin{bmatrix} f_{x_1} \\ \vdots \\ f_{x_n} \end{bmatrix}, \quad \text{and} \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I}).$$

Note that both \mathbf{y}_S and f_S are random vectors. Our goal is then to find a subset $S \subseteq \mathcal{X}$ of size n maximizing the information gain between our model f and \mathbf{y}_S .

This yields the maximization objective,

$$I(S) \doteq I(f_S; \mathbf{y}_S) = H[f_S] - H[f_S | \mathbf{y}_S]. \quad (8.20)$$

Here, $H[f_S]$ corresponds to the uncertainty about f_S before obtaining the observations \mathbf{y}_S and $H[f_S | \mathbf{y}_S]$ corresponds to the uncertainty about f_S , in expectation, after obtaining the observations \mathbf{y}_S .

Remark 8.7: Making optimal decisions with intrinsic rewards

Note that this objective function maps neatly onto our initial consideration of making optimal decisions under uncertainty from

Section 1.4. In fact, you can think of maximizing $I(S)$ simply as computing the optimal decision rule for the utility

$$r(\mathbf{y}_S, S) = H[f_S] - H[f_S \mid \mathbf{Y}_S = \mathbf{y}_S], \quad (8.21)$$

with $I(S) = \mathbb{E}_{\mathbf{y}_S}[r(\mathbf{y}_S, S)]$ measuring the expected utility of observations \mathbf{y}_S . Such a utility or reward function is often called an *intrinsic reward* since it does not measure an “objective” external quantity, but instead a “subjective” quantity that is internal to the model of f .

Observe that picking a subset of points $S \subseteq \mathcal{X}$ from the domain \mathcal{X} is a combinatorial problem. That is to say, we are optimizing a function over discrete sets. In general, such combinatorial optimization problems tend to be very difficult. It can be shown that maximizing mutual information is \mathcal{NP} -hard.

8.3 Submodularity of Mutual Information

We will look at optimizing mutual information in the following section. First, we want to introduce the notion of submodularity which is important in the analysis of discrete functions.

Definition 8.8 (Marginal gain). Given a (discrete) function $F : \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}$, the *marginal gain* of $x \in \mathcal{X}$ given $A \subseteq \mathcal{X}$ is defined as

$$\Delta_F(x \mid A) \doteq F(A \cup \{x\}) - F(A). \quad (8.22)$$

Intuitively, the marginal gain describes how much “adding” the additional x to A increases the value of F .

When maximizing mutual information, the marginal gain is $\textcircled{?}$

with $\mathbf{Y}_S = \mathbf{y}_S$ denoting an event

Problem 8.4

$$\Delta_I(x \mid A) = I(f_x; y_x \mid \mathbf{y}_A) \quad (8.23)$$

$$= H[y_x \mid \mathbf{y}_A] - H[\varepsilon_x]. \quad (8.24)$$

That is, when maximizing mutual information, the marginal gain corresponds to the difference between the uncertainty after observing y_A and the entropy of the noise $H[\varepsilon_x]$. Altogether, the marginal gain represents the reduction in uncertainty by observing $\{x\}$.

Definition 8.9 (Submodularity). A (discrete) function $F : \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}$ is *submodular* iff for any $x \in \mathcal{X}$ and any $A \subseteq B \subseteq \mathcal{X}$ it satisfies

$$F(A \cup \{x\}) - F(A) \geq F(B \cup \{x\}) - F(B). \quad (8.25)$$

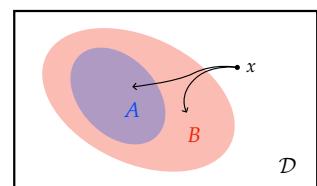


Figure 8.3: Monotone submodularity. The effect of “adding” x to the smaller set A is larger than the effect of adding x to the larger set B .

Equivalently, using our definition of marginal gain, we have that F is submodular iff for any $x \in \mathcal{X}$ and any $A \subseteq B \subseteq \mathcal{X}$,

$$\Delta_F(x | A) \geq \Delta_F(x | B). \quad (8.26)$$

That is, “adding” x to the smaller set A yields more marginal gain than adding x to the larger set B . In other words, the function F has “diminishing returns”. In this way, submodularity can be interpreted as a notion of “concavity” for discrete functions.

Definition 8.10 (Monotone submodularity). A function $F : \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}$ is called *monotone* iff for any $A \subseteq B \subseteq \mathcal{X}$ it satisfies

$$F(A) \leq F(B). \quad (8.27)$$

If F is also submodular, then F is called *monotone submodular*.

Theorem 8.11. *The objective I is monotone submodular.*

Proof. We fix arbitrary subsets $A \subseteq B \subseteq \mathcal{X}$ and any $x \in \mathcal{X}$. We have,

$$\begin{aligned} I \text{ is submodular} &\iff \Delta_I(x | A) \geq \Delta_I(x | B) \\ &\iff H[y_x | y_A] - H[\varepsilon_x] \geq H[y_x | y_B] - H[\varepsilon_x] \\ &\iff H[y_x | y_A] \geq H[y_x | y_B]. \end{aligned}$$

by submodularity in terms of marginal gain (8.26)

using Equation (8.24)

$H[\varepsilon_x]$ cancels

Due to the “information never hurts” principle (8.8) of entropy and as $A \subseteq B$, this is always true. Moreover,

$$\begin{aligned} I \text{ is monotone} &\iff I(A) \leq I(B) \\ &\iff I(f_A; y_A) \leq I(f_B; y_B) \\ &\iff I(f_B; y_A) \leq I(f_B; y_B) \\ &\iff H[f_B] - H[f_B | y_A] \leq H[f_B] - H[f_B | y_B] \\ &\iff H[f_B | y_A] \geq H[f_B | y_B], \end{aligned}$$

by the definition of monotonicity (8.27)

using the definition of I (8.20)

using $I(f_B; y_A) = I(f_A; y_A)$ as $y_A \perp f_B | f_A$

using the definition of MI (8.9)

$H[f_B]$ cancels

which is also satisfied due to the “information never hurts” principle (8.8). \square

The submodularity of I can be interpreted from the perspective of information theory. It turns out that submodularity is equivalent to the absence of synergy between observations $\textcircled{?}$. Intuitively, without synergies, acting greedily is enough to find a near-optimal solution. If there are synergies, then the combinatorial search problem is much harder, because single-step optimal actions do not necessarily lead us to the optimal solution. Consider the extreme case of having to solve a “needle in a haystack” problem, where only a single subset of \mathcal{X} with size k achieves objective value 1, with all other subsets achieving objective value 0. In this case, we can do nothing but exhaustively search through all $|\mathcal{X}|^k$ combinations to find the optimal solution.

Problem 8.5

8.4 Maximizing Mutual Information

As we cannot efficiently pick a set $S \subseteq \mathcal{X}$ to maximize mutual information but know that I is submodular, a natural approach is to maximize mutual information greedily. That is, we pick the locations x_1 through x_n individually by greedily finding the location with the maximal mutual information. The following general result for monotone submodular function maximization shows that, indeed, this greedy approach provides a good approximation.

Theorem 8.12 (Greedy submodular function maximization). *If the set function $F : \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}_{\geq 0}$ is monotone submodular, then greedily maximizing F is a $(1 - 1/e)$ -approximation:*²

$$F(S_n) \geq \left(1 - \frac{1}{e}\right) \max_{\substack{S \subseteq \mathcal{X} \\ |S|=n}} F(S). \quad (8.28)$$

Proof. Fix any $n \geq 1$. Let $S^* \in \arg \max \{F(S) \mid S \in \mathcal{X}, |S| \leq n\}$. We can assume $|S^*| = n$ due to the monotonicity (8.27) of F . We write, $\{x_1^*, \dots, x_n^*\} \doteq S^*$. We have,

$$\begin{aligned} F(S^*) &\leq F(S^* \cup S_t) \\ &= F(S_t) + \sum_{i=1}^n \Delta_F(x_i^* \mid S_t \cup \{x_1^*, \dots, x_{i-1}^*\}) \\ &\leq F(S_t) + \sum_{x^* \in S^*} \Delta_F(x^* \mid S_t) \\ &\leq F(S_t) + n(F(S_{t+1}) - F(S_t)). \end{aligned}$$

Let $\delta_t \doteq F(S^*) - F(S_t)$. Then,

$$\delta_t = F(S^*) - F(S_t) \leq n(F(S_{t+1}) - F(S_t)) = n(\delta_t - \delta_{t+1}).$$

This implies $\delta_{t+1} \leq (1 - 1/n)\delta_t$ and $\delta_n \leq (1 - 1/n)^n \delta_0 \leq \delta_0/e$, using the well-known inequality $1 - x \leq e^{-x}$ for all $x \in \mathbb{R}$.

Finally, observe that $\delta_0 = F(S^*) - F(\emptyset) \leq F(S^*)$ due to the non-negativity of F . We obtain,

$$\delta_n = F(S^*) - F(S_n) \leq \frac{\delta_0}{e} \leq \frac{F(S^*)}{e}.$$

Rearranging the terms yields the theorem. \square

Optional Readings

The original proof of greedy maximization for submodular functions was given by “An analysis of approximations for maximiz-

² $1 - 1/e \approx 0.632$

using monotonicity (8.27)

using the definition of marginal gain (8.22)

using submodularity (8.26)

using that $S_{t+1} = S_t \cup \{x\}$ is chosen such that $\Delta_F(x \mid S_t)$ is maximized (8.29)

ing submodular set functions" (Nemhauser et al., 1978).

For more background on maximizing submodular functions, see "Submodular function maximization" (Krause and Golovin, 2014).

Now that we have established that greedy maximization of mutual information is a decent approximation to maximizing the joint information of data, we will look at how this optimization problem can be solved in practice.

8.4.1 Uncertainty Sampling

When maximizing mutual information, at time t when we have already picked $S_t = \{x_1, \dots, x_t\}$, we need to solve the following optimization problem,

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} \Delta_I(x \mid S_t) \quad (8.29)$$

$$= \arg \max_{x \in \mathcal{X}} I(f_x; y_x \mid y_{S_t}). \quad (8.30) \quad \text{using Equation (8.23)}$$

Note that f_x and y_x are univariate random variables. Thus, using our formula for the mutual information of conditional linear Gaussians (8.13), we can simplify to,

$$= \arg \max_{x \in \mathcal{X}} \frac{1}{2} \log \left(1 + \frac{\sigma_t^2(x)}{\sigma_n^2} \right) \quad (8.31)$$

where $\sigma_t^2(x)$ is the (remaining) variance at x after observing S_t . Assuming the label noise is independent of x (i.e., homoscedastic),

$$= \arg \max_{x \in \mathcal{X}} \sigma_t^2(x). \quad (8.32)$$

Therefore, if f is modeled by a Gaussian and we assume homoscedastic noise, greedily maximizing mutual information corresponds to simply picking the point x with the largest variance. This algorithm is also called *uncertainty sampling*.

8.4.2 Heteroscedastic Noise

Uncertainty sampling is clearly problematic if the noise is heteroscedastic. If there are a particular set of inputs with a large aleatoric uncertainty dominating the epistemic uncertainty, uncertainty sampling will continuously choose those points even though the epistemic uncertainty will not be reduced substantially (cf. Figure 8.4).

Looking at Equation (8.31) suggests a natural fix. Instead of only considering the epistemic uncertainty $\sigma_t^2(x)$, we can also consider the ale-

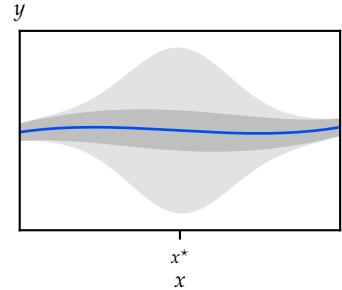


Figure 8.4: Uncertainty sampling with heteroscedastic noise. The epistemic uncertainty of the model is shown in a dark gray. The aleatoric uncertainty of the data is shown in a light gray. Uncertainty sampling would repeatedly pick points around x^* as they maximize the epistemic uncertainty, even though the aleatoric uncertainty at x^* is much larger than at the boundary.

atoric uncertainty $\sigma_n^2(x)$ by modeling heteroscedastic noise, yielding

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} \frac{1}{2} \log \left(1 + \frac{\sigma_t^2(x)}{\sigma_n^2(x)} \right) = \arg \max_{x \in \mathcal{X}} \frac{\sigma_t^2(x)}{\sigma_n^2(x)}. \quad (8.33)$$

Thus, we choose locations that trade large epistemic uncertainty with large aleatoric uncertainty. Ideally, we find a location where the epistemic uncertainty is large, and the aleatoric uncertainty is low, which promises a significant reduction of uncertainty around this location.

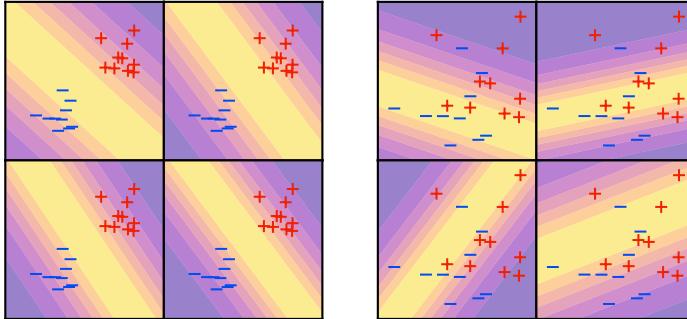
8.4.3 Classification

While we focused on regression, one can apply active learning also for other settings, such as (probabilistic) classification. In this setting, for any input x , a model produces a categorical distribution over labels y_x .³ Here, uncertainty sampling corresponds to selecting samples that maximize the entropy of the predicted label y_x ,

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} H[y_x | x_{1:t}, y_{1:t}]. \quad (8.34)$$

The entropy of a categorical distribution is simply a finite sum of weighted surprise terms.

³ see Section 1.3



This approach generally leads to sampling points that are close to the decision boundary. Often, the uncertainty is mainly dominated by label noise rather than epistemic uncertainty, and hence, we do not learn much from our observations. This is a similar problem to the one we encountered with uncertainty sampling in the setting of heteroscedastic noise.

This naturally suggests distinguishing between the aleatoric and epistemic uncertainty of the model f (parameterized by θ). To this end, mutual information can be used similarly as we have done with uncertainty sampling for regression,

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} I(\theta; y_x | x_{1:t}, y_{1:t}) \quad (8.35)$$

Figure 8.5: Uncertainty sampling in classification. The area with high uncertainty (as measured by entropy) is highlighted in yellow. The shown figures display each a sequence of model updates, each after one new observation. In the left figure, the classes are well-separated and uncertainty is dominated by epistemic uncertainty, whereas in the right figure the uncertainty is dominated by noise. In the latter case, if we mostly choose points x in the area of highest uncertainty (i.e., close to the decision boundary) to make observations, the label noise results in frequently changing models.

$$\begin{aligned}
&= \arg \max_{x \in \mathcal{X}} I(y_x; \theta | x_{1:t}, y_{1:t}) \\
&= \arg \max_{x \in \mathcal{X}} H[y_x | x_{1:t}, y_{1:t}] - H[y_x | \theta, x_{1:t}, y_{1:t}] \\
&= \arg \max_{x \in \mathcal{X}} H[y_x | x_{1:t}, y_{1:t}] - \mathbb{E}_{\theta|x_{1:t}, y_{1:t}} [H[y_x | \theta, x_{1:t}, y_{1:t}]] \quad (8.36) \\
&= \arg \max_{x \in \mathcal{X}} \underbrace{H[y_x | x_{1:t}, y_{1:t}]}_{\text{entropy of predictive posterior}} - \mathbb{E}_{\theta|x_{1:t}, y_{1:t}} \underbrace{H[y_x | \theta]}_{\text{entropy of likelihood}}. \quad (8.37)
\end{aligned}$$

using symmetry (8.11)

using the definition of mutual information (8.9)

using the definition of conditional entropy (8.2)

using the definition of entropy (5.27)
and assuming $y_x \perp x_{1:t}, y_{1:t} | \theta$

The first term measures the entropy of the averaged prediction while the second term measures the average entropy of predictions. Thus, the first term looks for points where the average prediction is not confident. In contrast, the second term penalizes points where many of the sampled models are not confident about their prediction, and thus looks for points where the models are confident in their predictions. This identifies those points x where the models *disagree* about the label y_x (that is, each model is “confident” but the models predict different labels). For this reason, this approach is known as *Bayesian active learning by disagreement* (BALD).

Note that the second term of the difference acts as a regularizer when compared to Equation (8.34). The second term mirrors our description of [aleatoric uncertainty](#) from Section 2.2. Recall that we interpreted aleatoric uncertainty as the average uncertainty for all models. Crucially, here we use entropy to “measure” uncertainty, whereas previously we have been using variance. Therefore, intuitively, Equation (8.36) subtracts the aleatoric uncertainty from the total uncertainty about the label.

Observe that both terms require approximate forms of the posterior distribution. In Chapters 5 and 6, we have seen various approaches from variational inference and MCMC methods which can be used here. The first term can be approximated by the predictive distribution of an approximated posterior which is obtained, for example, using variational inference. The nested entropy of the second term is typically easy to compute, as it corresponds to the entropy of the (discrete) likelihood distribution $p(y | x, \theta)$ of the model θ . The outer expectation of the second term may be approximated using (approximate) samples from the posterior distribution obtained via variational inference, MCMC, or some other method.

Optional Readings

- Gal, Islam, and Ghahramani (2017).
Deep Bayesian active learning with image data.

8.5 Learning Locally: Transductive Active Learning

So far we have explored how to select observations that provide us with the best predictor $f(\mathbf{x})$ across the entire domain $\mathbf{x} \in \mathcal{X}$. However, we typically utilize predictors to make predictions at a particular location \mathbf{x}^* . It is therefore a natural question to ask how to select observations that lead to the best individual prediction $f(\mathbf{x}^*)$ at the prediction target \mathbf{x}^* . The distinction between these two settings is closely related to the distinction between two general approaches to learning: *inductive learning* and *transductive learning*. Inductive learning aims to extract general rules from the data, that is, to extract and compress the most bits of information. In contrast, transductive learning aims to make the best prediction at a particular location, that is, to extract the most bits of information *relevant to the prediction* at \mathbf{x}^* . The concept of transduction was developed by Vladimir Vapnik in the 1980s who described its essence as follows:

“When solving a problem of interest, do not solve a more general problem as an intermediate step. Try to get the answer that you really need but not a more general one.” — Vladimir Vapnik

Remark 8.13: What are the prediction targets \mathbf{x}^* ?

Typically, in transductive learning, we cannot directly observe the value $f(\mathbf{x}^*)$ at the prediction target \mathbf{x}^* , for example, when learning from a fixed dataset $\mathcal{X}' \subset \mathcal{X}$. If we *could* observe $f(\mathbf{x}^*)$ directly (or perturbed by noise), solving the learning task would only require memorization. Instead, most interesting learning tasks require generalizing $f(\mathbf{x}^*)$ from the behavior of f at other locations. Therefore, transductive learning becomes interesting precisely when we cannot directly observe $f(\mathbf{x}^*)$.⁴

Note that this is similar to the inductive setting where we, in principle, could observe $f(\mathbf{x})$ at any location \mathbf{x} , but practically, we can only observe $f(\mathbf{x})$ at a finite number of locations. Since such an inductive model is then used to make predictions at any location \mathbf{x}^* , it also needs to generalize from the observations.

⁴ We will discuss an example of this kind in Problem 8.6.

Following the transductive paradigm, when already knowing that our goal is to predict $f(\mathbf{x}^*)$ our objective is to select observations that provide the most information about $f(\mathbf{x}^*)$. We can express this objective elegantly using the probabilistic framework from this chapter (MacKay, 1992; Hübotter et al., 2024):

$$\mathbf{x}_{t+1} \doteq \arg \max_{\mathbf{x} \in \mathcal{X}'} I(f_{\mathbf{x}^*}; y_{\mathbf{x}} \mid \mathbf{x}_{1:t}, y_{1:t}) = \arg \min_{\mathbf{x} \in \mathcal{X}'} H[f_{\mathbf{x}^*} \mid \mathbf{x}_{1:t}, y_{1:t}, y_{\mathbf{x}}]. \quad (8.38)$$

Hübotter et al. (2024) show that this objective leads to a remarkably different selection of observations compared to the inductive active learning objective we discussed earlier. Indeed, while the inductive objective focuses on selecting a *diverse* set of examples, the transductive objective also takes into account the *relevance* of the examples to the prediction target x^* . We can see this tradeoff between diversity and relevance by rewriting the transductive objective as

$$I(f_{x^*}; y_x | x_{1:t}, y_{1:t}) = \underbrace{I(f_{x^*}; y_x)}_{\text{relevance}} - \underbrace{I(f_{x^*}; y_x; y_{1:t} | x_{1:t})}_{\text{redundancy}} \quad (8.39)$$

using the definition of interaction information (8.18)

where the first term measures the information gain of y_x about f_{x^*} while the second term is the interaction information which measures the redundancy of the information in y_x and $y_{1:t}$ about f_{x^*} . In this way, transductive active learning describes a middle ground between traditional search and retrieval methods that focus on relevance on the one hand and inductive active learning which focuses on diversity on the other hand.

Optional Readings

- Hübotter, Sukhija, Treven, As, and Krause (2024).
Transductive active learning: Theory and applications.
- In modern machine learning, one often differentiates between a “pre-training” and a “fine-tuning” stage. During pre-training, a model is trained on a large dataset to extract general knowledge without a specific task in mind. Then, during fine-tuning, the model is adapted to a specific task by training on a smaller dataset. Whereas (inductive) active learning is closely linked to the pre-training stage, transductive active learning has been shown to be useful for task-specific fine-tuning:
- Hübotter, Bongni, Hakimi, and Krause (2025).
Efficiently Learning at Test-Time: Active Fine-Tuning of LLMs.
- Bagatella, Hübotter, Martius, and Krause (2024).
Active Fine-Tuning of Generalist Policies.

Discussion

We have discussed how to select the most informative data. Thereby, we focused mostly on *inductive* active learning which is applicable to “pre-training” when we aim to extract general knowledge from data, but also explored *transductive* active learning which is useful for “fine-tuning” when we aim to adapt a model to a specific task.

We focused on quantifying the “informativeness” of data by its information gain, which is a common approach, though many other viable

criteria exist:

Remark 8.14: Beyond mutual information

The problem of identifying which experiments to conduct in order to maximize learning is studied extensively in the field of *experimental design* where a set of observations S is commonly called a *design*. In the presence of a prior and likelihood, and a different possible posterior distribution for each design S , the field is also called *Bayesian experimental design* (Chaloner and Verdinelli, 1995; Rainforth et al., 2024; Mutný, 2024).

As we highlighted in the beginning of this chapter, how we measure the utility (i.e., the informativeness) of a design S is crucial. Such a measure is called a *design criterion* and a design which is optimal with respect to a design criterion is called an *optimal design*. The literature studies various design criteria beyond maximizing mutual information (i.e., minimizing posterior entropy). One popular alternative is to select the observations S that minimize the trace of the posterior covariance matrix, which corresponds to minimizing the posterior average variance.

Next, we will move to the topic of optimization and ask which data we should select to find the optimum of an unknown function as quickly as possible. In the following chapter, we will focus on “Bayesian optimization” (also called “bandit optimization”) where our aim is to *find and sample* the optimal point. A related task that is slightly more related to active learning is the “best-arm identification” problem where we aim only to *identify* the optimal point without sampling it. This problem is closely related transductive active learning (with the local task being defined by the location of the maximum) and so-called *entropy search* methods that minimize the entropy of the posterior distribution over the location or value of the maximum (akin to Equation (8.38)) are often used to solve this problem (Hennig and Schuler, 2012; Wang and Jegelka, 2017; Hvarfner et al., 2022).

Problems

8.1. Mutual information and KL divergence.

Show that by expanding the definition of mutual information,

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y}) &= \mathbb{E}_{(x,y) \sim p} \left[\log \frac{p(x,y)}{p(x)p(y)} \right] \\ &= \text{KL}(p(x,y) \| p(x)p(y)) \quad (8.40) \\ &= \mathbb{E}_{y \sim p} [\text{KL}(p(x | y) \| p(x))], \quad (8.41) \end{aligned}$$

where $p(\mathbf{x}, \mathbf{y})$ denotes the joint distribution and $p(\mathbf{x}), p(\mathbf{y})$ denote the marginal distributions of \mathbf{X} and \mathbf{Y} .

8.2. Non-monotonicity of cond. mutual information.

1. Show that if $\mathbf{X} \perp \mathbf{Z}$ then $I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) \geq I(\mathbf{X}; \mathbf{Y})$.
 2. Show that if $\mathbf{X} \perp \mathbf{Z} | \mathbf{Y}$ then $I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) \leq I(\mathbf{X}; \mathbf{Y})$.
 3. Note that the condition $\mathbf{X} \perp \mathbf{Z} | \mathbf{Y}$ is the Markov property, namely, $\{\mathbf{X}, \mathbf{Y}, \mathbf{Z}\}$ form a Markov chain with graphical model $\mathbf{X} \rightarrow \mathbf{Y} \rightarrow \mathbf{Z}$. This situation often occurs when data is processed sequentially.
- Prove

$$I(\mathbf{X}; \mathbf{Z}) \leq I(\mathbf{X}; \mathbf{Y}). \quad (8.42)$$

which is also known as the *data processing inequality*, and which says that processing cannot increase the information contained in a signal.

8.3. Interaction information.

1. Show that interaction information is symmetric.
2. Let $X_1, X_2 \sim \text{Bern}(p)$ for some $p \in (0, 1)$ and independent. We denote by $Y \doteq X_1 \oplus X_2$ their XOR. Compute $I(Y; X_1; X_2)$.

8.4. Marginal gain of maximizing mutual information.

Show that in the context of maximizing mutual information, the marginal gain is

$$\Delta_I(\mathbf{x} | A) = I(f_{\mathbf{x}}; y_{\mathbf{x}} | \mathbf{y}_A) = H[y_{\mathbf{x}} | \mathbf{y}_A] - H[\varepsilon_{\mathbf{x}}].$$

8.5. Submodularity means no synergy.

Show that submodularity is equivalent to the absence of synergy between observations. That is, show that for all $A \subseteq B$,

$$I(f_{\mathbf{x}}; y_{\mathbf{x}}; \mathbf{y}_{B \setminus A} | \mathbf{y}_A) \geq 0. \quad (8.43)$$

8.6. Transductive active learning.

Consider the prior distribution $X_i \sim \mathcal{N}(0, 1)$ with all X_i independent and consider the “output” variable

$$Z \doteq \sum_{i=1}^{100} i \cdot X_i.$$

Our observation when X_{i_t} is selected in round t is generated by

$$Y_t \doteq X_{i_t} + \varepsilon_t \quad \text{with } \varepsilon_t \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1).$$

For a set of inputs $S = \{i_1, \dots, i_t\} \subseteq \{1, \dots, 100\}$, we define the *information gain* of S as

$$F(S) \doteq I(Z; Y_{1:t}) = H[Z] - H[Z \mid Y_{1:t}]$$

where $H[Z]$ is the entropy of Z according to the prior and $H[Z \mid Y_{1:t}]$ is the conditional entropy after the observations Y_1, \dots, Y_t .

Note: The random vector (X_1, \dots, X_{100}, Z) is jointly Gaussian due to the closedness of Gaussians under linear transformations.

1. We define the *marginal information gain* $\Delta(j \mid S)$ for a new observation Y_j as

$$\Delta(j \mid S) \doteq F(S \cup \{j\}) - F(S).$$

Does $\Delta(j \mid S) \geq 0$ hold for all j and S ?

2. Is maximizing the marginal information gain $\Delta(i \mid S)$ equivalent to picking the point $i \in \{1, \dots, 100\}$ with maximum variance under the current posterior over X_i , that is, “equivalent to uncertainty sampling”?
3. Let us consider the alternative prior where $X_i \sim \text{Bern}(0.5)$ are fair coin flips which we observe directly (i.e., $Y_t = X_{i_t}$). For which of the following definitions of Z is the acquisition function $S \mapsto F(S)$ submodular?
 - (a) $Z \doteq X_1 \wedge \dots \wedge X_{100}$ with \wedge denoting the logical AND.
 - (b) $Z \doteq X_1 \vee \dots \vee X_{100}$ with \vee denoting the logical OR.
 - (c) $Z \doteq X_1 \oplus \dots \oplus X_{100}$ with \oplus denoting the logical XOR, i.e., the exclusive OR which returns 1 iff exactly one of the X_i is 1 and 0 otherwise.

9

Bayesian Optimization

Often, obtaining data is costly. In the previous chapter, this led us to investigate how we can optimally improve our understanding (i.e., reduce uncertainty) of the process we are trying to model. However, purely improving our understanding is often not good enough. In many cases, we want to use our improving understanding *simultaneously* to reach certain goals. This is a very common problem in artificial intelligence and will concern us for the rest of this manuscript. One common instance of this problem is the setting of optimization.

Given some function $f^* : \mathcal{X} \rightarrow \mathbb{R}$, suppose we want to find the

$$\arg \max_{x \in \mathcal{X}} f^*(x). \quad (9.1)$$

Now, contrary to classical optimization, we are interested in the setting where the function f^* is unknown to us (like a “black-box”). We are only able to obtain noisy observations of f^* ,

$$y_t = f^*(x_t) + \varepsilon_t. \quad (9.2)$$

Moreover, these noise-perturbed evaluations are costly to obtain. We will assume that similar alternatives yield similar results,¹ which is commonly encoded by placing a Gaussian process prior on f^* . This assumed correlation is fundamentally what will allow us to learn a model of f^* from relatively few samples.²

9.1 Exploration-Exploitation Dilemma

In Bayesian optimization, we want to learn a model of f^* and use this model to optimize f^* simultaneously. These goals are somewhat contrary. Learning a model of f^* requires us to explore the input space while using the model to optimize f^* requires us to focus on the most promising well-explored areas. This trade-off is commonly known as the *exploration-exploitation dilemma*, whereby

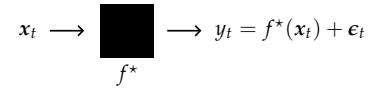


Figure 9.1: Illustration of Bayesian optimization. We pass an input x_t into the unknown function f^* to obtain noisy observations y_t .

¹ That is, f^* is “smooth”. We will be more precise in the subsequent parts of this chapter. If this were not the case, optimizing the function without evaluating it everywhere would not be possible. Fortunately, many interesting functions obey by this relatively weak assumption.

² There are countless examples of this problem in the “real world”. Instances are

- drug trials
- chemical engineering — the development of physical products
- recommender systems
- automatic machine learning — automatic tuning of model & hyperparameters
- and many more...

- *exploration* refers to choosing points that are “informative” with respect to the unknown function. For example, points that are far away from previously observed points (i.e., have high posterior variance);³ and
- *exploitation* refers to choosing promising points where we expect the function to have high values. For example, points that have a high posterior mean and a low posterior variance.

In other words, the exploration-exploitation dilemma refers to the challenge of learning enough to understand f^* , but not learning too much to lose track of the objective — optimizing f^* .

The exploration-exploitation dilemma is yet another example of the *principle of curiosity and conformity* which we introduced in Section 5.5.2 and encountered many times since in our study of approximate probabilistic inference. We will see in subsequent chapters that sequential decision-making is intimately related to probabilistic inference, and there we will also make this correspondence more precise.

9.2 Online Learning and Bandits

Bayesian optimization is closely related to a form of online learning. In *online learning* we are given a set of possible inputs \mathcal{X} and an unknown function $f^* : \mathcal{X} \rightarrow \mathbb{R}$. We are now asked to choose a sequence of inputs x_1, \dots, x_T online,⁴ and our goal is to maximize our cumulative reward $\sum_{t=1}^T f^*(x_t)$. Depending on what we observe about f^* , there are different variants of online learning. Bayesian optimization is closest to the so-called (stochastic) “bandit” setting.

9.2.1 Multi-Armed Bandits

The “*multi-armed bandits*” (MAB) problem is a classical, canonical formalization of the exploration-exploitation dilemma. In the MAB problem, we are provided with k possible actions (arms) and want to maximize our reward online within the time horizon T . We do not know the reward distributions of the actions in advance, however, so we need to trade learning the reward distribution with following the most promising action. Bayesian optimization can be interpreted as a variant of the MAB problem where there can be a potentially infinite number of actions (arms), but their rewards are correlated (because of the smoothness of the Gaussian process prior).

There exists a large body of work on this and similar problems in online decision-making. Much of this work develops theory on how to explore and exploit in the face of uncertainty. The shared prevalence of the exploration-exploitation dilemma signals a deep connection be-

³ We explored this topic (with strategies like uncertainty sampling) in the previous chapter.

⁴ *Online* is best translated as “sequential”. That is, we need to pick x_{t+1} based only on our prior observations y_1, \dots, y_t .

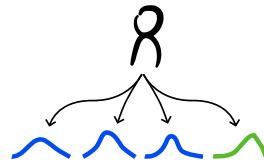


Figure 9.2: Illustration of a multi-armed bandit with four arms, each with a different reward distribution. The agent tries to identify the arm with the most beneficial reward distribution shown in green.

tween online learning and Bayesian optimization (and — as we will later come to see — reinforcement learning). Many of the approaches which we will encounter in the context of these topics are strongly related to methods in online learning.

One of the key principles of the theory on multi-armed bandits and reinforcement learning is the principle of “*optimism in the face of uncertainty*”, which suggests that it is a good guideline to explore where we can hope for the best outcomes. We will frequently come back to this general principle in our discussion of algorithms for Bayesian optimization and reinforcement learning.

9.2.2 Regret

The key performance metric in online learning is the regret.

Definition 9.1 (Regret). The (*cumulative*) *regret* for a time horizon T associated with choices $\{x_t\}_{t=1}^T$ is defined as

$$R_T \doteq \sum_{t=1}^T \underbrace{\left(\max_x f^*(x) - f^*(x_t) \right)}_{\text{instantaneous regret}} \quad (9.3)$$

$$= T \max_x f^*(x) - \sum_{t=1}^T f^*(x_t). \quad (9.4)$$

The regret can be interpreted as the additive loss with respect to the *static optimum* $\max_x f^*(x)$.

The goal is to find algorithms that achieve *sublinear* regret,

$$\lim_{T \rightarrow \infty} \frac{R_T}{T} = 0. \quad (9.5)$$

Importantly, if we use an algorithm which explores *forever*, e.g., by going to a random point \tilde{x} with a constant probability ϵ in each round, then the regret will grow linearly with time. This is because the instantaneous regret is at least $\epsilon(\max_x f^*(x) - f^*(\tilde{x}))$ and non-decreasing. Conversely, if we use an algorithm which *never* explores, then we might never find the static optimum, and hence, also incur constant instantaneous regret in each round, implying that regret grows linearly with time. Thus, achieving sublinear regret requires *balancing* exploration and exploitation.

Typically, online learning (and Bayesian optimization) consider stationary environments, hence the comparison to the static optimum. Dynamic environments are studied in *online algorithms* (see metrical task systems⁵, convex function chasing⁶, and generalizations of multi-armed bandits to changing reward distributions) and reinforcement

⁵ Metrical task systems are a classical example in online algorithms. Suppose we are moving in a (finite) decision space \mathcal{X} . In each round, we are given a “task” $f_t : \mathcal{X} \rightarrow \mathbb{R}$ which is more or less costly depending on our state $x_t \in \mathcal{X}$. In many contexts, it is natural to assume that it is also costly to move around in the decision space. This cost is modeled by a metric $d(\cdot, \cdot)$ on \mathcal{X} . In *metrical task systems*, we want to minimize our total cost,

$$\sum_{t=1}^T f_t(x_t) + d(x_t, x_{t-1}).$$

That is, we want to trade completing our tasks optimally with moving around in the state space. Crucially, we do not know the sequence of tasks f_t in advance. Due to the cost associated with moving in the decision space, previous choices affect the future!

⁶ Convex function chasing (or *convex body chasing*) generalize metrical task systems to continuous domains \mathcal{X} . To make any guarantees about the performance in these settings, one typically has to assume that the tasks f_t are convex. Note that this mirrors our assumption in Bayesian optimization that similar alternatives yield similar results.

learning. When operating in dynamic environments, other metrics such as the competitive ratio,⁷ which compares against the best *dynamic* choice, are useful. As we will later come to see in Section 13.1 in the context of reinforcement learning, operating in dynamic environments is deeply connected to a rich field of research called *control*.

9.3 Acquisition Functions

It is common to use a so-called *acquisition function* to greedily pick the next point to sample based on the current model.

Throughout our description of acquisition functions, we will focus on a setting where we model f^* using a Gaussian process which we denote by f . The methods generalize to other means of learning f^* such as Bayesian deep learning. The various acquisition functions F are used in the same way as is illustrated in Algorithm 9.2.

Algorithm 9.2: Bayesian optimization (with GPs)

```

1 initialize  $f \sim \mathcal{GP}(\mu_0, k_0)$ 
2 for  $t = 1$  to  $T$  do
3   choose  $x_t = \arg \max_{x \in \mathcal{X}} F(x; \mu_{t-1}, k_{t-1})$ 
4   observe  $y_t = f(x_t) + \epsilon_t$ 
5   perform a probabilistic update to obtain  $\mu_t$  and  $k_t$ 
```

Remark 9.3: Model selection

Selecting a model of f^* in sequential decision-making is much harder than in the i.i.d. data setting of supervised learning. There are mainly the following two dangers:

- the data sets collected in active learning and Bayesian optimization are *small*; and
- the data points are selected *dependently* on prior observations.

This leads to a specific danger of overfitting. In particular, due to feedback loops between the model and the acquisition function, one may end up sampling the same point repeatedly.

One approach to reduce the chance of overfitting is the use of hyperpriors which we mentioned previously in Section 4.4.2. Another approach that often works fairly well is to occasionally (according to some schedule) select points uniformly at random instead of using the acquisition function. This tends to prevent getting stuck in suboptimal parts of the state space.

⁷ To assess the performance in dynamic environments, we typically compare to a dynamic optimum. As these problems are difficult (we are usually not able to guarantee convergence to the dynamic optimum), one considers a multiplicative performance metric similar to the approximation ratio, the *competitive ratio*,

$$\text{cost(ALG)} \leq \alpha \cdot \text{cost(OPT)},$$

where OPT corresponds to the dynamic optimal choice (in hindsight).

One possible acquisition function is uncertainty sampling (8.31), which we discussed in the previous chapter. However, this acquisition function does not at all take into account the objective of maximizing f^* and focuses solely on exploration.

Suppose that our model f of f^* is well-calibrated, in the sense that the true function lies within its confidence bounds. Consider the best lower bound, that is, the maximum of the lower confidence bound. Now, if the true function is really contained in the confidence bounds, it must hold that the optimum is somewhere above this best lower bound. In particular, we can exclude all regions of the domain where the upper confidence bound (the optimistic estimate of the function value) is lower than the best lower bound. This is visualized in Figure 9.3.

Therefore, we only really care how the function looks like in the regions where the upper confidence bound is larger than the best lower bound. The key idea behind the methods that we will explore is to focus exploration on these plausible maximizers.

Note that it is crucial that our uncertainty about f reflects the “fit” of our model to the unknown function. If the model is not well calibrated or does not describe the underlying function at all, these methods will perform poorly. This is where we can use the Bayesian philosophy by imposing a prior belief that may be conservative.

9.3.1 Upper Confidence Bound

The principle of *optimism in the face of uncertainty* suggests picking the point where we can hope for the optimal outcome. In this setting, this corresponds to simply maximizing the *upper confidence bound* (UCB),

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} \mu_t(x) + \beta_{t+1} \sigma_t(x), \quad (9.6)$$

where $\sigma_t(x) \doteq \sqrt{k_t(x, x)}$ is the standard deviation at x and β_t regulates how confident we are about our model f (i.e., the choice of confidence interval).

This acquisition function naturally trades exploitation by preferring a large posterior mean with exploration by preferring a large posterior variance. Note that if $\beta_t = 0$ then UCB is purely exploitative, whereas, if $\beta_t \rightarrow \infty$, UCB recovers uncertainty sampling (i.e., is purely explorative).⁸ UCB is an example of an optimism-based method, as it greedily picks the point where we can hope for the best outcome.

As can be seen in Figure 9.4, the UCB acquisition function is generally non-convex. For selecting the next point, we can use approximate

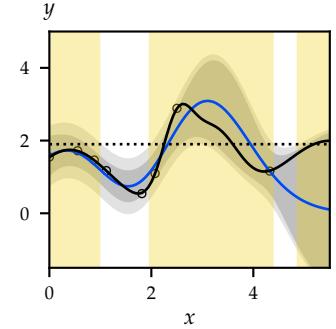


Figure 9.3: Optimism in Bayesian optimization. The **unknown function** is shown in black, our **model** in blue with gray confidence bounds. The dotted black line denotes the maximum lower bound. We can therefore focus our exploration to the yellow regions where the upper confidence bound is higher than the maximum lower bound.

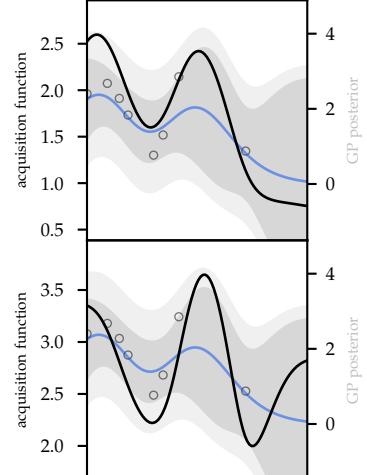


Figure 9.4: Plot of the UCB acquisition function for $\beta = 0.25$ and $\beta = 1$, respectively.

⁸ Due to the monotonicity of $(\cdot)^2$, it does not matter whether we optimize the variance or standard deviation at x .

global optimization techniques like Lipschitz optimization (in low dimensions) and gradient ascent with random initialization (in high dimensions). Another widely used technique is to sample some random points from the domain, score them according to this criterion, and simply take the best one.

The choice of β_t is crucial for the performance of UCB. Intuitively, for UCB to work even if the unknown function f^* is not contained in the confidence bounds, we use β_t to re-scale the confidence bounds to enclose f^* as shown in Figure 9.5. A theoretical analysis requires that β_t is chosen “correctly”. Formally, we say that the sequence β_t is chosen correctly if it leads to *well-calibrated confidence intervals*, that is, if with probability at least $1 - \delta$,

$$\forall t \geq 1, \forall x \in \mathcal{X} : f^*(x) \in \mathcal{C}_t(x) \doteq [\mu_{t-1}(x) \pm \beta_t(\delta) \cdot \sigma_{t-1}(x)]. \quad (9.7)$$

Bounds on $\beta_t(\delta)$ can be derived both in a “Bayesian” and in a “frequentist” setting. In the Bayesian setting, it is assumed that f^* is drawn from the prior GP, i.e., $f^* \sim \mathcal{GP}(\mu_0, k_0)$. However, in many cases this may be an unrealistic assumption. In the frequentist setting, it is assumed instead that f^* is a fixed element of a reproducing kernel Hilbert space $\mathcal{H}_k(\mathcal{X})$ which depending on the kernel k can encompass a large class of functions. We will discuss the Bayesian setting first and later return to the frequentist setting.

Theorem 9.4 (Bayesian confidence intervals, lemma 5.5 of Srinivas et al. (2010)). **?** Let $\delta \in (0, 1)$. Assuming $f^* \sim \mathcal{GP}(\mu_0, k_0)$ and Gaussian observation noise $\epsilon_t \sim \mathcal{N}(0, \sigma_n^2)$, the sequence

$$\beta_t(\delta) = O\left(\sqrt{\log(|\mathcal{X}|t/\delta)}\right) \quad (9.8)$$

satisfies $\mathbb{P}(\forall t \geq 1, x \in \mathcal{X} : f^*(x) \in \mathcal{C}_t(x)) \geq 1 - \delta$.

Under the assumption of well-calibrated confidence intervals, we can bound the regret of UCB.

Theorem 9.5 (Regret of GP-UCB, theorem 2 of Srinivas et al. (2010)). **?** If $\beta_t(\delta)$ is chosen “correctly” for a fixed $\delta \in (0, 1)$, with probability at least $1 - \delta$, greedily choosing the upper confidence bound yields cumulative regret

$$R_T = O\left(\beta_T(\delta)\sqrt{\gamma_T T}\right) \quad (9.9)$$

where

$$\gamma_T \doteq \max_{\substack{S \subseteq \mathcal{X} \\ |S|=T}} I(f_S; \mathbf{y}_S) = \max_{\substack{S \subseteq \mathcal{X} \\ |S|=T}} \frac{1}{2} \log \det(\mathbf{I} + \sigma_n^{-2} \mathbf{K}_{SS}) \quad (9.10)$$

is the maximum information gain after T rounds.

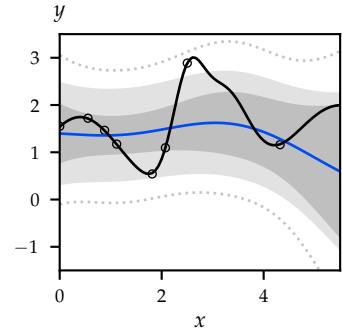


Figure 9.5: Re-scaling the confidence bounds. The dotted gray lines represent updated confidence bounds.

Problem 9.2

Problem 9.3

Observe that if the information gain is sublinear in T then we achieve sublinear regret and, in particular, converge to the true optimum. The information gain γ_T measures how much can be learned about f^* within T rounds. If the function is assumed to be smooth (perhaps even linear), then the information gain is smaller than if the function was assumed to be “rough”. Intuitively, the smoother the functions encoded by the prior, the smaller is the class of functions to choose from and the more can be learned from a single observation about “neighboring” points.

Theorem 9.6 (Information gain of common kernels, theorem 5 of Srinivas et al. (2010) and remark 2 of Vakili et al. (2021)). *Due to submodularity, we have the following bounds on the information gain of common kernels:*

- linear kernel

$$\gamma_T = O(d \log T), \quad (9.11)$$

- Gaussian kernel

$$\gamma_T = O((\log T)^{d+1}), \quad (9.12)$$

- Matérn kernel for $\nu > \frac{1}{2}$

$$\gamma_T = O\left(T^{\frac{d}{2\nu+d}} (\log T)^{\frac{2\nu}{2\nu+d}}\right). \quad (9.13)$$

The information gain of common kernels is illustrated in Figure 9.6. Notably, when all points in the domain are independent, the information gain is linear in T . This is because when the function f^* may be arbitrarily “rough”, we cannot generalize from a single observation to “neighboring” points, and as there are infinitely many points in the domain \mathcal{X} there are no diminishing returns. As one would expect, in this case, Theorem 9.5 does not yield sublinear regret. However, we can see from Theorem 9.6 that the information gain is sublinear for linear, Gaussian, and most Matérn kernels. Moreover, observe that unless the function is linear, the information gain grows exponentially with the dimension d . This is because the number of “neighboring” points (with respect to Euclidean geometry) decreases exponentially with the dimension which is also known as the *curse of dimensionality*.

As mentioned, the size of the confidence intervals can also be analyzed under a frequentist assumption on f^* .

Theorem 9.7 (Frequentist confidence intervals, theorem 2 of Chowdhury and Gopalan (2017)). *Let $\delta \in (0, 1)$. Assuming $f^* \in \mathcal{H}_k(\mathcal{X})$, we have that with probability at least $1 - \delta$, the sequence*

$$\beta_t(\delta) = \|f^*\|_k + \sigma_n \sqrt{2(\gamma_t + \log(1/\delta))} \quad (9.14)$$

satisfies $\mathbb{P}(\forall t \geq 1, \mathbf{x} \in \mathcal{X} : f^*(\mathbf{x}) \in \mathcal{C}_t(\mathbf{x})) \geq 1 - \delta$.

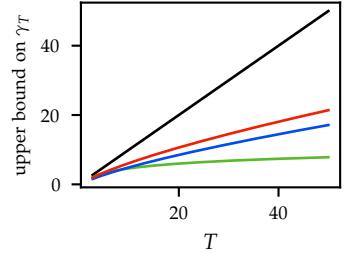


Figure 9.6: Information gain of **independent**, **linear**, **Gaussian**, and **Matérn** ($\nu \approx 0.5$) kernels with $d = 2$ (up to constant factors). The kernels with sublinear information gain have strong diminishing returns (due to their strong dependence between “close” points). In contrast, the independent kernel has no dependence between points in the domain, and therefore no diminishing returns. Intuitively, the “smoother” the class of functions modeled by the kernel, the stronger are the diminishing returns.

That is, β_t depends on the information gain of the kernel as well as on the “complexity” of f^* which is measured in terms of its norm in the underlying reproducing kernel Hilbert space $\mathcal{H}_k(\mathcal{X})$.

Remark 9.8: Bayesian vs frequentist assumption

Theorems 9.4 and 9.7 provide different bounds on $\beta_t(\delta)$ based on fundamentally different assumptions on the ground truth f^* : The Bayesian assumption is that f^* is drawn from the prior GP, whereas the frequentist assumption is that f^* is an element of a reproducing kernel Hilbert space $\mathcal{H}_k(\mathcal{X})$. The frequentist assumption holds uniformly for all functions f^* with $\|f^*\|_k < \infty$, whereas the Bayesian assumption holds only under the Bayesian “belief” that f^* is drawn from the prior GP.

Interestingly, neither assumption encompasses the other. This is because if $f \sim \mathcal{GP}(0, k)$ then it can be shown that almost surely $\|f\|_k = \infty$, which implies that $f \notin \mathcal{H}_k(\mathcal{X})$ (Srinivas et al., 2010).

We remark that Theorem 9.7 holds also under the looser assumption that observations are perturbed by σ_n -sub-Gaussian noise (cf. Equation (A.39)) instead of Gaussian noise. The bound on γ_T from Equation (9.13) for the Matérn kernel does not yield sublinear regret when combined with the standard regret bound from Theorem 9.5, however, Whitehouse et al. (2024) show that the regret of GP-UCB is sublinear also in this case provided σ_n^2 is chosen carefully.

This concludes our discussion of the UCB algorithm. We have seen that its regret can be analyzed under both Bayesian and frequentist assumptions on f^* .

9.3.2 Improvement

Another well-known family of methods is based on keeping track of a running optimum \hat{f}_t , and scoring points according to their improvement upon the running optimum. The *improvement* of x after round t is measured by

$$I_t(x) \doteq (f(x) - \hat{f}_t)_+ \quad (9.15)$$

where we use $(\cdot)_+$ to denote $\max\{0, \cdot\}$.

The *probability of improvement* (PI) picks the point that maximizes the probability to improve upon the running optimum,

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} \mathbb{P}(I_t(x) > 0 \mid x_{1:t}, y_{1:t}) \quad (9.16)$$

$$= \arg \max_{x \in \mathcal{X}} \mathbb{P}(f(x) > \hat{f}_t \mid x_{1:t}, y_{1:t}) \quad (9.17)$$

$$= \arg \max_{x \in \mathcal{X}} \Phi \left(\frac{\mu_t(x) - \hat{f}_t}{\sigma_t(x)} \right) \quad (9.18)$$

where Φ denotes the CDF of the standard normal distribution and we use that $f(x) | x_{1:t}, y_{1:t} \sim \mathcal{N}(\mu_t(x), \sigma_t^2(x))$. Probability of improvement tends to be biased in favor of exploitation, as it prefers points with large posterior mean and small posterior variance which is typically true “close” to the previously observed maximum \hat{f}_t .

Probability of improvement looks at *how likely* a point is to improve upon the running optimum. An alternative is to look at *how much* a point is expected to improve upon the running optimum. This acquisition function is called the *expected improvement* (EI),

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} \mathbb{E}[I_t(x) | x_{1:t}, y_{1:t}] . \quad (9.19)$$

Intuitively, EI seeks a large expected improvement (exploitation) while also preferring states with a large variance (exploration). Expected improvement yields the same regret bound as UCB (Nguyen et al., 2017).

The expected improvement acquisition function is often flat which makes it difficult to optimize in practice due to vanishing gradients. One approach addressing this is to instead optimize the logarithm of EI (Ament et al., 2024).

using linear transformations of Gaussians (1.78)

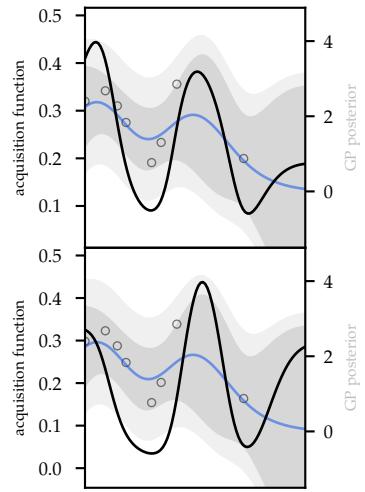


Figure 9.7: Plot of the PI and EI acquisition functions, respectively.

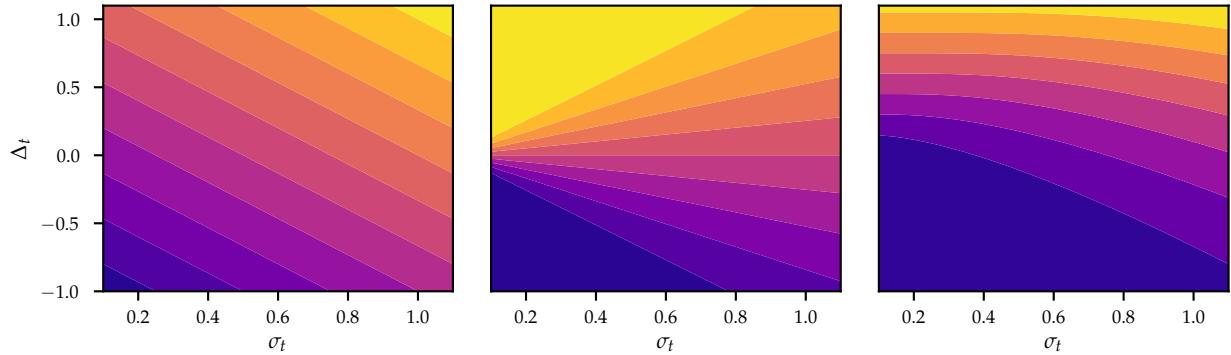


Figure 9.8: Contour lines of acquisition functions for varying $\Delta_t = \mu_t(x) - \hat{f}_t$ and σ_t . A brighter color corresponds to a larger acquisition function value. The first graph shows contour lines of UCB with $\beta_t = 0.75$, the second of PI, and the third of EI.

9.3.3 Thompson Sampling

We can also interpret the principle of *optimism in the face of uncertainty* in a slightly different way than we did with UCB (and EI). Suppose we select the next point according to the probability that it is optimal (assuming that the posterior distribution is an accurate representation of the uncertainty),

$$\pi(x | x_{1:t}, y_{1:t}) \doteq \mathbb{P}_{f|x_{1:t}, y_{1:t}} \left(f(x) = \max_{x'} f(x') \right) \quad (9.20)$$

$$\mathbf{x}_{t+1} \sim \pi(\cdot \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:t}). \quad (9.21)$$

This approach of sampling according to the *probability of maximality* π is called *probability matching*. Probability matching is exploratory as it prefers points with larger variance (as they automatically have a larger chance of being optimal), but at the same time exploitative as it effectively discards points with low posterior mean and low posterior variance. Unfortunately, it is generally difficult to compute π analytically given a posterior.

Instead, it is common to use a sampling-based approximation of π . Observe that the density π can be expressed as an expectation,

$$\pi(\mathbf{x} \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:t}) = \mathbb{E}_{f \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:t}} [\mathbb{1}\{f(\mathbf{x}) = \max_{\mathbf{x}'} f(\mathbf{x}')\}], \quad (9.22)$$

which we can approximate using Monte Carlo sampling (typically using a single sample),

$$\approx \mathbb{1}\{\tilde{f}_{t+1}(\mathbf{x}) = \max_{\mathbf{x}'} \tilde{f}_{t+1}(\mathbf{x}')\} \quad (9.23)$$

where $\tilde{f}_{t+1} \sim p(\cdot \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ is a sample from our posterior distribution. Observe that this approximation of π coincides with a point density at the maximizer of \tilde{f}_{t+1} .

The resulting algorithm is known as *Thompson sampling*. At time $t + 1$, we sample a function $\tilde{f}_{t+1} \sim p(\cdot \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ from our posterior distribution. Then, we simply maximize \tilde{f}_{t+1} ,

$$\mathbf{x}_{t+1} \doteq \arg \max_{\mathbf{x} \in \mathcal{X}} \tilde{f}_{t+1}(\mathbf{x}). \quad (9.24)$$

In many cases, the randomness in the realizations of \tilde{f}_{t+1} is already sufficient to effectively trade exploration and exploitation. Similar regret bounds to those of UCB can also be established for Thompson sampling (Russo and Van Roy, 2016; Kandasamy et al., 2018).

9.3.4 Information-Directed Sampling

After having looked at multiple methods that aim to balance exploitation of the current posterior distribution over f for immediate returns with exploration to reduce uncertainty about f for future returns, we will next discuss a method that makes this tradeoff explicit.

Denoting the instantaneous regret of choosing \mathbf{x} as $\Delta(\mathbf{x}) \doteq \max_{\mathbf{x}'} f^*(\mathbf{x}') - f^*(\mathbf{x})$ and by $I_t(\mathbf{x})$ some function capturing the “information gain” associated with observing \mathbf{x} in iteration $t + 1$, we can define the *information ratio*,

$$\Psi_t(\mathbf{x}) \doteq \frac{\Delta(\mathbf{x})^2}{I_t(\mathbf{x})}, \quad (9.25)$$

which was originally introduced by Russo and Van Roy (2016). Here exploitation reduces regret while exploration increases information gain, and hence, points \mathbf{x} that minimize the information ratio are those that most effectively balance exploration and exploitation. We can make the key observation that the regret $\Delta(\cdot)$ decreases when $I_t(\cdot)$ decreases, as a small $I_t(\cdot)$ implies that the algorithm has already learned a lot about the function f^* . The strength of this relationship is quantified by the information ratio:

Theorem 9.9 (Proposition 1 of Russo and Van Roy (2014) and theorem 8 of Kirschner and Krause (2018)). *For any iteration $T \geq 1$, let $\sum_{t=1}^T I_{t-1}(\mathbf{x}_t) \leq \gamma_T$ and suppose that $\Psi_{t-1}(\mathbf{x}_t) \leq \bar{\Psi}_T$ for all $t \in [T]$. Then, the cumulative regret is bounded by*

$$R_T \leq \sqrt{\gamma_T \bar{\Psi}_T T}. \quad (9.26)$$

Proof. By Equation (9.25), $r_t = \Delta(\mathbf{x}_t) = \sqrt{\Psi_{t-1}(\mathbf{x}_t) \cdot I_{t-1}(\mathbf{x}_t)}$. Hence,

$$\begin{aligned} R_T &= \sum_{t=1}^T r_t \\ &= \sum_{t=1}^T \sqrt{\Psi_{t-1}(\mathbf{x}_t) \cdot I_{t-1}(\mathbf{x}_t)} \\ &\leq \sqrt{\sum_{t=1}^T \Psi_{t-1}(\mathbf{x}_t) \cdot \sum_{t=1}^T I_{t-1}(\mathbf{x}_t)} \\ &\leq \sqrt{\gamma_T \bar{\Psi}_T T}. \end{aligned} \quad \square$$

using the Cauchy-Schwarz inequality

using the assumptions on $I_t(\cdot)$ and $\Psi_t(\cdot)$

Example 9.10: How to measure “information gain”?

One possibility of measuring the “information gain” is

$$I_t(\mathbf{x}) \doteq I(f_{\mathbf{x}}; \mathbf{y}_{\mathbf{x}} \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:t}) \quad (9.27)$$

which — as you may recall — is precisely the marginal gain of the utility $I(S) = I(f_S; \mathbf{y}_S)$ we were studying in Chapter 8. In this case,

$$\begin{aligned} \sum_{t=1}^T I_{t-1}(\mathbf{x}_t) &= \sum_{t=1}^T I(f_{\mathbf{x}_t}; \mathbf{y}_{\mathbf{x}_t} \mid \mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1}) \\ &= \sum_{t=1}^T I(f_{\mathbf{x}_{1:T}}; \mathbf{y}_{\mathbf{x}_t} \mid \mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1}) \\ &= I(f_{\mathbf{x}_{1:T}}; \mathbf{y}_{\mathbf{x}_{1:T}}) \\ &\leq \gamma_T. \end{aligned}$$

using $I(\mathbf{X}, \mathbf{Z}; \mathbf{Y}) \geq I(\mathbf{X}; \mathbf{Y})$ which follows from Equations (8.12) and (8.16) and is called *monotonicity of MI*
by repeated application of
Equation (8.16), also called the *chain rule*
of MI by definition of γ_T (9.10)

The regret bound from Theorem 9.9 suggests an algorithm which in each iteration chooses the point which minimizes the information ratio (9.25). However, this is not possible since $\Delta(\cdot)$ is unknown due to its dependence on f^* . Kirschner and Krause (2018) propose to use a surrogate to the regret which is based on the current model of f^* ,

$$\hat{\Delta}_t(\mathbf{x}) \doteq \max_{\mathbf{x}' \in \mathcal{X}} u_t(\mathbf{x}') - l_t(\mathbf{x}). \quad (9.28)$$

Here, $u_t(\mathbf{x}) \doteq \mu_t(\mathbf{x}) + \beta_{t+1}\sigma_t(\mathbf{x})$ and $l_t(\mathbf{x}) \doteq \mu_t(\mathbf{x}) - \beta_{t+1}\sigma_t(\mathbf{x})$ are the upper and lower confidence bounds of the confidence interval $\mathcal{C}_t(\mathbf{x})$ of $f^*(\mathbf{x})$, respectively. Similarly to our discussion of UCB, we make the assumption that the sequence β_t is chosen “correctly” (cf. Equation (9.7)) so that the confidence interval is well-calibrated and $\Delta(\mathbf{x}) \leq \hat{\Delta}_t(\mathbf{x})$ with high probability. The resulting algorithm

$$\mathbf{x}_{t+1} \doteq \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \hat{\Psi}_t(\mathbf{x}) \doteq \frac{\hat{\Delta}_t(\mathbf{x})^2}{I_t(\mathbf{x})} \right\} \quad (9.29)$$

is known as *information-directed sampling* (IDS).

Theorem 9.11 (Regret of IDS, lemma 8 of Kirschner and Krause (2018)).

② Let $\beta_t(\delta)$ be chosen “correctly” for a fixed $\delta \in (0, 1)$. Then, if the measure of information gain is $I_t(\mathbf{x}) = I(f_x; y_x | \mathbf{x}_{1:t}, y_{1:t})$, with probability at least $1 - \delta$, IDS has cumulative regret

$$R_T = O\left(\beta_T(\delta)\sqrt{\gamma_T T}\right). \quad (9.30)$$

Regret bounds such as Theorem 9.11 can be derived also for different measures of information gain. For example, the argument of Problem 9.6 also goes through for the “greedy” measure

$$I_t(\mathbf{x}) \doteq I(f_{\mathbf{x}_t^{\text{UCB}}}; y_x | \mathbf{x}_{1:t}, y_{1:t}) \quad (9.31)$$

which focuses exclusively on reducing the uncertainty at $\mathbf{x}_t^{\text{UCB}}$ rather than globally. We compare the two measures of information gain in Figure 9.9. Observe that the acquisition function depends critically on the choice of $I_t(\cdot)$ and is less sensitive to the scaling of confidence intervals.

IDS trades exploitation and exploration by balancing the (exploitative) regret surrogate with a measure of information gain (such as those studied in Chapter 8) that is purely explorative. In this way, IDS can account for kinds of information which are not addressed by alternative algorithms such as UCB or EI (Russo and Van Roy, 2014): Depending on the measure of information gain, IDS can select points to obtain *indirect information* about other points or *cumulating information* that does not immediately lead to a higher reward but only when combined with subsequent observations. Moreover, IDS avoids selecting points which yield *irrelevant information*.

Problem 9.6

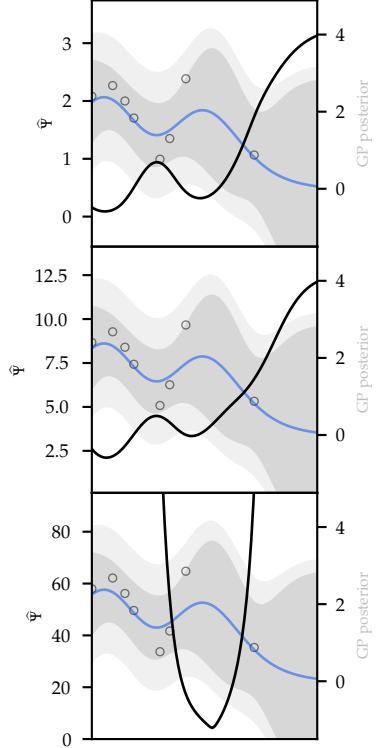


Figure 9.9: Plot of the surrogate information ratio $\hat{\Psi}$: IDS selects its minimizer. The first two plots use the “global” information gain measure from Example 9.10 with $\beta = 0.25$ and $\beta = 0.5$, respectively. The third plot uses the “greedy” information gain measure from Equation (9.31) and $\beta = 1$.

9.3.5 Probabilistic Inference

As we mentioned before, computing the probability of maximality π is generally intractable. You can think of computing π as attempting to fully solve the probabilistic inference problem associated with determining the optimum of $f \mid \mathbf{x}_{1:t}, y_{1:t}$. In many cases, it is useful to determine the probability that a point is optimal under the current posterior distribution, and we will consider one particular example in the following.

Example 9.12: Maximizing recall

In domains such as molecular design we often use machine learning to screen candidates for further manual testing. The goal here is to suggest a small set E from a large domain of molecules \mathcal{X} , so that the probability of E containing the optimal molecule, i.e., the *recall*, is maximized. Note that this task is quite different from online Bayesian optimization, for example, in BO we get sequential feedback that we can use to decide which inputs to query next. Nevertheless, we will see in this section that both tasks turn out to be closely related.

Let us assume that the maximizer is unique almost surely, which with GPs is automatically the case if there are no same-mean, perfectly-correlated entries. The recall task is then intimately related to the task of determining the probability of maximality, since

$$\begin{aligned} \arg \max_{E \subseteq \mathcal{X}: |E|=k} \mathbb{P}_{f \mid \mathbf{x}_{1:t}, y_{1:t}} \left(\max_{\mathbf{x} \in E} f(\mathbf{x}) = \max_{\mathbf{x}' \in \mathcal{X}} f(\mathbf{x}') \right) \\ = \arg \max_{E \subseteq \mathcal{X}: |E|=k} \sum_{\mathbf{x} \in E} \pi(\mathbf{x} \mid \mathbf{x}_{1:t}, y_{1:t}). \end{aligned} \quad (9.32)$$

Thus, to obtain the recall-optimal candidates for further testing, we need to find the probability of maximality. However, as mentioned, computing the probability of maximality π is generally intractable.

Figure 9.10 depicts an example of a recall task, where the black line shows the optimal recall achieved by knowing the probability of maximality π exactly. LITE (Menet et al., 2025) is an almost-linear time approximation of π , whereas TS selects points via Thompson sampling and MEANS selects the points with the highest posterior mean. Selecting E according to probability of maximality achieves much higher recall than the other intuitive heuristics.

by noting that for $\mathbf{x} \neq \mathbf{y}$, the events $\{f(\mathbf{x}) = \max_{\mathbf{x}' \in \mathcal{X}} f(\mathbf{x}')\}$ and $\{f(\mathbf{y}) = \max_{\mathbf{x}' \in \mathcal{X}} f(\mathbf{x}')\}$ are disjoint

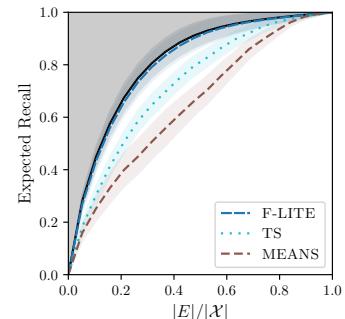


Figure 9.10: Menet et al. (2025) show in an experiment that selecting $E \subseteq \mathcal{X}$ according to the estimates of probability of maximality from LITE (here called F-LITE) is near-optimal. Intuitive heuristics such as Thompson sampling (TS) or selecting points with the highest posterior mean (MEANS) perform worse.

So how can we estimate the probability of maximality? LITE approximates π by

$$\begin{aligned}\pi(\mathbf{x} \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:t}) &= \mathbb{P}_{f|\mathbf{x}_{1:t}, \mathbf{y}_{1:t}}\left(f(\mathbf{x}) \geq \max_{\mathbf{x}'} f(\mathbf{x}')\right) \\ &\approx \mathbb{P}_{f|\mathbf{x}_{1:t}, \mathbf{y}_{1:t}}(f(\mathbf{x}) \geq \kappa^*)\end{aligned}\quad (9.33)$$

$$= \Phi\left(\frac{\mu_t(\mathbf{x}) - \kappa^*}{\sigma_t(\mathbf{x})}\right) \quad (9.34)$$

with Φ denoting the CDF of the standard normal distribution and κ^* chosen such that the approximation of π integrates to 1, so that it is a valid distribution.

Remarkably, LITE is intimately related to many of the BO methods we discussed in this chapter. First, and most obviously, Equation (9.33) looks similar to the PI acquisition function (9.17). But note that κ^* is not equal to the best observed value \hat{f}_t as in PI, but instead typically larger. If $\kappa^* > \hat{f}_t$, then Equation (9.33) is more exploratory than PI: it puts additional emphasis on points with large posterior variance and comparatively less emphasis on points with large posterior mean, which we illustrate in Figure 9.11.

An insightful interpretation of LITE is that it balances two different kinds of exploration: *optimism in the face of uncertainty* and *entropy regularization*. To see this, let us define the following variational objective:

$$\mathcal{W}(\pi) \doteq \sum_{\mathbf{x} \in \mathcal{X}} \pi(\mathbf{x}) \cdot \left(\underbrace{\mu_t(\mathbf{x})}_{\text{exploitation}} + \underbrace{\sqrt{2S'(\pi(\mathbf{x})) \cdot \sigma_t(\mathbf{x})}}_{\text{exploration}} \right) \quad (9.35)$$

with the “quasi-surprise” $S'(u) \doteq \frac{1}{2}(\phi(\Phi^{-1}(u))/u)^2$. The quasi-surprise $S'(\cdot)$ behaves similarly to the surprise $-\ln(\cdot)$. In fact, their asymptotics coincide:

$$S'(1) = 0 = -\ln(1) \quad \text{and} \quad S'(u) \rightarrow -\ln(u) \text{ as } u \rightarrow 0^+.$$

The objective \mathcal{W} is maximized for those probability distributions π that are concentrated around points with large mean $\mu_t(\mathbf{x})$ and points with large exploration bonus. We have seen that the uncertainty $\sigma_t(\mathbf{x})$ about $f(\mathbf{x})$ is the standard exploration bonus of UCB-style algorithms. In Equation (9.35), $\sigma_t(\mathbf{x})$ is weighted by the quasi-surprise, which acts as “entropy regularization”: it increases the entropy of π by uniformly pushing $\pi(\mathbf{x})$ away from zero. You can think of entropy regularization as encouraging a different kind of exploration than optimism by not making deterministic decisions like in UCB.

Menet et al. (2025) show that LITE (9.34) is the solution to the varia-

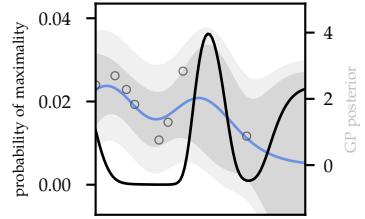


Figure 9.11: Plot of the probability of maximality as estimated by LITE.

tional problem (9.35) among valid probability distributions $\{\cdot\}$:

$$\arg \max_{\pi \in \Delta^{\mathcal{X}}} \mathcal{W}(\pi) = \Phi \left(\frac{\mu_t(\cdot) - \kappa^*}{\sigma_t(\cdot)} \right) \quad (9.36)$$

with κ^* such that the right-hand side sums to 1.⁹ This indicates that LITE and Thompson sampling, which samples from probability of maximality, achieve exploration through two means:

1. **Optimism:** by preferring points with large uncertainty $\sigma_t(x)$ about the reward value $f(x)$.
2. **Decision uncertainty:** by assigning some probability mass to all x , that is, by remaining uncertain about which x is the maximizer.

In our discussion of balancing exploration and exploitation in reinforcement learning, we will return to this dichotomy of exploration strategies.

Problem 9.7

⁹ $\Delta^{\mathcal{X}}$ is the probability simplex on \mathcal{X} .

Remark 9.13: But why is “exploration” useful for recall?

Let us pause for a moment and reflect on why this interpretation of LITE is remarkable. This interpretation shows that LITE is more “exploratory” than simply taking the highest posterior means. However, the recall task (9.32) differs from the standard BO task in that we do not collect any further observations, which we may use downstream to make better decisions. In the recall task, we are only interested in having the best shot at including the maximizer in the set E , without obtaining any further information or making any subsequent decisions. At first sight, this seems to suggest that we should be as “exploitative” as possible. Then how can it be that “exploration” is useful?

We will explore this question with the following example: You are observing a betting game. When placing a bet, players can either place a safe bet (“playing safe”) or a risky bet (“playing risky”). You now have to place a bet on their bets, and estimate which of the players will win the most money: those that play safe or those that play risky? Note that you do not care whether your guess ends up in 2nd place or last place among all players — you only care about whether your guess wins. That is, you are interested in *recalling* the best player, not in “ranking” the players.

Consider three players: one that plays safe and two that play risky. Suppose that the safe bet has payoff $S = 1$ while each risky bet has payoff $R \sim \mathcal{N}(0, 100)$. In expectation, the safe player will win the most money. However, one can see with just a little bit of algebra that the probability of either of the risky players winning the most money is $\approx 35\%$, whereas the safe player only wins with

probability $\approx 29\%$ (2). That is, it is in fact *optimal* to bet on either of the risky players since the best player might have vastly outperformed their expected winnings, and performed closer to their upper confidence bound. In summary, maximizing recall requires us to be “exploratory” since it is likely that the optimum among inputs is one that has performed better than expected, not simply the one with the highest expected performance.

Problem 9.8

Discussion

In this chapter, we have explored the exploration-exploitation dilemma in the context of optimizing black-box functions. To this end, we have explored various methods to balance exploration and exploitation. While computing the precise probability of maximality is generally intractable, we found that it can be understood approximately as balancing two sources of exploration: optimism in the face of uncertainty and entropy regularization.

In the following chapters, we will begin to discuss stateful settings where the black-box optimization task is known as “reinforcement learning”. Naturally, we will see that the exploration-exploitation dilemma is also a central challenge in reinforcement learning, and we will revisit many of the concepts we have discussed in this chapter.

Optional Readings

- Srinivas, Krause, Kakade, and Seeger (2010).
Gaussian process optimization in the bandit setting: No regret and experimental design.
- Golovin, Solnik, Moitra, Kochanski, Karro, and Sculley (2017).
Google vizier: A service for black-box optimization.
- Romero, Krause, and Arnold (2013).
Navigating the protein fitness landscape with Gaussian processes.
- Chowdhury and Gopalan (2017).
On kernelized multi-armed bandits.

Problems

9.1. Convergence to the static optimum.

Show that any algorithm where $\lim_{t \rightarrow \infty} f^*(x_t)$ exists achieves sublinear regret if and only if it converges to the static optimum, that is,

$$\lim_{t \rightarrow \infty} f^*(x_t) = \max_x f^*(x). \quad (9.37)$$

Hint: Use that if a sequence a_n converges to a as $n \rightarrow \infty$, then we have for the sequence

$$b_n \doteq \frac{1}{n} \sum_{i=1}^n a_i \quad (9.38)$$

that $\lim_{n \rightarrow \infty} b_n = a$. This is also known as the Cesàro mean.

9.2. Bayesian confidence intervals.

In this exercise, we derive Theorem 9.4.

1. For fixed $t \geq 1$ and $\mathbf{x} \in \mathcal{X}$, prove

$$\mathbb{P}(f^*(\mathbf{x}) \notin \mathcal{C}_t(\mathbf{x}) \mid \mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1}) \leq e^{-\beta_t^2/2}. \quad (9.39)$$

Hint: Bound $\mathbb{P}(Z > c)$ for $Z \sim \mathcal{N}(0, 1)$ and $c > 0$.

2. Prove Theorem 9.4.

9.3. Regret of GP-UCB.

To develop some intuition, we will derive Theorem 9.5.

1. Show that if Equation (9.7) holds, then for a fixed $t \geq 1$ the instantaneous regret r_t is bounded by $2\beta_t \sigma_{t-1}(\mathbf{x}_t)$.
2. Let $S_T \doteq \{\mathbf{x}_t\}_{t=1}^T$, and define $f_T \doteq f_{S_T}$ and $\mathbf{y}_T \doteq \mathbf{y}_{S_T}$. Prove

$$I(f_T; \mathbf{y}_T) = \frac{1}{2} \sum_{t=1}^T \log \left(1 + \frac{\sigma_{t-1}^2(\mathbf{x}_t)}{\sigma_n^2} \right). \quad (9.40)$$

3. Combine (1) and (2) to show Theorem 9.5. We assume w.l.o.g. that the sequence $\{\beta_t\}_t$ is monotonically increasing.

Hint: If $s \in [0, M]$ for some $M > 0$ then $s \leq C \cdot \log(1+s)$ with $C \doteq M / \log(1+M)$.

9.4. Sublinear regret of GP-UCB for a linear kernel.

Assume that $f^* \sim \mathcal{GP}(0, k)$ where k is the linear kernel

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'.$$

In addition, we assume that for any $\mathbf{x} \in \mathcal{X}$, $\|\mathbf{x}\|_2 \leq 1$. Moreover, recall that the points in a finite set $S \subseteq \mathcal{X}$ can be written in a matrix form (the “design matrix”) which we denote by $\mathbf{X}_S \in \mathbb{R}^{d \times |S|}$.

1. Prove that $\gamma_T = O(d \log T)$.
2. Deduce from (1) and Theorem 9.5 that $\lim_{T \rightarrow \infty} R_T/T = 0$.

9.5. Closed-form expected improvement.

Let us denote the acquisition function of EI from Equation (9.19) by $EI_t(\mathbf{x})$. In this exercise, we derive a closed-form expression.

1. Show that

$$\text{EI}_t(\mathbf{x}) = \int_{(\hat{f}_t - \mu_t(\mathbf{x}))/\sigma_t(\mathbf{x})}^{+\infty} (\mu_t(\mathbf{x}) + \sigma_t(\mathbf{x})\varepsilon - \hat{f}_t) \cdot \phi(\varepsilon) d\varepsilon \quad (9.41)$$

where ϕ is the PDF of the univariate standard normal distribution (1.5).

Hint: Reparameterize the posterior distribution

$$f(\mathbf{x}) \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:t} \sim \mathcal{N}(\mu_t(\mathbf{x}), \sigma_t^2(\mathbf{x}))$$

using a standard normal distribution.

2. Using the above expression, show that

$$\text{EI}_t(\mathbf{x}) = (\mu_t(\mathbf{x}) - \hat{f}_t)\Phi\left(\frac{\mu_t(\mathbf{x}) - \hat{f}_t}{\sigma_t(\mathbf{x})}\right) + \sigma_t(\mathbf{x})\phi\left(\frac{\mu_t(\mathbf{x}) - \hat{f}_t}{\sigma_t(\mathbf{x})}\right) \quad (9.42)$$

where $\Phi(u) \doteq \int_{-\infty}^u \phi(\varepsilon) d\varepsilon$ denotes the CDF of the standard normal distribution.

Note that the first term of Equation (9.42) encourages exploitation while the second term encourages exploration. EI can be seen as a special case of UCB where the confidence bounds are scaled depending on \mathbf{x} :

$$\beta_t = \phi(z_t(\mathbf{x}))/\Phi(z_t(\mathbf{x}))$$

for $z_t(\mathbf{x}) \doteq (\mu_t(\mathbf{x}) - \hat{f}_t)/\sigma_t(\mathbf{x})$.

9.6. Regret of IDS.

We derive Theorem 9.11.

1. Prove that for all $t \geq 1$, $\hat{\Delta}_t(\mathbf{x}_t^{\text{UCB}}) \leq 2\beta_{t+1}\sigma_t(\mathbf{x}_t^{\text{UCB}})$ where we denote by $\mathbf{x}_t^{\text{UCB}} \doteq \arg \max_{\mathbf{x} \in \mathcal{X}} u_t(\mathbf{x})$ the point maximizing the upper confidence bound after iteration t .
2. Using (1) and the assumption on I_t , bound $\hat{\Psi}_t(\mathbf{x}_t^{\text{UCB}})$.
Hint: You may find the hint of Problem 9.3 (3) useful.
3. Complete the proof of Theorem 9.11.

9.7. Variational form of LITE.

Let $|\mathcal{X}| < \infty$ be such that we can write \mathcal{W} as an objective function over the probability simplex $\Delta^{|\mathcal{X}|} \subset \mathbb{R}^{|\mathcal{X}|}$. Derive Equation (9.36).

Hint: First show that $\mathcal{W}(\cdot)$ is concave (i.e., minimizing $-\mathcal{W}(\cdot)$ is a convex optimization problem) and then use Lagrange multipliers to find the optimum.

9.8. Finding the winning player.

Consider the betting game from Remark 9.13 with two risky players, $R_1 \sim \mathcal{N}(0, 100)$ and $R_2 \sim \mathcal{N}(0, 100)$, and one safe player $S = 1$. Prove that the individual probability of any of the risky players winning the most money is larger than the winning probability of the safe player. We assume that players payoffs are mutually independent.

Hint: Compute the CDF of $R \doteq \max\{R_1, R_2\}$.

10

Markov Decision Processes

We will now turn to the topic of probabilistic planning. Planning deals with the problem of deciding which action an agent should play in a (stochastic) environment.¹ A key formalism for probabilistic planning in *known* environments are so-called Markov decision processes. Starting from the next chapter, we will look at reinforcement learning, which extends probabilistic planning to unknown environments.

Consider the setting where we have a sequence of states $(X_t)_{t \in \mathbb{N}_0}$ similarly to Markov chains. But now, the next state X_{t+1} of an agent does not only depend on the previous state X_t but also depends on the last action A_t of this agent.

Definition 10.1 ((Finite) Markov decision process, MDP). A *(finite) Markov decision process* is specified by

- a (finite) set of *states* $X \doteq \{1, \dots, n\}$,
- a (finite) set of *actions* $A \doteq \{1, \dots, m\}$,
- *transition probabilities*

$$p(x' | x, a) \doteq \mathbb{P}(X_{t+1} = x' | X_t = x, A_t = a) \quad (10.1)$$

which is also called the *dynamics model*, and

- a *reward function* $r : X \times A \rightarrow \mathbb{R}$ which maps the current state x and an action a to some reward.

The reward function may also depend on the next state x' , however, we stick to the above model for simplicity. Also, the reward function can be random with mean r . Observe that r induces the sequence of rewards $(R_t)_{t \in \mathbb{N}_0}$, where

$$R_t \doteq r(X_t, A_t), \quad (10.2)$$

which is sometimes used in the literature instead of r .

Crucially, we assume the dynamics model p and the reward function r to be known. That is, we operate in a known environment. For now,

¹ An environment is *stochastic* as opposed to deterministic, when the outcome of actions is random.

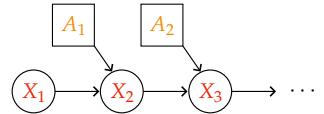


Figure 10.1: Directed graphical model of a Markov decision process with hidden states X_t and actions A_t .

we also assume that the environment is *fully observable*. In other words, we assume that our agent knows its current state. In Section 10.4, we discuss how this method can be extended to the partially observable setting.

Our fundamental objective is to learn how the agent should behave to optimize its reward. In other words, given its current state, the agent should decide (optimally) on the action to play. Such a decision map — whether optimal or not — is called a policy.

Definition 10.2 (Policy). A *policy* is a function that maps each state $x \in X$ to a probability distribution over the actions. That is, for any $t > 0$,

$$\pi(a | x) \doteq \mathbb{P}(A_t = a | X_t = x). \quad (10.3)$$

In other words, a policy assigns to each action $a \in A$, a probability of being played given the current state $x \in X$.

We assume that policies are stationary, that is, do not change over time.

Remark 10.3: Stochastic policies

We will see later in this chapter that in fully observable environments optimal policies are always deterministic. Thus, there is no need to consider stochastic policies in the context of Markov decision processes. For this chapter, you can think of a policy π simply as a deterministic mapping $\pi : X \rightarrow A$ from current state to played action. In the context of reinforcement learning, we will later see in Section 12.5.1 that randomized policies are important in trading exploration and exploitation.

Observe that a policy induces a Markov chain $(X_t^\pi)_{t \in \mathbb{N}_0}$ with transition probabilities,

$$p^\pi(x' | x) \doteq \mathbb{P}(X_{t+1}^\pi = x' | X_t^\pi = x) = \sum_{a \in A} \pi(a | x) p(x' | x, a). \quad (10.4)$$

This is crucial: if our agent follows a fixed policy (i.e., decision-making protocol) then the evolution of the process is described fully by a Markov chain.

As mentioned, we want to maximize the reward. There are many models of calculating a score from the infinite sequence of rewards $(R_t)_{t \in \mathbb{N}_0}$. For the purpose of our discussion of Markov decision processes and reinforcement learning, we will focus on a very common reward called discounted payoff.

Definition 10.4 (Discounted payoff). The *discounted payoff* (also called

discounted total reward) from time t is defined as the random variable,

$$G_t \doteq \sum_{m=0}^{\infty} \gamma^m R_{t+m} \quad (10.5)$$

where $\gamma \in [0, 1)$ is the *discount factor*.

Remark 10.5: Other reward models

Other well-known methods for combining rewards into a score are

$$\begin{array}{lll} (\text{instantaneous}) & (\text{finite-horizon}) & (\text{mean payoff}) \\ G_t \doteq R_t, & G_t \doteq \sum_{m=0}^{T-1} R_{t+m}, & \text{and} \quad G_t \doteq \liminf_{T \rightarrow \infty} \frac{1}{T} \sum_{m=0}^T R_{t+m}. \end{array}$$

The methods that we will discuss can also be analyzed using these or other alternative reward models.

We now want to understand the effect of the starting state and initial action on our optimization objective G_t . To analyze this, it is common to use the following two functions:

Definition 10.6 (State value function). The *state value function*,²

$$v_t^\pi(x) \doteq \mathbb{E}_\pi[G_t \mid X_t = x], \quad (10.7)$$

measures the average discounted payoff from time t starting from state $x \in X$.

Definition 10.7 (State-action value function). The *state-action value function* (also called *Q-function*),

$$q_t^\pi(x, a) \doteq \mathbb{E}_\pi[G_t \mid X_t = x, A_t = a] \quad (10.8)$$

$$= r(x, a) + \gamma \sum_{x' \in X} p(x' \mid x, a) \cdot v_{t+1}^\pi(x'), \quad (10.9)$$

measures the average discounted payoff from time t starting from state $x \in X$ and with playing action $a \in A$. In other words, it combines the immediate return with the value of the next states.

Note that both $v_t^\pi(x)$ and $q_t^\pi(x, a)$ are deterministic scalar-valued functions. Because we assumed stationary dynamics, rewards, and policies, the discounted payoff starting from a given state x will be independent of the start time t . Thus, we write $v^\pi(x) \doteq v_0^\pi(x)$ and $q^\pi(x, a) \doteq q_0^\pi(x, a)$ without loss of generality.

10.1 Bellman Expectation Equation

Let us now see how we can compute the value function,

$$v^\pi(x) = \mathbb{E}_\pi[G_0 \mid X_0 = x]$$

² Recall that following a fixed policy π induces a Markov chain $(X_t^\pi)_{t \in \mathbb{N}_0}$. We define

$$\mathbb{E}_\pi[\cdot] \doteq \mathbb{E}_{(X_t^\pi)_{t \in \mathbb{N}_0}}[\cdot] \quad (10.6)$$

as an expectation over all possible sequences of states $(x_t)_{t \in \mathbb{N}_0}$ within this Markov chain.

by expanding the defition of the discounted payoff (10.5); corresponds to one step in the induced Markov chain

using the definition of the value function (10.7)

$$\begin{aligned}
&= \mathbb{E}_\pi \left[\sum_{m=0}^{\infty} \gamma^m R_m \mid X_0 = x \right] \\
&= \mathbb{E}_\pi \left[\gamma^0 R_0 \mid X_0 = x \right] + \gamma \mathbb{E}_\pi \left[\sum_{m=0}^{\infty} \gamma^m R_{m+1} \mid X_0 = x \right] \\
&= r(x, \pi(x)) + \gamma \mathbb{E}_{x'} \left[\mathbb{E}_\pi \left[\sum_{m=0}^{\infty} \gamma^m R_{m+1} \mid X_1 = x' \right] \mid X_0 = x \right] \\
&= r(x, \pi(x)) + \gamma \sum_{x' \in X} p(x' \mid x, \pi(x)) \mathbb{E}_\pi \left[\sum_{m=0}^{\infty} \gamma^m R_{m+1} \mid X_1 = x' \right] \\
&= r(x, \pi(x)) + \gamma \sum_{x' \in X} p(x' \mid x, \pi(x)) \mathbb{E}_\pi \left[\sum_{m=0}^{\infty} \gamma^m R_m \mid X_0 = x' \right] \\
&= r(x, \pi(x)) + \gamma \sum_{x' \in X} p(x' \mid x, \pi(x)) \mathbb{E}_\pi [G_0 \mid X_0 = x'] \\
&= r(x, \pi(x)) + \gamma \sum_{x' \in X} p(x' \mid x, \pi(x)) \cdot v^\pi(x'). \tag{10.10}
\end{aligned}$$

using the definition of the discounted payoff (10.5)

using linearity of expectation (1.20)

by simplifying the first expectation and conditioning the second expectation on X_1

expanding the expectation on X_1 and using conditional independence of the discounted payoff of X_0 given X_1

shifting the start time of the discounted payoff using stationarity

using the definition of the discounted payoff (10.5)

using the definition of the value function (10.7)

interpreting the sum as an expectation

This equation is known as the *Bellman expectation equation*, and it shows a recursive dependence of the value function on itself. The intuition is clear: the value of the current state corresponds to the reward from the next action plus the discounted sum of all future rewards obtained from the subsequent states.

For stochastic policies, the above calculation can be extended to yield,

$$v^\pi(x) = \sum_{a \in A} \pi(a \mid x) \left(r(x, a) + \gamma \sum_{x' \in X} p(x' \mid x, a) v^\pi(x') \right) \tag{10.12}$$

$$= \mathbb{E}_{a \sim \pi(x)} [r(x, a) + \gamma \mathbb{E}_{x' \mid x, a} [v^\pi(x')]] \tag{10.13}$$

$$= \mathbb{E}_{a \sim \pi(x)} [q^\pi(x, a)]. \tag{10.14}$$

For stochastic policies, by also conditioning on the first action, one can obtain an analogous equation for the state-action value function,

$$q^\pi(x, a) = r(x, a) + \gamma \sum_{x' \in X} p(x' \mid x, a) \sum_{a' \in A} \pi(a' \mid x') q^\pi(x', a') \tag{10.15}$$

$$= r(x, a) + \gamma \mathbb{E}_{x' \mid x, a} \mathbb{E}_{a' \sim \pi(x')} [q^\pi(x', a')]. \tag{10.16}$$

Note that it does not make sense to consider a similar recursive formula for the state-action value function in the setting of deterministic policies as the action played when in state $x \in X$ is uniquely determined as $\pi(x)$. In particular,

$$v^\pi(x) = q^\pi(x, \pi(x)). \tag{10.17}$$

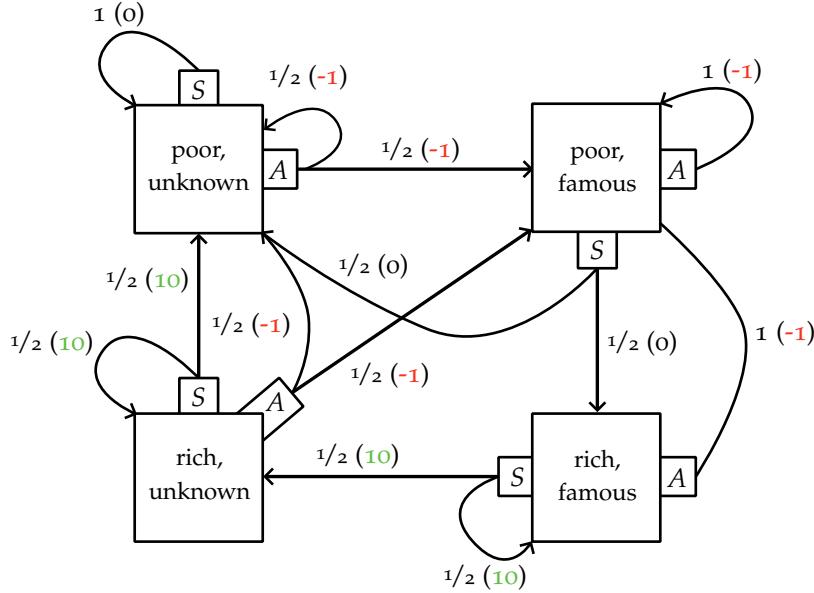


Figure 10.2: Example of an MDP, which we study in Problem 10.1. Suppose you are building a company. The shown MDP models “how to become rich and famous”. Here, the action S is short for *saving* and the action A is short for *advertising*.

Suppose you begin by being “poor and unknown”. Then, the greedy action (i.e., the action maximizing instantaneous reward) is to save. However, within this simplified environment, saving when you are poor and unknown means that you will remain poor and unknown forever. As the potential rewards in other states are substantially larger, this simple example illustrates that following the greedy choice is generally not optimal.

The example is adapted from Andrew Moore’s lecture notes on MDPs (Moore, 2002).

10.2 Policy Evaluation

Bellman’s expectation equation tells us how we can find the value function v^π of a fixed policy π using a system of linear equations! Using,

$$\begin{aligned} v^\pi &\doteq \begin{bmatrix} v^\pi(1) \\ \vdots \\ v^\pi(n) \end{bmatrix}, \quad r^\pi \doteq \begin{bmatrix} r(1, \pi(1)) \\ \vdots \\ r(n, \pi(n)) \end{bmatrix}, \quad \text{and} \\ P^\pi &\doteq \begin{bmatrix} p(1 | 1, \pi(1)) & \cdots & p(n | 1, \pi(1)) \\ \vdots & \ddots & \vdots \\ p(1 | n, \pi(n)) & \cdots & p(n | n, \pi(n)) \end{bmatrix} \end{aligned} \tag{10.18}$$

and a little bit of linear algebra, the Bellman expectation equation (10.10) is equivalent to

$$v^\pi = r^\pi + \gamma P^\pi v^\pi \tag{10.19}$$

$$\iff (I - \gamma P^\pi)v^\pi = r^\pi$$

$$\iff v^\pi = (I - \gamma P^\pi)^{-1}r^\pi. \tag{10.20}$$

Solving this linear system of equations (i.e., performing matrix inversion) takes cubic time in the size of the state space.

10.2.1 Fixed-point Iteration

To obtain an (approximate) solution of v^π , we can use that it is the unique fixed-point of the affine mapping $B^\pi : \mathbb{R}^n \rightarrow \mathbb{R}^n$,

$$B^\pi v \doteq r^\pi + \gamma P^\pi v. \quad (10.21)$$

Using this fact (which we will prove in just a moment), we can use fixed-point iteration of B^π .

Algorithm 10.8: Fixed-point iteration

```

1 initialize  $v^\pi$  (e.g., as  $\mathbf{0}$ )
2 for  $t = 1$  to  $T$  do
3    $v^\pi \leftarrow B^\pi v^\pi = r^\pi + \gamma P^\pi v^\pi$ 

```

Fixed-point iteration has computational advantages, for example, for sparse transitions.

Theorem 10.9. v^π is the unique fixed-point of B^π .

Proof. It is immediate from Bellman's expectation equation (10.19) and the definition of B^π (10.21) that v^π is a fixed-point of B^π . To prove uniqueness, we will show that B^π is a contraction.

Remark 10.10: Contractions

A *contraction* is a concept from topology. In a Banach space $(\mathcal{X}, \|\cdot\|)$ (a metric space with a norm), $f : \mathcal{X} \rightarrow \mathcal{X}$ is a contraction iff there exists a $k < 1$ such that

$$\|f(x) - f(y)\| \leq k \cdot \|x - y\| \quad (10.22)$$

for any $x, y \in \mathcal{X}$. By the *Banach fixed-point theorem*, a contraction admits a unique fixed-point. Intuitively, by iterating the function f , the distance to any fixed-point shrinks by a factor k in each iteration, hence, converges to 0. As we cannot converge to multiple fixed-points simultaneously, the fixed-point of a contraction f must be unique.

Let $v \in \mathbb{R}^n$ and $v' \in \mathbb{R}^n$ be arbitrary initial guesses. We use the L_∞ space,³

$$\begin{aligned} \|B^\pi v - B^\pi v'\|_\infty &= \|r^\pi + \gamma P^\pi v - r^\pi - \gamma P^\pi v'\|_\infty \\ &= \gamma \|P^\pi(v - v')\|_\infty \\ &\leq \gamma \max_{x \in X} \sum_{x' \in X} p(x' | x, \pi(x)) \cdot |v(x') - v'(x')|. \end{aligned}$$

³ The L_∞ norm (also called supremum norm) is defined as

$$\|x\|_\infty \doteq \max_i |x(i)|. \quad (10.23)$$

using the definition of B^π (10.21)

using the definition of the L_∞ norm (10.23), expanding the multiplication, and using $|\sum_i a_i| \leq \sum_i |a_i|$

$$\leq \gamma \|\mathbf{v} - \mathbf{v}'\|_\infty. \quad (10.24)$$

Thus, by Equation (10.22), \mathbf{B}^π is a contraction and by Banach's fixed-point theorem \mathbf{v}^π is its unique fixed-point. \square

Let \mathbf{v}_t^π be the value function estimate after t iterations. Then, we have for the convergence of fixed-point iteration,

$$\begin{aligned} \|\mathbf{v}_t^\pi - \mathbf{v}^\pi\|_\infty &= \|\mathbf{B}^\pi \mathbf{v}_{t-1}^\pi - \mathbf{B}^\pi \mathbf{v}^\pi\|_\infty \\ &\leq \gamma \|\mathbf{v}_{t-1}^\pi - \mathbf{v}^\pi\|_\infty \\ &= \gamma^t \|\mathbf{v}_0^\pi - \mathbf{v}^\pi\|_\infty. \end{aligned} \quad (10.25)$$

using the update rule of fixed-point iteration and $\mathbf{B}^\pi \mathbf{v}^\pi = \mathbf{v}^\pi$
using (10.24)
by induction

This shows that fixed-point iteration converges to \mathbf{v}^π exponentially fast.

10.3 Policy Optimization

Recall that our goal was to find an optimal policy,

$$\pi^* \doteq \arg \max_{\pi} \mathbb{E}_{\pi}[G_0]. \quad (10.26)$$

We can alternatively characterize an optimal policy as follows: We define a partial ordering over policies by

$$\pi \geq \pi' \iff v^\pi(x) \geq v^{\pi'}(x) \quad (\forall x \in X). \quad (10.27)$$

π^* is then simply a policy which is maximal according to this partial ordering.

It follows that all optimal policies have identical value functions. Subsequently, we use $v^* \doteq v^{\pi^*}$ and $q^* \doteq q^{\pi^*}$ to denote the state value function and state-action value function arising from an optimal policy, respectively. As an optimal policy maximizes the value of each state, we have that

$$v^*(x) = \max_{\pi} v^\pi(x), \quad q^*(x, a) = \max_{\pi} q^\pi(x, a). \quad (10.28)$$

Simply optimizing over each policy is not a good idea as there are m^n deterministic policies in total. It turns out that we can do much better.

10.3.1 Greedy Policies

Consider a policy that acts greedily according to the immediate return. It is fairly obvious that this policy will not perform well because the agent might never get to high-reward states. But what if someone could tell us not just the immediate return, but the long-term value of

the states our agent can reach in a single step? If we knew the value of each state our agent can reach, then we can simply pick the action that maximizes the expected value. We will make this approach precise in the next section.

This thought experiment suggests the definition of a greedy policy with respect to a value function.

Definition 10.11 (Greedy policy). The *greedy policy* with respect to a state-action value function q is defined as

$$\pi_q(x) \doteq \arg \max_{a \in A} q(x, a). \quad (10.29)$$

Analogously, we define the *greedy policy* with respect to a state value function v ,

$$\pi_v(x) \doteq \arg \max_{a \in A} r(x, a) + \gamma \sum_{x' \in X} p(x' | x, a) \cdot v(x'). \quad (10.30)$$

We can use v and q interchangeably $\textcircled{?}$.

Problem 10.2

10.3.2 Bellman Optimality Equation

Observe that following the greedy policy π_v , will lead us to a new value function v^{π_v} . With respect to this value function, we can again obtain a greedy policy, of which we can then obtain a new value function. In this way, the correspondence between greedy policies and value functions induces a cyclic dependency, which is visualized in Figure 10.3.

It turns out that the optimal policy π^* is a fixed-point of this dependency. This is made precise by the following theorem.

Theorem 10.12 (Bellman's theorem). *A policy π^* is optimal iff it is greedy with respect to its own value function. In other words, π^* is optimal iff $\pi^*(x)$ is a distribution over the set $\arg \max_{a \in A} q^*(x, a)$.*

In particular, if for every state there is a unique action that maximizes the state-action value function, the policy π^* is deterministic and unique,

$$\pi^*(x) = \arg \max_{a \in A} q^*(x, a). \quad (10.31)$$

Proof. It is a direct consequence of Equation (10.28) that a policy is optimal iff it is greedy with respect to q^* . \square

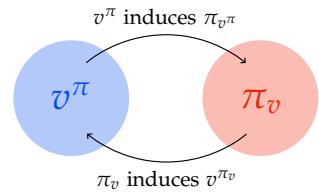


Figure 10.3: Cyclic dependency between **value function** and **greedy policy**.

This theorem confirms our intuition from the previous section that greedily following an optimal value function is itself optimal. In particular, Bellman's theorem shows that there always exists an optimal policy which is deterministic and stationary.

We have seen, that π^* is a fixed-point of greedily picking the best action according to its state-action value function. The converse is also true:

Corollary 10.13. *The optimal value functions v^* and q^* are a fixed-point of the so-called Bellman update,*

$$v^*(x) = \max_{a \in A} q^*(x, a), \quad (10.32)$$

$$= \max_{a \in A} r(x, a) + \gamma \mathbb{E}_{x' | x, a} [v^*(x')] \quad (10.33)$$

$$q^*(x, a) = r(x, a) + \gamma \mathbb{E}_{x' | x, a} \left[\max_{a' \in A} q^*(x', a') \right]. \quad (10.34)$$

using the definition of the q -function (10.9)

Proof. It follows from Equation (10.14) that

$$v^*(x) = \mathbb{E}_{a \sim \pi^*(x)} [q^*(x, a)]. \quad (10.35)$$

Thus, as π^* is greedy with respect to q^* , $v^*(x) = \max_{a \in A} q^*(x, a)$.

Equation (10.34) follows analogously from Equation (10.15). \square

These equations are also called the *Bellman optimality equations*. Intuitively, the Bellman optimality equations express that the value of a state under an optimal policy must equal the expected return for the best action from that state. Bellman's theorem is also known as *Bellman's optimality principle*, which is a more general concept.

The two perspectives of Bellman's theorem naturally suggest two separate ways of finding the optimal policy. Policy iteration uses the perspective from Equation (10.31) of π^* as a fixed-point of the dependency between greedy policy and value function. In contrast, value iteration uses the perspective from Equation (10.32) of v^* as the fixed-point of the Bellman update. Another approach which we will not discuss here is to use a linear program where the Bellman update is interpreted as a set of linear inequalities.

Bellman's optimality principle Bellman's optimality equations for MDPs are one of the main settings of Bellman's optimality principle. However, Bellman's optimality principle has many other important applications, for example in dynamic programming. Broadly speaking, Bellman's optimality principle says that optimal solutions to decision problems can be decomposed into optimal solutions to sub-problems.

10.3.3 Policy Iteration

Starting from an arbitrary initial policy, policy iteration as shown in Algorithm 10.14 uses the Bellman expectation equation to compute the value function of that policy (as we have discussed in Section 10.2) and then chooses the greedy policy with respect to that value function as its next iterate.

Algorithm 10.14: Policy iteration

```

1 initialize  $\pi$  (arbitrarily)
2 repeat
3   | compute  $v^\pi$ 
4   | compute  $\pi_{v^\pi}$ 
5   |  $\pi \leftarrow \pi_{v^\pi}$ 
6 until converged

```

Let π_t be the policy after t iterations. We will now show that policy iteration converges to the optimal policy. The proof is split into two parts. First, we show that policy iteration improves policies monotonically. Then, we will use this fact to show that policy iteration converges.

Lemma 10.15 (Monotonic improvement of policy iteration). *We have,*

- $v^{\pi_{t+1}}(x) \geq v^{\pi_t}(x)$ for all $x \in X$; and
- $v^{\pi_{t+1}}(x) > v^{\pi_t}(x)$ for at least one $x \in X$, unless $v^{\pi_t} \equiv v^*$.

Proof. We consider the Bellman update from (10.32) as the mapping $B^* : \mathbb{R}^n \rightarrow \mathbb{R}^n$,

$$(B^*v)(x) \doteq \max_{a \in A} q(x, a), \quad (10.36)$$

where q is the state-action value function corresponding to the state value function $v \in \mathbb{R}^n$. Recall that after obtaining v^{π_t} , policy iteration first computes the greedy policy w.r.t. v^{π_t} , $\pi_{t+1} \doteq \pi_{v^{\pi_t}}$, and then computes its value function $v^{\pi_{t+1}}$.

To establish the (weak) monotonic improvement of policy iteration, we consider a fixed-point iteration (cf. Algorithm 10.8) of $v^{\pi_{t+1}}$ initialized by v^{π_t} . We denote the iterates by \tilde{v}_τ , in particular, we have that $\tilde{v}_0 = v^{\pi_t}$ and $\lim_{\tau \rightarrow \infty} \tilde{v}_\tau = v^{\pi_{t+1}}$.⁴ First, observe that for the first iteration of fixed-point iteration,

$$\begin{aligned} \tilde{v}_1(x) &= (B^*v^{\pi_t})(x) \\ &= \max_{a \in A} q^{\pi_t}(x, a) \\ &\geq q^{\pi_t}(x, \pi_t(x)) \\ &= v^{\pi_t}(x) \\ &= \tilde{v}_0(x). \end{aligned}$$

Let us now consider a single iteration of fixed-point iteration. We have,

$$\tilde{v}_{\tau+1}(x) = r(x, \pi_{t+1}(x)) + \gamma \sum_{x' \in X} p(x' | x, \pi_{t+1}(x)) \cdot \tilde{v}_\tau(x').$$

⁴ using the convergence of fixed-point iteration (10.25)

using that π_{t+1} is greedy wrt. v^{π_t}

using the definition of the Bellman update (10.36)

using (10.17)

using the definition of $\tilde{v}_{\tau+1}$ (10.21)

Using an induction on τ , we conclude,

$$\begin{aligned} &\geq r(x, \pi_{t+1}(x)) + \gamma \sum_{x' \in X} p(x' | x, \pi_{t+1}(x)) \cdot \tilde{v}_{\tau-1}(x') \\ &= \tilde{v}_\tau(x). \end{aligned}$$

using the induction hypothesis,
 $\tilde{v}_\tau(x') \geq \tilde{v}_{\tau-1}(x')$

This establishes the first claim,

$$v^{\pi_{t+1}} = \lim_{\tau \rightarrow \infty} \tilde{v}_\tau \geq \tilde{v}_0 = v^{\pi_t}. \quad (10.37)$$

For the second claim, recall from Bellman's theorem (10.32) that v^* is a (unique) fixed-point of the Bellman update B^* .⁵ In particular, we have $v^{\pi_{t+1}} \equiv v^{\pi_t}$ if and only if $v^{\pi_{t+1}} \equiv v^{\pi_t} \equiv v^*$. In other words, if $v^{\pi_t} \not\equiv v^*$ then Equation (10.37) is strict for at least one $x \in X$ and $v^{\pi_{t+1}} \not\equiv v^{\pi_t}$. This proves the strict monotonic improvement of policy iteration. \square

⁵ We will show in Equation (10.38) that B^* is a contraction, implying that v^* is the unique fixed-point of B^* .

Theorem 10.16 (Convergence of policy iteration). *For finite Markov decision processes, policy iteration converges to an optimal policy.*

Proof. Finite Markov decision processes only have a finite number of deterministic policies (albeit exponentially many). Observe that policy iteration only considers deterministic policies, and recall that there is an optimal policy that is deterministic. As the value of policies strictly increase in each iteration until an optimal policy is found, policy iteration must converge in finite time. \square

It can be shown that policy iteration converges to an exact solution in a polynomial number of iterations (Ye, 2011). Each iteration of policy iteration requires computing the value function, which we have seen to be of cubic complexity in the number of states.

10.3.4 Value Iteration

As we have mentioned, another natural approach of finding the optimal policy is to interpret v^* as the fixed point of the Bellman update. Recall our definition of the Bellman update from Equation (10.36),

$$(B^*v)(x) = \max_{a \in A} q(x, a),$$

where q was the state-action value function associated with the state value function v . The value iteration algorithm is shown in Algorithm 10.17.

We will now prove the convergence of value iteration using the fixed-point interpretation.

Theorem 10.18 (Convergence of value iteration). *Value iteration converges asymptotically to an optimal policy.*

Algorithm 10.17: Value iteration

-
- 1 initialize $v(x) \leftarrow \max_{a \in A} r(x, a)$ for each $x \in X$
 - 2 **for** $t = 1$ **to** ∞ **do**
 - 3 $v(x) \leftarrow (B^*v)(x) = \max_{a \in A} q(x, a)$ for each $x \in X$
 - 4 choose π_v
-

Proof. Clearly, value iteration converges if v^* is the unique fixed-point of B^* . We already know from Bellman's theorem (10.32) that v^* is a fixed-point of B^* . It remains to show that it is indeed the unique fixed-point.

Analogously to our proof of the convergence of fixed-point iteration to the value function v^π , we show that B^* is a contraction. Fix arbitrary $v, v' \in \mathbb{R}^n$, then

$$\begin{aligned}
\|B^*v - B^*v'\|_\infty &= \max_{x \in X} |(B^*v)(x) - (B^*v')(x)| \\
&= \max_{x \in X} \left| \max_{a \in A} q(x, a) - \max_{a \in A} q'(x, a) \right| \\
&\leq \max_{x \in X} \max_{a \in A} |q(x, a) - q'(x, a)| \\
&\leq \gamma \max_{x \in X} \max_{a \in A} \sum_{x' \in X} p(x' | x, a) |v(x') - v'(x')| \\
&\leq \gamma \|v - v'\|_\infty
\end{aligned} \tag{10.38}$$

where q and q' are the state-action value functions associated with v and v' , respectively. By Equation (10.22), B^* is a contraction and by Banach's fixed-point theorem v^* is its unique fixed-point. \square

Remark 10.19: Value iteration as a dynamic program

Let us denote by v_t the value function estimate after the t -th iteration. Observe that $v_t(x)$ corresponds to the maximum expected reward when starting in state x and the “world ends” after t time steps. In particular, v_0 corresponds to the maximum immediate reward. This suggests a different perspective on value iteration (akin to dynamic programming) where in each iteration we extend the time horizon of our approximation by one time step.

For any $\epsilon > 0$, value iteration converges to an ϵ -optimal solution in polynomial time. However, unlike policy iteration, value iteration does not generally reach the *exact* optimum in a finite number of iterations. Recalling the update rule of value iteration, its main benefit is that each iteration only requires a sum over all possible actions a in state

using the definition of the L_∞ norm
(10.23)

using the definition of the Bellman update (10.36)

using $|\max_x f(x) - \max_x g(x)| \leq \max_x |f(x) - g(x)|$
using the definition of the Q-function (10.9) and $|\sum_i a_i| \leq \sum_i |a_i|$

using $\sum_{x' \in X} p(x' | x, a) = 1$ and
 $|v(x') - v'(x')| \leq \|v - v'\|_\infty$

x and a sum over all reachable states x' from x . In sparse Markov decision processes,⁶ an iteration of value iteration can be performed in (virtually) constant time.

10.4 Partial Observability

So far we have focused on the fully observable setting. That is, at any time, our agent knows its current state. We have seen that we can efficiently find the optimal policy (as long as the Markov decision process is finite).

We have already encountered the partially observable setting in Chapter 3, where we discussed filtering. In this section, we consider how Markov decision processes can be extended to a partially observable setting where the agent can only access noisy observations Y_t of its state X_t .

Definition 10.20 (Partially observable Markov decision process, POMDP).

Similarly to a Markov decision process, a *partially observable Markov decision process* is specified by

- a set of *states* X ,
- a set of *actions* A ,
- *transition probabilities* $p(x' | x, a)$, and
- a *reward function* $r : X \times A \rightarrow \mathbb{R}$.

Additionally, it is specified by

- a set of *observations* Y , and
- *observation probabilities*

$$o(y | x) \doteq \mathbb{P}(Y_t = y | X_t = x). \quad (10.39)$$

Whereas MDPs are controlled Markov chains, POMDPs are controlled hidden Markov models.

Remark 10.21: Hidden Markov models

A hidden Markov model is a Markovian process with unobservable states X_t and observations Y_t that depend on X_t in a known way.

Definition 10.22 (Hidden Markov model, HMM). A *hidden Markov model* is specified by

- a set of *states* X ,
- *transition probabilities* $p(x' | x) \doteq \mathbb{P}(X_{t+1} = x' | X_t = x)$ (also called *motion model*), and
- a *sensor model* $o(y | x) \doteq \mathbb{P}(Y_t = y | X_t = x)$.

Following from its directed graphical model shown in Figure 10.5,

⁶Sparsity refers to the interconnectivity of the state space. When only few states are reachable from any state, we call an MDP sparse.

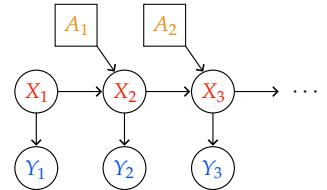


Figure 10.4: Directed graphical model of a partially observable Markov decision process with hidden states X_t , observables Y_t , and actions A_t .

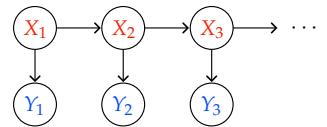


Figure 10.5: Directed graphical model of a hidden Markov model with hidden states X_t and observables Y_t .

its joint probability distribution factorizes into

$$\mathbb{P}(x_{1:t}, y_{1:t}) = \mathbb{P}(x_1) \cdot o(y_1 | x_1) \cdot \prod_{i=2}^t p(x_i | x_{i-1}) \cdot o(y_i | x_i). \quad (10.40)$$

Observe that a Kalman filter can be viewed as a hidden Markov model with conditional linear Gaussian motion and sensor models and a Gaussian prior on the initial state. In particular, the tasks of *filtering*, *smoothing*, and *predicting* which we discussed extensively in Chapter 3 are also of interest for hidden Markov models.

A widely used application of hidden Markov models is to find the most likely sequence (also called *most likely explanation*) of hidden states $x_{1:t}$ given a series of observations $y_{1:t}$,⁷ that is, to find

$$\arg \max_{x_{1:t}} \mathbb{P}(x_{1:t} | y_{1:t}). \quad (10.41)$$

This task can be solved in linear time by a simple backtracking algorithm known as the *Viterbi algorithm*.

POMDPs are a very powerful model, but very hard to solve in general. POMDPs can be reduced to a Markov decision process with an enlarged state space. The key insight is to consider an MDP whose states are the *beliefs*,

$$b_t(x) \doteq \mathbb{P}(X_t = x | y_{1:t}, a_{1:t-1}), \quad (10.42)$$

about the current state in the POMDP. In other words, the states of the MDP are probability distributions over the states of the POMDP. We will make this more precise in the following.

Let us assume that our prior belief about the state of our agent is given by $b_0(x) \doteq \mathbb{P}(X_0 = x)$. Keeping track of how beliefs change over time is known as *filtering*, which we already encountered in Section 3.1. Given a prior belief b_t , an action taken a_t , and a new observation y_{t+1} , the belief state can be updated as follows,

$$\begin{aligned} b_{t+1}(x) &= \mathbb{P}(X_{t+1} = x | y_{1:t+1}, a_{1:t}) \\ &= \frac{1}{Z} \mathbb{P}(y_{t+1} | X_{t+1} = x) \mathbb{P}(X_{t+1} = x | y_{1:t}, a_{1:t}) \\ &= \frac{1}{Z} o(y_{t+1} | x) \mathbb{P}(X_{t+1} = x | y_{1:t}, a_{1:t}) \\ &= \frac{1}{Z} o(y_{t+1} | x) \sum_{x' \in X} p(x | x', a_t) \mathbb{P}(X_t = x' | y_{1:t}, a_{1:t-1}) \\ &= \frac{1}{Z} o(y_{t+1} | x) \sum_{x' \in X} p(x | x', a_t) b_t(x') \end{aligned} \quad (10.43)$$

⁷ This is useful in many applications such as speech recognition, decoding data that was transmitted over a noisy channel, beat detection, and many more.

by the definition of beliefs (10.42)

using Bayes' rule (1.45)

using the definition of observation probabilities (10.39)

by conditioning on the previous state x' , noting a_t does not influence X_t

using the definition of beliefs (10.42)

where

$$Z \doteq \sum_{x \in X} o(y_{t+1} | x) \sum_{x' \in X} p(x | x', a_t) b_t(x'). \quad (10.44)$$

Thus, the updated belief state is a deterministic mapping from the previous belief state depending only on the (random) observation y_{t+1} and the taken action a_t . Note that this obeys a Markovian structure of transition probabilities with respect to the beliefs b_t .

The sequence of belief-states defines the sequence of random variables $(B_t)_{t \in \mathbb{N}_0}$,

$$B_t \doteq X_t | y_{1:t}, a_{1:t-1}, \quad (10.45)$$

where the (state-)space of all beliefs is the (infinite) space of all probability distributions over X ,⁸

$$\mathcal{B} \doteq \Delta^X \doteq \left\{ \mathbf{b} \in \mathbb{R}^{|X|} : \mathbf{b} \geq \mathbf{0}, \sum_{i=1}^{|X|} \mathbf{b}(i) = 1 \right\}. \quad (10.46)$$

A Markov decision process, where every belief corresponds to a state is called a belief-state MDP.

Definition 10.23 (Belief-state Markov decision process). Given a POMDP, the corresponding *belief-state Markov decision process* is a Markov decision process specified by

- the *belief space* $\mathcal{B} \doteq \Delta^X$ depending on the *hidden states* X ,
- the set of *actions* A ,
- *transition probabilities*

$$\tau(b' | b, a) \doteq \mathbb{P}(B_{t+1} = b' | B_t = b, A_t = a), \quad (10.47)$$

- and *rewards*

$$\rho(b, a) \doteq \mathbb{E}_{x \sim b}[r(x, a)] = \sum_{x \in X} b(x)r(x, a). \quad (10.48)$$

It remains to derive the transition probabilities τ in terms of the original POMDP. We have,

$$\begin{aligned} \tau(b_{t+1} | b_t, a_t) &= \mathbb{P}(b_{t+1} | b_t, a_t) \\ &= \sum_{y_{t+1} \in Y} \mathbb{P}(b_{t+1} | b_t, a_t, y_{t+1}) \mathbb{P}(y_{t+1} | b_t, a_t). \end{aligned} \quad (10.49) \quad \text{by conditioning on } y_{t+1} \in Y$$

Using the Markovian structure of the belief updates, we naturally set,

$$\mathbb{P}(b_{t+1} | b_t, a_t, y_{t+1}) \doteq \begin{cases} 1 & \text{if } b_{t+1} \text{ matches the belief update} \\ & \text{of Equation (10.42) given } b_t, a_t, \text{ and} \\ & y_{t+1}, \\ 0 & \text{otherwise.} \end{cases} \quad (10.50)$$

⁸This definition naturally extends to continuous state spaces X .

The final missing piece is the likelihood of an observation y_{t+1} given the prior belief b_t and action a_t , which using our interpretation of beliefs corresponds to

$$\begin{aligned}\mathbb{P}(y_{t+1} | b_t, a_t) &= \mathbb{E}_{x \sim b_t} \left[\mathbb{E}_{x' | x, a_t} [\mathbb{P}(y_{t+1} | X_{t+1} = x')] \right] \\ &= \mathbb{E}_{x \sim b_t} \left[\mathbb{E}_{x' | x, a_t} [o(y_{t+1} | x')] \right] \\ &= \sum_{x \in X} b_t(x) \sum_{x' \in X} p(x' | x, a_t) \cdot o(y_{t+1} | x').\end{aligned}\quad (10.51)$$

using the definition of observation probabilities (10.39)

In principle, we can now apply arbitrary algorithms for planning in MDPs to POMDPs. Of course, the problem is that there are infinitely many beliefs, even for a finite state space X .⁹ The belief-state MDP has therefore an infinitely large belief space \mathcal{B} . Even when only planning over finite horizons, exponentially many beliefs can be reached. So the belief space blows-up very quickly.

We will study MDPs with large state spaces (where transition dynamics and rewards are unknown) in Chapters 12 and 13. Similar methods can also be used to approximately solve POMDPs.

A key idea in approximate solutions to POMDPs is that most belief states are never reached. A common approach is to discretize the belief space by sampling or by applying a dimensionality reduction. Examples are *point-based value iteration* (PBVI) and *point-based policy iteration* (PBPI) (Shani et al., 2013).

Discussion

Even though we focus on the fully observed setting throughout this manuscript, the partially observed setting can be reduced to the fully observed setting with very large state spaces. In the next chapter, we will consider learning and planning in unknown Markov decision processes (i.e., reinforcement learning) for small state spaces. The setting of small state and action spaces is also known as the *tabular setting*. Then, in the final two chapters, we will consider approximate methods for large state and action spaces. In particular, in Section 13.1, we will revisit the problem of probabilistic planning in known Markov decision processes, but with continuous state and action spaces.

⁹ You can think of an $|X|$ -dimensional space. Here, all points whose coordinates sum to 1 correspond to probability distributions (i.e., beliefs) over the hidden states X . The convex hull of these points is also known as the $(|X| - 1)$ -dimensional *probability simplex* (cf. Appendix A.1.2). Now, by definition of the $(|X| - 1)$ -dimensional probability simplex as a polytope in $|X| - 1$ dimensions, we can conclude that its boundary consists of infinitely many points in $|X|$ dimensions. Noting that these points corresponded to the probability distributions on $|X|$, we conclude that there are infinitely many such distributions.

Problems

10.1. Value functions.

Recall the example of “becoming rich and famous” from Figure 10.2. Consider the policy, $\pi \equiv S$ (i.e., to always *save*) and let $\gamma = 1/2$. Show that the (rounded) state-action value function q^π is as follows:

	save	advertise
poor, unknown	0	0.1
poor, famous	4.4	1.2
rich, famous	17.8	1.2
rich, unknown	13.3	0.1

Shown in bold is the state value function v^π .

10.2. Greedy policies.

Show that if q and v arise from the same policy, that is, q is defined in terms of v as per Equation (10.9), then

$$\pi_v \equiv \pi_q. \quad (10.52)$$

This implies that we can use v and q interchangeably.

10.3. Optimal policies.

Again, recall the example of “becoming rich and famous” from Figure 10.2.

1. Show that the policy $\pi \equiv S$, which we considered in Problem 10.1, is not optimal.
2. Instead, consider the policy

$$\pi' \equiv \begin{cases} A & \text{if poor and unknown} \\ S & \text{otherwise} \end{cases}$$

and let $\gamma = 1/2$. Show that the (rounded) state-action value function $q^{\pi'}$ is as follows:

	save	advertise
poor, unknown	0.8	1.6
poor, famous	4.5	1.2
rich, famous	17.8	1.2
rich, unknown	13.4	0.2

Shown in bold is the state value function $v^{\pi'}$.

3. Is the policy π' optimal?

10.4. Linear convergence of policy iteration.

Denote by π_t the policy obtained by policy iteration after t iterations. Use that the Bellman operator B^* is a contraction with the unique fixed-point v^* to show that

$$\|v^{\pi_t} - v^*\|_\infty \leq \gamma^t \|v^{\pi_0} - v^*\|_\infty \quad (10.53)$$

where v^π and v^* are vector representations of the functions v^π and v^* , respectively.

*Hint: Recall from Lemma 10.15 that $v^{\pi_{t+1}} \geq B^*v^{\pi_t} \geq v^{\pi_t}$.*

10.5. Reward modification.

A key technique for solving sequential decision problems is the modification of reward functions that leaves the optimal policy unchanged while improving sample efficiency or convergence rates. This exercise looks at simple ways of modifying rewards and understanding how these modifications affect the optimal policy.

Consider two Markov decision processes $M \doteq (X, A, p, r)$ and $M' \doteq (X, A, p, r')$ where the reward function r is modified to obtain r' , and the rewards are bounded and discounted by the discount factor $\gamma \in [0, 1]$. Let π_M^* be the optimal policy for M .

1. Suppose $r'(x) = \alpha r(x)$, where $\alpha > 0$. Show that the optimal policy π^* of M is also an optimal policy of M' .
2. Given a modification of the form $r'(x) = r(x) + c$, where $c > 0$ is a constant scalar, show that the optimal policy π_M^* can be different from $\pi_{M'}^*$.
3. Another way of modifying the reward function is through *reward shaping* where one supplies additional rewards to the agent to guide the learning process. When one has no knowledge of the underlying transition dynamics p , a commonly used transformation is $r'(x, x') = r(x, x') + f(x, x')$ where f is a *potential-based* shaping function defined as

$$f(x, x') \doteq \gamma\phi(x') - \phi(x), \quad \phi : X \rightarrow \mathbb{R}. \quad (10.54)$$

Show that the optimal policy remains unchanged under this definition of f .

10.6. A partially observable fishing problem.

We model an angler's decisions while fishing, where the states are partially observable. There are two states: (1) **Fish** (F): A fish is hooked on the line. (2) **No fish** (\bar{F}): No fish is hooked on the line.

The angler can choose between two actions:

- **Pull up the rod (P):** If there is a fish on the line (F), there is a 90% chance of catching it (reward +10, transitioning to \bar{F}) and a 10% chance of it escaping (reward -1, transitioning to \bar{F}). If there is no fish (\bar{F}), pulling up the rod results in no catch, staying in \bar{F} with a reward of -5.
- **Waiting (W):** All waiting actions result in a reward of -1. In state F , there is a 60% chance of the fish staying (remaining in F) and a 40% chance of it escaping (transitioning to \bar{F}). In state \bar{F} , there is a 50% chance of a fish biting (transitioning to F) and a 50% chance of no change (remaining in \bar{F}).

Suggestion: Draw the MDP transition diagram. Draw each transition with action, associated probability, and associated reward.

Since the angler cannot directly observe whether there is a fish on the line, they receive a noisy observation about the state. This observation can be:

- o_1 : The signal suggests that a fish might be on the line.
- o_2 : The signal suggests that there is no fish on the line.

The observation model, which defines the probability of receiving each observation given the true state is as follows:

	$\mathbb{P}(o_1 \mid \cdot)$	$\mathbb{P}(o_2 \mid \cdot)$
F	0.8	0.2
\bar{F}	0.3	0.7

The angler's goal is to choose actions that maximize their overall reward, balancing the chances of catching a fish against the cost of waiting and unsuccessful pulls.

1. Given an initial belief $b_0(F) = b_0(\bar{F}) = 0.5$, the angler chooses to wait and observes o_1 . Compute the updated belief b_1 using the observation model and belief update equation (10.42).
2. Given belief $b_1(F) \approx 0.765$ and $b_1(\bar{F}) \approx 0.235$, compute the updated belief b_2 for both actions P (pull) and W (wait), both in the case where you observe o_1 (fish likely) and o_2 (fish unlikely).

11

Tabular Reinforcement Learning

11.1 The Reinforcement Learning Problem

Reinforcement learning is concerned with probabilistic planning in unknown environments. This extends our study of known environments in the previous chapter. Those environments are still modeled by Markov decision processes, but in reinforcement learning, we do not know the dynamics p and rewards r in advance. Hence, reinforcement learning is at the intersection of the theories of probabilistic planning (i.e., Markov decision processes) and learning (e.g., multi-armed bandits), which we covered extensively in the previous chapters.

We will continue to focus on the fully observed setting, where the agent knows its current state. As we have seen in the previous section, the partially observed setting corresponds to a fully observed setting with an enlarged state space. In this chapter, we will begin by considering reinforcement learning with small state and action spaces. This setting is often called the *tabular setting*, as the value functions can be computed exhaustively for all states and stored in a table.

Clearly, the agent needs to trade exploring and learning about the environment with exploiting its knowledge to maximize rewards. Thus, the exploration-exploitation dilemma, which was at the core of Bayesian optimization (see Section 9.1), also plays a crucial role in reinforcement learning. In fact, Bayesian optimization can be viewed as reinforcement learning with a fixed state: In each round, the agent plays an action, aiming to find the action that maximizes the reward. However, playing the same action multiple times yields the same reward, implying that we remain in a single state. In the context of Bayesian optimization, we used “regret” as performance metric: in the jargon of planning, minimizing regret corresponds to maximizing the cumulative reward.

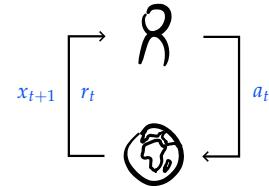


Figure 11.1: In reinforcement learning, an agent interacts with its environment in a sequence of rounds. After playing an action a_t , it observes rewards r_t and its new state x_{t+1} . The agent then uses this information to learn how to act to maximize reward.

Another key challenge of reinforcement learning is that the observed data is dependent on the played actions. This is in contrast to the setting of supervised learning that we have been considering in earlier chapters, where the data is sampled independently.

11.1.1 Trajectories

The data that the agent collects is modeled using so-called trajectories.

Definition 11.1 (Trajectory). A *trajectory* τ is a (possibly infinite) sequence,

$$\tau \doteq (\tau_0, \tau_1, \tau_2, \dots), \quad (11.1)$$

of *transitions*,

$$\tau_i \doteq (x_i, a_i, r_i, x_{i+1}), \quad (11.2)$$

where $x_i \in X$ is the starting state, $a_i \in A$ is the played action, $r_i \in \mathbb{R}$ is the attained reward, and $x_{i+1} \in X$ is the ending state.

In the context of learning a dynamics and rewards model, x_i and a_i can be understood as inputs, and r_i and x_{i+1} can be understood as labels of a regression problem.

Crucially, the newly observed states x_{t+1} and the rewards r_t (across multiple transitions) are conditionally independent given the previous states x_t and actions a_t . This follows directly from the Markovian structure of the underlying Markov decision process.¹ Formally, we have,

$$X_{t+1} \perp X_{t'+1} \mid X_t, X_{t'}, A_t, A_{t'}, \quad (11.3a)$$

$$R_t \perp R_{t'} \mid X_t, X_{t'}, A_t, A_{t'}, \quad (11.3b)$$

for any $t, t' \in \mathbb{N}_0$. In particular, if $x_t = x_{t'}$ and $a_t = a_{t'}$, then x_{t+1} and $x_{t'+1}$ are independent samples according to the transition model $p(X_{t+1} \mid x_t, a_t)$. Analogously, if $x_t = x_{t'}$ and $a_t = a_{t'}$, then r_t and $r_{t'}$ are independent samples of the reward model $r(x_t, a_t)$. As we will see later in this chapter and especially in Chapter 13, this independence property is crucial for being able to learn about the underlying Markov decision process. Notably, this implies that we can apply the law of large numbers (A.36) and Hoeffding's inequality (A.41) to our estimators of both quantities.

The collection of data is commonly classified into two settings. In the *episodic setting*, the agent performs a sequence of “training” rounds (called *episodes*). In the beginning of each episode, the agent is reset to some initial state. In contrast, in the *continuous setting* (or non-episodic /

¹ Recall the Markov property (6.6), which assumes that in the underlying Markov decision process (i.e., in our environment) the future state of an agent is independent of past states given the agent’s current state. This is commonly called a Markovian structure. From this Markovian structure, we gather that repeated encounters of state-action pairs result in independent trials of the transition model and rewards.

online setting), the agent learns online. Especially, every action, every reward, and every state transition counts.

The episodic setting is more applicable to an agent playing a computer game. That is, the agent is performing in a simulated environment that is easy to reset. The continuous setting is akin to an agent that is deployed to the “real world”. In principle, real-world agents can be trained in simulated environments before being deployed. However, this bears the risk of learning to exploit or rely on features of the simulated environment that are not present in the real environment. Sometimes, using a simulated environment for training is downright impossible, as the real environment is too complex.

11.1.2 *On-policy and Off-policy Methods*

Another important distinction in how data is collected, is the distinction between on-policy and off-policy methods. As the names suggest, *on-policy* methods are used when the agent has control over its own actions, in other words, the agent can freely choose to follow any policy. Being able to follow a policy is helpful, for example because it allows the agent to experiment with trading exploration and exploitation.

In contrast, *off-policy* methods can be used even when the agent cannot freely choose its actions. Off-policy methods are therefore able to make use of purely observational data. This might be data that was collected by another agent, a fixed policy, or during a previous episode. Off-policy methods are therefore more *sample-efficient* than on-policy methods. This is crucial, especially in settings where conducting experiments (i.e., collecting new data) is expensive.

11.2 *Model-based Approaches*

Approaches to reinforcement learning are largely categorized into two classes. *Model-based* approaches aim to learn the underlying Markov decision process. More concretely, they learn models of the dynamics p and rewards r . They then use these models to perform planning (i.e., policy optimization) in the underlying Markov decision process. In contrast, *model-free* approaches learn the value function directly. We begin by discussing model-based approaches to the tabular setting. In Section 11.4, we cover model-free approaches.

11.2.1 *Learning the Underlying Markov Decision Process*

Recall that the underlying Markov decision process was specified by its dynamics $p(x' | x, a)$ that correspond to the probability of entering

state $x' \in X$ when playing action $a \in A$ from state $x \in X$, and its rewards $r(x, a)$ for playing action $a \in A$ in state $x \in X$. A natural first idea is to use maximum likelihood estimation to approximate these quantities.

We can think of each transition $x' | x, a$ as sampling from a categorical random variable of which we want to estimate the success probabilities for landing in each of the states. Therefore, as we have seen in Example A.7, the MLE of the dynamics model coincides with the sample mean,

$$\hat{p}(x' | x, a) = \frac{N(x' | x, a)}{N(a | x)} \quad (11.4)$$

where $N(x' | x, a)$ counts the number of transitions from state x to state x' when playing action a and $N(a | x)$ counts the number of transitions that start in state x and play action a (regardless of the next state). Similarly, for the rewards model, we obtain the following maximum likelihood estimate (i.e., sample mean),

$$\hat{r}(x, a) = \frac{1}{N(a | x)} \sum_{\substack{t=0 \\ x_t=x \\ a_t=a}}^{\infty} r_t. \quad (11.5)$$

It is immediate that both estimates are unbiased as both correspond to a sample mean.

Still, for the models of our environment to become accurate, our agent needs to visit *each* state-action pair (x, a) numerous times. Note that our estimators for dynamics and rewards are only well-defined when we visit the corresponding state-action pair at least once. However, in a stochastic environment, a single visit will likely not result in an accurate model. We can use Hoeffding's inequality (A.41) to gauge how accurate the estimates are after only a limited number of visits.

11.3 Balancing Exploration and Exploitation

The next natural question is how to use our current model of the environment to pick actions such that exploration and exploitation are traded effectively. This is what we will consider next.

Given the estimated MDP given by \hat{p} and \hat{r} , we can compute the optimal policy using either policy iteration or value iteration. For example, using value iteration, we can compute the optimal state-action value function Q^* within the *estimated* MDP, and then employ the greedy policy

$$\pi(x) = \arg \max_{a \in A} Q^*(x, a). \quad (11.6)$$

Recall from Equation (10.31) that this corresponds to always picking the best action under the *current* model (that is, π is the optimal policy). But since the model is inaccurate, while potentially quickly generating some reward, we will likely get stuck in a suboptimal state.

11.3.1 ε -greedy

Consider the other extreme: If we always pick a random action, we will eventually(!) estimate the dynamics and rewards correctly, yet we will do extremely poorly in terms of maximizing rewards along the way. To trade exploration and exploitation, a natural idea is to balance these two extremes.

Arguably, the simplest idea is the following: At each time step, throw a biased coin. If this coin lands heads, we pick an action uniformly at random among all actions. If the coin lands tails, we pick the best action under our current model. This algorithm is called ε -greedy, where the probability of a coin landing heads at time t is ε_t .

Algorithm 11.2: ε -greedy

```

1 for  $t = 0$  to  $\infty$  do
2   sample  $u \in \text{Unif}([0, 1])$ 
3   if  $u \leq \varepsilon_t$  then pick action uniformly at random among all actions
4   else pick best action under the current model

```

The ε -greedy algorithm provides a general framework for addressing the exploration-exploitation dilemma. When the underlying MDP is learned using Monte Carlo estimation as we discussed in Section 11.2.1, the resulting algorithm is known as *Monte Carlo control*. However, the same framework can also be used in the model-free setting where we pick the best action without estimating the full underlying MDP. We discuss this approach in greater detail in Section 11.4.

Amazingly, this simple algorithm already works quite well. Nevertheless, it can clearly be improved. The key problem of ε -greedy is that it explores the state space in an uninformed manner. In other words, it explores ignoring all past experience. It thus does not eliminate clearly suboptimal actions. This is a problem, especially as we typically have many state-action pairs and recalling that we have to explore each such pair many times to learn an accurate model.

Remark 11.3: Asymptotic convergence

It can be shown that Monte Carlo control converges to an optimal policy (albeit slowly) almost surely when the learned policy is “greedy in the limit with infinite exploration”.

Definition 11.4 (Greedy in the limit with infinite exploration, GLIE).

A sequence of policies π_t is said to be *greedy in the limit with infinite exploration* if

1. all state-action pairs are explored infinitely many times,²

$$\lim_{t \rightarrow \infty} N_t(x, a) = \infty \quad \text{and} \quad (11.7)$$

2. the policy converges to a greedy policy,

$$\lim_{t \rightarrow \infty} \pi_t(a | x) = \mathbb{1}\{a = \arg \max_{a' \in A} Q_t^*(x, a')\} \quad (11.8)$$

where we denote by $N_t(x, a)$ the number of transitions from state x playing action a until time t , and Q_t^* is the optimal state-action value function in the estimated MDP at time t .

Note that ε -greedy is GLIE with probability 1 if the sequence $(\varepsilon_t)_{t \in \mathbb{N}_0}$ satisfies the Robbins-Monro (RM) conditions (A.60),

$$\varepsilon_t \geq 0 \quad \forall t, \quad \sum_{t=0}^{\infty} \varepsilon_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} \varepsilon_t^2 < \infty.$$

The RM-conditions are satisfied, for example, if $\varepsilon_t = 1/t$.

Theorem 11.5 (Convergence of Monte Carlo control). *GLIE Monte Carlo control converges to an optimal policy with probability 1.*

Intuitively, the probability of exploration converges to zero, and hence, the policy will “eventually coincide” with the greedy policy. Moreover, the greedy policy will “eventually coincide” with the optimal policy due to an argument akin to the convergence of policy iteration,³ and using that each state-action pair is visited infinitely often.

² That all state-action pairs are chosen is a fundamental requirement. There is no reason why any algorithm would converge to the true value function for *all* states when it only sees some state-action pairs finitely many times, or even not at all.

11.3.2 Softmax Exploration

An alternative to using ε -greedy for trading between greedy exploitation and uniform exploration is the so-called *softmax exploration* or *Boltzmann exploration*. Given the agent is in state x , we pick action a with probability,

$$\pi_\lambda(a | x) \propto \exp\left(\frac{1}{\lambda} Q^*(x, a)\right), \quad (11.9)$$

³ see Lemma 10.15 and Theorem 10.16

which is the Gibbs distribution with temperature parameter $\lambda > 0$. Observe that for $\lambda \rightarrow 0$, softmax exploration corresponds to greedily maximizing the Q-function (i.e., greedy exploitation), whereas for $\lambda \rightarrow \infty$, softmax exploration explores uniformly at random. This can outperform ε -greedy as the exploration is directed towards actions with larger estimated value.

11.3.3 Optimism

Recall from our discussion of multi-armed bandits in Section 9.2.1 that a key principle in effectively trading exploration and exploitation is *optimism in the face of uncertainty*. Let us apply this principle to the reinforcement learning setting. The key idea is to assume that the dynamics and rewards model “work in our favor” until we have learned “good estimates” of the true dynamics and rewards.

More formally, if $r(x, a)$ is unknown, we set $\hat{r}(x, a) = R_{\max}$, where R_{\max} is the maximum reward our agent can attain during a single transition. Similarly, if $p(x' | x, a)$ is unknown, we set $\hat{p}(x^* | x, a) = 1$, where x^* is a “fairy-tale state”. The fairy-tale state corresponds to everything our agent could wish for, that is,

$$\hat{p}(x^* | x^*, a) = 1 \quad \forall a \in A, \quad (11.10)$$

$$\hat{r}(x^*, a) = R_{\max} \quad \forall a \in A. \quad (11.11)$$

In practice, the decision of when to assume that the learned dynamics and reward models are “good enough” has to be tuned.

In using these optimistic estimates of p and r , we obtain an optimistic underlying Markov decision process that exhibits a bias towards exploration. In particular, the rewards attained in this MDP, are an upper bound of the true reward. The resulting algorithm is known as the R_{\max} algorithm.

How many transitions are “enough”? We can use Hoeffding’s inequality (A.41) to get a rough idea! The key here, is our observation from Equation (11.3) that the transitions and rewards are conditionally independent given the state-action pairs since, as we have discussed in Section 6.1.4 on the ergodic theorem, Hoeffding’s inequality does not hold for dependent samples. In this case, Hoeffding’s inequality tells us that for the absolute approximation error to be below ϵ with probability at least $1 - \delta$, we need

$$N(a | x) \geq \frac{R_{\max}^2}{2\epsilon^2} \log \frac{2}{\delta}. \quad (11.12) \quad \text{see (A.42)}$$

Lemma 11.7 (Exploration and exploitation of R_{\max}). *Every T time steps, with high probability, R_{\max} either*

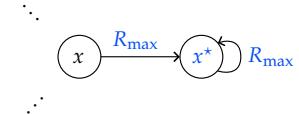


Figure 11.2: Illustration of the fairy-tale state of R_{\max} . If in doubt, the agent believes actions from the state x to lead to the fairy-tale state x^* with maximal rewards. This encourages the exploration of unknown states.

Algorithm 11.6: R_{\max} algorithm

```

1 add the fairy-tale state  $x^*$  to the Markov decision process
2 set  $\hat{r}(x, a) = R_{\max}$  for all  $x \in X$  and  $a \in A$ 
3 set  $\hat{p}(x^* | x, a) = 1$  for all  $x \in X$  and  $a \in A$ 
4 compute the optimal policy  $\hat{\pi}$  for  $\hat{r}$  and  $\hat{p}$ 
5 for  $t = 0$  to  $\infty$  do
6   execute policy  $\hat{\pi}$  (for some number of steps)
7   for each visited state-action pair  $(x, a)$ , update  $\hat{r}(x, a)$ 
8   estimate transition probabilities  $\hat{p}(x' | x, a)$ 
9   after observing “enough” transitions and rewards, recompute the
      optimal policy  $\hat{\pi}$  according the current model  $\hat{p}$  and  $\hat{r}$ .

```

- obtains near-optimal reward; or
- visits at least one unknown state-action pair.⁴

Here, T depends on the mixing time of the Markov chain induced by the optimal policy.

Theorem 11.8 (Convergence of R_{\max} , Brafman and Tennenholz (2002)).
With probability at least $1 - \delta$, R_{\max} reaches an ϵ -optimal policy in a number of steps that is polynomial in $|X|$, $|A|$, T , $1/\epsilon$, $1/\delta$, and R_{\max} .

11.3.4 Challenges of Model-based Approaches

We have seen that the R_{\max} algorithm performs remarkably well in the tabular setting. However, there are important computational limitations to the model-based approaches that we discussed so far.

First, observe that the (tabular) model-based approach requires us to store $\hat{p}(x' | x, a)$ and $\hat{r}(x, a)$ in a table. This table already has $O(n^2m)$ entries. Even though polynomial in the size of the state and action spaces, this quickly becomes unmanageable.

Second, the model-based approach requires us to “solve” the learned Markov decision processes to obtain the optimal policy (using policy or value iteration). As we continue to learn over time, we need to find the optimal policy many times. R_{\max} recomputes the policy after each state-action pair is observed sufficiently often, so $O(nm)$ times.

11.4 Model-free Approaches

In the previous section, we have seen that learning and remembering the model as well as planning within the estimated model can potentially be quite expensive in the model-based approach. We there-

⁴Note that in the tabular setting, there are “only” polynomially many state-action pairs.

fore turn to model-free methods that estimate the value function directly. Thus, they require neither remembering the full model nor planning (i.e., policy optimization) in the underlying Markov decision process. We will, however, return to model-based methods in Chapter 13 to see that promise lies in combining methods from model-based reinforcement learning with methods from model-free reinforcement learning.

A significant benefit to model-based reinforcement learning is that it is inherently off-policy. That is, any trajectory regardless of the policy used to obtain it can be used to improve the model of the underlying Markov decision process. In the model-free setting, this is not necessarily true. By default, estimating the value function according to the data from a trajectory, will yield an estimate of the value function corresponding to the policy that was used to sample the data.

We will start by discussing on-policy methods and later see how the value function can be estimated off-policy.

11.4.1 On-policy Value Estimation

Let us suppose, our agent follows a fixed policy π . Then, the corresponding value function v^π is given as

$$\begin{aligned} v^\pi(x) &= r(x, \pi(x)) + \gamma \sum_{x' \in X} p(x' | x, \pi(x)) \cdot v^\pi(x') \\ &= \mathbb{E}_{R_0, X_1}[R_0 + \gamma v^\pi(X_1) | X_0 = x, A_0 = \pi(x)] \end{aligned} \quad (11.13)$$

Our first instinct might be to use a Monte Carlo estimate of this expectation. Due to the conditional independence of the transitions (11.3), Monte Carlo approximation does yield an unbiased estimate,

$$\approx r + \gamma v^\pi(x'), \quad (11.14)$$

where the agent observed the transition (x, a, r, x') . Note that to estimate this expectation we use a single(!) sample,⁵ unlike our previous applications of Monte Carlo sampling where we usually averaged over m samples. However, there is one significant problem in this approximation. Our approximation of v^π does in turn depend on the (unknown) true value of v^π !

The key idea is to use a bootstrapping estimate of the value function instead. That is, in place of the true value function v^π , we will use a “running estimate” V^π . In other words, whenever observing a new transition, we use our previous best estimate of v^π to obtain a new estimate V^π . We already encountered bootstrapping briefly in Section 7.3.4 in the context of probabilistic ensembles in Bayesian deep learning. More generally, *bootstrapping* refers to approximating a true

using the definition of the value function (10.7)

interpreting the above expression as an expectation over the random quantities R_0 and X_1

⁵ The idea is that we will use this approximation repeatedly as our agent collects new data to achieve the same effect as initially averaging over multiple samples.

quantity (e.g., v^π) by using an empirical quantity (e.g., V^π), which itself is constructed using samples from the true quantity that is to be approximated.

Due to its use in estimating the value function, bootstrapping is a core concept to model-free reinforcement learning. Crucially, using a bootstrapping estimate generally results in biased estimates of the value function. Moreover, due to relying on a single sample, the estimates from Equation (11.14) tend to have very large variance.

The variance of the estimate is typically reduced by mixing new estimates of the value function with previous estimates using a learning rate α_t . This yields the *temporal-difference learning* algorithm.

Algorithm 11.9: Temporal-difference (TD) learning

```

1 initialize  $V^\pi$  arbitrarily (e.g., as  $\mathbf{0}$ )
2 for  $t = 0$  to  $\infty$  do
3   follow policy  $\pi$  to obtain the transition  $(x, a, r, x')$ 
4    $V^\pi(x) \leftarrow (1 - \alpha_t)V^\pi(x) + \alpha_t(r + \gamma V^\pi(x'))$            // (11.15)

```

The update rule is sometimes written equivalently as

$$V^\pi(x) \leftarrow V^\pi(x) + \alpha_t(r + \gamma V^\pi(x') - V^\pi(x)). \quad (11.16)$$

Thus, the update to $V^\pi(x)$ is proportional to the learning rate and the difference between the previous estimate and the renewed estimate using the new observation.

Theorem 11.10 (Convergence of TD-learning, Jaakkola et al. (1993)).
If $(\alpha_t)_{t \in \mathbb{N}_0}$ satisfies the RM-conditions (A.60) and all state-action pairs are chosen infinitely often, then V^π converges to v^π with probability 1.

Importantly, note that due to the Monte Carlo approximation of Equation (11.13) with respect to transitions attained by following policy π , TD-learning is fundamentally on-policy. That is, for the estimates V^π to converge to the true value function v^π , the transitions that are used for the estimation must follow policy π .

11.4.2 SARSA: On-policy Control

TD-learning merely estimates the value function of a fixed policy π . To find the optimal policy π^* , we can use an analogue of policy iteration (see Algorithm 10.14). Here, it is more convenient to use an estimate of the state-action value function q^π which can be obtained

analogously to the bootstrapping estimate of v^π (11.14),

$$\begin{aligned} q^\pi(x, a) &= r(x, a) + \gamma \sum_{x' \in X} p(x' | x, a) \sum_{a' \in A} \pi(a' | x') q^\pi(x', a') \\ &= \mathbb{E}_{R_0, X_1, A_1} [R_0 + \gamma q^\pi(X_1, A_1) | X_0 = x, A_0 = a] \end{aligned} \quad (11.17)$$

$$\approx r + \gamma q^\pi(x', a'), \quad (11.18)$$

where the agent observed transitions (x, a, r, x') and (x', a', r', x'') .

The update rule from TD-learning is therefore adapted to⁶

$$Q^\pi(x, a) \leftarrow (1 - \alpha_t) Q^\pi(x, a) + \alpha_t (r + \gamma Q^\pi(x', a')). \quad (11.19)$$

This algorithm is known as *SARSA* (short for *state-action-reward-state-action*). Similar convergence guarantees to those of TD-learning can also be derived for SARSA.

Theorem 11.11 (Convergence of SARSA, Singh et al. (2000)). *If $(\alpha_t)_{t \in \mathbb{N}_0}$ satisfies the RM-conditions (A.6o) and all state-action pairs are chosen infinitely often, then Q^π converges to q^π with probability 1.*

The policy iteration scheme to identify the optimal policy can be outlined as follows: In each iteration t , we estimate the value function q^{π_t} of policy π_t with the estimate Q^{π_t} obtained from SARSA. We then choose the greedy policy with respect to Q^{π_t} as the next policy π_{t+1} . However, due to the on-policy nature of SARSA, we cannot reuse any data between the iterations. Moreover, it turns out that in practice, when using only finitely many samples, this form of greedily optimizing Markov decision processes does not explore enough. At least partially, this can be compensated for by injecting noise when choosing the next action, e.g., by following an ϵ -greedy policy or using softmax exploration.

11.4.3 Off-policy Value Estimation

Consider the following slight adaptation of the derivation of SARSA (11.18),

$$\begin{aligned} q^\pi(x, a) &= r(x, a) + \gamma \sum_{x' \in X} p(x' | x, a) \sum_{a' \in A} \pi(a' | x') q^\pi(x', a') \\ &= \mathbb{E}_{R_0, X_1} \left[R_0 + \gamma \sum_{a' \in A} \pi(a' | X_1) q^\pi(X_1, a') \mid X_0 = x, A_0 = a \right] \end{aligned} \quad (11.20)$$

$$\approx r + \gamma \sum_{a' \in A} \pi(a' | x') q^\pi(x', a'), \quad (11.21)$$

using Bellman's expectation equation
(10.15)

interpreting the above expression as an expectation over R_0, X_1 and A_1
Monte Carlo approximation with a single sample

⁶ Note that for deterministic policies π , $Q^\pi(x', a') = Q^\pi(x', \pi(x')) = V^\pi(x')$ if the transitions are obtained by following policy π .

Monte Carlo approximation with a single sample

where the agent observed the transition (x, a, r, x') . This yields the update rule,

$$Q^\pi(x, a) \leftarrow (1 - \alpha_t)Q^\pi(x, a) + \alpha_t \left(r + \gamma \sum_{a' \in A} \pi(a' | x') Q^\pi(x', a') \right). \quad (11.22)$$

This adapted update rule *explicitly* chooses the subsequent action a' according to policy π whereas SARSA absorbs this choice into the Monte Carlo approximation. The algorithm has analogous convergence guarantees to those of SARSA.

Crucially, this algorithm is off-policy. That is, we can use transitions that were obtained according to *any* policy to estimate the value of a fixed policy π , which we may have never used! Perhaps this seems contradictory at first, but it is not. As noted, the key difference to the on-policy TD-learning and SARSA is that our estimate of the Q-function explicitly keeps track of the next-performed action. It does so for any action in any state. Moreover, note that the transitions that are due to the dynamics model and rewards are unaffected by the used policy. They merely depend on the originating state-action pair. We can therefore use the instances where other policies played action $\pi(x)$ in state x to estimate the performance of π .

11.4.4 Q-learning: Off-policy Control

It turns out that there is a way to estimate the value function of the optimal policy directly. Recall from Bellman's theorem (10.31) that the optimal policy π^* can be characterized in terms of the optimal state-action value function q^* ,

$$\pi^*(x) = \arg \max_{a \in A} q^*(x, a).$$

π^* corresponds to greedily maximizing the value function.

Analogously to our derivation of SARSA (11.18), only using Bellman's theorem (10.32) in place of Bellman's expectation equation (10.15), we obtain,

$$\begin{aligned} q^*(x, a) &= r(x, a) + \gamma \sum_{x' \in X} p(x' | x, a) \max_{a' \in A} q^*(x', a') \\ &= \mathbb{E}_{R_0, X_1} \left[R_0 + \gamma \max_{a' \in A} q^*(X_1, a') \mid X_0 = x, A_0 = a \right] \quad (11.23) \end{aligned}$$

$$\approx r + \gamma \max_{a' \in A} q^*(x', a'), \quad (11.24)$$

where the agent observed the transition (x, a, r, x') . Using a bootstrapping estimate Q^* for q^* , we obtain a structurally similar algorithm to

using that the Q-function is a fixed-point of the Bellman update, see Bellman's theorem (10.32)
interpreting the above expression as an expectation over R_0 and X_1
Monte Carlo approximation with a single sample

TD-learning and SARSA — only for estimating the optimal Q-function directly! This algorithm is known as *Q-learning*. Whereas we have seen that the optimal policy can be found using SARSA in a policy-iteration-like scheme, Q-learning is conceptually similar to value iteration.

Algorithm 11.12: Q-learning

```

1 initialize  $Q^*(x, a)$  arbitrarily (e.g., as 0)
2 for  $t = 0$  to  $\infty$  do
3   observe the transition  $(x, a, r, x')$ 
4    $Q^*(x, a) \leftarrow (1 - \alpha_t)Q^*(x, a) + \alpha_t(r + \gamma \max_{a' \in A} Q^*(x', a'))$ 
      // (11.25)

```

Similarly to TD-learning, the update rule can also be expressed as

$$Q^*(x, a) \leftarrow Q^*(x, a) + \alpha_t \left(r + \gamma \max_{a' \in A} Q^*(x', a') - Q^*(x, a) \right). \quad (11.26)$$

Crucially, the Monte Carlo approximation of Equation (11.23) does not depend on the policy. Thus, Q-learning is an off-policy method.

Theorem 11.13 (Convergence of Q-learning, Jaakkola et al. (1993)). *If $(\alpha_t)_{t \in \mathbb{N}_0}$ satisfies the RM-conditions (A.60) and all state-action pairs are chosen infinitely often (that is, the sequence of policies used to obtain the transitions is GLIE), then Q^* converges to q^* with probability 1.*

It can be shown that with probability at least $1 - \delta$, Q-learning converges to an ϵ -optimal policy in a number of steps that is polynomial in $\log |X|$, $\log |A|$, $1/\epsilon$ and $\log 1/\delta$ (Even-Dar et al., 2003).

11.4.5 Optimistic Q-learning

The next natural question is how to effectively trade exploration and exploitation to both visit all state-action pairs many times, but also attain a high reward.

However, as we have seen in Section 11.3, random “uninformed” exploration like ϵ -greedy and softmax exploration explores the state space very slowly. We therefore return to the principle of *optimism in the face of uncertainty*, which already led us to the R_{\max} algorithm in the model-based setting. We will now additionally assume that the rewards are non-negative, that is, $0 \leq r(x, a) \leq R_{\max}$ ($\forall x \in X, a \in A$). It turns out that a similar algorithm to R_{\max} also exists for (model-free) Q-learning: it is called *optimistic Q-learning* and shown in Algorithm 11.14.

Algorithm 11.14: Optimistic Q-learning

```

1 initialize  $Q^*(x, a) = V_{\max} \prod_{t=1}^{T_{\text{init}}} (1 - \alpha_t)^{-1}$ 
2 for  $t = 0$  to  $\infty$  do
3   pick action  $a_t = \arg \max_{a \in A} Q^*(x, a)$  and observe the transition
       $(x, a_t, r, x')$ 
4    $Q^*(x, a_t) \leftarrow (1 - \alpha_t)Q^*(x, a_t) + \alpha_t(r + \gamma \max_{a' \in A} Q^*(x', a'))$ 
      // (11.27)

```

Here,

$$V_{\max} \doteq \frac{R_{\max}}{1 - \gamma} \geq \max_{x \in X, a \in A} q^*(x, a),$$

is an upper bound on the discounted return and T_{init} is some initialization time. Intuitively, the initialization of Q^* corresponds to the best-case long-term reward, assuming that all individual rewards are upper bounded by R_{\max} . This is shown by the following lemma.

Lemma 11.15. Denote by Q_t^* , the approximation of q^* attained in the t -th iteration of optimistic Q-learning. Then, for any state-action pair (x, a) and iteration t such that $N(a | x) \leq T_{\text{init}}$,

$$Q_t^*(x, a) \geq V_{\max} \geq q^*(x, a). \quad (11.28)$$

⁷ $N(a | x)$ is the number of times action a is performed in state x .

Proof. We write $\beta_\tau \doteq \prod_{i=1}^\tau (1 - \alpha_i)$ and $\eta_{i,\tau} \doteq \alpha_i \prod_{j=i+1}^\tau (1 - \alpha_j)$. Using the update rule of optimistic Q-learning (11.27), we have

$$Q_t^*(x, a) = \beta_{N(a|x)} Q_0^*(x, a) + \sum_{i=1}^{N(a|x)} \eta_{i,N(a|x)} (r + \gamma \max_{a_i \in A} Q_{t_i}^*(x_i, a_i)) \quad (11.29)$$

where x_i is the next state arrived at time t_i when action a is performed the i -th time in state x .

Using the assumption that the rewards are non-negative, from Equation (11.29) and $Q_0^*(x, a) = V_{\max}/\beta_{T_{\text{init}}}$, we immediately have

$$\begin{aligned} Q_t^*(x, a) &\geq \frac{\beta_{N(a|x)}}{\beta_{T_{\text{init}}}} V_{\max} \\ &\geq V_{\max}. \end{aligned} \quad \square \quad \text{using } N(a | x) \leq T_{\text{init}}$$

Now, if T_{init} is chosen large enough, it can be shown that optimistic Q-learning converges quickly to an optimal policy.

Theorem 11.16 (Convergence of optimistic Q-learning, Even-Dar and Mansour (2001)). *With probability at least $1 - \delta$, optimistic Q-learning obtains an ϵ -optimal policy after a number of steps that is polynomial in $|X|$,*

$|A|$, $1/\epsilon$, $\log 1/\delta$, and R_{\max} where the initialization time T_{init} is upper bounded by a polynomial in the same coefficients.

Note that for Q-learning, we still need to store $Q^*(x, a)$ for any state-action pair in memory. Thus, Q-learning requires $O(nm)$ memory. During each transition, we need to compute

$$\max_{a \in A} Q^*(x', a')$$

once. If we run Q-learning for T iterations, this yields a time complexity of $O(Tm)$. Crucially, for sparse Markov decision processes where, in most states, only few actions are permitted, each iteration of Q-learning can be performed in (virtually) constant time. This is a big improvement of the quadratic (in the number of states) performance of the model-based R_{\max} algorithm.

Discussion

We have seen that both the model-based R_{\max} algorithm and the model-free Q-learning take time polynomial in the number of states $|X|$ and the number of actions $|A|$ to converge. While this is acceptable in small grid worlds, this is completely unacceptable for large state and action spaces.

Often, domains are continuous, for example when modeling beliefs about states in a partially observable environment. Also, in many structured domains (e.g., chess or multiagent planning), the size of the state and action space is exponential in the size of the input. In the final two chapters, we will therefore explore how model-free and model-based methods can be used (approximately) in such large domains.

Problems

11.1. Q-learning.

Assume the following grid world, where from state A the agent can go to the right and down, and from state B to the left and down. From states G_1 and G_2 the only action is to exit. The agent receives a reward (+10 or +1) only when exiting.

Rewards		States	
0	0	A	B
+10	+1	G_1	G_2

We assume the discount factor $\gamma = 1$ and that all actions are deterministic.

1. We observe the following two episodes:

		Episode 1		Episode 2	
x	a	x'	r	x	a
A	\downarrow	G_1	0	B	\leftarrow
G_1	exit		10	A	\downarrow
				G_1	0
				G_1	exit
					10

Assume $\alpha = 0.3$, and that Q-values of all non-terminal states are initialized to 0.5. What are the Q-values

$$Q^*(A, \downarrow), \quad Q^*(G_1, \text{exit}), \quad Q^*(G_2, \text{exit})$$

learned by executing Q-learning with the above episodes?

2. Will Q-learning converge to q^* for all state-action pairs (x, a) if we repeat episode 1 and episode 2 infinitely often? If not, design a sequence of episodes that leads to convergence.
3. How does the choice of initial Q-values influence the convergence of the Q-learning algorithm when episodes are obtained off-policy?
4. Determine v^* for all states.

12

Model-free Reinforcement Learning

In the previous chapter, we have seen methods for tabular settings. Our goal now is to extend the model-free methods like TD-learning and Q-learning to large state-action spaces \mathcal{X} and \mathcal{A} . We have seen that a crucial bottleneck of these methods is the parameterization of the value function. If we want to store the value function in a table, we need at least $O(|\mathcal{X}|)$ space. If we learn the Q function, we even need $O(|\mathcal{X}| \cdot |\mathcal{A}|)$ space. Also, for large state-action spaces, the time required to compute the value function for every state-action pair exactly will grow polynomially in the size of the state-action space. Hence, a natural idea is to learn *approximations* of these functions with a low-dimensional parameterization. Such approximations were the focus of the first few chapters and are, in fact, the key idea behind methods for reinforcement learning in large state-action spaces.

12.1 Tabular Reinforcement Learning as Optimization

To begin with, let us reinterpret the model-free methods from the previous section, TD-learning and Q-learning, as solving an optimization problem, where each iteration corresponds to a single gradient update. We will focus on TD-learning here, but the same interpretation applies to Q-learning. Recall the update rule of TD-learning (11.15),

$$V^\pi(x) \leftarrow (1 - \alpha_t)V^\pi(x) + \alpha_t(r + \gamma V^\pi(x')).$$

Note that this looks just like the update rule of an optimization algorithm! We can parameterize our estimates V^π with parameters θ that are updated according to the gradient of some loss function, assuming fixed bootstrapping estimates. In particular, in the tabular setting (i.e., over a finite domain), we can parameterize the value function exactly by learning a separate parameter for each state,

$$\theta \doteq [\theta(1), \dots, \theta(n)], \quad V^\pi(x; \theta) \doteq \theta(x). \quad (12.1)$$

To re-derive the above update rule as a gradient update, let us consider the following loss function,

$$\bar{\ell}(\boldsymbol{\theta}; x, r) \doteq \frac{1}{2}(v^\pi(x) - \boldsymbol{\theta}(x))^2 \quad (12.2)$$

$$= \frac{1}{2}\left(r + \gamma \mathbb{E}_{x'|x, \pi(x)}[v^\pi(x')] - \boldsymbol{\theta}(x)\right)^2 \quad (12.3)$$

Note that this loss corresponds to a standard squared loss of the difference between the parameter $\boldsymbol{\theta}(x)$ and the label $v^\pi(x)$ we want to learn.

We can now find the gradient of this loss. Elementary calculations yield,

$$\nabla_{\boldsymbol{\theta}(x)} \bar{\ell}(\boldsymbol{\theta}; x, r) = \boldsymbol{\theta}(x) - \left(r + \gamma \mathbb{E}_{x'|x, \pi(x)}[v^\pi(x')]\right). \quad (12.4)$$

Now, we cannot compute this derivative because we cannot compute the expectation. Firstly, the expectation is over the true value function which is unknown to us. Secondly, the expectation is over the transition model which we are trying to avoid in model-free methods.

To resolve the first issue, analogously to TD-learning, instead of learning the true value function v^π which is unknown, we learn the bootstrapping estimate V^π . Recall that the core principle behind bootstrapping as discussed in Section 11.4.1 is that this bootstrapping estimate V^π is *treated* as if it were independent of the current estimate of the value function $\boldsymbol{\theta}$. To emphasize this, we write $V^\pi(x; \boldsymbol{\theta}^{\text{old}}) \approx v^\pi(x)$ where $\boldsymbol{\theta}^{\text{old}} = \boldsymbol{\theta}$ but $\boldsymbol{\theta}^{\text{old}}$ is treated as a constant with respect to $\boldsymbol{\theta}$.¹

To resolve the second issue, analogously to the introduction of TD-learning in the previous chapter, we will use a Monte Carlo estimate using a single sample. Recall that this is only possible because the transitions are conditionally independent given the state-action pair.

using Bellman's expectation equation
(10.11)

Remark 12.1: Sample (in)efficiency of model-free methods

Taking these two shortcuts are two of the main reasons why model-free methods such as TD-learning and Q-learning are usually *sample inefficient*. This is because using a bootstrapping estimate leads to “(initially) incorrect” and “unstable” targets of the optimization problem,² and Monte Carlo estimation with a single sample leads to a large variance. Recall that the theoretical guarantees for model-free methods in the tabular setting therefore required that all state-action pairs are visited infinitely often.

Using the aforementioned shortcuts, let us define the loss ℓ after observing the single transition (x, a, r, x') ,

$$\ell(\boldsymbol{\theta}; x, r, x') \doteq \frac{1}{2}\left(r + \gamma \boldsymbol{\theta}^{\text{old}}(x') - \boldsymbol{\theta}(x)\right)^2. \quad (12.5)$$

¹ That is, the bootstrapping estimate $V^\pi(x; \boldsymbol{\theta}^{\text{old}})$ is assumed to be constant with respect to $\boldsymbol{\theta}(x)$ in the same way that $v^\pi(x)$ is constant with respect to $\boldsymbol{\theta}(x)$.

If we were not using the bootstrapping estimate, the following derivation of the gradient of the loss would not be this simple.

² We explore this in some more capacity in Section 12.2.1 where we cover heuristic approaches to alleviate this problem to some degree.

We define the gradient of this loss with respect to $\theta(x)$ as

$$\begin{aligned}\delta_{\text{TD}} &\doteq \nabla_{\theta(x)} \ell(\theta; x, r, x') \\ &= \theta(x) - (r + \gamma \theta^{\text{old}}(x')).\end{aligned}\quad (12.6)$$

This error term is also called *temporal-difference (TD) error*. The temporal difference error compares the previous estimate of the value function to the bootstrapping estimate of the value function. We know from the law of large numbers (A.36) that Monte Carlo averages are unbiased.³ We therefore have,

$$\mathbb{E}_{x'|x, \pi(x)}[\delta_{\text{TD}}] \approx \nabla_{\theta(x)} \bar{\ell}(\theta; x, r). \quad (12.7)$$

Naturally, we can use these unbiased gradient estimates with respect to the loss $\bar{\ell}$ to perform stochastic gradient descent. This yields the update rule,

$$V^\pi(x; \theta) = \theta(x) \leftarrow \theta(x) - \alpha_t \delta_{\text{TD}} \quad (12.8)$$

$$\begin{aligned}&= (1 - \alpha_t) \theta(x) + \alpha_t (r + \gamma \theta^{\text{old}}(x')) \\ &= (1 - \alpha_t) V^\pi(x; \theta) + \alpha_t (r + \gamma V^\pi(x'; \theta^{\text{old}})).\end{aligned}\quad (12.9)$$

Observe that this gradient update coincides with the update rule of TD-learning (11.15). Therefore, TD-learning is essentially performing stochastic gradient descent using the TD-error as an unbiased gradient estimate.⁴ Crucially, TD-learning performs stochastic gradient descent with respect to the bootstrapping estimate of the value function V^π and not the true value function v^π ! Stochastic gradient descent with a bootstrapping estimate is also called *stochastic semi-gradient descent*. Importantly, the optimization target $r + \gamma V^\pi(x'; \theta^{\text{old}})$ from the loss ℓ is now *moving* between iterations which introduces some practical challenges we will discuss in Section 12.2.1. We have seen in the previous chapter that using a bootstrapping estimate still guarantees (asymptotic) convergence to the true value function.

12.2 Value Function Approximation

To scale to large state spaces, it is natural to approximate the value function using a parameterized model, $V(x; \theta)$ or $Q(x, a; \theta)$. You may think of this as a regression problem where we map state(-action) pairs to a real number. Recall from the previous section that this is a strict generalization of the tabular setting, as we could use a separate parameter to learn the value function for each individual state-action pair. Our goal for large state-action spaces is to exploit the smoothness properties⁵ of the value function to condense the representation.

³Crucially, the samples are unbiased with respect to the approximate label in terms of the bootstrapping estimate only. Due to bootstrapping the value function, the estimates are not unbiased with respect to the true value function. Moreover, the variance of each individual estimation of the gradient is large, as we only consider a single transition.

using stochastic gradient descent with learning rate α_t , see Algorithm A.29

using the definition of the temporal difference error (12.6)

substituting $V^\pi(x; \theta)$ for $\theta(x)$

⁴An alternative interpretation is that TD-learning performs gradient descent with respect to the loss ℓ .

⁵That is, the value function takes similar values in “similar” states.

A straightforward approach is to use a linear function approximation with the hand-designed feature map ϕ ,

$$Q(x, a; \theta) \doteq \theta^\top \phi(x, a). \quad (12.10)$$

A common alternative is to use a deep neural network to learn these features instead. Doing so is also known as *deep reinforcement learning*.⁶

We will now apply the derivation from the previous section directly to Q-learning. For Q-learning, after observing the transition (x, a, r, x') , the loss function is given as

$$\ell(\theta; x, a, r, x') \doteq \frac{1}{2} \left(r + \gamma \max_{a' \in \mathcal{A}} Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2. \quad (12.11)$$

Here, we simply use Bellman's optimality equation (10.32) to estimate $q^*(x, a)$, instead of the estimation of $v^\pi(x)$ using Bellman's expectation equation for TD-learning. The difference between the current approximation and the optimization target,

$$\delta_B \doteq r + \gamma \max_{a' \in \mathcal{A}} Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta), \quad (12.12)$$

is called the *Bellman error*. Analogously to TD-learning,⁷ we obtain the gradient update,

$$\theta \leftarrow \theta - \alpha_t \nabla_\theta \ell(\theta; x, a, r, x') \quad (12.13)$$

$$\begin{aligned} &= \theta - \alpha_t \nabla_\theta \frac{1}{2} \left(r + \gamma \max_{a' \in \mathcal{A}} Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2 \\ &= \theta + \alpha_t \delta_B \nabla_\theta Q^*(x, a; \theta). \end{aligned} \quad (12.14)$$

When using a neural network to learn Q^* , we can use automatic differentiation to obtain unbiased gradient estimates. In the case of linear function approximation, we can compute the gradient exactly,

$$\begin{aligned} &= \theta + \alpha_t \delta_B \nabla_\theta \theta^\top \phi(x, a) \\ &= \theta + \alpha_t \delta_B \phi(x, a). \end{aligned} \quad (12.15)$$

In the tabular setting, this algorithm is identical to Q-learning and, in particular, converges to the true Q-function q^* .⁸ There are few such results in the approximate setting. Usage in practice indicates that using an approximation of the value function "should be fine" when a "rich-enough" class of functions is used.

12.2.1 Heuristics

The vanilla stochastic semi-gradient descent is very slow. In this subsection, we will discuss some "tricks of the trade" to improve its performance.

⁶ Note that often non-Bayesian deep learning (i.e., point estimates of the weights of a neural network) is applied. In the final chapter, Chapter 13, we will explore the benefits of using Bayesian deep learning.

⁷ compare to Equation (12.8)

using the definition of ℓ (12.11)

using the chain rule

using the linear approximation of Q^* (12.10)

⁸ see Theorem 11.13

Stabilizing optimization targets: There are mainly two problems. The first problem is that, as mentioned previously, the bootstrapping estimate changes after each iteration. As we are trying to learn an approximate value function that depends on the bootstrapping estimate, this means that the optimization target is “moving” between iterations. In practice, moving targets lead to stability issues. The first family of techniques we discuss here aim to “stabilize” the optimization targets.

One such technique is called *neural fitted Q-iteration* or *deep Q-networks* (DQN) (Mnih et al., 2015). DQN updates the neural network used for the approximate bootstrapping estimate infrequently to maintain a constant optimization target across multiple episodes. How this is implemented exactly varies. One approach is to clone the neural network and maintain one changing neural network (“online network”) for the most recent estimate of the Q-function which is parameterized by θ , and one fixed neural network (“target network”) used as the target which is parameterized by θ^{old} and which is updated infrequently.

This can be implemented by maintaining a data set \mathcal{D} of observed transitions (the so-called *replay buffer*) and then “every once in a while” (e.g., once $|\mathcal{D}|$ is large enough) solving a regression problem, where the labels are determined by the target network. This yields a loss term where the target is fixed across all transitions in the replay buffer \mathcal{D} ,

$$\ell_{\text{DQN}}(\theta; \mathcal{D}) \doteq \frac{1}{2} \sum_{(x, a, r, x') \in \mathcal{D}} \left(r + \gamma \max_{a' \in \mathcal{A}} Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2. \quad (12.16)$$

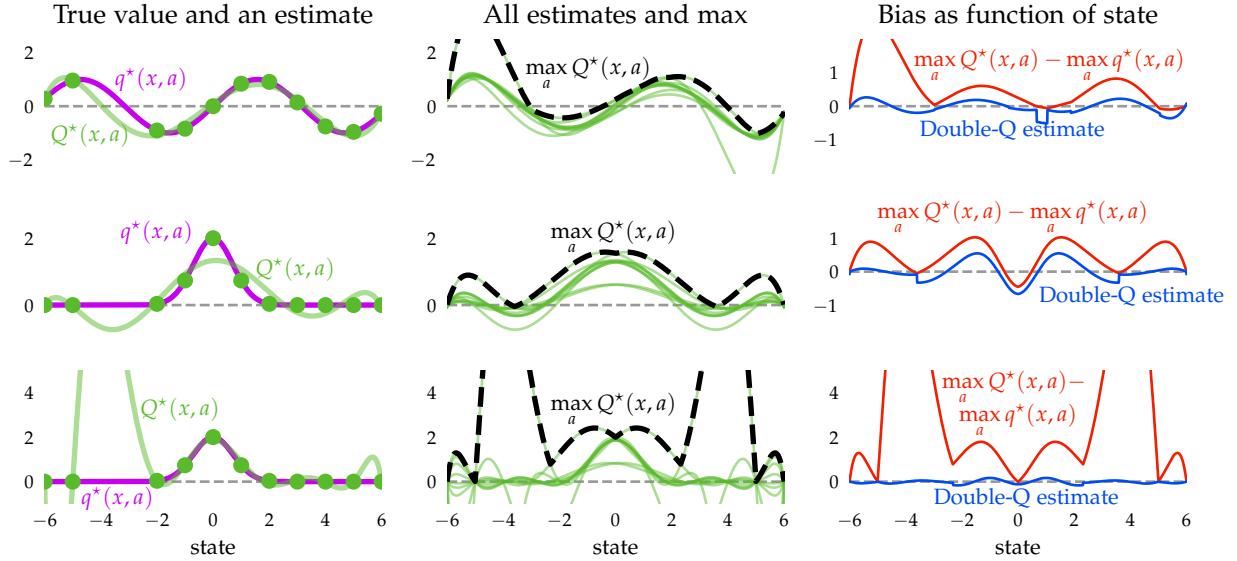
compare to the Q-learning loss (12.11)

The loss can also be interpreted (in an online sense) as performing regular Q-learning with the modification that the target network θ^{old} is not updated to θ after every observed transition, but instead only after observing $|\mathcal{D}|$ -many transitions. This technique is known as *experience replay*. Another approach is *Polyak averaging* where the target network is gradually “nudged” by the neural network used to estimate the Q-function.

Maximization bias: Now, observe that the estimates Q^* are noisy estimates of q^* and consider the term,

$$\max_{a' \in \mathcal{A}} q^*(x', a') \approx \max_{a' \in \mathcal{A}} Q^*(x', a'; \theta^{\text{old}}),$$

from the loss function (12.16). This term maximizes a *noisy* estimate of q^* , which leads to a *biased* estimate of $\max q^*$ as can be seen in Figure 12.1. The fact that the update rules of Q-learning and DQN are affected by inaccuracies (i.e., noise in the estimates) of the learned Q-function is known as the “*maximization bias*”.



Double DQN (DDQN) is an algorithm that addresses this maximization bias (Van Hasselt et al., 2016). Instead of picking the optimal action with respect to the old network, it picks the optimal action with respect to the new network,

$$\ell_{\text{DDQN}}(\theta; \mathcal{D}) \doteq \frac{1}{2} \sum_{(x, a, r, x') \in \mathcal{D}} \left(r + \gamma Q^*(x', a^*(x'; \theta); \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2 \quad (12.17)$$

$$\text{where } a^*(x'; \theta) \doteq \arg \max_{a' \in \mathcal{A}} Q^*(x', a'; \theta). \quad (12.18)$$

Intuitively, this change ensures that the evaluation of the target network is consistent with the updated Q-function, which makes the algorithm more robust to noise.

Similarly to DQN, this can also be interpreted as the online update,

$$\theta \leftarrow \theta + \alpha_t \left(r + \gamma Q^*(x', a^*(x'; \theta); \theta^{\text{old}}) - Q^*(x, a; \theta) \right) \nabla_{\theta} Q^*(x, a; \theta) \quad (12.19)$$

after observing a single transition (x, a, r, x') where while differentiating, $a^*(x'; \theta)$ is treated as constant with respect to θ . θ^{old} is then updated to θ after observing $|\mathcal{D}|$ -many transitions.

12.3 Policy Approximation

Q-learning defines a policy implicitly by

$$\pi^*(x) \doteq \arg \max_{a \in \mathcal{A}} Q^*(x, a). \quad (12.20)$$

Figure 12.1: Illustration of overestimation during learning. In each state (x -axis), there are 10 actions. The left column shows the true values $v^*(x)$ (purple line). All true action values are defined by $q^*(x, a) \doteq v^*(x)$. The green line shows estimated values $Q^*(x, a)$ for one action as a function of state, fitted to the true value at several sampled states (green dots). The middle column plots show all the estimated values (green), and the maximum of these values (dashed black). The maximum is higher than the true value (purple, left plot) almost everywhere. The right column plots show the difference in red. The blue line in the right plots is the estimate used by Double Q-learning with a second set of samples for each state. The blue line is much closer to zero, indicating less bias. The three rows correspond to different true functions (left, purple) or capacities of the fitted function (left, green). Reproduced with permission from “Deep reinforcement learning with double q-learning” (Van Hasselt et al., 2016).

Q-learning also maximizes over the set of all actions in its update step while learning the Q-function. This is intractable for large and, in particular, continuous action spaces. A natural idea to escape this limitation is to immediately learn an approximate parameterized policy,

$$\pi^*(x) \approx \pi(x; \varphi) \doteq \pi_\varphi(x). \quad (12.21)$$

Methods that find an approximate policy are also called *policy search methods* or *policy gradient methods*.

Whereas with Q-learning, exploration can be encouraged by using an ε -greedy policy, softmax exploration, or an optimistic initialization, we will see later that policy gradient methods fundamentally rely on randomized policies for exploration.

Remark 12.2: Notation

We refer to deterministic policies π in bold, as they can be interpreted as vector-valued functions from \mathcal{X} to \mathcal{A} . We still refer to randomized policies by π , as for each state $x \in \mathcal{X}$ they are represented as a PDF over actions \mathcal{A} .

In particular, we denote by $\pi_\varphi(a | x)$ the probability (density) of playing action a when in state x according to π_φ .

12.3.1 Estimating Policy Values

We will begin by attributing a “value” to a policy. Recall the definition of the discounted payoff G_t from time t , which we are aiming to maximize,

$$G_t = \sum_{m=0}^{\infty} \gamma^m R_{t+m}. \quad \text{see Equation (10.5)}$$

We define $G_{t:T}$ to be the *bounded discounted payoff* until time T ,

$$G_{t:T} \doteq \sum_{m=0}^{T-1-t} \gamma^m R_{t+m}. \quad (12.22)$$

Based on these two random variables, we can define the policy value function:

Definition 12.3 (Policy value function). The *policy value function*,

$$J(\pi) \doteq \mathbb{E}_\pi[G_0] = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \right], \quad (12.23)$$

measures the expected discounted payoff of policy π .⁹ We also define the bounded variant,

$$J_T(\pi) \doteq \mathbb{E}_\pi[G_{0:T}] = \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} \gamma^t R_t \right]. \quad (12.24)$$

⁹ We neglect here that implicitly one also averages over the initial state if this state is not fixed.

For simplicity, we will abbreviate $J(\varphi) \doteq J(\pi_\varphi)$.

Remark 12.4: Notation of policy value function

We adopt the more common notation $J(\pi)$ for the policy value function, as opposed to $j(\pi)$, which would be consistent with our notation of the (true) value functions v^π, q^π . So don't be confused by this: just like the value functions v^π, q^π , the policy value function $J(\pi)$ is a deterministic object, measuring the mean discounted payoff. We will use $\hat{J}(\pi)$ to refer to our estimates of the policy value function.

Naturally, we want to maximize $J(\varphi)$. That is, we want to solve

$$\varphi^* \doteq \arg \max_{\varphi} J(\varphi) \quad (12.25)$$

which is a non-convex optimization problem. Let us see how $J(\varphi)$ can be evaluated to understand the optimization problem better. We will again use a Monte Carlo estimate. Recall that a fixed φ induces a unique Markov chain, which can be simulated. In the episodic setting, each episode (also called *rollout*) of length T yields an independently sampled trajectory,

$$\tau^{(i)} \doteq ((x_0^{(i)}, a_0^{(i)}, r_0^{(i)}, x_1^{(i)}), (x_1^{(i)}, a_1^{(i)}, r_1^{(i)}, x_2^{(i)}), \dots) \quad (12.26)$$

Simulating m rollouts yields the samples $\tau^{(1)}, \dots, \tau^{(m)} \sim \Pi_\varphi$, where Π_φ is the distribution of all trajectories (i.e., rollouts) of the Markov chain induced by policy π_φ . We denote the (bounded) discounted payoff of the i -th rollout by

$$g_{0:T}^{(i)} \doteq \sum_{t=0}^{T-1} \gamma^t r_t^{(i)} \quad (12.27)$$

where $r_t^{(i)}$ is the reward at time t of the i -th rollout. Using a Monte Carlo approximation, we can then estimate $J_T(\varphi)$. Moreover, due to the exponential discounting of future rewards, it is reasonable to approximate the policy value function using bounded trajectories,

$$J(\varphi) \approx J_T(\varphi) \approx \hat{J}_T(\varphi) \doteq \frac{1}{m} \sum_{i=1}^m g_{0:T}^{(i)}. \quad (12.28)$$

12.3.2 Policy Gradient

Policy gradient methods solve the above optimization problem (12.25) by stochastic gradient ascent on the policy parameter φ :

$$\varphi \leftarrow \varphi + \eta \nabla_\varphi J(\varphi). \quad (12.29)$$

How can we compute the policy gradient? Let us first formally define the distribution over trajectories Π_φ that we introduced in the previous section. We can specify the probability of a specific trajectory τ under a policy π_φ by

$$\Pi_\varphi(\tau) = p(x_0) \prod_{t=0}^{T-1} \pi_\varphi(a_t | x_t) p(x_{t+1} | x_t, a_t). \quad (12.30)$$

For optimizing $J(\varphi)$ we need to obtain unbiased gradient estimates:

$$\nabla_\varphi J(\varphi) \approx \nabla_\varphi J_T(\varphi) = \nabla_\varphi \mathbb{E}_{\tau \sim \Pi_\varphi}[G_{0:T}]. \quad (12.31)$$

Note that the expectation integrates over the measure Π_φ , which depends on the parameter φ . Thus, we cannot move the gradient operator inside the expectation as we have often done previously (cf. Appendix A.1.5). This should remind you of the reparameterization trick (see Equation (5.62)) that we used to solve a similar gradient in the context of variational inference. In this context, however, we cannot apply the reparameterization trick.¹⁰ Fortunately, there is another way of estimating this gradient.

Theorem 12.5 (Score gradient estimator). *Under some regularity assumptions, we have*

$$\nabla_\varphi \mathbb{E}_{\tau \sim \Pi_\varphi}[G_0] = \mathbb{E}_{\tau \sim \Pi_\varphi}[G_0 \nabla_\varphi \log \Pi_\varphi(\tau)]. \quad (12.32)$$

This estimator of the gradient is called the score gradient estimator.

Proof. To begin with, let us look at the so-called *score function* of the distribution Π_φ , $\nabla_\varphi \log \Pi_\varphi(\tau)$. Using the chain rule, the score function can be expressed as

$$\nabla_\varphi \log \Pi_\varphi(\tau) = \frac{\nabla_\varphi \Pi_\varphi(\tau)}{\Pi_\varphi(\tau)} \quad (12.33)$$

and by rearranging the terms, we obtain

$$\nabla_\varphi \Pi_\varphi(\tau) = \Pi_\varphi(\tau) \nabla_\varphi \log \Pi_\varphi(\tau). \quad (12.34)$$

This is called the *score function trick*.¹¹

Now, assuming that state and action spaces are continuous, we obtain

$$\begin{aligned} \nabla_\varphi \mathbb{E}_{\tau \sim \Pi_\varphi}[G_0] &= \nabla_\varphi \int \Pi_\varphi(\tau) \cdot G_0 d\tau \\ &= \int \nabla_\varphi \Pi_\varphi(\tau) \cdot G_0 d\tau \\ &= \int G_0 \cdot \Pi_\varphi(\tau) \nabla_\varphi \log \Pi_\varphi(\tau) d\tau \\ &= \mathbb{E}_{\tau \sim \Pi_\varphi}[G_0 \nabla_\varphi \log \Pi_\varphi(\tau)]. \end{aligned} \quad \square$$

¹⁰This is because the distribution Π_φ is generally not reparameterizable. We will, however, see that reparameterization gradients are also useful in reinforcement learning. See, e.g., Sections 12.5.1 and 13.1.2.

¹¹We have already applied this “trick” in Problem 6.9.

using the definition of expectation (1.19)

using the regularity assumptions to swap gradient and integral

using the score function trick (12.34)

interpreting the integral as an expectation over Π_φ

Intuitively, maximizing $J(\varphi)$ increases the probability of policies with high returns and decreases the probability of policies with low returns.

To use the score gradient estimator for estimating the gradient, we need to compute $\nabla_\varphi \log \Pi_\varphi(\tau)$.

$$\begin{aligned} & \nabla_\varphi \log \Pi_\varphi(\tau) \\ &= \nabla_\varphi \left(\log p(x_0) + \sum_{t=0}^{T-1} \log \pi_\varphi(a_t | x_t) + \sum_{t=0}^{T-1} \log p(x_{t+1} | x_t, a_t) \right) \\ &= \nabla_\varphi \log p(x_0) + \sum_{t=0}^{T-1} \nabla_\varphi \log \pi_\varphi(a_t | x_t) + \sum_{t=0}^{T-1} \nabla_\varphi \log p(x_{t+1} | x_t, a_t) \\ &= \sum_{t=0}^{T-1} \nabla_\varphi \log \pi_\varphi(a_t | x_t). \end{aligned} \quad (12.35)$$

using the definition of the distribution over trajectories Π_φ

using that the first and third term are independent of φ

When using a neural network for the parameterization of the policy π , we can use automatic differentiation to compute the gradients.

The expectation of the score gradient estimator (12.32) can be approximated using Monte Carlo sampling,

$$\nabla_\varphi J_T(\varphi) \approx \nabla_\varphi \hat{J}_T(\varphi) \doteq \frac{1}{m} \sum_{i=1}^m g_{0:T}^{(i)} \sum_{t=0}^{T-1} \nabla_\varphi \log \pi_\varphi(a_t^{(i)} | x_t^{(i)}). \quad (12.36)$$

However, typically the variance of these estimates is very large. Using so-called *baselines*, we can reduce the variance dramatically (7).

Problem 12.3

Lemma 12.6 (Score gradients with baselines). *We have,*

$$\mathbb{E}_{\tau \sim \Pi_\varphi} [G_0 \nabla_\varphi \log \Pi_\varphi(\tau)] = \mathbb{E}_{\tau \sim \Pi_\varphi} [(G_0 - b) \nabla_\varphi \log \Pi_\varphi(\tau)]. \quad (12.37)$$

Here, $b \in \mathbb{R}$ is called a baseline.

Proof. For the term to the right, we have due to linearity of expectation (1.20),

$$\begin{aligned} \mathbb{E}_{\tau \sim \Pi_\varphi} [(G_0 - b) \nabla_\varphi \log \Pi_\varphi(\tau)] &= \mathbb{E}_{\tau \sim \Pi_\varphi} [G_0 \nabla_\varphi \log \Pi_\varphi(\tau)] \\ &\quad - \mathbb{E}_{\tau \sim \Pi_\varphi} [b \cdot \nabla_\varphi \log \Pi_\varphi(\tau)]. \end{aligned}$$

Thus, it remains to show that the second term is zero,

$$\begin{aligned} \mathbb{E}_{\tau \sim \Pi_\varphi} [b \cdot \nabla_\varphi \log \Pi_\varphi(\tau)] &= b \cdot \int \Pi_\varphi(\tau) \nabla_\varphi \log \Pi_\varphi(\tau) d\tau \\ &= b \cdot \int \Pi_\varphi(\tau) \frac{\nabla_\varphi \Pi_\varphi(\tau)}{\Pi_\varphi(\tau)} d\tau \\ &= b \cdot \int \nabla_\varphi \Pi_\varphi(\tau) d\tau \\ &= b \cdot \nabla_\varphi \int \Pi_\varphi(\tau) d\tau \\ &= b \cdot \nabla_\varphi 1 = 0. \end{aligned}$$

using the definition of expectation (1.19)

substituting the score function (12.33), “undoing the score function trick”

$\Pi_\varphi(\tau)$ cancels

integrating a PDF over its domain is 1 and the derivative of a constant is 0

One can even show, that we can subtract arbitrary baselines depending on *previous states* (?)

Problem 12.4

Example 12.7: Downstream returns

A commonly used state-dependent baseline is

$$b(\tau_{0:t-1}) \doteq \sum_{m=0}^{t-1} \gamma^m r_m. \quad (12.38)$$

This baseline subtracts the returns of all actions before time t . Intuitively, using this baseline, the score gradient only considers downstream returns. Recall from Equation (12.22) that we defined $G_{t:T}$ as the bounded discounted payoff from time t . It is also commonly called the (bounded) *downstream return* (or *reward to go*) beginning at time t .

For a fixed trajectory τ that is bounded at time T , we have

$$G_0 - b(\tau_{0:t-1}) = \gamma^t G_{t:T}, \quad (12.39)$$

yielding the gradient estimator,

$$\begin{aligned} \nabla_\varphi J(\varphi) &\approx \nabla_\varphi J_T(\varphi) = \mathbb{E}_{\tau \sim \Pi_\varphi} [G_0 \nabla_\varphi \log \Pi_\varphi(\tau)] \\ &= \mathbb{E}_{\tau \sim \Pi_\varphi} \left[\sum_{t=0}^{T-1} \gamma^t G_{t:T} \nabla_\varphi \log \pi_\varphi(a_t | x_t) \right]. \end{aligned} \quad (12.40)$$

using the score gradient estimator
(12.32)

using a state-dependent baseline
(12.102)

Performing stochastic gradient descent with the score gradient estimator and downstream returns is known as the *REINFORCE algorithm* (Williams, 1992) which is shown in Algorithm 12.8.

Algorithm 12.8: REINFORCE algorithm

```

1 initialize policy weights  $\varphi$ 
2 repeat
3   generate an episode (i.e., rollout) to obtain trajectory  $\tau$ 
4   for  $t = 0$  to  $T - 1$  do
5     set  $g_{t:T}$  to the downstream return from time  $t$ 
6      $\varphi \leftarrow \varphi + \eta \gamma^t g_{t:T} \nabla_\varphi \log \pi_\varphi(a_t | x_t)$            // (12.41)
7 until converged

```

The variance of REINFORCE can be reduced further. A common tech-

nique is to subtract a term b_t from the downstream returns,

$$\nabla_{\varphi} J(\varphi) = \mathbb{E}_{\tau \sim \Pi_{\varphi}} \left[\sum_{t=0}^T \gamma^t (G_{t:T} - b_t) \nabla_{\varphi} \log \pi_{\varphi}(a_t | x_t) \right]. \quad (12.42)$$

For example, we can subtract the t -independent mean reward to go,

$$b_t \doteq b = \frac{1}{T} \sum_{t'=0}^{T-1} G_{t':T}. \quad (12.43)$$

The main advantage of policy gradient methods such as REINFORCE is that they can be used in continuous action spaces. However, REINFORCE is not guaranteed to find an optimal policy. Even when operating in very small domains, REINFORCE can get stuck in local optima.

Typically, policy gradient methods are slow due to the large variance in the score gradient estimates. Because of this, they need to take small steps and require many rollouts of a Markov chain. Moreover, we cannot reuse data from previous rollouts, as policy gradient methods are fundamentally on-policy.¹²

Next, we will combine value approximation techniques like Q-learning and policy gradient methods, leading to an often more practical family of methods called actor-critic methods.

¹² This is because the score gradient estimator is used to obtain gradients of the policy value function with respect to the *current* policy.

12.4 On-policy Actor-Critics

Actor-Critic methods reduce the variance of policy gradient estimates by using ideas from value function approximation. They use function approximation both to approximate value functions *and* to approximate policies. The goal for these algorithms is to scale to reinforcement learning problems, where we both have large state spaces and large action spaces.

12.4.1 Advantage Function

A key concept of actor-critic methods is the advantage function.

Definition 12.9 (Advantage function). Given a policy π , the *advantage function*,

$$a^{\pi}(x, a) \doteq q^{\pi}(x, a) - v^{\pi}(x) \quad (12.44)$$

$$= q^{\pi}(x, a) - \mathbb{E}_{a' \sim \pi(x)} [q^{\pi}(x, a')], \quad (12.45)$$

using Equation (10.14)

measures the advantage of picking action $a \in \mathcal{A}$ when in state $x \in \mathcal{X}$ over simply following policy π .

It follows immediately from Equation (12.45) that for any policy π and state $x \in \mathcal{X}$, there exists an action $a \in \mathcal{A}$ such that $a^\pi(x, a)$ is non-negative,

$$\max_{a \in \mathcal{A}} a^\pi(x, a) \geq 0. \quad (12.46)$$

Moreover, it follows directly from Bellman's theorem (10.31) that

$$\pi \text{ is optimal} \iff \forall x \in \mathcal{X}, a \in \mathcal{A} : a^\pi(x, a) \leq 0. \quad (12.47)$$

In other words, quite intuitively, π is optimal if and only if there is no action that has an advantage in any state over the action that is played by π .

Finally, we can re-define the greedy policy π_q with respect to the state-action value function q as

$$\pi_q(x) \doteq \arg \max_{a \in \mathcal{A}} a(x, a) \quad (12.48)$$

since

$$\arg \max_{a \in \mathcal{A}} a(x, a) = \arg \max_{a \in \mathcal{A}} q(x, a) - v(x) = \arg \max_{a \in \mathcal{A}} q(x, a),$$

as $v(x)$ is independent of a . This coincides with our initial definition of greedy policies in Equation (10.29). Intuitively, the advantage function is a shifted version of the state-action value function q that is relative to 0. Using this quantity rather than q often has numerical advantages.

12.4.2 Policy Gradient Theorem

Recall the score gradient estimator (12.40) that we had introduced in the previous section,

$$\nabla_\varphi J_T(\varphi) = \mathbb{E}_{\tau \sim \Pi_\varphi} \left[\sum_{t=0}^{T-1} \gamma^t G_{t:T} \nabla_\varphi \log \pi_\varphi(a_t | x_t) \right].$$

Previously, we have approximated the policy value function $J(\varphi)$ by the bounded policy value function $J_T(\varphi)$. We said that this approximation was “reasonable” due to the diminishing returns. Essentially, we have “cut off the tails” of the policy value function. Let us now reinterpret score gradients while taking into account the tails of $J(\varphi)$.

$$\begin{aligned} \nabla_\varphi J(\varphi) &= \lim_{T \rightarrow \infty} \nabla_\varphi J_T(\varphi) \\ &= \sum_{t=0}^{\infty} \mathbb{E}_{\tau \sim \Pi_\varphi} [\gamma^t G_t \nabla_\varphi \log \pi_\varphi(a_t | x_t)]. \end{aligned} \quad (12.49)$$

substituting the score gradient estimator with downstream returns (12.40) and using linearity of expectation (1.20)

Observe that because the expectations only consider downstream returns, we can disregard all data from the trajectory prior to time t . Let us define

$$\tau_{t:\infty} \doteq ((\mathbf{x}_t, \mathbf{a}_t, r_t, \mathbf{x}_{t+1}), (\mathbf{x}_{t+1}, \mathbf{a}_{t+1}, r_{t+1}, \mathbf{x}_{t+2}), \dots), \quad (12.50)$$

as the trajectory from time step t . Then,

$$= \sum_{t=0}^{\infty} \mathbb{E}_{\tau_{t:\infty} \sim \Pi_{\varphi}} [\gamma^t G_t \nabla_{\varphi} \log \pi_{\varphi}(\mathbf{a}_t | \mathbf{x}_t)].$$

We now condition on \mathbf{x}_t and \mathbf{a}_t ,

$$= \sum_{t=0}^{\infty} \mathbb{E}_{\mathbf{x}_t, \mathbf{a}_t} [\gamma^t \mathbb{E}_{r_t, \tau_{t+1:\infty}} [G_t | \mathbf{x}_t, \mathbf{a}_t] \nabla_{\varphi} \log \pi_{\varphi}(\mathbf{a}_t | \mathbf{x}_t)].$$

using that $\pi_{\varphi}(\mathbf{a}_t | \mathbf{x}_t)$ is a constant given \mathbf{x}_t and \mathbf{a}_t

Observe that averaging over the trajectories $\mathbb{E}_{\tau \sim \Pi_{\varphi}} [\cdot]$ that are sampled according to policy π_{φ} is equivalent to our shorthand notation $\mathbb{E}_{\pi_{\varphi}} [\cdot]$ from Equation (10.6),

$$\begin{aligned} &= \sum_{t=0}^{\infty} \mathbb{E}_{\mathbf{x}_t, \mathbf{a}_t} [\gamma^t \mathbb{E}_{\pi_{\varphi}} [G_t | \mathbf{x}_t, \mathbf{a}_t] \nabla_{\varphi} \log \pi_{\varphi}(\mathbf{a}_t | \mathbf{x}_t)] \\ &= \sum_{t=0}^{\infty} \mathbb{E}_{\mathbf{x}_t, \mathbf{a}_t} [\gamma^t q^{\pi_{\varphi}}(\mathbf{x}_t, \mathbf{a}_t) \nabla_{\varphi} \log \pi_{\varphi}(\mathbf{a}_t | \mathbf{x}_t)]. \end{aligned} \quad (12.51)$$

using the definition of the Q-function (10.8)

It turns out that $\mathbb{E}_{\mathbf{x}_t, \mathbf{a}_t} [q^{\pi_{\varphi}}(\mathbf{x}_t, \mathbf{a}_t)]$ exhibits much less variance than our previous estimator $\mathbb{E}_{\mathbf{x}_t, \mathbf{a}_t} \mathbb{E}_{\pi_{\varphi}} [G_t | \mathbf{x}_t, \mathbf{a}_t]$. Equation (12.51) is known as the *policy gradient theorem*.

Often, the policy gradient theorem is stated in a slightly rephrased form in terms of the *discounted state occupancy measure*,¹³

$$\rho_{\varphi}^{\infty}(\mathbf{x}) \doteq (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t p_{\mathbf{X}_t}(\mathbf{x}). \quad (12.52)$$

The factor $(1 - \gamma)$ ensures that ρ_{φ}^{∞} is a probability density as

$$\int \rho_{\varphi}^{\infty}(\mathbf{x}) d\mathbf{x} = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \int p_{\mathbf{X}_t}(\mathbf{x}) d\mathbf{x} = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t = 1.$$

Intuitively, $\rho_{\varphi}^{\infty}(\mathbf{x})$ measures how often we visit state \mathbf{x} when following policy π_{φ} . It can be thought of as a “discounted frequency”.

Theorem 12.10 (Policy gradient theorem in terms of ρ_{φ}^{∞}). *Policy gradients can be represented in terms of the Q-function,*

$$\nabla_{\varphi} J(\varphi) \propto \mathbb{E}_{\mathbf{x} \sim \rho_{\varphi}^{\infty}} \mathbb{E}_{\mathbf{a} \sim \pi_{\varphi}(\cdot | \mathbf{x})} [q^{\pi_{\varphi}}(\mathbf{x}, \mathbf{a}) \nabla_{\varphi} \log \pi_{\varphi}(\mathbf{a} | \mathbf{x})]. \quad (12.53)$$

Proof. The right hand side of Equation (12.51) can be expressed as

$$\sum_{t=0}^{\infty} \int p_{\mathbf{X}_t}(\mathbf{x}) \cdot \mathbb{E}_{\mathbf{a}_t \sim \pi_{\varphi}(\cdot | \mathbf{x})} [\gamma^t q^{\pi_{\varphi}}(\mathbf{x}, \mathbf{a}_t) \nabla_{\varphi} \log \pi_{\varphi}(\mathbf{a}_t | \mathbf{x})] d\mathbf{x}$$

¹³ Depending on the reward setting, there exist various variations of the policy gradient theorem. We derived the variant for infinite-horizon discounted payoffs. “Reinforcement learning: An introduction” (Sutton and Barto, 2018) derive the variant for undiscounted average rewards.

$$= \frac{1}{1-\gamma} \int \rho_\varphi^\infty(x) \cdot \mathbb{E}_{a \sim \pi_\varphi(\cdot|x)} [q^{\pi_\varphi}(x, a) \nabla_\varphi \log \pi_\varphi(a | x)] dx$$

where we swapped the order of sum and integral and reorganized terms. \square

Matching our intuition, according to the policy gradient theorem, maximizing $J(\varphi)$ corresponds to increasing the probability of actions with a large value and decreasing the probability of actions with a small value, taking into account how often the resulting policy visits certain states.

Observe that we cannot use the policy gradient to calculate the gradient exactly, as we do not know q^{π_φ} . Instead, we will use bootstrapping estimates Q^{π_φ} of q^{π_φ} .

12.4.3 A First Actor-Critic

Actor-Critic methods consist of two components:

- a parameterized policy, $\pi(a | x; \varphi) \doteq \pi_\varphi$, which is called *actor*; and
- a value function approximation, $q^{\pi_\varphi}(x, a) \approx Q^{\pi_\varphi}(x, a; \theta)$, which is called *critic*. In the following, we will abbreviate Q^{π_φ} by Q . (12.55)

In deep reinforcement learning, neural networks are used to parameterize both actor and critic. Therefore, in principle, the actor-critic framework allows scaling to both large state spaces and large action spaces. We begin by discussing on-policy actor-critics.

One approach in the online setting (i.e., non-episodic setting), is to simply use SARSA for learning the critic. To learn the actor, we use stochastic gradient descent with gradients obtained using single samples from

$$\nabla_\varphi J(\varphi) \approx \nabla_\varphi \hat{J}(\varphi) \doteq \sum_{t=0}^{\infty} \mathbb{E}_{(x_t, a_t) \sim \pi_\varphi} [\gamma^t Q(x_t, a_t; \theta) \nabla_\varphi \log \pi_\varphi(a_t | x_t)] \quad (12.56)$$

where Q is a bootstrapping estimate of q^{π_φ} . This algorithm is known as *online actor-critic* or *Q actor-critic* and shown in Algorithm 12.11.

Comparing to the derivation for TD-learning from Equation (12.8), we observe that Equation (12.58) corresponds to the SARSA update rule.¹⁴ Due to the use of SARSA for learning the critic, this algorithm is fundamentally on-policy.

Crucially, by neglecting the dependence of the bootstrapping estimate Q on the policy parameters φ , we introduce bias in the gradient estimates. In other words, using the bootstrapping estimate Q means that

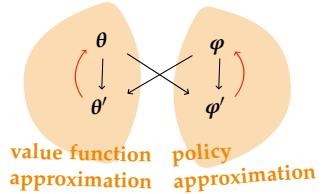


Figure 12.2: Illustration of one iteration of actor-critic methods. The dependencies between the actors and critics are shown as arrows. Methods differ in the exact order in which actor and critic are updated.

see Equation (12.51)

¹⁴ The gradient with respect to θ_Q appears analogously to our derivation of approximate Q-learning (12.15).

Algorithm 12.11: Online actor-critic

```

1 initialize parameters  $\varphi$  and  $\theta$ 
2 repeat
3   use  $\pi_\varphi$  to obtain transition  $(x, a, r, x')$ 
4    $\delta = r + \gamma Q(x', \pi_\varphi(x'); \theta) - Q(x, a; \theta)$ 
5   // actor update
6    $\varphi \leftarrow \varphi + \eta \gamma^t Q(x, a; \theta) \nabla_\varphi \log \pi_\varphi(a | x)$            // (12.57)
7   // critic update
8    $\theta \leftarrow \theta + \eta \delta \nabla_\theta Q(x, a; \theta)$            // (12.58)
9 until converged

```

the resulting gradient direction might not be a valid ascent direction. In particular, the actor is not guaranteed to improve. Still, it turns out that under strong so-called “compatibility conditions” that are rarely satisfied in practice, a valid ascent direction can be guaranteed.

12.4.4 Improved Actor-Critics

Reducing variance: To further reduce the variance of the gradient estimates, it turns out that a similar approach to the baselines we discussed in the previous section on policy gradient methods is useful. A common approach is to subtract the state value function from estimates of the Q-function,

$$\varphi \leftarrow \varphi + \eta_t \gamma^t (Q(x, a; \theta) - V(x; \theta)) \nabla_\varphi \log \pi_\varphi(a | x) \quad (12.59)$$

$$= \varphi + \eta_t \gamma^t A(x, a; \theta) \nabla_\varphi \log \pi_\varphi(a | x) \quad (12.60)$$

where $A(x, a; \theta)$ is a bootstrapped estimate of the advantage function a^{π_φ} . This algorithm is known as *advantage actor-critic* (A2C) (Mnih et al., 2016). Recall that the Q-function is an absolute quantity, whereas the advantage function is a relative quantity, where the sign is informative for the gradient direction. Intuitively, an absolute value is harder to estimate than the sign. Actor-Critic methods are therefore often implemented with respect to the advantage function rather than the Q-function.

Taking a step back, observe that policy gradient methods such as REINFORCE generally have *high variance* in their gradient estimates. However, due to using Monte Carlo estimates of G_t , the gradient estimates are *unbiased*. In contrast, using a bootstrapped Q-function to obtain gradient estimates yields estimates with a *smaller variance*, but those estimates are *biased*. We are faced with a *bias-variance tradeoff*. A natural approach is therefore to blend both gradient estimates to allow

using the definition of the advantage function (12.44)

for effectively trading bias and variance. This leads to algorithms such as *generalized advantage estimation* (GAE) (Schulman et al., 2016).

Exploration: Similarly to REINFORCE, actor-critic methods typically rely on randomization in the policy to encourage exploration, the idea being that if the policy is stochastic, then the agent will visit a diverse set of states. The inherent stochasticity of the policy is, however, often insufficient. A common problem is that the policy quickly “collapses” to a deterministic policy since the objective function is greedily exploitative. A common workaround is to use an ε -greedy policy (cf. Section 11.3.1) or to explicitly encourage the policy to exhibit uncertainty by adding an entropy term to the objective function (more on this in Section 12.6). However, note that for on-policy methods, changing the policy also changes the value function learned by the critic.

Improving sample efficiency: Actor-Critic methods typically suffer from low sample efficiency. When additionally using an on-policy method, actor-critics often need an extremely large number of interactions before learning a near-optimal policy, because they cannot reuse past data. Allowing to reuse past data is a major advantage of off-policy methods like Q-learning.

One well-known variant that slightly improves the sample efficiency is *trust-region policy optimization* (TRPO) (Schulman et al., 2015). TRPO uses multiple iterations, where in each iteration a fixed critic is used to optimize the policy.¹⁵ During iteration k , we select

$$\boldsymbol{\varphi}_{k+1} \leftarrow \arg \max_{\boldsymbol{\varphi}} \hat{J}(\boldsymbol{\varphi}) \quad \text{subject to } \mathbb{E}_{x \sim \rho_{\boldsymbol{\varphi}_k}^{\infty}} \text{KL}(\pi_{\boldsymbol{\varphi}_k}(\cdot | x) \| \pi_{\boldsymbol{\varphi}}(\cdot | x)) \leq \delta$$

(12.61)

for some fixed $\delta > 0$ and where

$$\hat{J}(\boldsymbol{\varphi}) \doteq \mathbb{E}_{x \sim \rho_{\boldsymbol{\varphi}_k}^{\infty}, a \sim \pi_{\boldsymbol{\varphi}_k}(\cdot | x)} [w_k(\boldsymbol{\varphi}; x, a) A^{\pi_{\boldsymbol{\varphi}_k}}(x, a)]. \quad (12.62)$$

Notably, \hat{J} is an expectation with respect to the *previous* policy $\pi_{\boldsymbol{\varphi}_k}$ and the previous critic $A^{\pi_{\boldsymbol{\varphi}_k}}$. TRPO uses *importance sampling* where the importance weights (called “likelihood ratios”),

$$w_k(\boldsymbol{\varphi}; x, a) \doteq \frac{\pi_{\boldsymbol{\varphi}}(a | x)}{\pi_{\boldsymbol{\varphi}_k}(a | x)},$$

are used to correct for taking the expectation over the previous policy. When $w_k \approx 1$ the policies $\pi_{\boldsymbol{\varphi}}$ and $\pi_{\boldsymbol{\varphi}_k}$ are similar, whereas when $w_k \ll 1$ or $w_k \gg 1$, the policies differ significantly. To be able to assume that the fixed critic is a good approximation within a certain “trust region” (i.e., one iteration), we impose the constraint

$$\mathbb{E}_{x \sim \rho_{\boldsymbol{\varphi}_k}^{\infty}} \text{KL}(\pi_{\boldsymbol{\varphi}_k}(\cdot | x) \| \pi_{\boldsymbol{\varphi}}(\cdot | x)) \leq \delta$$

¹⁵ Intuitively, each iteration performs a collection of gradient ascent steps.

to optimize only in the “neighborhood” of the current policy. This constraint is also necessary for the importance weights not to blow up.

Remark 12.12: Estimating the KL-divergence

Instead of naive computation with $\mathbb{E}_{a \sim \pi_{\varphi_k}(\cdot|x)}[-\log w_k(\varphi; x, a)]$, the KL-divergence is commonly estimated by Monte Carlo samples of

$$\begin{aligned} & \text{KL}(\pi_{\varphi_k}(\cdot|x) \| \pi_{\varphi}(\cdot|x)) \\ &= \mathbb{E}_{a \sim \pi_{\varphi_k}(\cdot|x)}[w_k(\varphi; x, a) - 1 - \log w_k(\varphi; x, a)], \end{aligned}$$

which adds the “baseline” $w_k(\varphi; x, a) - 1$ with mean 0. Observe that this estimator is unbiased, always non-negative since $\log(x) \leq x - 1$ for all x , while having a lower variance than the naive estimator.

Taking the expectation with respect to the previous policy π_{φ_k} means that we can reuse data from rollouts within the same iteration. That is, TRPO allows reusing past data as long as it can still be “trusted”. This makes TRPO “somewhat” off-policy. Fundamentally, though, TRPO is still an on-policy method.

Proximal policy optimization (PPO) is a family of heuristic variants of TRPO which replace the constrained optimization problem of Equation (12.61) by the unconstrained optimization of a regularized objective (Schulman et al., 2017; Wang et al., 2020). PPO algorithms often work well in practice. One canonical PPO method uses the modified objective

$$\varphi_{k+1} \leftarrow \arg \max_{\varphi} \hat{J}(\varphi) - \lambda \mathbb{E}_{x \sim \rho_{\varphi_k}^{\infty}} \text{KL}(\pi_{\varphi_k}(\cdot|x) \| \pi_{\varphi}(\cdot|x)) \quad (12.63)$$

with some $\lambda > 0$, which regularizes towards the trust region. Another common variant of PPO is based on controlling the importance weights directly rather than regularizing by the KL-divergence. PPO is used, for example, to train large-scale language models such as GPT (Stiennon et al., 2020; OpenAI, 2023) which we will discuss in more detail in Section 12.7. There we will also see that the objective from Equation (12.63) can be cast as performing probabilistic inference.

Improving computational efficiency: A practical problem with the above methods is that the estimation of the advantage function $A(x, a; \theta)$ requires training a separate critic, next to the policy (i.e., the actor) parameterized by φ . This can be computationally expensive. In particular, when both models are large neural networks (think multiple billions of parameters each), training both models is computationally

prohibitive. Now, recall that we introduced critics in the first place to reduce the variance of the policy gradient estimates.¹⁶ *Group relative policy optimization* (GRPO) replaces the critic in PPO with simple Monte Carlo estimates of the advantage function (Shao et al., 2024):

$$\begin{aligned}\hat{J}(\boldsymbol{\varphi}) &\doteq \mathbb{E}_{\{\tau^{(i)}\}_{i=1}^m \sim \Pi_{\boldsymbol{\varphi}_k}(\cdot|\mathbf{x})} \left[\frac{1}{m} \sum_{i=1}^m \sum_{t=1}^T w_k(\boldsymbol{\varphi}; \mathbf{a}^{(i)}) \hat{A}_{t,i}^{\pi_{\boldsymbol{\varphi}_k}} \right], \\ \text{where } \hat{A}_{t,i}^{\pi_{\boldsymbol{\varphi}_k}} &\doteq \frac{g_{t:T}^{(i)} - \text{mean}(\{\tau^{(i)}\})}{\text{std}(\{\tau^{(i)}\})}\end{aligned}\quad (12.64)$$

estimates the advantage of action $\mathbf{a}^{(i)}$ at time t by comparing to the mean reward and normalizing by the standard deviation of rewards from all trajectories $\tau^{(i)}$.¹⁷ GRPO combines Monte Carlo sampling and baselines for variance reduction with the trust-region optimization of PPO, leading to a method that is more sample efficient than naive REINFORCE while being computationally more efficient than PPO.

12.5 Off-policy Actor-Critics

In many applications, sample efficiency is crucial. Either because requiring too many interactions is computationally prohibitive or because obtaining sufficiently many samples for learning a near-optimal policy is simply impossible. We therefore now want to look at a separate family of actor-critic methods, which are off-policy, and hence, allow for the reuse of past data. These algorithms use the reparameterization gradient estimates, which we encountered before in the context of variational inference,¹⁸ instead of score gradient estimators.

The on-policy methods that we discussed in the previous section can be understood as performing a variant of *policy iteration*, where we use an estimate of the state-action value function of the current policy and then try to improve that policy by acting greedily with respect to this estimate. They mostly vary in how improving the policy is traded with improving the estimate of its value. Fundamentally, these methods rely on policy evaluation.¹⁹

The techniques that we will introduce in this section are much more closely related to value iteration, essentially making use of Bellman's optimality principle to learn the optimal value function directly which characterizes the optimal policy.

To begin with, let us assume that the policy π is deterministic. We will later lift this restriction in Section 12.5.1. Recall that our initial motivation to consider policy gradient methods and then actor-critic

¹⁶ That is, we moved from the REINFORCE policy update (12.41) to the actor-critic policy update (12.57).

¹⁷ In Equation (12.36), we have already seen that using a Monte Carlo estimate of returns is a simple approach to reduce variance without needing to learn a critic.

¹⁸ see Section 5.5.1

¹⁹ Policy evaluation is at the core of policy iteration. See Algorithm 10.14 for the definition of policy iteration and Section 10.2 for a summary of policy evaluation in the context of Markov decision processes.

methods was the intractability of the DQN loss

$$\ell_{\text{DQN}}(\theta; \mathcal{D}) = \frac{1}{2} \sum_{(x, a, r, x') \in \mathcal{D}} \left(r + \gamma \max_{a' \in \mathcal{A}} Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2$$

see Equation (12.16)

when the action space \mathcal{A} is large. What if we simply replace the exact maximum over actions by a parameterized policy?

$$\ell_{\text{DQN}}(\theta; \mathcal{D}) \approx \frac{1}{2} \sum_{(x, a, r, x') \in \mathcal{D}} \left(r + \gamma Q^*(x', \pi_\varphi(x'); \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2. \quad (12.65)$$

We want to train our parameterized policy to learn the maximization over actions, that is, to approximate the greedy policy²⁰

$$\pi_\varphi(x) \approx \pi^*(x) = \arg \max_{a \in \mathcal{A}} Q^*(x, a; \theta). \quad (12.66)$$

The key idea is that if we use a “rich-enough” parameterization of policies, selecting the greedy policy with respect to Q^* is equivalent to

$$\varphi^* = \arg \max_{\varphi} \mathbb{E}_{x \sim \mu} [Q^*(x, \pi_\varphi(x); \theta)] \quad (12.67)$$

where $\mu(x) > 0$ is an *exploration distribution* over states with full support.²¹ We refer to this expectation by

$$\hat{J}_\mu(\varphi; \theta) \doteq \mathbb{E}_{x \sim \mu} [Q^*(x, \pi_\varphi(x); \theta)]. \quad (12.68)$$

²⁰ Here, we already apply the improvement of DDQN to use the most-recent estimate of the Q -function for action selection (see Equation (12.17)).

Commonly, the exploration distribution μ is taken to be the distribution that samples states uniformly from a replay buffer. Note that we can easily obtain unbiased gradient estimates of \hat{J}_μ with respect to φ :

$$\nabla_\varphi \hat{J}_\mu(\varphi; \theta) = \mathbb{E}_{x \sim \mu} [\nabla_\varphi Q^*(x, \pi_\varphi(x); \theta)]. \quad (12.69)$$

see Appendix A.1.5

Analogously to on-policy actor-critics (see Equation (12.56)), we use a bootstrapping estimate of Q^* . That is, we neglect the dependence of the critic Q^* on the actor π_φ , and in particular, the policy parameters φ . We have seen that bootstrapping works with Q-learning, so there is reason to hope that it will work in this context too. This then allows us to use the chain rule to compute the gradient,

$$\nabla_\varphi Q^*(x, \pi_\varphi(x); \theta) = D_\varphi \pi_\varphi(x) \cdot \nabla_a Q^*(x, a; \theta)|_{a=\pi_\varphi(x)}. \quad (12.70)$$

This corresponds to evaluating the bootstrapping estimate of the Q -function at $\pi_\varphi(x)$ and obtaining a gradient estimate of the policy estimate (e.g., through automatic differentiation). Note that as π_φ is vector-valued, $D_\varphi \pi_\varphi(x)$ is the Jacobian of π_φ evaluated at x .

Exploration: Now that we have estimates of the gradient of our optimization target \hat{J}_μ , it is natural to ask how we should select actions (based on π_φ) to trade exploration and exploitation. As we have seen, policy gradient techniques rely on the randomness in the policy to explore, but here we consider deterministic policies. As our method is off-policy, a simple idea in continuous action spaces is to add Gaussian noise to the action selected by π_φ — also known as *Gaussian noise “dithering”*.²² This corresponds to an algorithm called *deep deterministic policy gradients* (Lillicrap et al., 2016) shown in Algorithm 12.13. This algorithm is essentially equivalent to Q-learning with function approximation (e.g., DQN),²³ with the only exception that we replace the maximization over actions with the learned policy π_φ .

²² Intuitively, this adds “additional randomness” to the policy π_φ .

²³ see Equation (12.16)

Algorithm 12.13: Deep deterministic policy gradients, DDPG

```

1 initialize  $\varphi, \theta$ , a (possibly non-empty) replay buffer  $\mathcal{D} = \emptyset$ 
2 set  $\varphi^{\text{old}} = \varphi$  and  $\theta^{\text{old}} = \theta$ 
3 for  $t = 0$  to  $\infty$  do
4   observe state  $x$ , pick action  $a = \pi_\varphi(x) + \varepsilon$  for  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \lambda I)$ 
5   execute  $a$ , observe  $r$  and  $x'$ 
6   add  $(x, a, r, x')$  to the replay buffer  $\mathcal{D}$ 
7   if collected “enough” data then
8     // policy improvement step
9     for some iterations do
10       sample a mini-batch  $B$  of  $\mathcal{D}$ 
11       for each transition in  $B$ , compute the label
12          $y = r + \gamma Q^*(x', \pi(x'; \varphi^{\text{old}}); \theta^{\text{old}})$ 
13         // critic update
14          $\theta \leftarrow \theta - \eta \nabla_\theta \frac{1}{B} \sum_{(x, a, r, x', y) \in B} (y - Q^*(x, a; \theta))^2$ 
15         // actor update
16          $\varphi \leftarrow \varphi + \eta \nabla_\varphi \frac{1}{B} \sum_{(x, a, r, x', y) \in B} Q^*(x, \pi(x; \varphi); \theta)$ 
17          $\theta^{\text{old}} \leftarrow (1 - \rho)\theta^{\text{old}} + \rho\theta$ 
18          $\varphi^{\text{old}} \leftarrow (1 - \rho)\varphi^{\text{old}} + \rho\varphi$ 

```

Twin delayed DDPG (TD3) is an extension of DDPG that uses two separate critic networks for predicting the maximum action and evaluating the policy (Fujimoto et al., 2018). This addresses the maximization bias akin to Double-DQN. TD3 also applies delayed updates to the actor network, which increases stability.

12.5.1 Randomized Policies

We have seen that randomized policies naturally encourage exploration. With deterministic actor-critic methods like DDPG, we had to inject Gaussian noise to enforce sufficient exploration. A natural question is therefore whether we can also handle randomized policies in this framework of off-policy actor-critics.

The key idea is to replace the squared loss of the critic,

$$\ell_{\text{DQN}}(\theta; \mathcal{D}) \approx \frac{1}{2} \left(r + \gamma Q^*(x', \pi(x'; \varphi); \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2,$$

which only considers the fixed action $\pi(x'; \varphi^{\text{old}})$ with an expected squared loss,

$$\ell_{\text{DQN}}(\theta; \mathcal{D}) \approx \mathbb{E}_{a' \sim \pi(x'; \varphi)} \left[\frac{1}{2} \left(r + \gamma Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2 \right], \quad (12.71)$$

which considers a distribution over actions.

It turns out that we can still compute gradients of this expectation.

$$\begin{aligned} & \nabla_{\theta} \mathbb{E}_{a' \sim \pi(x'; \varphi)} \left[\frac{1}{2} \left(r + \gamma Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2 \right] \\ &= \mathbb{E}_{a' \sim \pi(x'; \varphi)} \left[\nabla_{\theta} \frac{1}{2} \left(r + \gamma Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2 \right]. \end{aligned}$$

refer to the squared loss of Q-learning
(12.11)

see Appendix A.1.5

Similarly to our definition of the Bellman error (12.12), we define by

$$\delta_B(a') \doteq r + \gamma Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta), \quad (12.72)$$

the *Bellman error* for a fixed action a' . Using the chain rule, we obtain

$$= \mathbb{E}_{a' \sim \pi(x'; \varphi)} [\delta_B(a') \nabla_{\theta} Q^*(x, a; \theta)]. \quad (12.73)$$

Note that this is identical to the gradient in DQN (12.14), except that now we have an expectation over actions. As we have done many times already, we can use automatic differentiation to obtain gradient estimates of $\nabla_{\theta} Q^*(x, a; \theta)$. This provides us with a method of obtaining unbiased gradient estimates for the critic.

We also need to reconsider the actor update. When using a randomized policy, the objective function changes to

$$\hat{J}_{\mu}(\varphi; \theta) \doteq \mathbb{E}_{x \sim \mu} \mathbb{E}_{a \sim \pi(x; \varphi)} [Q^*(x, a; \theta)].$$

of which we can obtain gradients via

$$\nabla_{\varphi} \hat{J}_{\mu}(\varphi; \theta) = \mathbb{E}_{x \sim \mu} \nabla_{\varphi} \mathbb{E}_{a \sim \pi(x; \varphi)} [Q^*(x, a; \theta)]. \quad (12.74) \quad \text{see Appendix A.1.5}$$

Note that the inner expectation is with respect to a measure that depends on the parameters φ , which we are trying to optimize. We therefore cannot move the gradient operator inside the expectation. This is a problem that we have already encountered several times. In the previous section on policy gradients, we used the score gradient estimator.²⁴ Earlier, in Chapter 5 on variational inference we have already seen reparameterization gradients.²⁵ Here, if our policy is reparameterizable, we can use the *reparameterization trick* from Theorem 5.19!

Example 12.14: Reparameterization gradients for Gaussians

Suppose we use a Gaussian parameterization of policies,

$$\pi(x; \varphi) \doteq \mathcal{N}(\mu(x; \varphi), \Sigma(x; \varphi)).$$

Then, using conditional linear Gaussians, our action a is given by

$$a = g(\varepsilon; x, \varphi) \doteq \Sigma^{1/2}(x; \varphi)\varepsilon + \mu(x; \varphi), \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, I) \quad (12.75)$$

where $\Sigma^{1/2}(x; \varphi)$ is the square root of $\Sigma(x; \varphi)$. This coincides with our earlier application of the reparameterization trick to Gaussians in Example 5.20.

As we have seen, not only Gaussians are reparameterizable. In general, we call a distribution (in this context, a policy) reparameterizable iff $a \sim \pi(x; \varphi)$ is such that $a = g(\varepsilon; x, \varphi)$, where $\varepsilon \sim \phi$ is an independent random variable.

Then, we have,

$$\begin{aligned} & \nabla_{\varphi} \mathbb{E}_{a \sim \pi(x; \varphi)} [Q^*(x, a; \theta)] \\ &= \mathbb{E}_{\varepsilon \sim \phi} [\nabla_{\varphi} Q^*(x, g(\varepsilon; x, \varphi); \theta)] \end{aligned} \quad (12.76)$$

$$= \mathbb{E}_{\varepsilon \sim \phi} [\nabla_a Q^*(x, a; \theta)|_{a=g(\varepsilon; x, \varphi)} \cdot D_{\varphi}g(\varepsilon; x, \varphi)]. \quad (12.77)$$

using the reparameterization trick (5.63)

using the chain rule analogously to Equation (12.70)

In this way, we can obtain unbiased gradient estimates for reparameterizable policies. This general technique does not only apply to continuous action spaces. For discrete action spaces, there is the analogous so-called *Gumbel-max trick*, which we will not discuss in greater detail here.

The algorithm that uses Equation (12.73) to obtain gradients for the critic and reparameterization gradients for the actor is called *stochastic value gradients* (SVG) (Heess et al., 2015).

²⁴ see Equation (12.32)

²⁵ see (5.63)

12.6 Maximum Entropy Reinforcement Learning

In practice, algorithms like SVG often do not explore enough. A key issue with relying on randomized policies for exploration is that they might collapse to deterministic policies. That is, the algorithm might quickly reach a local optimum, where all mass is placed on a single action.

A simple trick that encourages a little bit of extra exploration is to regularize the randomized policies “away” from putting all mass on a single action. In other words, we want to encourage the policies to exhibit some uncertainty. A natural measure of uncertainty is entropy, which we have already seen several times.²⁶ This approach is known as *entropy regularization* or *maximum entropy reinforcement learning* (MERL). Canonically, entropy regularization is applied to finite-horizon rewards (cf. Remark 10.5), yielding the optimization problem of maximizing

$$J_\lambda(\boldsymbol{\varphi}) \doteq J(\boldsymbol{\varphi}) + \lambda H[\Pi_{\boldsymbol{\varphi}}] \quad (12.78)$$

$$= \sum_{t=1}^T \mathbb{E}_{(x_t, a_t) \sim \Pi_{\boldsymbol{\varphi}}} [r(x_t, a_t) + \lambda H[\pi_{\boldsymbol{\varphi}}(\cdot | x_t)]], \quad (12.79)$$

²⁶ see Section 5.4

where we have a preference for entropy in the actor distribution to encourage exploration which is regulated by the temperature parameter λ . As $\lambda \rightarrow 0$, we recover the “standard” reinforcement learning objective (here for finite-horizon rewards):

$$J(\boldsymbol{\varphi}) = \sum_{t=1}^T \mathbb{E}_{(x_t, a_t) \sim \Pi_{\boldsymbol{\varphi}}} [r(x_t, a_t)]. \quad (12.80)$$

Here, for notational convenience, we begin the sum with $t = 1$ rather than $t = 0$.

12.6.1 Entropy Regularization as Probabilistic Inference

The entropy-regularized objective from Equation (12.79) leads us to a remarkable interpretation of reinforcement learning and, more generally, decision-making under uncertainty as solving an inference problem akin to variational inference. The framing of “control as inference” will lead us to contemporary algorithms for reinforcement learning as well as paint a path for decision-making under uncertainty beyond stationary MDPs.

Let us denote by Π_* the distribution over trajectories τ under the optimal policy π^* . By framing the problem of optimal control as an inference problem in a hidden Markov model with hidden “optimality variables” $O_t \in \{0, 1\}$ indicating whether the played action a_t was

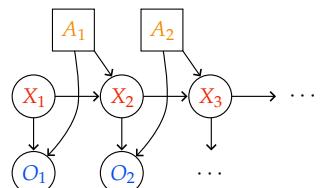


Figure 12.3: Directed graphical model of the underlying hidden Markov model with hidden states X_t , optimality variables O_t , and actions A_t .

optimal we can derive Π_* analytically. That is to say, when $O_t = 1$ and $O_{t+1:T} \equiv 1$ the policy from time t onwards was optimal. To simplify the notation, we will denote the event $O_t = 1$ by \mathcal{O}_t .

We consider the HMM defined by the Gibbs distribution

$$p(\mathcal{O}_t | \mathbf{x}_t, \mathbf{a}_t) \propto \exp\left(\frac{1}{\lambda} r(\mathbf{x}_t, \mathbf{a}_t)\right), \quad \text{with } \lambda > 0 \quad (12.81)$$

which is a natural choice as we have seen in Problem 6.7 that the Gibbs distribution maximizes entropy subject to $\mathbb{E}[O_t \cdot r(\mathbf{x}_t, \mathbf{a}_t) | \mathbf{x}_t, \mathbf{a}_t] < \infty$.

The distribution over trajectories conditioned on optimality of actions (i.e., conditioned on $\mathcal{O}_{1:T}$) is given by

$$\Pi_*(\tau) \doteq p(\tau | \mathcal{O}_{1:T}) = p(\mathbf{x}_1) \prod_{t=1}^T p(\mathbf{a}_t | \mathbf{x}_t, \mathcal{O}_t) p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{a}_t). \quad (12.82)$$

It remains to determine $p(\mathbf{a}_t | \mathbf{x}_t, \mathcal{O}_t)$ which corresponds to the optimal policy $\pi^*(\mathbf{a}_t | \mathbf{x}_t)$. It is generally useful to think of the situation where the prior policy $p(\mathbf{a}_t | \mathbf{x}_t)$ is uniform on \mathcal{A} ,²⁷ in which case by Bayes' rule (1.45), $p(\mathbf{a}_t | \mathbf{x}_t, \mathcal{O}_t) \propto p(\mathcal{O}_t | \mathbf{x}_t, \mathbf{a}_t)$, so

$$\Pi_*(\tau) \propto \left[p(\mathbf{x}_1) \prod_{t=1}^T p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{a}_t) \right] \exp\left(\frac{1}{\lambda} \sum_{t=1}^T r(\mathbf{x}_t, \mathbf{a}_t)\right). \quad (12.83)$$

Recall that our fundamental goal is to approximate Π_* with a distribution over trajectories Π_φ under the parameterized policy π_φ . It is therefore a natural idea to minimize $\text{KL}(\Pi_\varphi \| \Pi_*)$:²⁸

$$\begin{aligned} & \arg \min_{\varphi} \text{KL}(\Pi_\varphi \| \Pi_*) \\ &= \arg \min_{\varphi} H[\Pi_\varphi \| \Pi_*] - H[\Pi_\varphi] \\ &= \arg \max_{\varphi} \mathbb{E}_{\tau \sim \Pi_\varphi} [\log \Pi_*(\tau) - \log \Pi_\varphi(\tau)] \\ &= \arg \max_{\varphi} \mathbb{E}_{\tau \sim \Pi_\varphi} \left[\sum_{t=1}^T r(\mathbf{x}_t, \mathbf{a}_t) - \lambda \log \pi_\varphi(\mathbf{a}_t | \mathbf{x}_t) \right] \\ &= \arg \max_{\varphi} \sum_{t=1}^T \mathbb{E}_{(\mathbf{x}_t, \mathbf{a}_t) \sim \Pi_\varphi} [r(\mathbf{x}_t, \mathbf{a}_t) + \lambda H[\pi_\varphi(\cdot | \mathbf{x}_t)]]. \end{aligned} \quad (12.84)$$

That is, entropy regularization is equivalent to minimizing the KL-divergence from Π_* to Π_φ . This highlights a very natural tradeoff between exploration and exploitation, wherein $H[\Pi_\varphi \| \Pi_*]$ encourages exploitation and $H[\Pi_\varphi]$ encourages exploration.

It can be shown that a “softmax” version of the Bellman optimality equation (10.34) can be obtained for Equation (12.84) (?):

$$q^*(\mathbf{x}, \mathbf{a}) = \frac{1}{\lambda} r(\mathbf{x}, \mathbf{a}) + \mathbb{E}_{\mathbf{x}' \sim \mathbf{x}, \mathbf{a}} \left[\log \int_{\mathcal{A}} \exp(q^*(\mathbf{x}', \mathbf{a}')) d\mathbf{a}' \right] \quad (12.85)$$

Using Equation (12.30). We assume here that the dynamics and initial state distribution are “fixed”, that is, we assume $p(\mathbf{x}_1 | \mathcal{O}_{1:T}) = p(\mathbf{x}_1)$ and $p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{a}_t, \mathcal{O}_{1:T}) = p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{a}_t)$.

²⁷ This is not a restriction as any informative prior can be pushed into Equation (12.81).

²⁸ Observe that we cannot easily minimize forward-KL as we cannot sample from Π_* . In the context of RL, it can be argued that the mode-seeking behavior of reverse-KL is preferable over the moment-matching behavior of forward-KL (Levine, 2018).

using the definition of KL-divergence (5.34)

using the definition of cross-entropy (5.32) and entropy (5.27)

using Equations (12.30) and (12.83) and simplifying

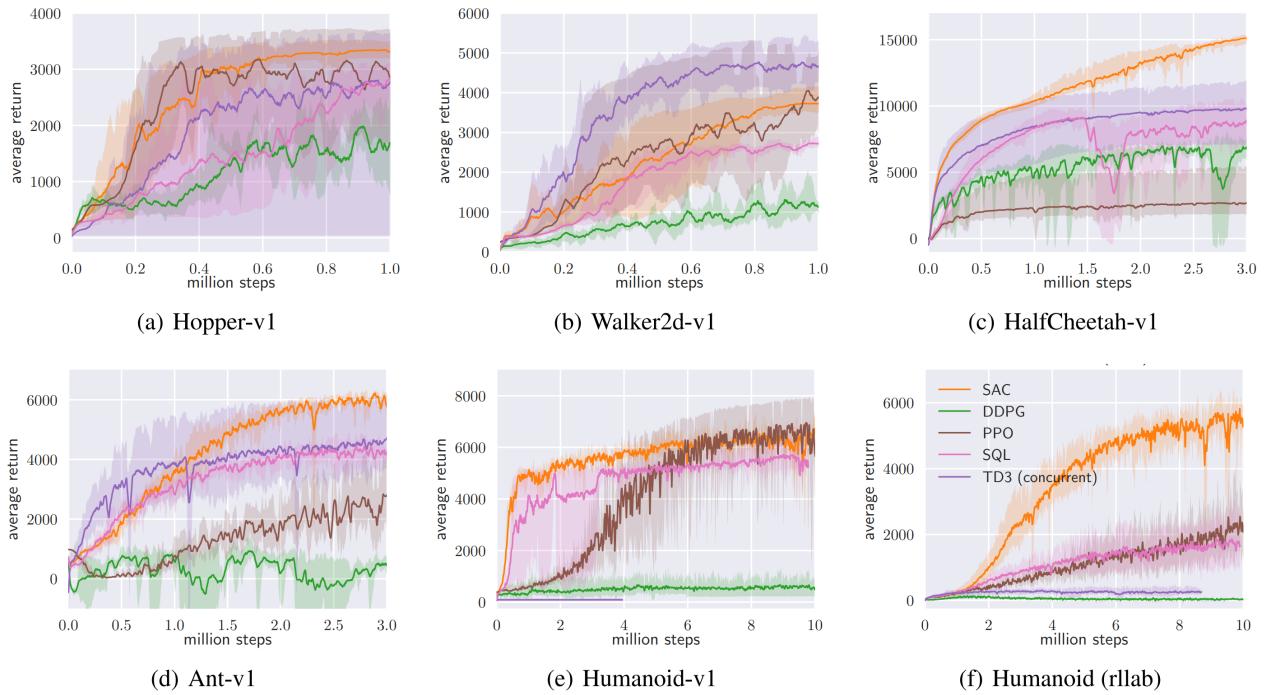
using the definition of entropy (5.27) and linearity of expectation (1.20)

Problem 12.7

with the convention that $q^*(x_T, a) = 0$ for all a .²⁹ Here, q^* is called a *soft value function*. As we will see in Problem 12.7, the optimal policy has the form $\pi^*(a | x) \propto \exp(q^*(x, a))$, that is, it simply corresponds to performing softmax exploration (11.9) with the soft value function. The second term of Equation (12.85) quantifies downstream rewards. In comparison to the “standard” Bellman optimality equation (10.34), the soft value function is less greedy which tends to encourage robustness.

Analogously to Q-learning, the soft value function q^* can be approximated via a bootstrapped “critic” Q^* which is called *soft Q-learning* (Levine, 2018). Note that computing the optimal policy requires computing an integral over the action space, which is typically intractable for continuous action spaces. As discussed in Sections 12.3 to 12.5 and analogously to actor-critic methods such as DDPG and SVG, we can learn a parameterized policy (i.e., an “actor”) π_ϕ to approximate the optimal policy π^* . The resulting algorithm, *soft actor critic* (SAC) (Haarnoja et al., 2018a,b), is widely used. Due to its off-policy nature, it is also relatively sample efficient.

²⁹ Note that Equation (10.34) was derived in the infinite horizon setting with discounted rewards, whereas here we study the finite horizon setting. Equation (12.85) is the natural extension of the standard Bellman optimality equation in the finite horizon setting where the downstream rewards are measured by a softmax rather than the greedy policy.



We can also express this optimization in terms of an evidence lower

Figure 12.4: Comparison of training curves of a selection of on-policy and off-policy policy gradient methods. Reproduced with permission from “Soft actor-critic algorithms and applications” (Haarnoja et al., 2018b).

bound. The evidence lower bound for the observations $\mathcal{O}_{1:T}$ is

$$L(\Pi_\varphi, \Pi_\star; \mathcal{O}_{1:T}) = \mathbb{E}_{\tau \sim \Pi_\varphi} [\log p(\mathcal{O}_{1:T} | \tau) + \log \Pi(\tau) - \log \Pi_\varphi(\tau)] \quad (12.86)$$

using the definition of the ELBO (5.55b)

where Π denotes the distribution over trajectories (12.30) under the prior policy $p(a_t | x_t)$. This is commonly written as the variational free energy

$$-L(\Pi_\varphi, \Pi_\star; \mathcal{O}_{1:T}) = \mathbb{E}_{\tau \sim \Pi_\varphi} [S[p(\mathcal{O}_{1:T} | \tau)]] + KL(\Pi_\varphi \| \Pi) \quad (12.87)$$

$$= \underbrace{S[p(\mathcal{O}_{1:T})]}_{\text{"extrinsic" value}} + \underbrace{KL(\Pi_\varphi \| \Pi_\star)}_{\text{"epistemic" value}}. \quad (12.88)$$

using that $\Pi_\star(\tau) = p(\tau | \mathcal{O}_{1:T})$

which we already encountered in Section 5.5.2 in the context of variational inference. Here, the “extrinsic” value is independent of the variational distribution Π_φ and can be thought of as a fixed “problem cost”, whereas the “epistemic” value can be interpreted as the approximation error or “solution cost”. To summarize, we have seen that

$$\arg \max_{\varphi} J_\lambda(\varphi) = \arg \min_{\varphi} KL(\Pi_\varphi \| \Pi_\star) = \arg \max_{\varphi} L(\Pi_\varphi, \Pi_\star; \mathcal{O}_{1:T}).$$

Recall that the free energy $-L(\Pi_\varphi, \Pi_\star; \mathcal{O}_{1:T})$ is a variational upper bound to the surprise about observations $S[\Pi_\star(\mathcal{O}_{1:T})]$ when following an optimal policy.³⁰ For example, undesirable states incur a low reward while desirable states yield a high reward, and thus, if we expect optimality of actions (i.e., $\mathcal{O}_{1:T}$) paths leading to such states have high and low surprise, respectively.³¹

An agent which acts to minimize free energy with $\mathcal{O}_{1:T}$ can be thought of as hallucinating to perform optimally, and acting to minimize the surprise about having played suboptimal actions. Think, for example, about the robotics task of moving an arm to a new position. Intuitively, minimizing free energy solves this task by “hallucinating that the arm is at the goal position”, and then minimizing the surprise with respect to this perturbed world model. In this way, MERL can be understood as identifying paths of least surprise akin to the *free energy principle*.

³⁰ see Section 5.5.2

³¹ This is because, a path leading to a low reward state will include sub-optimal actions.

Remark 12.15: Towards active inference

Maximum entropy reinforcement learning makes one fundamental assumption, namely, that the “biased” distribution about observations specifying the underlying HMM is

$$p(\mathcal{O}_t | x_t, a_t) \propto \exp\left(\frac{1}{\lambda} r(x_t, a_t)\right).$$

This assumption is key to the framing of optimal control (in an

unknown MDP) with a certain reward function as an inference problem. One could conceive other HMMs. For example, Fellows et al. (2019) propose an HMM defined in terms of the current value function of the agent:

$$p(\mathcal{O}_t | \mathbf{x}_t, \mathbf{a}_t) \propto \exp\left(\frac{1}{\lambda} Q(\mathbf{x}_t, \mathbf{a}_t; \boldsymbol{\theta})\right).$$

Crucially, the choice of $p(\mathcal{O}_t | \mathbf{x}_t, \mathbf{a}_t)$ is the only place where the reward enters the inference problem, and one can conceive of settings where a “stationary” (i.e., time-independent) reward is not assumed to exist. The general approach to decision-making as probabilistic inference presented in this section (but for possibly reward-independent and non-stationary HMMs) is known as *active inference* (Friston et al., 2015; Millidge et al., 2020, 2021; Parr et al., 2022).

12.7 Learning from Preferences

So far, we have been assuming that the agent is presented with a reward signal after every played action. This is a natural assumption in domains such as games and robotics — even though it often requires substantial “reward engineering” to break down complex tasks with sparse rewards to more manageable tasks (cf. Section 13.3). In many other domains such as an agent learning to drive a car or a chatbot, it is unclear how one can even quantify the reward associated with an action or a sequence of actions. For example, in the context of autonomous driving it is typically desired that agents behave “human-like” even though a different driving behavior may also reach the destination safely.

The task of “aligning” the behavior of an agent to human expectations is difficult in complex domains such as the physical world and language, yet crucial for their practical use. To this end, one can conceive of alternative ways for presenting “feedback” to the agent:

- The classical feedback in reinforcement learning is a numerical score. Consider, for example, a recommender system for movies. The feedback is obtained after a movie was recommended to a user by a user-rating on a given scale (often 1 to 10). This rating is informative as it corresponds to an *absolute value assessment*, allowing to place the recommendation in a complete ranking of all previous recommendations. However, numerical feedback of this type can be error-prone as it is scale-dependent (different users may ascribe different value to a recommendation rated a 7).

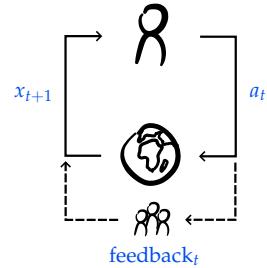


Figure 12.5: We generalize the perspective of reinforcement learning from Figure 11.1 by allowing the feedback to come from either the environment or an evaluation by other agents (e.g., humans), and by allowing the feedback to come in other forms than a numerical reward.

- An alternative feedback mechanism is comparison-based. The user is presented with k alternative actions and selects their preferred action (or alternatively returns a ranking of actions). This feedback is typically easy to obtain as humans are fast in making “this-or-that” decisions. However, in contrast to numerical rewards, the feedback provides information only on the user’s *relative preferences*. That is, such preference feedback encodes fewer bits of information than score feedback, and it therefore often takes longer to learn complex behavior from this feedback.

Remark 12.16: Context-dependent feedback

We neglect here that feedback is often context-dependent (Lindner and El-Assady, 2022; Casper et al., 2023). For example, if someone is asked whether they prefer “ice cream” over “pizza” the answer may depend on whether they are hungry and the weather.

12.7.1 Language Models as Agents

In the following, we discuss approaches to learning from preference feedback in the context of autoregressive³² large language models and chatbots. A chatbot is an agent (often based on a transformer architecture, Vaswani et al. (2017)) parameterized by φ that given a *prompt* x returns a (stochastic) *response* y . The autoregressive generation of the response can be understood in terms of a policy $\pi_\varphi(y_{t+1} | x, y_{1:t})$ which generates the next token given the prompt and all previous tokens.³³ We denote the policy over complete responses (i.e., the chatbot) by

$$\Pi_\varphi(y | x) = \prod_{t=0}^{T-1} \pi_\varphi(y_{t+1} | x, y_{1:t}). \quad (12.89)$$

In RL jargon, the agents action corresponds to the choice of next token y_{t+1} and the deterministic dynamics add this token to the current (incomplete) response $y_{1:t}$. A full trajectory y consists sequentially of all tokens $y_{1:T}$ comprising a response, and the prompt x can be interpreted as a *context* to this trajectory. Observe that Equation (12.89) is derived from the general representation of the distribution over trajectories (12.30), noting that prior and dynamics are deterministic.

The standard pipeline for applying pre-trained³⁴ large language models such as GPT (OpenAI, 2023) to downstream tasks consists of two main steps which are illustrated in Figure 12.6 (Stiennon et al., 2020): (1) supervised fine-tuning and (2) post-training using preference feedback.

The first step is to fine-tune the language model with supervised learning on high-quality data for the downstream task of interest. For ex-

³² An autoregressive model predicts/-generates the next “token” as a function of previous tokens.

³³ In large language models, a “token” is usually taken to be a letter, word, or something in-between. A special token is used to terminate the response.

³⁴ The pre-trained language model is usually obtained by self-supervised training on a large corpus of text. *Self-supervised learning* generates labeled training data from an unlabeled data source by selectively “masking-out” parts of the data. When training an autoregressive language model, labeled training data can be obtained by repeatedly “masking-out” the next word in a sentence, and training the language model to predict this word. Such large models that can be fine-tuned and “post-trained” to various downstream tasks are also called *foundation models*.

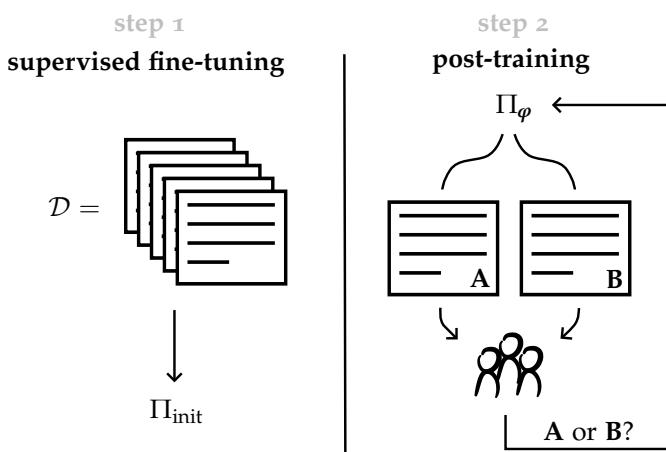


Figure 12.6: Illustration of the learning pipeline of a large language model.

ample, when the goal is to build a chatbot, this data may consist of desirable responses to some exemplary prompts. We will denote the parameters of the fine-tuned language model by φ^{init} , and its associated policy by Π_{init} .

The second step is then to “post-train” the language model Π_{init} from the first step using human feedback. Here, it is important that Π_{init} is already capable of producing sensible output (i.e., with correct spelling and grammar). Learning this from scratch using only preference feedback would take far too long. Instead, the post-training step is used to align the agent to the task and user preferences.

In each iteration of post-training, the model is prompted with prompts x to produce pairs of answers $(y_A, y_B) \sim \Pi_\varphi(\cdot | x)$. These answers are presented to human labelers who express their preference for one of the answers, denoted $y_A \succ y_B | x$. A popular choice for modeling preferences is the *Bradley-Terry model* which stipulates that the human preference distribution is given by

$$p(y_A \succ y_B | x, r) = \frac{\exp(r(y_A | x))}{\exp(r(y_A | x)) + \exp(r(y_B | x))} \quad (12.90)$$

for some unknown latent reward model $r(y | x)$ (Bradley and Terry, 1952). This can be written in terms of the logistic function σ (5.9):

$$= \sigma(r(y_A | x) - r(y_B | x)). \quad (12.91)$$

as seen in Problem 7.1 this is the Gibb's distribution with energy $-r(y | x)$ in a binary classification problem

Remark 12.17: Outcome rewards

Note that the Bradley-Terry model attributes reward only to “complete” responses. We call such a reward an *outcome reward*.³⁵ While the following discussion is on outcome rewards (which is most common in the context of language models), everything translates to individual per-step rewards.

³⁵ Framing this in terms of an individual “per-step” rewards which we have seen so far, this corresponds to a (sparse) reward which is zero until the final action.

The aggregated human feedback $\mathcal{D} = \{\mathbf{y}_A^{(i)} \succ \mathbf{y}_B^{(i)} \mid \mathbf{x}^{(i)}\}_{i=1}^n$ across n different prompts is then used to update the language model π_φ . In the next two sections, we discuss two standard approaches to post-training: reinforcement learning from human feedback (RLHF) and direct preference optimization (DPO).

12.7.2 Reinforcement Learning from Human Feedback

RLHF separates the post-training step into two stages (Stiennon et al., 2020). First, the human feedback is used to learn an approximate reward model r_θ . This reward model is then used in the second stage to determine a refined policy Π_φ .

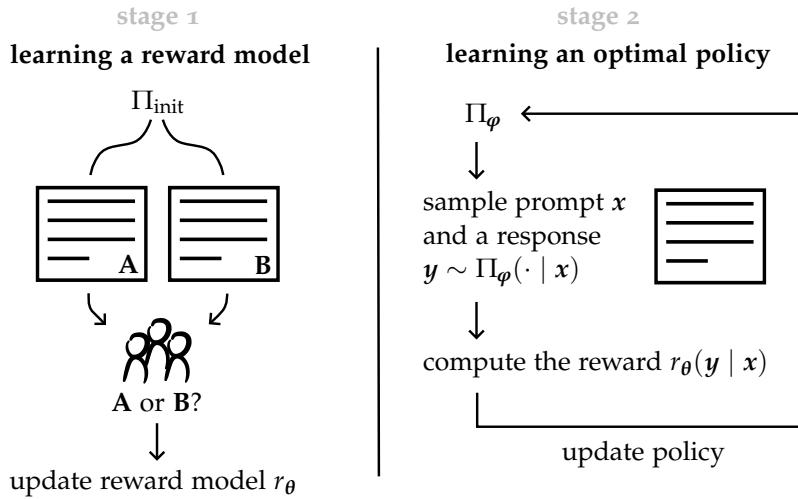


Figure 12.7: Illustration of the post-training process of RLHF.

Learning a reward model: During the first stage, the initial policy obtained by supervised fine-tuning Π_{init} is used to generate proposals $(\mathbf{y}_A, \mathbf{y}_B)$ for some exemplary prompts \mathbf{x} , which are then ranked according to the preference of human labelers. This preference data can be used to learn a reward model r_θ by maximum likelihood esti-

mation (or equivalently minimizing cross-entropy):

$$\begin{aligned} & \arg \max_{\theta} p(\mathcal{D} \mid r_{\theta}) \\ &= \arg \max_{\theta} \mathbb{E}_{(y_A \succ y_B \mid x) \sim \mathcal{D}} [\log \sigma(r_{\theta}(y_A \mid x) - r_{\theta}(y_B \mid x))] \end{aligned} \quad (12.92)$$

using Equation (12.91) and SGD

This is analogous to the standard maximum likelihood estimation of the reward model in model-based RL with score feedback which we discussed in Section 11.2.1. The reward model r_{θ} is often initialized from the initial policy π_{init} by placing a linear layer producing a scalar output on top of the final transformer layer.

Learning an optimal policy: One can now employ the methods from this and previous chapters to determine the optimal policy for the approximate reward r_{θ} . Due to the use of an *approximate* reward, however, simply maximizing r_{θ} surfaces the so-called “reward gaming” problem which is illustrated in Figure 12.8. As the responses generated by the learned policy π_{ϕ} deviate from the distribution of \mathcal{D} (i.e., the distribution induced by the initial policy π_{init}), the approximate reward model becomes inaccurate. The approximate reward may severely overestimate the true reward in regions far away from the training data. A common approach to address this problem is to regularize the policy search towards policies whose distribution does not stray away “too far” from the training distribution.

Analogously to the trust-region methods we discussed in Section 12.4.4, the deviation from the initial policy is typically controlled by maximizing the regularized objective

$$J_{\lambda}(\phi; \phi^{\text{init}} \mid x) \doteq \mathbb{E}_{y \sim \Pi_{\phi}(\cdot \mid x)} [r(y \mid x)] - \lambda \text{KL}(\Pi_{\phi}(\cdot \mid x) \parallel \Pi_{\text{init}}(\cdot \mid x))$$

in expectation over prompts x sampled uniformly at random from the dataset \mathcal{D} . Note that this coincides with the PPO objective from Equation (12.63) with an outcome reward. We can expand the regularization term to obtain

$$\begin{aligned} &= \underbrace{\mathbb{E}_{y \sim \Pi_{\phi}(\cdot \mid x)} [r(y \mid x)] + \lambda H[\Pi_{\phi}(\cdot \mid x)]}_{\text{entropy-regularized RL}} \quad (12.93) \\ &\quad - \lambda H[\Pi_{\phi}(\cdot \mid x) \parallel \Pi_{\text{init}}(\cdot \mid x)] \end{aligned}$$

which indicates an intimate relationship to entropy-regularized RL.³⁶

The optimal policy maximizing $J_{\lambda}(\phi; \phi^{\text{init}} \mid x)$ is $\textcircled{?}$

$$\Pi_{\star}(y \mid x) \propto \Pi_{\text{init}}(y \mid x) \exp\left(\frac{1}{\lambda} r(y \mid x)\right) \quad (12.94)$$

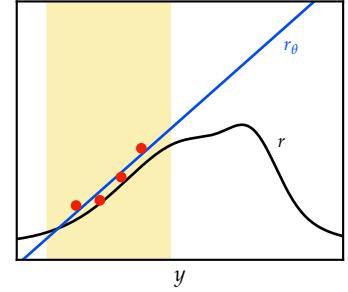


Figure 12.8: Illustration of “reward gaming”. Shown in black is the true reward $r(y \mid x)$ for a fixed prompt x . Shown in blue is the approximation based on the feedback \mathcal{D} to the responses shown in red. The yellow region symbolizes responses y where the approximate reward can still be “trusted”.

using the definition of KL-divergence (5.34)

³⁶ compare to Equation (12.79)

Problem 12.8

and can be interpreted as a probabilistic update to the prior Π_{init} where $\exp(\frac{1}{\lambda}r(\mathbf{y} \mid \mathbf{x}))$ corresponds to the “likelihood of optimality”. As $\lambda \rightarrow \infty$, $\Pi_{\star} \rightarrow \Pi_{\text{init}}$, and as $\lambda \rightarrow 0$, the optimal policy reduces to deterministically picking the response with the highest reward.

In practice, sampling from Π_{\star} explicitly is intractable, but note that any of the approximate inference methods discussed in Part I are applicable here. In the context of chatbots, it is important that the sampling from the resulting policy is efficient (i.e., the chatbot should respond quickly to prompts). Most commonly, the optimization problem of maximizing $j_{\varphi}(\varphi^{\text{init}} \mid \mathbf{x}[\lambda])$ (with the estimated reward r_{θ}) is solved approximately within a parameterized family of policies. This is typically done using policy gradient methods such as PPO (Stiennon et al., 2020) or GRPO (Guo et al., 2025).

Remark 12.18: Other (non-preference) reward models

Note that the post-training pipeline we described here is agnostic to the choice of reward model. That is, instead of post-training our language model to comply with human preferences, we could have also post-trained it to maximize any other reward signal. More recently, works have explored the use of reward models for challenging “reasoning problems” such as mathematical calculations, where the reward is usually based on the correctness of the answer.³⁷ One prominent example for this are “reasoning” models such as the DeepSeek-R1 model (Guo et al., 2025), which was trained with GRPO.

12.7.3 Direct Preference Optimization

Observe that the reward model can be expressed in terms of its associated optimal policy:

$$r(\mathbf{y} \mid \mathbf{x}) = \lambda \log \frac{\Pi_{\star}(\mathbf{y} \mid \mathbf{x})}{\Pi_{\text{init}}(\mathbf{y} \mid \mathbf{x})} + \text{const.} \quad (12.95)$$

In particular, it follows that given a fixed prior Π_{init} , *any* policy Π_{φ} has a family of associated reward models with respect to which it is optimal! We denote by

$$r_{[\varphi]}(\mathbf{y} \mid \mathbf{x}) \doteq \lambda \log \frac{\Pi_{\varphi}(\mathbf{y} \mid \mathbf{x})}{\Pi_{\text{init}}(\mathbf{y} \mid \mathbf{x})} \quad (12.96)$$

the “simplest” of these reward models. Remembering the characterization of the optimal policy from Equation (12.94) it follows immediately that Π_{φ} is optimal with respect to $r_{[\varphi]}$.

Instead of first learning an approximate reward model and then finding the associated optimal policy, DPO exploits the relationship of

³⁷ On training data where the answer to problem \mathbf{x} is known to be $y^*(\mathbf{x})$, the reward is simply $r(\mathbf{y} \mid \mathbf{x}) = \mathbb{1}\{y^*(\mathbf{x}) \in \mathbf{y}\}$. That is, the reward is 1 if the response \mathbf{y} contains the correct answer and 0 otherwise.

by reorganizing the terms of Equation (12.94)

Equation (12.95) to learn the optimal policy directly (Rafailov et al., 2023). Substituting r_θ in the maximum likelihood estimation from Equation (12.92), yields the objective

$$\begin{aligned} & \mathbb{E}_{(y_A \succ y_B | x) \sim \mathcal{D}} \left[\log \sigma \left(r_{[\varphi]}(y_A | x) - r_{[\varphi]}(y_B | x) \right) \right] \\ &= \mathbb{E}_{(y_A \succ y_B | x) \sim \mathcal{D}} \left[\log \sigma \left(\lambda \log \frac{\Pi_\varphi(y_A | x)}{\Pi_{\text{init}}(y_A | x)} - \lambda \log \frac{\Pi_\varphi(y_B | x)}{\Pi_{\text{init}}(y_B | x)} \right) \right]. \end{aligned} \tag{12.97}$$

using Equation (12.95)

Gradients can be computed via automatic differentiation:

$$\begin{aligned} & \lambda \mathbb{E}_{(y_A \succ y_B | x) \sim \mathcal{D}} \left[\underbrace{\sigma(r_{[\varphi]}(y_B | x) - r_{[\varphi]}(y_A | x))}_{\text{weight according to error of reward estimate}} \right. \\ & \quad \left. \left[\underbrace{\nabla_\varphi \log \Pi_\varphi(y_A | x)}_{\text{increase likelihood of } y_A} - \underbrace{\nabla_\varphi \log \Pi_\varphi(y_B | x)}_{\text{decrease likelihood of } y_B} \right] \right]. \end{aligned} \tag{12.98}$$

Intuitively, DPO successively increases the likelihood of preferred responses y_A and decreases the likelihood of dispreferred responses y_B . Examples $(y_A \succ y_B | x)$ are weighted by the strength of regularization λ and by the degree to which the implicit reward model incorrectly orders the responses.

Discussion

In this chapter, we studied central ideas in actor-critic methods. We have seen two main approaches to use policy-gradient methods. We began, in Section 12.3, by introducing the REINFORCE algorithm which uses policy gradients and Monte Carlo estimation, but suffered from large variance in the gradient estimates of the policy value function. In Section 12.4, we have then seen a number of actor-critic methods such as A2C and GAE behaving similarly to policy iteration that exhibit less variance, but are very sample inefficient due to their on-policy nature. TRPO improves the sample efficiency slightly, but not fundamentally.

In Section 12.5, we discussed a second family of policy gradient techniques that generalize Q-learning and are akin to value iteration. For reparameterizable policies, this led us to algorithms such as DDPG, TD3, SVG. Importantly, these algorithms are significantly more sample efficient than on-policy policy gradient methods, which often results in much faster learning of a near-optimal policy. In Section 12.6, we discussed entropy regularization which frames reinforcement learning as probabilistic inference, and we derived the SAC algorithm which is widely used and works quite well in practice.

Finally, in Section 12.7, we studied two canonical approaches to learning from preference feedback: RLHF which separately learns reward model and policy and DPO which learns a policy directly. We have seen that RLHF is akin to model-based RL as it explicitly learns a reward model through maximum likelihood estimation. In contrast, DPO is more closely related to model-free RL and policy gradient methods as it learns the optimal policy directly.

Optional Readings

- **A3C:** Mnih, Badia, Mirza, Graves, Lillicrap, Harley, Silver, and Kavukcuoglu (2016).
Asynchronous methods for deep reinforcement learning.
- **GAE:** Schulman, Moritz, Levine, Jordan, and Abbeel (2016).
High-dimensional continuous control using generalized advantage estimation.
- **TRPO:** Schulman, Levine, Abbeel, Jordan, and Moritz (2015).
Trust region policy optimization.
- **PPO:** Schulman, Wolski, Dhariwal, Radford, and Klimov (2017).
Proximal policy optimization algorithms.
- **DDPG:** Lillicrap, Hunt, Pritzel, Heess, Erez, Tassa, Silver, and Wierstra (2016).
Continuous control with deep reinforcement learning.
- **TD3:** Fujimoto, Hoof, and Meger (2018).
Addressing function approximation error in actor-critic methods.
- **SVG:** Heess, Wayne, Silver, Lillicrap, Erez, and Tassa (2015).
Learning continuous control policies by stochastic value gradients.
- **SAC:** Haarnoja, Zhou, Abbeel, and Levine (2018a).
Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor.
- **DPO:** Rafailov, Sharma, Mitchell, Ermon, Manning, and Finn (2023).
Direct preference optimization: Your language model is secretly a reward model.

Problems

12.1. Q-learning and function approximation.

Consider the MDP of Figure 12.9 and set $\gamma = 1$.

1. Using Bellman's theorem, prove that $v^*(x) = -|x - 4|$ is the optimal value function.
2. Suppose we observe the following episode:

x	a	x'	r
3	-1	2	-1
2	1	3	-1
3	1	4	-1
4	1	4	0

We initialize all Q-values to 0. Compute the updated Q-values using Q-learning with learning rate $\alpha = 1/2$.

3. We will now approximate the Q-function with a linear function. We let

$$Q(x, a; \mathbf{w}) \doteq xw_0 + aw_1 + w_2$$

where $\mathbf{w} = [w_0 \ w_1 \ w_2]^\top \in \mathbb{R}^3$.

Suppose we have $\mathbf{w}^{\text{old}} = [1 \ -1 \ -2]^\top$ and $\mathbf{w} = [-1 \ 1 \ 1]^\top$, and we observe the transition $\tau = (2, -1, -1, 1)$. Use the learning rate $\alpha = 1/2$ to compute $\nabla_{\mathbf{w}} \ell(\mathbf{w}; \tau)$ and the updated weights $\mathbf{w}' = \mathbf{w} - \alpha \nabla_{\mathbf{w}} \ell(\mathbf{w}; \tau)$.

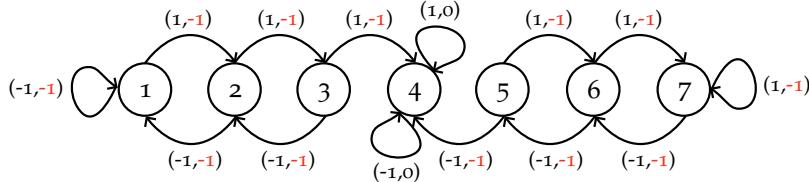


Figure 12.9: MDP studied in Problem 12.1. Each arrow marks a (deterministic) transition and is labeled with (action, reward).

12.2. Eligibility vector.

The vector $\nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(a_t | x_t)$ is commonly called *eligibility vector*. In the following, we assume that the action space \mathcal{A} is finite and denote it by A .

If we parameterize $\pi_{\boldsymbol{\varphi}}$ as a softmax distribution

$$\pi_{\boldsymbol{\varphi}}(a | x) \doteq \frac{\exp(h(x, a, \boldsymbol{\varphi}))}{\sum_{b \in A} \exp(h(x, b, \boldsymbol{\varphi}))} \quad (12.99)$$

with linear preferences $h(x, a, \boldsymbol{\varphi}) \doteq \boldsymbol{\varphi}^\top \boldsymbol{\varphi}(x, a)$ where $\boldsymbol{\varphi}(x, a)$ is some feature vector, what is the form of the eligibility vector?

12.3. Variance of score gradients with baselines.

In this exercise, we will see a sufficient condition for baselines to reduce the variance of score gradient estimators.

- Suppose for a random vector \mathbf{X} , we want to estimate $\mathbb{E}[f(\mathbf{X})]$ for some function f . Assume that you are given a function g and also its expectation $\mathbb{E}[g(\mathbf{X})]$. Instead of estimating $\mathbb{E}[f(\mathbf{X})]$ directly, we will instead estimate $\mathbb{E}[f(\mathbf{X}) - g(\mathbf{X})]$ as we know from linearity of expectation (1.20) that

$$\mathbb{E}[f(\mathbf{X})] = \mathbb{E}[f(\mathbf{X}) - g(\mathbf{X})] + \mathbb{E}[g(\mathbf{X})].$$

Prove that if $\frac{1}{2}\text{Var}[g(\mathbf{X})] \leq \text{Cov}[f(\mathbf{X}), g(\mathbf{X})]$, then

$$\text{Var}[f(\mathbf{X}) - g(\mathbf{X})] \leq \text{Var}[f(\mathbf{X})]. \quad (12.100)$$

- Consider estimating $\nabla_{\boldsymbol{\varphi}} J(\boldsymbol{\varphi})$. Prove that if $b^2 \leq 2b \cdot r(\mathbf{x}, \mathbf{a})$ for every state $\mathbf{x} \in \mathcal{X}$ and action $\mathbf{a} \in \mathcal{A}$, then

$$\text{Var}_{\tau \sim \Pi_{\boldsymbol{\varphi}}}[(G_0 - b)\nabla_{\boldsymbol{\varphi}} \log \Pi_{\boldsymbol{\varphi}}(\tau)] \leq \text{Var}_{\tau \sim \Pi_{\boldsymbol{\varphi}}}[G_0 \nabla_{\boldsymbol{\varphi}} \log \Pi_{\boldsymbol{\varphi}}(\tau)]. \quad (12.101)$$

12.4. Score gradients with state-dependent baselines.

For a sequence of *state-dependent baselines* $\{b(\tau_{0:t-1})\}_{t=1}^T$ where

$$\tau_{0:t-1} \doteq (\tau_0, \tau_1, \dots, \tau_{t-1}),$$

show that

$$\begin{aligned} & \mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}}[G_0 \nabla_{\boldsymbol{\varphi}} \log \Pi_{\boldsymbol{\varphi}}(\tau)] \\ &= \mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}}\left[\sum_{t=0}^{T-1}(G_0 - b(\tau_{0:t-1}))\nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t \mid \mathbf{x}_t)\right] \end{aligned} \quad (12.102)$$

where we write $b(\tau_{0:-1}) = 0$.

12.5. Policy gradients with downstream returns.

Suppose we are training an agent to solve a computer game. There are only two possible actions, specifically:

- do nothing*; and
- move*.

Each episode lasts for four ($T = 3$) time steps. The policy π_{θ} is completely determined by the parameter $\theta \in [0, 1]$. Here, for simplicity, we have assumed that the policy is independent of the current state. The probability of moving (action 2) is equal to θ and the probability of doing nothing (action 1) is $1 - \theta$.

Initially, $\theta = 0.5$. One episode is played with this initial policy and the results are

$$\text{actions} = (1, 0, 1, 0), \quad \text{rewards} = (1, 0, 1, 1).$$

Compute the policy gradient estimate with downstream returns, discount factor $\gamma = 1/2$, and the provided *single* sample $\tau \sim \Pi_\theta$.

12.6. Policy gradient with an exponential family.

- Suppose, we can choose between two actions $a \in \{0, 1\}$ in each state. A natural stochastic policy is induced by a Bernoulli distribution,

$$a \sim \text{Bern}(\sigma(f_\varphi(x))), \quad (12.103)$$

where σ is the logistic function from Equation (5.9). First, write down the expression for $\pi_\varphi(a | x)$. Then, derive the expression for $\nabla_\varphi J(\varphi)$ in terms of q^{π_φ} , $\sigma(f_\varphi(x))$, and $\nabla_\varphi f_\varphi(x)$ using the policy gradient theorem.

- The Bernoulli distribution is part of a family of distributions that allows for a much simpler derivation of the gradient than was necessary in (1). A univariate *exponential family* is a family of distributions whose PDF (or PMF) can be expressed in canonical form as

$$\pi_\varphi(a | x) = h(a) \exp(a f_\varphi(x) - A(f_\varphi(x))) \quad (12.104)$$

where h , f , and A are known functions.³⁸ Derive the expression of the policy gradient $\nabla_\varphi J(\varphi)$ for such a distribution.

- Can you relate the results of the previous two exercises (1) and (2)? What are h and A in case of the Bernoulli distribution?
- The Gaussian distribution with unit variance $\mathcal{N}(f_\varphi(x), 1)$ is of the same canonical form with

$$A(f_\varphi(x)) = \frac{f_\varphi(x)^2}{2}. \quad (12.105)$$

Determine the policy gradient $\nabla_\varphi J(\varphi)$.

- For a Gaussian policy, can we instead apply the reparameterization trick (5.65) that we have seen in the context of variational inference? If yes, how? If not, why?

12.7. Soft value function.

In this exercise, we derive the optimal policy solving Equation (12.84) and the soft value function from Equation (12.85).

- We let $\beta(a_t | x_t) \doteq \exp(\frac{1}{\lambda} r(x_t, a_t))$, $Z(x) \doteq \int_{\mathcal{A}} \beta(a | x) da$, and denote by $\hat{\pi}(\cdot | x)$ the policy $\beta(\cdot | x)/Z(x)$. Show that

$$\begin{aligned} & \text{KL}(\Pi_\varphi \| \Pi_\star) \\ &= \sum_{t=1}^T \mathbb{E}_{x_t \sim \Pi_\varphi} [\text{KL}(\pi_\varphi(\cdot | x_t) \| \hat{\pi}(\cdot | x_t)) - \log Z(x_t)]. \end{aligned} \quad (12.106)$$

- Show that if the space of policies parameterized by φ is sufficiently expressive, $\pi^*(a | x) \propto \exp(q^*(x, a))$ solves Equation (12.84).

³⁸ Observe that this form is equivalent to the form introduced in Equation (5.48) where $f_\varphi(x)$ is the natural parameter, and we let $A(f) = \log Z(f)$.

12.8. PPO as probabilistic inference.

1. Consider the same generative model as was introduced in Section 12.6.1 when we interpreted entropy-regularized RL as probabilistic inference. Only now, assume that $T = 1$ ³⁹ and suppose that the prior over actions is $p(\mathbf{y} | \mathbf{x}) = \Pi_{\text{init}}(\mathbf{y} | \mathbf{x})$ rather than uniform. Define the distribution over optimal trajectories $\Pi_{\star}(\mathbf{y} | \mathbf{x})$ as before in Equation (12.82). Show that for any context \mathbf{x} ,

$$\arg \min_{\boldsymbol{\varphi}} \text{KL}(\Pi_{\boldsymbol{\varphi}}(\cdot | \mathbf{x}) \| \Pi_{\star}(\cdot | \mathbf{x})) = \arg \max_{\boldsymbol{\varphi}} J_{\lambda}(\boldsymbol{\varphi}; \boldsymbol{\varphi}^{\text{init}} | \mathbf{x}). \quad (12.107)$$

2. Conclude that the policy maximizing $J_{\lambda}(\boldsymbol{\varphi}; \boldsymbol{\varphi}^{\text{init}} | \mathbf{x})$ is

$$\Pi_{\star}(\mathbf{y} | \mathbf{x}) \propto \Pi_{\text{init}}(\mathbf{y} | \mathbf{x}) \exp\left(\frac{1}{\lambda} r(\mathbf{y} | \mathbf{x})\right). \quad (12.108)$$

³⁹ since, in this section, we have been considering outcome rewards

13

Model-based Reinforcement Learning

In this final chapter, we will revisit the model-based approach to reinforcement learning. We will see some advantages it offers over model-free methods. In particular, we will use the machinery of probabilistic inference, which we developed in the first chapters, to quantify uncertainty about our model and use this uncertainty for planning, exploration, and reliable decision-making.

To recap, in Chapter 11, we began by discussing model-based reinforcement learning which attempts to learn the underlying Markov decision process and then use it for planning. We have seen that in the tabular setting, computing and storing the entire model is computationally expensive. This led us to consider the family of model-free approaches, which learn the value function directly, and as such can be considered more economical in the amount of data that they store.

In Chapter 12, we saw that using function approximation, we were able to scale model-free methods to very large state and action spaces. We will now explore similar ideas in the model-based framework. Namely, we will use function approximation to condense the representation of our model of the environment. More concretely, we will learn an approximate dynamics model $f \approx p$ and approximate rewards r , which is also called a *world model*.

There are a few ways in which the model-based approach is advantageous. First, if we have an accurate model of the environment, we can use it for planning — ideally also for interacting safely with the environment. However, in practice, we will rarely have such a model. In fact, the accuracy of our model will depend on the past experience of our agent and the region of the state-action space. Understanding the uncertainty in our model of the environment is crucial for planning. In particular, quantifying uncertainty is necessary to drive safe(!) exploration and avoid undesired states.

Moreover, modeling uncertainty in our model of the environment can be extremely useful in deciding where to explore. Learning a model can therefore help to dramatically reduce the sample complexity over model-free techniques. Often times this is crucial when developing agents for real-world use as in such settings, exploration is costly and potentially dangerous.

Algorithm 13.1 describes the general approach to model-based reinforcement learning.

Algorithm 13.1: Model-based reinforcement learning (outline)

- 1 start with an initial policy π and no (or some) initial data \mathcal{D}
 - 2 **for** several episodes **do**
 - 3 roll out policy π to collect data
 - 4 learn a model of the dynamics f and rewards r from data
 - 5 plan a new policy π based on the estimated models
-

We face three main challenges in model-based reinforcement learning. First, given a fixed model, we need to perform planning to decide on which actions to play. Second, we need to learn models f and r accurately and efficiently. Third, we need to effectively trade exploration and exploitation. We will discuss these three topics in the following.

13.1 Planning

There exists a large literature on planning in various settings. These settings can be mainly characterized as

- discrete or continuous action spaces;
- fully- or partially observable state spaces;
- constrained or unconstrained; and
- linear or nonlinear dynamics.

In Chapter 10, we have already seen algorithms such as policy iteration and value iteration, which can be used to solve planning exactly in tabular settings. In the following, we will now focus on the setting of continuous state and action spaces, fully observable state spaces, no constraints, and nonlinear dynamics.

13.1.1 Deterministic Dynamics

To begin with, let us assume that our dynamics model is deterministic and known. That is, given a state-action pair, we know the subsequent state,

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{a}_t). \quad (13.1)$$

We continue to focus on the setting of infinite-horizon discounted returns (10.5), which we have been considering throughout our discussion of reinforcement learning. This yields the objective,

$$\max_{\mathbf{a}_{0:\infty}} \sum_{t=0}^{\infty} \gamma^t r(\mathbf{x}_t, \mathbf{a}_t) \quad \text{such that } \mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{a}_t). \quad (13.2)$$

Now, because we are optimizing over an infinite time horizon, we cannot solve this optimization problem directly. This problem is studied ubiquitously in the area of *optimal control*. We will discuss one central idea from optimal control that is widely used in model-based reinforcement learning, and will later return to using this idea for learning parametric policies in Section 13.1.3.

Planning over finite horizons: The key idea of a classical algorithm from optimal control called *receding horizon control* (RHC) or *model predictive control* (MPC) is to iteratively plan over finite horizons. That is, in each round, we plan over a finite time horizon H and carry out the first action.

Algorithm 13.2: Model predictive control, MPC

```

1 for  $t = 0$  to  $\infty$  do
2   observe  $\mathbf{x}_t$ 
3   plan over a finite horizon  $H$ ,
4     
$$\max_{\mathbf{a}_{t:t+H-1}} \sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} r(\mathbf{x}_\tau, \mathbf{a}_\tau) \quad \text{such that } \mathbf{x}_{\tau+1} = f(\mathbf{x}_\tau, \mathbf{a}_\tau) \quad (13.3)$$

5   carry out action  $\mathbf{a}_t$ 

```

Observe that the state \mathbf{x}_τ can be interpreted as a deterministic function $\mathbf{x}_\tau(\mathbf{a}_{t:\tau-1})$, which depends on all actions from time t to time $\tau - 1$ and the state \mathbf{x}_t . To solve the optimization problem of a single iteration, we therefore need to maximize,

$$J_H(\mathbf{a}_{t:t+H-1}) \doteq \sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} r(\mathbf{x}_\tau(\mathbf{a}_{t:\tau-1}), \mathbf{a}_\tau). \quad (13.4)$$

This optimization problem is in general non-convex. If the actions are continuous and the dynamics and reward models are differentiable, we can nevertheless obtain analytic gradients of \hat{J}_H . This can be done using the chain rule and “backpropagating” through time, analogously to backpropagation in neural networks.¹ Especially for large horizons H , this optimization problem becomes difficult to solve exactly due to local optima and vanishing/exploding gradients.

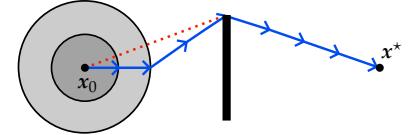


Figure 13.1: Illustration of model predictive control in a deterministic transition model. The agent starts in position x_0 and wants to reach x^* despite the black obstacle. We use the reward function $r(\mathbf{x}) = -\|\mathbf{x} - \mathbf{x}^*\|$. The gray concentric circles represent the length of a single step. We plan with a time horizon of $H = 2$. Initially, the agent does not “see” the black obstacle, and therefore moves straight towards the goal. As soon as the agent sees the obstacle, the optimal trajectory is “replanned”. The dotted red line corresponds to the optimal trajectory, the agent’s steps are shown in blue.

¹ see Section 7.1.4

Tree search: Often, heuristic global optimization methods (also called “search methods”) are used to optimize \hat{J}_H . An example are *random shooting methods*, which find the optimal choice among a set of random proposals. Of course, obtaining a “good” set of randomly proposed action sequences is crucial. A naive way of generating the proposals is to pick them uniformly at random. This strategy, however, usually does not perform very well as it corresponds to suggesting random walks of the state space. Alternatives are to sample from a Gaussian distribution or using the *cross-entropy method* which gradually adapts the sampling distribution by reweighing samples according to the rewards they produce.

Algorithm 13.3: Random shooting methods

- 1 generate m sets of random samples, $\mathbf{a}_{t:t+H-1}^{(i)}$
 - 2 pick the sequence of actions $\mathbf{a}_{t:t+H-1}^{(i^*)}$ where
- $$i^* = \arg \max_{i \in [m]} J_H(\mathbf{a}_{t:t+H-1}^{(i)}) \quad (13.5)$$
-

The evolution of the state can be visualized as a tree where — if dynamics are deterministic — the branching is fully determined by the played actions. For this reason, classical tree search algorithms can be employed, such as *alpha-beta pruning* or *Monte Carlo tree search* (MCTS), which was used for example by “AlphaZero” (Silver et al., 2016, 2017) and can be viewed as an advanced variant of a shooting method.

Finite-horizon planning with sparse rewards: A common problem of finite-horizon methods is that in the setting of sparse rewards, there is often no signal that can be followed. You can think of an agent that operates in a kitchen and tries to find a box of candy. Yet, to get to this box, it needs to perform a large number of actions. In particular, if this number of actions is larger than the horizon H , then the local optimization problem of MPC will not take into account the reward for choosing this long sequence of actions. Thus, the box of candy will likely never be found by our agent.

A solution to this problem is to amend a long-term value estimate to the finite-horizon sum. The idea is to not only consider the rewards attained *while* following the actions $\mathbf{a}_{t:t_H-1}$, but to also consider the value of the final state \mathbf{x}_{t+H} , which estimates the discounted sum of

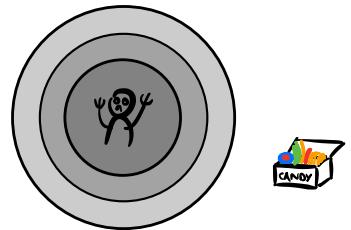


Figure 13.2: Illustration of finite-horizon planning with sparse rewards. When the finite time horizon does not suffice to “reach” a reward, the agent has no signal to follow.

future rewards.

$$\hat{J}_H(\mathbf{a}_{t:t+H-1}) \doteq \underbrace{\sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} r(\mathbf{x}_\tau(\mathbf{a}_{t:\tau-1}), \mathbf{a}_\tau)}_{\text{short-term}} + \underbrace{\gamma^H V(\mathbf{x}_{t+H})}_{\text{long-term}}. \quad (13.6)$$

Intuitively, $\gamma^H V(\mathbf{x}_{t+H})$ is estimating the tail of the infinite sum.

Remark 13.4: Planning generalizes model-free methods!

Observe that for $H = 1$, when we use the value function estimate associated with this MPC controller, maximizing J_H coincides with using the greedy policy π_V . That is,

$$\mathbf{a}_t = \arg \max_{\mathbf{a} \in \mathcal{A}} \hat{J}_1(\mathbf{a}) = \pi_V(\mathbf{x}_t). \quad (13.7)$$

Thus, by looking ahead for a single time step, we recover the approaches from the model-free setting in this model-based setting! Essentially, if we do not plan long-term and only consider the value estimate, the model-based setting reduces to the model-free setting. However, in the model-based setting, we are now able to use our model of the transition dynamics to anticipate the downstream effects of picking a particular action \mathbf{a}_t . This is one of the fundamental reasons for why model-based approaches are typically severely more sample efficient than model-free methods.

To obtain the value estimates, we can use the approaches we discussed in detail in Section 11.4, such as TD-learning for on-policy value estimates and Q-learning for off-policy value estimates. For large state-action spaces, we can use their approximate variants, which we discussed in Section 12.1 and Section 12.2. To improve value estimates, we can obtain “artificial” data by rolling out policies within our model. This is a key advantage over model-free methods as, once a sufficiently accurate model has been learned, data for value estimation can be generated efficiently in simulation.

13.1.2 Stochastic Dynamics

How can we extend this approach to planning to a stochastic transition model? A natural extension of model predictive control is to do what is called *stochastic average approximation* (SAA) or *trajectory sampling* (Chua et al., 2018). Like in MPC, we still optimize over a deterministic sequence of actions, but now we will average over all resulting trajectories.

Algorithm 13.5: Trajectory sampling

```

1 for  $t = 0$  to  $\infty$  do
2   observe  $x_t$ 
3   optimize expected performance over a finite horizon  $H$ ,

$$\max_{\mathbf{a}_{t:t+H-1}} \mathbb{E}_{\mathbf{x}_{t+1:t+H}} \left[ \sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} r_\tau + \gamma^H V(\mathbf{x}_{t+H}) \mid \mathbf{a}_{t:t+H-1}, f \right] \quad (13.8)$$

4   carry out action  $a_t$ 

```

Intuitively, trajectory sampling optimizes over a much simpler object — namely a deterministic sequence of actions of length H — than finding a policy, which corresponds to finding an optimal decision tree mapping states to actions. Of course, using trajectory sampling (from an arbitrary starting state) implies such a policy. However, trajectory sampling never computes this policy explicitly, and rather, in each step, only plans over a finite horizon.

Computing the expectation exactly is typically not possible as this involves solving a high-dimensional integral of nonlinear functions. Instead, a common approach is to use Monte Carlo estimates of this expectation. This approach is known as *Monte Carlo trajectory sampling*. The key issue with using sampling based estimation is that the sampled trajectory (i.e., sampled sequence of states) we obtain, depends on the actions we pick. In other words, the measure we average over depends on the decision variables — the actions. This is a problem that we have seen several times already! It naturally suggests using the reparameterization trick.²

Previously, we have used the reparameterization trick to reparameterize variational distributions (see Section 5.5.1) and to reparameterize policies (see Section 12.5.1). It turns out that we can use the exact same approach for reparameterizing the transition model. We say that a (stochastic) transition model f is *reparameterizable* iff $x_{t+1} \sim f(x_t, a_t)$ is such that $x_{t+1} = g(\varepsilon; x_t, a_t)$, where $\varepsilon \sim \phi$ is a random variable that is independent of x_t and a_t . We have already seen in Example 12.14 (in the context of stochastic policies) how a Gaussian transition model can be reparameterized.

In this case, x_τ is determined recursively by $a_{t:\tau-1}$ and $\varepsilon_{t:\tau-1}$,

$$\begin{aligned} x_\tau &= x_\tau(\varepsilon_{t:\tau-1}; a_{t:\tau-1}) \\ &\doteq g(\varepsilon_{\tau-1}; g(\dots; (\varepsilon_{t+1}; g(\varepsilon_t; x_t, a_t), a_{t+1}), \dots), a_{\tau-1}). \end{aligned} \quad (13.9)$$



Figure 13.3: Illustration of trajectory sampling. High-reward states are shown in brighter colors. The agent iteratively plans a finite number of time steps into the future and picks the best initial action.

² see Theorem 5.19 and Equation (12.76)

This allows us to obtain unbiased estimates of J_H using Monte Carlo approximation,

$$\begin{aligned}\widehat{J}_H(\mathbf{a}_{t:t+H-1}) \approx \frac{1}{m} \sum_{i=1}^m & \left(\sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} r(\mathbf{x}_\tau(\boldsymbol{\varepsilon}_{t:\tau-1}^{(i)}; \mathbf{a}_{t:\tau-1}), \mathbf{a}_\tau) \right. \\ & \left. + \gamma^H V(\mathbf{x}_{t+H}(\boldsymbol{\varepsilon}_{t:t+H-1}^{(i)}; \mathbf{a}_{t:t+H-1})) \right) \quad (13.10)\end{aligned}$$

where $\boldsymbol{\varepsilon}_{t:t+H-1}^{(i)} \stackrel{\text{iid}}{\sim} \phi$ are independent samples. To optimize this approximation we can again compute analytic gradients or use shooting methods as we have discussed in Section 13.1.1 for deterministic dynamics.

13.1.3 Parametric Policies

When using algorithms such as model predictive control for planning, planning needs to be done online before each time we take an action. This is called *closed-loop control* and can be expensive. Especially when the time horizon is large, or we encounter similar states many times (leading to “repeated optimization problems”), it can be beneficial to “store” the planning decision in a (deterministic) policy,

$$\mathbf{a}_t = \pi(\mathbf{x}_t; \boldsymbol{\varphi}) \doteq \pi_{\boldsymbol{\varphi}}(\mathbf{x}_t). \quad (13.11)$$

This policy can then be trained offline and evaluated cheaply online, which is known as *open-loop control*.

This is akin to a problem that we have discussed in detail in the previous chapter when extending Q-learning to large action spaces. There, this led us to discuss policy gradient and actor-critic methods. Recall that in Q-learning, we seek to follow the greedy policy,

$$\pi^*(\mathbf{x}) = \arg \max_{\mathbf{a} \in \mathcal{A}} Q^*(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta}), \quad \text{see Equation (12.20)}$$

and therefore had to solve an optimization problem over all actions. We accelerated this by learning an approximate policy that “mimicked” this optimization,

$$\boldsymbol{\varphi}^* = \arg \max_{\boldsymbol{\varphi}} \underbrace{\mathbb{E}_{\mathbf{x} \sim \mu} [Q^*(\mathbf{x}, \pi_{\boldsymbol{\varphi}}(\mathbf{x}); \boldsymbol{\theta})]}_{\widehat{J}_\mu(\boldsymbol{\varphi}; \boldsymbol{\theta})} \quad \text{see Equation (12.67)}$$

where $\mu(\mathbf{x}) > 0$ was some exploration distribution that has full support and thus leads to the exploration of all states. The key idea was that if we use a differentiable approximation Q and a differentiable parameterization of policies, which is “rich enough”, then both optimizations are equivalent, and we can use the chain rule to obtain

gradient estimates of the second expression. We then used this to derive the *deep deterministic policy gradients* (DDPG) and *stochastic value gradients* (SVG) algorithms. It turns out that there is a very natural analogue to DDPG/SVG for model-based reinforcement learning.

Instead of maximizing the Q-function directly, we use finite-horizon planning to estimate the immediate value of the policy within the next H time steps and simply use the Q-function to approximate the terminal value (i.e., the tails of the infinite sum). Then, our objective becomes,

$$\hat{J}_{\mu,H}(\boldsymbol{\varphi}; \boldsymbol{\theta}) \doteq \mathbb{E}_{x_0 \sim \mu, x_{1:H} | \pi_{\boldsymbol{\varphi}}, f} \left[\sum_{\tau=0}^{H-1} \gamma^\tau r_\tau + \gamma^H Q^*(x_H, \pi_{\boldsymbol{\varphi}}(x_H); \boldsymbol{\theta}) \right] \quad (13.12)$$

This approach naturally extends to randomized policies using reparameterization gradients, which we have discussed in Section 12.5.1. Analogously to Remark 13.4, for $H = 0$, this coincides with the DDPG objective! For larger time horizons, the look-ahead takes into account the transition model for planning next time steps. This tends to help dramatically in improving policies *much more rapidly* between episodes. Instead of just gradually improving policies a little by slightly adapting the policy to the learned value function estimates (as in model-free RL), we use the model to anticipate the consequences of actions multiple time steps ahead. This is at the heart of model-based reinforcement learning.

Essentially, we are using methods such as Q-learning and DDPG/SVG as subroutines within the framework of model predictive control to do much bigger steps in policy improvement than to slightly improve the next picked action. To encourage exploration, it is common to extend the objective in Equation (13.12) by an additional entropy regularization term as seen in Section 12.6.

Remark 13.6: Planning as inference

We have been using *Monte Carlo rollouts* to estimate the expectation of Equation (13.12). That is, we have been using a Monte Carlo approximation (e.g., a sample mean) using samples obtained by “rolling out” the induced Markov chain of a fixed policy.

It would certainly be preferable to compute the expectation exactly, however, this is generally not possible as this involves solving a high dimensional integral. Recall that we faced the same problem when studying inference in the first half of the manuscript. In both problems, we need to approximate high-dimensional integrals (i.e., expectations). This suggests a deep connection between

the problems of planning and inference. It is therefore not surprising that many techniques for approximate inference that we have seen earlier can also be applied to planning.

- *Monte Carlo approximation* which we have been focusing on during our discussion of planning is a very simple inference algorithm — approximating an expectation by sampling from the distribution that is averaged over. This allowed us to obtain unbiased gradient estimates (which may have high variance).
- An alternative approach is *moment matching* (cf. Section 5.4.6). Instead of approximating the expectation, here, we approximate the distribution over trajectories using a tractable distribution (e.g., a Gaussian) and “matching” their moments. This then allows us to analytically compute gradients of the expectation in Equation (13.12). A prominent example of this approach is *probabilistic inference for learning control* (PILCO).
- In Section 12.6, we have used *variational inference* for planning, and seen that it coincides with entropy regularization as implemented by the *soft actor critic* (SAC) algorithm.

13.2 Learning

Thus far, we have considered known environments. That is, we assumed that the transition model f and the rewards r are known. In reinforcement learning, f and r are (of course!) not known. Instead, we have to estimate them from data. This will also be crucial in our later discussion of exploration in Section 13.3 where we explore methods of driving data collection to learn what we need to learn about the world.

First, let us revisit one of our key observations when we first introduced the reinforcement learning problem. Namely, that the observed transitions x' and rewards r are conditionally independent given the state-action pairs (x, a) .³ This is due to the Markovian structure of the underlying Markov decision process.

³ see Equation (11.3)

This is the key observation that allows us to treat the estimation of the dynamics and rewards as a simple regression problem (or a density estimation problem when the quantities are stochastic rather than deterministic). More concretely, we can estimate the dynamics and rewards off-policy using the standard supervised learning techniques we discussed in earlier chapters, from a replay buffer

$$\mathcal{D} = \{(\underbrace{x_t, a_t}_{\text{“input”}}, \underbrace{r_t, x_{t+1}}_{\text{“label”}})\}_t. \quad (13.13)$$

Here, x_t and a_t are the “inputs”, and r_t and x_{t+1} are the “labels” of the

regression problem. Due to the conditional independence of the labels given the inputs, we have independent label noise (i.e., “independent training data”) which is the basic assumption that we have been making throughout our discussion of techniques for probabilistic machine learning in Part I.

The key difference to supervised learning is that the set of inputs depends on how we act. That is, the current inputs arise from previous policies, and the inputs which we will observe in the future will depend on the model (and policy) obtained from the current data: we have feedback loops! We will come back to this aspect of reinforcement learning in the next section on exploration. For now, recall we only assume that we have used an arbitrary policy to collect some data, which we then stored in a replay buffer, and which we now want to use to learn the “best-possible” model of our environment.

13.2.1 Probabilistic Inference

In the following, we will discuss how we can use the techniques from probabilistic inference, which we have seen in Part I, to learn the dynamics and reward models. Thereby, we will focus on learning the transition model f as learning the reward model r is completely analogous. For learning deterministic dynamics or rewards, we can use for example Gaussian processes (cf. Chapter 4) or deep neural networks (cf. Chapter 7). We will now focus on the setting where the dynamics are stochastic, that is,

$$\mathbf{x}_{t+1} \sim f(\mathbf{x}_t, \mathbf{a}_t; \boldsymbol{\psi}). \quad (13.14)$$

Example 13.7: Conditional Gaussian dynamics

We could, for example, use a conditional Gaussian for the transition model,

$$\mathbf{x}_{t+1} \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_t, \mathbf{a}_t; \boldsymbol{\psi}), \boldsymbol{\Sigma}(\mathbf{x}_t, \mathbf{a}_t; \boldsymbol{\psi})). \quad (13.15)$$

As we have seen in Equation (A.7), we can rewrite the covariance matrix $\boldsymbol{\Sigma}$ as a product of a lower-triangular matrix \mathcal{L} and its transpose using the Cholesky decomposition $\boldsymbol{\Sigma} = \mathcal{L}\mathcal{L}^\top$ of $\boldsymbol{\Sigma}$. This allows us to represent the model by only $n(n + 1)/2$ parameters. Moreover, we have learned that Gaussians are reparameterizable, which we have seen to be useful for planning.

Note that this model reduces to a deterministic model if the covariance is zero. So the stochastic transition models encompass all deterministic models. Moreover, in many applications (such as robotics), it is often useful to use stochastic models to attribute

slight inaccuracies and measurement noise to a small uncertainty in the model.

A first approach might be to obtain a point estimate for f , either through maximum likelihood estimation (which we have seen to overfit easily) or through maximum a posteriori estimation. If, for example, f is represented as a deep neural network, we have already seen how to find the MAP estimate of its weights in Section 7.2.1.

Remark 13.8: The key pitfall of point estimates

However, using point estimates leads to a *key pitfall* of model-based reinforcement learning. Using a point estimate of the model for planning — even if this point estimate is very accurate — often performs *very poorly*. The reason is that planning is very good at overfitting (i.e., exploiting) small errors in the transition model. Moreover, the errors in the model estimate compound over time when using a longer time horizon H . The key to remedy this pitfall lies in being robust to misestimated models. This naturally suggests quantifying the uncertainty in our model estimate and taking it into account during planning.⁴ In the following section, we will rediscover that estimates of epistemic uncertainty are also extremely useful for driving (safe) exploration — something that we have already encountered in our discussion of Bayesian optimization.

In the following, we will differentiate between the epistemic uncertainty and the aleatoric uncertainty. Recall from Section 2.2 that epistemic uncertainty corresponds to our uncertainty about the model, $p(f | \mathcal{D})$, while aleatoric uncertainty corresponds to the uncertainty of the transitions in the underlying Markov decision process (which can be thought of as “irreducible” noise), $p(x_{t+1} | f, x_t, a_t)$.

Intuitively, probabilistic inference of dynamics models corresponds to learning a distribution over possible models f and r given prior *beliefs*, where f and r characterize the underlying Markov decision process. This goes to show another benefit of the model-based over the model-free approach to reinforcement learning. Namely, that it is much easier to encode prior knowledge about the transition and rewards model.

Example 13.9: Inference with conditional Gaussian dynamics

Let us revisit inference with our conditional Gaussian dynamics

⁴ Quantifying the uncertainty of an estimate is a problem that we have spent the first few chapters exploring. Notably, refer to

- Section 2.2 for a description of epistemic and aleatoric uncertainty;
- Chapter 4 for our use of uncertainty estimates in the context of Gaussian processes;
- Chapter 7 for our use of uncertainty estimates in the context of Bayesian deep learning; and
- Chapters 8 and 9 for our use of epistemic uncertainty estimates to drive exploration.

model from Equation (13.15),

$$\mathbf{x}_{t+1} \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_t, \mathbf{a}_t; \boldsymbol{\psi}), \boldsymbol{\Sigma}(\mathbf{x}_t, \mathbf{a}_t; \boldsymbol{\psi})).$$

Recall that in the setting of Bayesian deep learning, most approximate inference techniques represented the approximate posterior using some form of a mixture of Gaussians.⁵

$$p(\mathbf{x}_{t+1} | \mathcal{D}, \mathbf{x}_t, \mathbf{a}_t) \approx \frac{1}{m} \sum_{i=1}^m \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_t, \mathbf{a}_t; \boldsymbol{\psi}^{(i)}), \boldsymbol{\Sigma}(\mathbf{x}_t, \mathbf{a}_t; \boldsymbol{\psi}^{(i)})). \quad (13.16)$$

Hereby, the epistemic uncertainty is represented by the variance between mixture components, and the aleatoric uncertainty by the average variance within the components.⁶

In supervised learning, we often conflated the notions of epistemic and aleatoric uncertainty. In the context of planning, there is an important consequence of the decomposition into epistemic and aleatoric uncertainty. Recall that the epistemic uncertainty corresponds to a distribution over Markov decision processes f , whereas the aleatoric uncertainty corresponds to the randomness in the transitions within one such MDP f . Crucially, this randomness in the transitions must be consistent within a single MDP! That is, once we selected a single MDP for planning, we should disregard the epistemic uncertainty and solely focus on the randomness of the transitions. Then, to take into account epistemic uncertainty, we should average our plan across the different realizations of f . This yields the following Monte Carlo estimate of our reward \hat{J}_H ,

$$\hat{J}_H(\mathbf{a}_{t:t+H-1}) \approx \frac{1}{m} \sum_{i=1}^m \hat{J}_H(\mathbf{a}_{t:t+H-1}; \mathbf{f}^{(i)}) \quad \text{where} \quad (13.17)$$

$$\hat{J}_H(\mathbf{a}_{t:t+H-1}; f) \doteq \sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} r(\mathbf{x}_\tau(\boldsymbol{\varepsilon}_{t:\tau-1}^{(i)}; \mathbf{a}_{t:\tau-1}, f), \mathbf{a}_\tau) + \gamma^H V(\mathbf{x}_{t+H}). \quad (13.18)$$

Here, $\mathbf{f}^{(i)} \stackrel{\text{iid}}{\sim} p(f | \mathcal{D})$ are independent samples of the transition model, and $\boldsymbol{\varepsilon}_{t:t+H-1}^{(i)} \stackrel{\text{iid}}{\sim} \boldsymbol{\phi}$ parameterizes the dynamics analogously to Equation (13.9):

$$\begin{aligned} & \mathbf{x}_\tau(\boldsymbol{\varepsilon}_{t:\tau-1}; \mathbf{a}_{t:\tau-1}, f) \\ & \doteq f(\boldsymbol{\varepsilon}_{\tau-1}; f(\dots; (\mathbf{x}_t; \mathbf{a}_t), \mathbf{a}_{t+1}), \dots), \mathbf{a}_{\tau-1}). \end{aligned} \quad (13.19)$$

Observe that the epistemic and aleatoric uncertainty are treated differently. Within a particular MDP f , we ensure that randomness (i.e., ale-

⁵ see

- Equation (7.18) for variational inference;
- Equation (7.21a) for Markov chain Monte Carlo;
- Equation (7.28) for dropout regularization; and
- Equation (7.29) for probabilistic ensembles.

⁶ see Section 7.3.1

atoric uncertainty) is simulated consistently using our previous framework from our discussion of planning.⁷ The Monte Carlo samples of f take into account the epistemic uncertainty about the transition model. In our previous discussion of planning, we assumed the Markov decision process f to be fixed. Essentially, in Equation (13.17) we are now using Monte Carlo trajectory sampling as a subroutine and average over an “ensemble” of Markov decision processes.

⁷ see Equation (13.10)

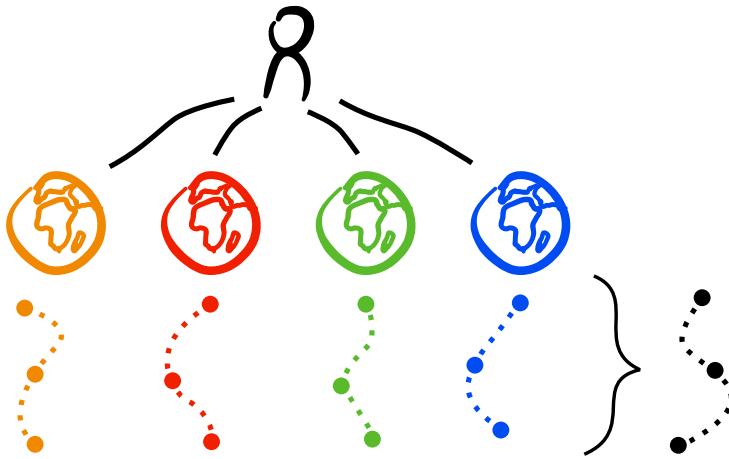


Figure 13.4: Illustration of planning with epistemic uncertainty and Monte Carlo sampling. The agent considers m alternative “worlds”. Within each world, it plans a sequence of actions over a finite time horizon. Then, the agent averages all optimal initial actions from all worlds. Crucially, each world by itself is *consistent*. That is, its transition model (i.e., the aleatoric uncertainty of the model) is constant.

The same approach that we have seen in Section 13.1.3 can be used to “compile” these plans into a parametric policy that can be trained offline, in which case, we write $\hat{J}_H(\pi)$ instead of $\hat{J}_H(a_{t:t+H-1})$.

This leads us to a first greedily exploitative algorithm for model-based reinforcement learning, which is shown in Algorithm 13.10. This algorithm is purely exploitative as it greedily maximizes the expected reward with respect to the transition model, taking into account epistemic uncertainty.

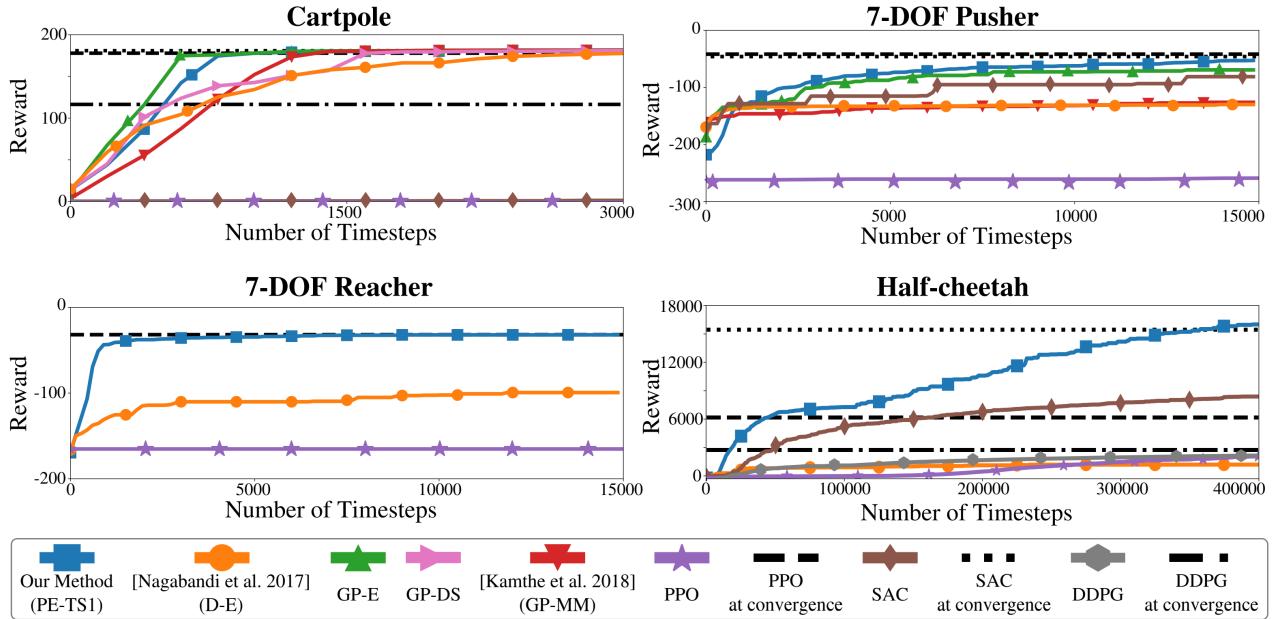
Algorithm 13.10: Greedy exploitation for model-based RL

- 1 start with (possibly empty) data $\mathcal{D} = \emptyset$ and a prior $p(f) = p(f | \mathcal{D})$
 - 2 **for** several episodes **do**
 - 3 | plan a new policy π to (approximately) maximize,

$$\max_{\pi} \mathbb{E}_{f \sim p(\cdot | \mathcal{D})} [\hat{J}_H(\pi; f)] \quad (13.20)$$
 - 4 | roll out policy π to collect more data
 - 5 | update posterior $p(f | \mathcal{D})$
-

In the context of Gaussian process models, this algorithm is called

probabilistic inference for learning control (PILCO) (Deisenroth and Rasmussen, 2011), which was the first demonstration of how sample efficient model-based reinforcement learning can be. As was mentioned in Remark 13.6, the originally proposed PILCO uses moment matching instead of Monte Carlo averaging.



In the context of neural networks, this algorithm is called *probabilistic ensembles with trajectory sampling* (PETS) (Chua et al., 2018), which is one of the state-of-the-art algorithms. PETS uses an ensemble of conditional Gaussian distributions over weights, trajectory sampling for evaluating the performance and model predictive control for planning.

Notably, PETS does not explicitly explore. Exploration only happens due to the uncertainty in the model, which already drives exploration to some extent. In many settings, however, this incentive is not sufficient for exploration — especially when rewards are sparse.

Figure 13.5: Sample efficiency of model-free and model-based reinforcement learning. DDPG and SAC are shown as constant (black) lines, because they take an order of magnitude more time steps before learning a good model. They also compare the PETS algorithm (blue) to deterministic ensembles (orange), where the transition models are assumed to be deterministic (or to have covariance 0). Reproduced with permission from “Deep reinforcement learning in a handful of trials using probabilistic dynamics models” (Chua et al., 2018).

Optional Readings

- **PILCO:** Deisenroth and Rasmussen (2011).
PILCO: A model-based and data-efficient approach to policy search.
- **PETS:** Chua, Calandra, McAllister, and Levine (2018).
Deep reinforcement learning in a handful of trials using probabilistic dynamics models.

13.2.2 Partial Observability

Depending on the task, it may be difficult to learn the observed dynamics directly. For example, when learning an artificial player for a computer game based only on the games' visual interface it may be difficult to predict the next frame in pixel space. Instead, a common approach is to treat the visual interface as an observation in a POMDP with a hidden latent space (cf. Section 10.4), and to learn the dynamics within the latent space and the observation model separately.

An example of this approach is the *deep planning network* (PlaNet) algorithm which learns such a POMDP via variational inference and uses the cross-entropy method for closed-loop planning of action sequences (Hafner et al., 2019). The *Dreamer* algorithm replaces the closed-loop planning of PlaNet by open-loop planning with entropy regularization (Hafner et al., 2020, 2021). Notably, neither PlaNet nor Dreamer explicitly take into account epistemic model uncertainty during planning but rather use point estimates.

13.3 Exploration

Exploration is critical when interacting with unknown environments, such as in reinforcement learning. We first encountered the exploration-exploitation dilemma in our discussion of Bayesian optimization, where we aimed at maximizing an unknown function in as little time as possible by making noisy observations.⁸ Within the framework of Bayesian optimization, we used so-called acquisition functions for selecting the next point at which to observe the unknown function. Observe that these acquisition functions are analogous to policies in the setting of reinforcement learning. In particular, the policy that uses greedy exploitation like we have seen in the previous section is analogous to simply picking the point that maximizes the mean of the posterior distribution. In the context of Bayesian optimization, we have already seen that this is insufficient for exploration and can easily get stuck in locally optimal solutions. As reinforcement learning is a strict generalization of Bayesian optimization, there is no reason why such a strategy should be sufficient now.

⁸ see Chapter 9

Recall from our discussion of Bayesian optimization that we “solved” this problem by using the epistemic uncertainty in our model of the unknown function to pick the next point to explore. This is what we will now explore in the context of reinforcement learning.

One simple strategy that we already investigated is the addition of some *exploration noise*. In other words, one follows a greedy exploitative strategy, but every once in a while, one chooses a random action

(like in ε -greedy); or one adds additional noise to the selected actions (known as Gaussian noise “dithering”). We have already seen that in difficult exploration tasks, these strategies are not systematic enough.

Two other approaches that we have considered before are Thompson sampling (cf. Section 9.3.3) and, more generally, the principle of *optimism in the face of uncertainty* (cf. Section 9.2.1).

13.3.1 Thompson Sampling

Recall from Section 9.3.3 the main idea behind Thompson sampling: namely, that the randomness in the realizations of f from the posterior distribution is already enough to drive exploration. That is, instead of picking the action that performs best on average across all realizations of f , Thompson sampling picks the action that performs best for a *single realization* of f . The epistemic uncertainty in the realizations of f leads to variance in the picked actions and provides an additional incentive for exploration. This yields Algorithm 13.11 which is an immediate adaptation of greedy exploitation and straightforward to implement.

Algorithm 13.11: Thompson sampling

- 1 start with (possibly empty) data $\mathcal{D} = \emptyset$ and a prior $p(f) = p(f | \mathcal{D})$
 - 2 **for** several episodes **do**
 - 3 sample a model $f \sim p(\cdot | \mathcal{D})$
 - 4 plan a new policy π to (approximately) maximize,
 - 5
$$\max_{\pi} \hat{f}_H(\pi; f) \quad (13.21)$$
 - 6 roll out policy π to collect more data
 - 6 update posterior $p(f | \mathcal{D})$
-

13.3.2 Optimistic Exploration

We have already seen in the context of Bayesian optimization and tabular reinforcement learning that optimism is a central pillar for exploration. But how can we explore optimistically in model-based reinforcement learning?

Let us consider a set $\mathcal{M}(\mathcal{D})$ of *plausible models* given some data \mathcal{D} . Optimistic exploration would then optimize for the most advantageous model among all models that are plausible given the seen data.

Example 13.12: Plausible conditional Gaussians

In the context of conditional Gaussians, we can consider the set of all models such that the prediction of a single dimension i lies in some confidence interval,

$$\begin{aligned} f_i(\mathbf{x}, \mathbf{a}) \in \mathcal{C}_i &\doteq [\mu_i(\mathbf{x}, \mathbf{a} \mid \mathcal{D}) - \beta_i \sigma_i(\mathbf{x}, \mathbf{a} \mid \mathcal{D}), \\ &\quad \mu_i(\mathbf{x}, \mathbf{a} \mid \mathcal{D}) + \beta_i \sigma_i(\mathbf{x}, \mathbf{a} \mid \mathcal{D})], \end{aligned} \quad (13.22)$$

where β_i tunes the width of the confidence interval, analogously to Section 9.3.1. The set of plausible models is then given as

$$\mathcal{M}(\mathcal{D}) \doteq \{f \mid \forall \mathbf{x} \in \mathcal{X}, \mathbf{a} \in \mathcal{A}, i \in [d] : f_i(\mathbf{x}, \mathbf{a}) \in \mathcal{C}_i\}. \quad (13.23)$$

When compared to greedy exploitation, instead of taking the optimal step on average with respect to all realizations of the transition model f , optimistic exploration as shown in Algorithm 13.13 takes the optimal step with respect to the most optimistic model among all transition models that are consistent with the data.

Algorithm 13.13: Optimistic exploration

-
- 1 start with (possibly empty) data $\mathcal{D} = \emptyset$ and a prior $p(f) = p(f \mid \mathcal{D})$
 - 2 **for** several episodes **do**
 - 3 plan a new policy π to (approximately) maximize,

$$\max_{\pi} \max_{f \in \mathcal{M}(\mathcal{D})} \hat{J}_H(\pi; f) \quad (13.24)$$
 - 4 roll out policy π to collect more data
 - 5 update posterior $p(f \mid \mathcal{D})$
-

In general, the joint maximization over π and f is very difficult to solve computationally. Yet, remarkably, it turns out that this complex optimization problem can be reduced to standard model-based reinforcement learning with a fixed model. The key idea is to consider an agent that can control its “luck”. In other words, we assume that the agent believes it can control the outcome of its actions — or rather choose which of the plausible dynamics it follows. The “luck” of the agent can be considered as additional decision variables. Consider the optimization problem,

$$\pi_{t+1} \doteq \arg \max_{\pi} \max_{\eta(\cdot) \in [-1,1]^d} \hat{J}_H(\pi; \hat{f}_t) \quad (13.25)$$

with “optimistic” dynamics

$$\hat{f}_{t,i}(\mathbf{x}, \mathbf{a}) \doteq \mu_{t,i}(\mathbf{x}, \mathbf{a}) + \beta_{t,i} \eta_i(\mathbf{x}, \mathbf{a}) \sigma_{t,i}(\mathbf{x}, \mathbf{a}). \quad (13.26)$$