

Andreas Krause

Probabilistic Artificial Intelligence

Fall 2023



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Institute for Machine Learning
Department of Computer Science

Contents. Bayesian learning and inference. Bayesian linear regression, Gaussian processes, Bayesian Deep Learning. Variational inference, Markov chain Monte Carlo methods. Active learning, Bayesian optimization and the exploration-exploitation dilemma. Markov decision processes, and Reinforcement learning in tabular and large state-action spaces with model-based and model-free approaches.

ACKNOWLEDGEMENT

These notes were written and edited by **Jonas Hübötter**.

- Many figures were inspired by the accompanying Jupyter notebooks created by Sebastian Curi.
- Parts of section 1.1 were adapted from earlier notes by Mohammad Reza Karimi.
- Figure 12.2 was contributed by Hado van Hasselt.
- Some exercises are adapted from previous iterations of the course.
- We thank Zhiyuan Hu and Shyam Sundhar Ramesh for proofreading.

CONTRIBUTING

You are encouraged to raise issues and suggest fixes for anything you think can be improved.

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

REPOSITORY: <https://gitlab.inf.ethz.ch/OU-KRAUSE/pai-script>

JUPYTER NOTEBOOKS WITH EXAMPLES

<https://gitlab.inf.ethz.ch/OU-KRAUSE/pai-demos>

SYLLABUS

1. *Introduction* chapter 1

Part I: Probabilistic Machine Learning

2. *Bayesian Linear Regression & Filtering* chapter 2 through to section 2.4 and chapter 3
3. *Gaussian Processes 1* rest of chapter 2 and chapter 4 through to section 4.3
4. *Gaussian Processes 2* rest of chapter 4
5. *Variational Inference* chapter 5 up to section 5.5.1
6. *Markov Chain Monte Carlo* section 5.5.1 and chapter 6
7. *Bayesian Deep Learning* chapter 7

Part II: Sequential Decision-Making

8. *Active Learning & Bayesian Optimization* chapter 8 and chapter 9
9. *Markov Decision Processes* chapter 10
10. *Reinforcement Learning 1* chapter 11
11. *Reinforcement Learning 2* chapter 12 through to section 12.3
12. *Reinforcement Learning 3* section 12.4 through to section 12.5
13. *Reinforcement Learning 4* chapter 13

COMPILATION DATE: **December 2, 2023**

This set of notes was written for the course PROBABILISTIC ARTIFICIAL INTELLIGENCE (263-5210-00L) at ETH Zürich. Distribution of these notes without the permission of the authors is prohibited.

© 2023 ETH Zürich. All rights reserved.

Contents

Preface vii

Summary of Notation ix

1	<i>Fundamentals</i>	1
1.1	<i>Probability</i>	1
1.2	<i>Supervised Learning</i>	18
1.3	<i>Bayesian Learning and Inference</i>	19
1.4	<i>Parameter Estimation</i>	22
1.5	<i>Optimization</i>	30
1.6	<i>Multivariate Normal Distribution</i>	37
	 <i>I Probabilistic Machine Learning</i>	 45
2	<i>Bayesian Linear Regression</i>	47
2.1	<i>Linear Regression</i>	47
2.2	<i>Weight-space View</i>	50
2.3	<i>Aleatoric and Epistemic Uncertainty</i>	54
2.4	<i>Recursive Bayesian Updates</i>	54
2.5	<i>Non-linear Regression</i>	56
2.6	<i>Function-space View</i>	56
3	<i>Kalman Filters</i>	61
3.1	<i>Bayesian Filtering</i>	63
3.2	<i>Kalman Update</i>	65

3.3	<i>Predicting</i>	66
4	<i>Gaussian Processes</i>	69
4.1	<i>Learning and Inference</i>	70
4.2	<i>Sampling</i>	71
4.3	<i>Kernel Functions</i>	71
4.4	<i>Model Selection</i>	78
4.5	<i>Approximations</i>	83
5	<i>Variational Inference</i>	91
5.1	<i>Laplace Approximation</i>	91
5.2	<i>Inference with a Variational Posterior</i>	95
5.3	<i>Blueprint of Variational Inference</i>	98
5.4	<i>Information Theory</i>	99
5.5	<i>Evidence Lower Bound</i>	111
6	<i>Markov Chain Monte Carlo Methods</i>	121
6.1	<i>Markov Chains</i>	122
6.2	<i>Elementary Sampling Methods</i>	130
6.3	<i>Sampling as Optimization</i>	134
7	<i>Bayesian Deep Learning</i>	147
7.1	<i>Artificial Neural Networks</i>	147
7.2	<i>Bayesian Neural Networks</i>	151
7.3	<i>Approximate Inference</i>	152
7.4	<i>Calibration</i>	157
<i>II Sequential Decision-Making</i>		161
8	<i>Active Learning</i>	163
8.1	<i>Conditional Entropy</i>	163
8.2	<i>Mutual Information</i>	165
8.3	<i>Submodularity of Mutual Information</i>	168
8.4	<i>Maximizing Mutual Information</i>	170
8.5	<i>Beyond Mutual Information: Experimental Design</i>	174

8.6	<i>Beyond Global Learning</i>	175
9	<i>Bayesian Optimization</i>	179
9.1	<i>Exploration-Exploitation Dilemma</i>	179
9.2	<i>Online Learning and Bandits</i>	180
9.3	<i>Acquisition Functions</i>	182
9.4	<i>Model Selection</i>	192
10	<i>Markov Decision Processes</i>	195
10.1	<i>Bellman Expectation Equation</i>	198
10.2	<i>Policy Evaluation</i>	199
10.3	<i>Policy Optimization</i>	201
10.4	<i>Partial Observability</i>	209
11	<i>Tabular Reinforcement Learning</i>	215
11.1	<i>The Reinforcement Learning Problem</i>	215
11.2	<i>Model-based Approaches</i>	217
11.3	<i>Balancing Exploration and Exploitation</i>	218
11.4	<i>Model-free Approaches</i>	222
12	<i>Model-free Reinforcement Learning</i>	231
12.1	<i>Tabular Reinforcement Learning as Optimization</i>	231
12.2	<i>Value Function Approximation</i>	233
12.3	<i>Policy Approximation</i>	238
12.4	<i>Actor-Critic Methods</i>	245
12.5	<i>Maximum Entropy Reinforcement Learning</i>	256
13	<i>Model-based Reinforcement Learning</i>	263
13.1	<i>Planning</i>	264
13.2	<i>Learning</i>	271
13.3	<i>Exploration</i>	276
A	<i>Mathematical Background</i>	287
A.1	<i>Useful Matrix Identities and Inequalities</i>	287

B Solutions 289

Acronyms 351

Index 353

Errata 361

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 the course “Probabilistic Artificial Intelligence” is about. The first part covers probabilistic approaches to machine learning. We discuss how Bayesian approaches to learning allow to differentiate 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 Bayesian learning, 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 course 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.

Summary of Notation

We follow these general rules:

- uppercase italic for constants N
- lowercase italic for indices i and scalar variables x
- lowercase italic bold for vectors x , entries are denoted $x(i)$
- uppercase italic bold for matrices M , entries are denoted $M(i, j)$
- uppercase italic for random variables X
- uppercase bold for random vectors \mathbf{X} , entries are denoted $\mathbf{X}(i)$
- uppercase italic for sets A
- uppercase calligraphy for spaces (usually infinite sets) \mathcal{A}

\doteq	equality by definition
\approx	approximately equals
\propto	proportional to (up to multiplicative constants), $f \propto g$ iff $\exists k. \forall x. f(x) = k \cdot g(x)$
const	an (additive) constant
\mathbb{N}	set of natural numbers $\{1, 2, \dots\}$
\mathbb{N}_0	set of natural numbers, including 0, $\mathbb{N} \cup \{0\}$
\mathbb{R}	set of real numbers
$[m]$	set of natural numbers from 1 to m , $\{1, 2, \dots, m-1, m\}$
$i : j$	subset of natural numbers between i and j , $\{i, i+1, \dots, j-1, j\}$
$(a, b]$	real interval between a and b including b but not including a
$f : A \rightarrow B$	function f from elements of set A to elements of set B
$f \circ g$	function composition, $f(g(\cdot))$
$(\cdot)_+$	$\max\{0, \cdot\}$
log	logarithm with base e
$\mathcal{P}(A)$	power set (set of all subsets) of A
$\mathbb{1}\{predicate\}$	indicator function ($\mathbb{1}\{predicate\} \doteq 1$ if the <i>predicate</i> is true, else 0)
\odot	Hadamard (element-wise) product
\leftarrow	assignment

$\nabla f(x) \in \mathbb{R}^n$	gradient of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at a point $x \in \mathbb{R}^n$
$Dg(x) \in \mathbb{R}^{m \times n}$	Jacobian of a function $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at a point $x \in \mathbb{R}^n$
$Hf(x) \in \mathbb{R}^{n \times n}$	Hessian of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at a point $x \in \mathbb{R}^n$
$\nabla \cdot \mathbf{F}$	divergence operation on vector field \mathbf{F}
Δf	Laplacian of a scalar field $f : \mathbb{R}^n \rightarrow \mathbb{R}$
$f \in \mathcal{O}(g)$	f grows at most as fast as g (up to constant factors), $0 \leq \limsup_{n \rightarrow \infty} \left \frac{f(n)}{g(n)} \right < \infty$
$f \in \tilde{\mathcal{O}}(g)$	f grows at most as fast as g up to constant and logarithmic factors
$\ \cdot\ _\alpha$	α -norm
$\ \cdot\ _A$	Mahalanobis norm induced by matrix A

LINEAR ALGEBRA

I	identity matrix
A^\top	transpose of matrix A
A^{-1}	inverse of invertible matrix A
$A^{1/2}$	square root of a symmetric and positive semi-definite matrix A
$\det(A)$	determinant of A
$\text{tr}(A)$	trace of A , $\sum_i A(i, i)$
$\text{diag}_{i \in I} \{a_i\}$	diagonal matrix with elements a_i , indexed according to the set I

PROBABILITY

Ω	sample space
\mathcal{A}	event space
\mathbb{P}	probability measure
$X \sim P$	random variable X follows the distribution P
$X_{1:n} \stackrel{\text{iid}}{\sim} P$	random variables $X_{1:n}$ are independent and identically distributed according to distribution P
$x \sim P$	value x is sampled according to distribution P
P_X	cumulative distribution function of a random variable X
\bar{P}_X	tail distribution function of a random variable X
P_X^{-1}	quantile function of a random variable X
p_X	probability mass function (if discrete) or probability density function (if continuous) of a random variable X
$\Delta^{\mathcal{A}}$	set of all probability distributions over the set \mathcal{A}
δ_α	Dirac delta function, point density at α
$g\sharp p$	pushforward of a density p under perturbation g
$X \perp Y$	random variable X is independent of random variable Y
$X \perp Y \mid Z$	random variable X is conditionally independent of random variable Y given random variable Z

$\mathbb{E}[X]$	expected value of random variable X
$\mathbb{E}_{x \sim X}[f(x)]$	expected value of the random variable $f(X)$, $\mathbb{E}[f(X)]$
$\mathbb{E}[X Y]$	conditional expectation of random variable X given random variable Y
$\text{Cov}[X, Y]$	covariance of random variable X and random variable Y
$\text{Cor}[X, Y]$	correlation of random variable X and random variable Y
$\text{Var}[X]$	variance of random variable X
$\text{Var}[X Y]$	conditional variance of random variable X given random variable Y
Σ_X	covariance matrix of random vector \mathbf{X}
Λ_X	precision matrix of random vector \mathbf{X}
$\text{MSE}(X)$	mean squared error of random variable X
\bar{X}_n	sample mean of random variable X with n samples
S_n^2	sample variance of random variable X with n samples
$X_n \xrightarrow{\text{a.s.}} X$	the sequence of random variables X_n converges almost surely to X
$X_n \xrightarrow{\mathbb{P}} X$	the sequence of random variables X_n converges to X in probability
$X_n \xrightarrow{\mathcal{D}} X$	the sequence of random variables X_n converges to X in distribution
$S[u]$	surprise associated with an event of probability u
$H[p], H[X]$	entropy of distribution p (or random variable X)
$H[p q]$	cross-entropy of distribution q with respect to distribution p
$\text{KL}(p q)$	KL-divergence of distribution q with respect to distribution p
$J(p q)$	relative Fisher information of distribution q with respect to distribution p
$H[X Y]$	conditional entropy of random variable X given random variable Y
$H[X, Y]$	joint entropy of random variables X and Y
$I(X; Y)$	mutual information of random variables X and Y
$I(X; Y Z)$	conditional mutual information of random variables X and Y given random variable Z
$\mathcal{N}(\mu, \Sigma)$	normal distribution with mean μ covariance Σ
$\text{Laplace}(\mu, h)$	Laplace distribution with mean μ scale h
$\text{Unif}(S)$	uniform distribution on the set S
$\text{Bern}(p)$	Bernoulli distribution with success probability p
$\text{Bin}(n, p)$	binomial distribution with n trials and success probability p
$\text{Beta}(\alpha, \beta)$	beta distribution with shape parameters α and β
$\text{Gamma}(\alpha, \beta)$	gamma distribution with shape α and rate β
$\text{Cauchy}(m, \tau)$	Cauchy distribution with location m and scale τ
$\text{Pareto}(c, \alpha)$	Pareto distribution with cutoff threshold c and shape α

SUPERVISED LEARNING

θ	parameterization of a model
\mathcal{X}	input space
\mathcal{Y}	label space
$x \in \mathcal{X}$	input

$\epsilon(x)$	zero-mean noise, sometimes assumed to be independent of x
$y \in \mathcal{Y}$	(noisy) label, $f(x) + \epsilon(x)$ where f is unknown
$\mathcal{D} \subseteq \mathcal{X} \times \mathcal{Y}$	labeled training data, $\{(x_i, y_i)\}_{i=1}^n$
$\mathbf{X} \in \mathbb{R}^{n \times d}$	design matrix when $\mathcal{X} = \mathbb{R}^d$
$\Phi \in \mathbb{R}^{n \times e}$	design matrix in feature space \mathbb{R}^e
$\mathbf{y} \in \mathbb{R}^n$	label vector when $\mathcal{Y} = \mathbb{R}$
$p(\theta)$	prior belief about θ
$p(\theta \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$	posterior belief about θ given training data
$p(y_{1:n} \mid \mathbf{x}_{1:n}, \theta)$	likelihood of training data under the model parameterized by θ
$p(y_{1:n} \mid \mathbf{x}_{1:n})$	marginal likelihood of training data
$\hat{\theta}_{\text{MLE}}$	maximum likelihood estimate of θ
$\hat{\theta}_{\text{MAP}}$	maximum a posteriori estimate of θ
$\ell_{\text{nl}}(\theta; \mathcal{D})$	negative log-likelihood of the training data \mathcal{D} under model θ

BAYESIAN LINEAR MODELS

$\mathbf{w} \in \mathbb{R}^d$	weights of linear function $f(x; \mathbf{w}) = \mathbf{w}^\top \mathbf{x}$
$\hat{\mathbf{w}}_{\text{ls}}$	least squares estimate of \mathbf{w}
$\hat{\mathbf{w}}_{\text{ridge}}$	ridge estimate of \mathbf{w}
$\hat{\mathbf{w}}_{\text{lasso}}$	lasso estimate of \mathbf{w}
$\mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$	prior
$\mathcal{N}(\mathbf{w}^\top \mathbf{x}, \sigma_n^2)$	likelihood
$\boldsymbol{\mu} \in \mathbb{R}^d$	posterior mean, $\sigma_n^{-2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y}$
$\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$	posterior covariance matrix, $(\sigma_n^{-2} \mathbf{X}^\top \mathbf{X} + \sigma_p^{-2} \mathbf{I})^{-1}$
$\mathbf{K} \in \mathbb{R}^{n \times n}$	kernel matrix, $\sigma_p^2 \mathbf{X} \mathbf{X}^\top$
σ	logistic function
$\text{Bern}(\sigma(\mathbf{w}^\top \mathbf{x}))$	logistic likelihood
$\ell_{\log}(\cdot; \mathbf{x}, y)$	logistic loss of a single training example (\mathbf{x}, y)

KALMAN FILTERS

\mathbf{X}_t	sequence of hidden states in \mathbb{R}^d
\mathbf{Y}_t	sequence of observations in \mathbb{R}^m
$\mathbf{F} \in \mathbb{R}^{d \times d}$	motion model
$\mathbf{H} \in \mathbb{R}^{m \times d}$	sensor model
$\boldsymbol{\epsilon}_t$	zero-mean motion noise with covariance matrix $\boldsymbol{\Sigma}_x$
$\boldsymbol{\eta}_t$	zero-mean sensor noise with covariance matrix $\boldsymbol{\Sigma}_y$
$\mathbf{K}_t \in \mathbb{R}^{d \times m}$	Kalman gain

GAUSSIAN PROCESSES

$\mu : \mathcal{X} \rightarrow \mathbb{R}$	mean function
$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$	kernel function / covariance function
$f \sim \mathcal{GP}(\mu, k)$	f is a Gaussian process with mean function μ and kernel function k
$\mathcal{H}_k(\mathcal{X})$	reproducing kernel Hilbert space associated with kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

DEEP MODELS

$W_l \in \mathbb{R}^{n_l \times n_{l-1}}$	weight matrix of layer l
$v^{(l)} \in \mathbb{R}^{n_l}$	activations of layer l
φ	activation function
Tanh	hyperbolic tangent activation function
ReLU	rectified linear unit activation function
$\sigma_i(f)$	softmax function computing the probability mass of class i given outputs f

VARIATIONAL INFERENCE

\mathcal{Q}	variational family
$\lambda \in \Lambda$	variational parameters
q_λ	variational posterior parameterized by λ
$L(q, p; \mathcal{D})$	evidence lower bound for data \mathcal{D} of variational posterior q and true posterior $p(\cdot \mid \mathcal{D})$

MARKOV CHAINS

S	set of n states
X_t	sequence of states
$p(x' \mid x)$	transition function, probability of going from state x to state x'
$p^{(t)}(x' \mid x)$	probability of reaching x' from x in exactly t steps
$P \in \mathbb{R}^{n \times n}$	transition matrix
q_t	distribution over states at time t
π	stationary distribution
$\ \mu - \nu\ _{\text{TV}}$	total variation distance between two distributions μ and ν
τ_{TV}	mixing time with respect to total variation distance

MARKOV CHAIN MONTE CARLO METHODS

$r(x' \mid x)$	proposal distribution, probability of proposing x' when in x
----------------	--

$\alpha(x' x)$	acceptance distribution, probability of accepting the proposal x' when in x
f	energy function

ACTIVE LEARNING

$S \subseteq \mathcal{X}$	set of observations
$I(S)$	maximization objective, quantifying the “information value” of S
$\Delta_I(x A)$	marginal gain of observation x with respect to objective I given prior observations A

BAYESIAN OPTIMIZATION

R_T	cumulative regret for time horizon T
$F(x; \mu, \sigma)$	acquisition function
$\mathcal{C}_t(x)$	confidence interval of $f^*(x)$ after round t
$\beta_t(\delta)$	scale of confidence interval to achieve confidence level δ
γ_T	maximum information gain after T rounds

REINFORCEMENT LEARNING

X, \mathcal{X}	set of states
A, \mathcal{A}	set of actions
$p(x' x, a)$	dynamics model, probability of transitioning from state x to state x' when playing action a
r	reward function
X_t	sequence of states
A_t	sequence of actions
R_t	sequence of rewards
$\pi(a x)$	policy, probability of playing action a when in state x
G_t	discounted payoff from time t
γ	discount factor
$v_t^\pi(x)$	state value function, average discounted payoff from time t starting from state x
$q_t^\pi(x, a)$	state-action value function, average discounted payoff from time t starting from state x playing action a
$a_t^\pi(x, a)$	advantage function, $q_t^\pi(x, a) - v_t^\pi(x)$
$j(\pi)$	policy value function, expected reward of policy π

1

Fundamentals

1.1 Probability

Probability is the natural extension of Boolean logic (where every statement is either true or false) to reasoning under uncertainty (Jaynes, 2003). 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. Often, probabilities are used to describe *beliefs* that one has about the outcome of a trial. While it is natural to consider the relative frequency of the outcome in repeated experiments as our belief, 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 we can conduct multiple independent experiments. In these settings, our notion of probability is simply a (subjective) measure of uncertainty about occurrence of events and does not require any experiments. This notion is commonly understood as *Bayesian reasoning* or the Bayesian interpretation of probability.¹

¹ In the early 20th century, Bruno De Finetti has done foundational work to formalize this notion.

Readings

In this section, we will briefly recall the fundamentals of probability theory, which we will need in the subsequent chapters.

For a more thorough introduction, read chapter 6 of “Mathematics for machine learning” (Deisenroth et al., 2020).

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.² This set of events is required to be a σ -algebra over the sample space.

² We use $\mathcal{P}(\Omega)$ to denote the *power set* (set of all subsets) of Ω .

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
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 $\{A_i \in \mathcal{A}\}_i$.³

³ We say that a set of sets $\{A_i\}_i$ is disjoint if for all $i \neq j$ we have $A_i \cap A_j = \emptyset$.

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} .

Exercise 1.5: 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 $A \in \mathcal{A}$, $\mathbb{P}(\bar{A}) = 1 - \mathbb{P}(A)$.
3. For any countable set of events $\{A_i \in \mathcal{A}\}_i$,

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) \leq \sum_{i=1}^{\infty} \mathbb{P}(A_i). \quad (1.1)$$

which is called a *union bound*.

▷ *Solution*

Example 1.6: 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 1.7 (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,⁴ and where X respects the information available in the σ -algebra \mathcal{A} . That is,⁵

$$\forall S \subseteq \mathcal{T} : \{\omega \in \Omega : X(\omega) \in S\} \in \mathcal{A}. \quad (1.2)$$

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

$$\mathbb{P}(X \in S) = \mathbb{P}(\{\omega \in \Omega : X(\omega) \in S\}). \quad (1.3)$$

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

$$p_X(x) \doteq \mathbb{P}(X = x) \quad (\text{in the discrete setting}), \quad (1.4)$$

$$P_X(x) \doteq \mathbb{P}(X \leq x). \quad (1.5)$$

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.

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.

Remark 1.8: Notation

We use $X \sim P$ to denote that a random variable X is distributed according to the distribution P . Due to the unique characterization of random variables through distributions, the two notions are often conflated. Sometimes we use P as a random variable (but

⁴For a random variable that maps a graph to its number of edges, $\mathcal{T} = \mathbb{N}_0$. For our purposes, you can generally assume $\mathcal{T} \subseteq \mathbb{R}$.

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

refer to a random variable that is distributed according to P).

This notation is not to be confused with $x \sim P$, which we use to denote that a value x is sampled according to a distribution P .

Further, we often abbreviate $\mathbb{P}(X = x)$ by $\mathbb{P}(x)$ if the correspondence between X and x is clear from the context. Sometimes, we also use $p(x)$ equivalently (to denote the PMF or PDF of X).

Example 1.9: Common discrete distributions

- *Bernoulli* $\text{Bern}(p)$ describes (biased) coin flips. Its domain is $\Omega = \{0, 1\}$ where 1 corresponds to heads and 0 corresponds to tails. $p \in [0, 1]$ is the probability of the coin landing heads, that is, $\mathbb{P}(X = 1) = p$ and $\mathbb{P}(X = 0) = 1 - p$.
- *Binomial* $\text{Bin}(n, p)$ counts the number of heads in n independent Bernoulli trials, each with probability of heads p .
- *Categorical* $\text{Cat}(m, p_1, \dots, p_m)$ is a generalization of the Bernoulli distribution and represents throwing a (biased) m -sided die. Its domain is $\Omega = [m]$ and we have $\mathbb{P}(X = i) = p_i$. We require $p_i \geq 0$ and $\sum_{i \in [m]} p_i = 1$.
- *Multinomial* $\text{Mult}(n, m, p_1, \dots, p_m)$ counts the number of outcomes of each side in n independent Categorical trials.
- *Uniform* $\text{Unif}(S)$ assigns identical probability mass to all values in the set S . That is, $\mathbb{P}(X = x) = \frac{1}{|S|}$ ($\forall x \in S$).

Remark 1.10: Probability simplex

We denote by Δ^m the set of all categorical distributions on m classes. Observe that Δ^m is an $(m - 1)$ -dimensional convex polytope.

To see this, let us consider the m -dimensional space of probabilities $[0, 1]^m$. It follows from our characterization of the categorical distribution in example 1.9 that there is a one-to-one correspondence between probability distributions over m classes and points in the space $[0, 1]^m$ where all coordinates sum to one.

This $(m - 1)$ -dimensional subspace of $[0, 1]^m$ is also known as the *probability simplex*.

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 1.11: Normal distribution / Gaussian

An 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

$$\mathcal{N}(x; \mu, \sigma^2) \doteq \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right). \quad (1.6)$$

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.

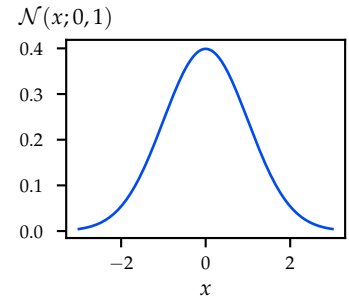


Figure 1.1: PDF of the standard normal distribution. Observe that the PDF is symmetric around the mode.

Much of the relevance of the normal distribution stems from the *central limit theorem*, which states that the sum of any i.i.d. (independent and identically distributed) samples tends to a normal distribution as the sample size goes to infinity (even if the samples themselves are not normally distributed).

Fact 1.12 (Central limit theorem by Lindeberg-Lévy). *Given independent and identically distributed random variables X_1, \dots, X_n with finite expectation μ and variance σ^2 . Then,*

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu) = \sqrt{n}(\bar{X}_n - \mu) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma^2) \quad (1.7)$$

as $n \rightarrow \infty$.

Here, \bar{X}_n denotes the sample mean of X_1, \dots, X_n (see example 1.36) and $\xrightarrow{\mathcal{D}}$ denotes convergence in distribution (see eq. (1.80)).

1.1.5 Universality of the Uniform and Sampling

Sampling from a continuous random variable is a common and non-trivial problem. A family of techniques is motivated by the following general property of the uniform distribution, colloquially known as the “universality of the uniform”.

First, given a random variable $X \sim P$, we call

$$P^{-1}(u) \doteq \min\{x \mid P(x) \geq u\} \quad \text{for all } 0 < u < 1. \quad (1.8)$$

the *quantile function* of X . That is, $P^{-1}(u)$ corresponds to the value x such that the probability of X being at most x is u . If the CDF P is invertible, then P^{-1} coincides with the inverse of P .

Theorem 1.13 (Universality of the uniform). *If $U \sim \text{Unif}([0, 1])$ and P is invertible, then $P^{-1}(U) \sim P$.*

Proof. Let $U \sim \text{Unif}([0, 1])$. Then,

$$\mathbb{P}(P^{-1}(U) \leq x) = \mathbb{P}(U \leq P(x)) = P(x). \quad \square$$

This implies that if we are able to sample from $\text{Unif}([0, 1])$,⁶ then we are able to sample from any distribution with invertible CDF. This method is known as *inverse transform sampling*.

In the case of Gaussians, we learned that the CDF cannot be expressed in closed-form (and hence, is not invertible), however, for practical purposes, the quantile function of Gaussians can be approximated well.

Uniform distribution The (continuous) *uniform distribution* $\text{Unif}([a, b])$ is the only distribution that assigns constant density to all points in the support $[a, b]$. That is, it has PDF

$$p(u) = \begin{cases} \frac{1}{b-a} & \text{if } u \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

and CDF

$$P(u) = \begin{cases} \frac{u-a}{b-a} & \text{if } u \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

using $\mathbb{P}(U \leq u) = u$ if $U \sim \text{Unif}([0, 1])$

⁶ This is done in practice using so-called *pseudo-random number generators*.

1.1.6 Joint Probability

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

$$\mathbb{P}(A, B) \doteq \mathbb{P}(A \cap B). \quad (1.9)$$

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 . 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 1.14 (Sum rule). *We have that*

$$p(x_{1:i-1}, x_{i+1:n}) = \sum_{x_i \in X_i(\Omega)} p(x_{1:i-1}, x_i, x_{i+1:n}), \quad (1.10a)$$

$$p(x_{1:i-1}, x_{i+1:n}) = \int_{X_i(\Omega)} p(x_{1:i-1}, x_i, x_{i+1:n}) dx_i. \quad (1.10b)$$

if X_i is discrete and continuous respectively.

1.1.7 Conditional Probability

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

Definition 1.15 (Conditional probability). Given two events A and B such that $\mathbb{P}(B) > 0$, the probability of A conditioned on B is given as

$$\mathbb{P}(A | B) \doteq \frac{\mathbb{P}(A, B)}{\mathbb{P}(B)}. \quad (1.11)$$

Simply rearranging the terms yields,

$$\mathbb{P}(A, B) = \mathbb{P}(A | B) \cdot \mathbb{P}(B) = \mathbb{P}(B | A) \cdot \mathbb{P}(A). \quad (1.12)$$

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.

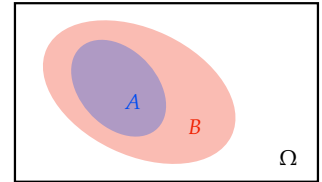


Figure 1.2: Conditioning an event A on another event B can be understood as replacing the universe of all possible outcomes Ω by the observed outcomes B . Then, the conditional probability is simply expressing the likelihood of A given that B occurred.

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

$$p_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} \mid \mathbf{y}) \doteq \frac{p_{\mathbf{X},\mathbf{Y}}(\mathbf{x}, \mathbf{y})}{p_{\mathbf{Y}}(\mathbf{y})}. \quad (1.13)$$

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

Remark 1.16: Notation

When $\mathbf{X} \sim p$, we write $p(\cdot \mid \mathbf{y})$ to denote the distribution of the random vector $\mathbf{X} \mid \mathbf{y}$. In the literature, $p(\mathbf{X} \mid \mathbf{y})$ is sometimes used to denote equivalent quantities.

When sampling $\mathbf{x} \sim p(\cdot \mid \mathbf{y})$ and the distribution p is clear from context, we abbreviate $\mathbf{x} \mid \mathbf{y}$.

As noted previously, if the correspondence between the random vectors \mathbf{X} and \mathbf{Y} and their states \mathbf{x} and \mathbf{y} is clear from context, we often omit the indices of the PDF p and simply write $p(\mathbf{x} \mid \mathbf{y})$, $p(\mathbf{x}, \mathbf{y})$ and $p(\mathbf{y})$. However, this does *not* mean that they refer to the same distribution p .

Extending eq. (1.12) to arbitrary random vectors yields the product rule (also called *chain rule of probability*).

Fact 1.17 (Product rule). *Given random variables $X_{1:n}$,*

$$p(x_{1:n}) = p(x_1) \cdot \prod_{i=2}^n p(x_i \mid x_{1:i-1}). \quad (1.14)$$

Combining sum rule and product rule, we can compute marginal probabilities too. If \mathbf{Y} is continuous, then

$$p(\mathbf{x}) = \int_{\mathbf{Y}(\Omega)} p(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \int_{\mathbf{Y}(\Omega)} p(\mathbf{x} \mid \mathbf{y}) \cdot p(\mathbf{y}) d\mathbf{y} \quad (1.15)$$

first using the sum rule (1.10) then the product rule (1.14)

Analogously, one can condition on a discrete random variable by replacing the integral with a sum. The above is called the *law of total probability* (LOTP). Colloquially, this is often referred to as to “condition” 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 $p(\mathbf{x} \mid \mathbf{y})$ and $p(\mathbf{y})$ are straightforward to understand.

Example 1.18: Point densities and Dirac’s delta function

The *Dirac delta function* δ_α is a function satisfying $\delta_\alpha(x) = 0$ for any $x \neq \alpha$ and $\int \delta_\alpha(x) dx = 1$. δ_α is also called a *point density* at α .

As an example, let us consider the random variable $Y = g(X)$, which is defined in terms of another random variable X and a continuously differentiable function $g : \mathbb{R} \rightarrow \mathbb{R}$. Using the sum

rule and product rule, we can express the PDF of Y as

$$p(y) = \int p(y | x) \cdot p(x) dx = \int \delta_{g(x)}(y) \cdot p(x) dx. \quad (1.16)$$

In section 1.1.12, we discuss how one can obtain an explicit representation of p_Y using the “change of variables” formula.

Exercise 1.19: 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 vertex v eventually?

▷ *Solution*

1.1.8 Independence

Two random vectors \mathbf{X} and \mathbf{Y} are *independent* (denoted $\mathbf{X} \perp \mathbf{Y}$) iff 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

$$P_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}) = P_{\mathbf{X}}(\mathbf{x}), \quad p_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}) = p_{\mathbf{X}}(\mathbf{x}). \quad (1.17)$$

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

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:

$$P_{\mathbf{X},\mathbf{Y}}(\mathbf{x}, \mathbf{y}) = P_{\mathbf{X}}(\mathbf{x}) \cdot P_{\mathbf{Y}}(\mathbf{y}), \quad p_{\mathbf{X},\mathbf{Y}}(\mathbf{x}, \mathbf{y}) = p_{\mathbf{X}}(\mathbf{x}) \cdot p_{\mathbf{Y}}(\mathbf{y}). \quad (1.18)$$

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

We often say that a set of random variables $\{X_i\}_i$ is i.i.d. to abbreviate that the X_i are “(mutually) independent and identically distributed”, and we write $X_i \stackrel{\text{iid}}{\sim} p$.

A “weaker” notion of independence is conditional independence.⁷ 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

⁷ It is explained in remark 1.20 how “weaker” is to be interpreted in this context.

value of one random vector \mathbf{Y} does not affect the distribution of the other random vector \mathbf{X} , namely if

$$P_{\mathbf{X}|\mathbf{Y},\mathbf{Z}}(\mathbf{x} | \mathbf{y}, \mathbf{z}) = P_{\mathbf{X}|\mathbf{Z}}(\mathbf{x} | \mathbf{z}), \quad (1.19a)$$

$$p_{\mathbf{X}|\mathbf{Y},\mathbf{Z}}(\mathbf{x} | \mathbf{y}, \mathbf{z}) = p_{\mathbf{X}|\mathbf{Z}}(\mathbf{x} | \mathbf{z}). \quad (1.19b)$$

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:

$$P_{\mathbf{X},\mathbf{Y}|\mathbf{Z}}(\mathbf{x}, \mathbf{y} | \mathbf{z}) = P_{\mathbf{X}|\mathbf{Z}}(\mathbf{x} | \mathbf{z}) \cdot P_{\mathbf{Y}|\mathbf{Z}}(\mathbf{y} | \mathbf{z}), \quad (1.20a)$$

$$p_{\mathbf{X},\mathbf{Y}|\mathbf{Z}}(\mathbf{x}, \mathbf{y} | \mathbf{z}) = p_{\mathbf{X}|\mathbf{Z}}(\mathbf{x} | \mathbf{z}) \cdot p_{\mathbf{Y}|\mathbf{Z}}(\mathbf{y} | \mathbf{z}). \quad (1.20b)$$

Remark 1.20: 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 $X \perp X | X \not\Rightarrow X \perp X$.⁸ Neither does independence imply conditional independence: for example, $X \perp Y \not\Rightarrow X \perp Y | X + Y$.⁹

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 $X \perp Y | Z$.

The independence and conditional independence of events is defined analogously.

1.1.9 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 (as a directed acyclic graph) of the factorization of the joint distribution into a product of conditional distributions. Given the sequence of random variables $\{X_i\}_{i=1}^n$, their joint distribution can be expressed as

$$\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n \mathbb{P}(X_i = x_i | \text{parents}(X_i)) \quad (1.21)$$

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

⁸ $X \perp X | X$ is true trivially.

⁹ Knowing X and $X + Y$ already implies the value of Y , and hence, $X \not\perp Y | X + Y$.

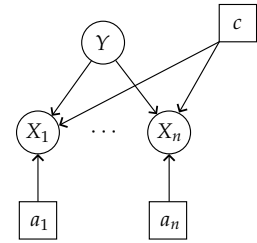


Figure 1.3: Example of a directed graphical model. The random variables X_1, \dots, X_n are mutually independent given the random variable Y . The squared rectangular nodes are used to represent dependencies on parameters c, a_1, \dots, a_n .

a conditional independence of a random variable X with a random variable Y given their parents,¹⁰

$$X \perp Y \mid \text{parents}(X), \text{parents}(Y). \quad (1.22)$$

Thus, eq. (1.21) simply uses the product rule and the conditional independence relationships to factorize the joint distribution.

An example of a directed graphical model is given in fig. 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 fig. 1.4.

1.1.10 Expectation

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

$$\mathbb{E}[\mathbf{X}] \doteq \sum_{\mathbf{x} \in \mathbf{X}(\Omega)} \mathbf{x} \cdot p(\mathbf{x}) \quad (1.23a)$$

$$\mathbb{E}[\mathbf{X}] \doteq \int_{\mathbf{X}(\Omega)} \mathbf{x} \cdot p(\mathbf{x}) d\mathbf{x} \quad (1.23b)$$

if \mathbf{X} is discrete or continuous, respectively.¹¹ A very special and often used property of expectations is their *linearity*. For any random vectors \mathbf{X} and \mathbf{Y} in \mathbb{R}^n and any $\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^m$ we have

$$\mathbb{E}[\mathbf{A}\mathbf{X} + \mathbf{b}] = \mathbf{A}\mathbb{E}[\mathbf{X}] + \mathbf{b} \quad \text{and} \quad (1.24)$$

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

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

$$\mathbb{E}[\mathbf{X}\mathbf{Y}^\top] = \mathbb{E}[\mathbf{X}] \cdot \mathbb{E}[\mathbf{Y}]^\top. \quad (1.26)$$

¹⁰ 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.

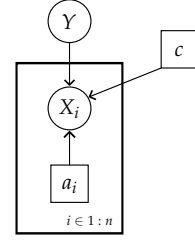


Figure 1.4: The same directed graphical model as in fig. 1.3 using plate notation.

¹¹ In infinite probability spaces, absolute convergence of $\mathbb{E}[\mathbf{X}]$ is necessary for the existence of $\mathbb{E}[\mathbf{X}]$.

Remark 1.21: Notation of summary statistics

Often, when taking the expectation, the brackets around the random vector are omitted: $\mathbb{E}\mathbf{X}$. For notational convenience, we often use the state \mathbf{x} of the random vector \mathbf{X} within the brackets: $\mathbb{E}_{\mathbf{x} \sim \mathbf{X}}[\mathbf{x}]$. If the correspondence between \mathbf{x} and \mathbf{X} is clear from

context, we often abbreviate this expectation with $\mathbb{E}_{\mathbf{X}}[x]$ or $\mathbb{E}[x]$. The analogous notation is used for covariance and variance (cf. section 1.1.11).

Remark 1.22: Gradients of expectations

We will often encounter gradients of expectations, $\nabla_{\theta} \mathbb{E}_{\mathbf{X}}[f(\mathbf{X}, \theta)]$.

Fact 1.23 (Differentiation under the integral sign). *Let $f(x, t)$ be differentiable in x , integrable in t , and be such that*

$$\left| \frac{\partial f(x, t)}{\partial t} \right| \leq g(x)$$

for some integrable function g . Then, if the distribution of \mathbf{X} is not parameterized by t ,¹²

$$\frac{\partial \mathbb{E}[f(\mathbf{X}, t)]}{\partial t} = \mathbb{E} \left[\frac{\partial f(\mathbf{X}, t)}{\partial t} \right]. \quad (1.27)$$

Therefore, if the distribution of \mathbf{X} is not parameterized by θ ,

$$\nabla_{\theta} \mathbb{E}[f(\mathbf{X}, \theta)] = \mathbb{E}[\nabla_{\theta} f(\mathbf{X}, \theta)]. \quad (1.28)$$

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

Fact 1.24 (Law of the unconscious statistician, LOTUS).

$$\mathbb{E}[g(\mathbf{X})] = \int_{\mathbf{X}(\Omega)} g(x) \cdot p(x) dx \quad (1.30)$$

where $g : \mathbf{X}(\Omega) \rightarrow \mathbb{R}^n$ is a “nice” function¹³ 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 discussed in section 1.1.12.

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

$$\mathbb{E}[\mathbf{X} \mid \mathbf{Y} = \mathbf{y}] \doteq \int_{\mathbf{X}(\Omega)} x \cdot p_{\mathbf{X}|\mathbf{Y}}(x \mid \mathbf{y}) dx. \quad (1.31)$$

The definition is analogous for discrete random variables. Note that $\mathbb{E}[\mathbf{X} \mid \mathbf{Y} = \cdot]$ defines a deterministic mapping from \mathbf{y} to $\mathbb{E}[\mathbf{X} \mid \mathbf{Y} = \mathbf{y}]$.

¹² That is, $\mathbb{E}[f(\mathbf{X}, t)] = \int p(x) \cdot f(x, t) dx$ and $p(x)$ does not depend on t .

Gradient The *gradient* of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at a point $x \in \mathbb{R}^n$ is

$$\nabla f(x) \doteq \left[\frac{\partial f(x)}{\partial x(1)} \quad \cdots \quad \frac{\partial f(x)}{\partial x(n)} \right]^{\top}. \quad (1.29)$$

¹³ g being a differentiable function whose inverse is monotonic (this is satisfied in most cases) is sufficient.

Therefore, $\mathbb{E}[\mathbf{X} \mid \mathbf{Y}]$ is itself a random vector:

$$\mathbb{E}[\mathbf{X} \mid \mathbf{Y}](\omega) = \mathbb{E}[\mathbf{X} \mid \mathbf{Y} = \mathbf{Y}(\omega)] \quad (1.32)$$

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

Theorem 1.25 (Tower rule). *Given random vectors \mathbf{X} and \mathbf{Y} , we have*

$$\mathbb{E}_{\mathbf{Y}}[\mathbb{E}_{\mathbf{X}}[\mathbf{X} \mid \mathbf{Y}]] = \mathbb{E}[\mathbf{X}]. \quad (1.33)$$

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

$$\mathbb{E}[\mathbf{X} \mid \mathbf{Y} = \mathbf{y}] = \int \mathbf{x} \cdot p(\mathbf{x} \mid \mathbf{y}) d\mathbf{x}$$

and hence

$$\begin{aligned} \mathbb{E}[\mathbb{E}[\mathbf{X} \mid \mathbf{Y}]] &= \int \left(\int \mathbf{x} \cdot p(\mathbf{x} \mid \mathbf{y}) d\mathbf{x} \right) p(\mathbf{y}) d\mathbf{y} \\ &= \int \int \mathbf{x} \cdot p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} && \text{by definition of conditional densities (1.13)} \\ &= \int \mathbf{x} \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} d\mathbf{x} && \text{by Fubini's theorem} \\ &= \int \mathbf{x} \cdot p(\mathbf{x}) d\mathbf{x} && \text{using the sum rule (1.10)} \\ &= \mathbb{E}[\mathbf{X}]. \end{aligned} \quad \square$$

Note that the tower rule can be used analogously to the law of total probability (1.15) to “condition” on a random vector \mathbf{Y} . For this reason, the tower rule is also known as the *law of total expectation* (LOTE).

Exercise 1.26: Law of total expectation

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

$$\mathbb{E}[\mathbf{X}] = \sum_{i=1}^k \mathbb{E}[\mathbf{X} \mid A_i] \cdot \mathbb{P}(A_i). \quad (1.34)$$

▷ *Solution*

1.1.11 Covariance and Variance

Covariance measures the linear dependence between two random vectors. Given two random vectors \mathbf{X} in \mathbb{R}^n and \mathbf{Y} in \mathbb{R}^m , their *covariance* is defined as

$$\text{Cov}[\mathbf{X}, \mathbf{Y}] \doteq \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{Y} - \mathbb{E}[\mathbf{Y}])^\top] \quad (1.35)$$

$$= \mathbb{E}[\mathbf{XY}^\top] - \mathbb{E}[\mathbf{X}] \cdot \mathbb{E}[\mathbf{Y}]^\top \quad (1.36)$$

$$= \text{Cov}[\mathbf{Y}, \mathbf{X}]^\top \in \mathbb{R}^{n \times m}. \quad (1.37)$$

A direct consequence of the definition of covariance (1.35) is that given linear maps $\mathbf{A} \in \mathbb{R}^{n' \times n}$, $\mathbf{B} \in \mathbb{R}^{m' \times m}$, 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

$$\text{Cov}[\mathbf{AX} + \mathbf{c}, \mathbf{BY} + \mathbf{d}] = \mathbf{A} \text{Cov}[\mathbf{X}, \mathbf{Y}] \mathbf{B}^\top. \quad (1.38)$$

Two random vectors \mathbf{X} and \mathbf{Y} are said to be *uncorrelated* iff $\text{Cov}[\mathbf{X}, \mathbf{Y}] = \mathbf{0}$. Note that if \mathbf{X} and \mathbf{Y} are independent, then eq. (1.26) implies that \mathbf{X} and \mathbf{Y} are uncorrelated. The reverse does not hold in general.

Remark 1.27: Correlation

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

$$\text{Cor}[\mathbf{X}, \mathbf{Y}](i, j) \doteq \frac{\text{Cov}[X_i, Y_j]}{\sqrt{\text{Var}[X_i] \text{Var}[Y_j]}} \in [-1, 1]. \quad (1.39)$$

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

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.¹⁴

The length of a random variable X in this inner product space is called its *standard deviation*,

$$\|X\| = \sqrt{\text{Cov}[X, X]} = \sqrt{\text{Var}[X]} \doteq \sigma[X]. \quad (1.40)$$

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

$$\cos \theta = \frac{\text{Cov}[X, Y]}{\|X\| \|Y\|} = \text{Cor}[X, Y]. \quad (1.41)$$

Thus,

$$\theta = \arccos \text{Cor}[X, Y]. \quad (1.42)$$

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”).

¹⁴ 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, $\text{Cov}[X, Y] = \|X\| \|Y\| \cos \theta$

The variance is a measure of uncertainty about the value of a random vector. Given the random vector \mathbf{X} in \mathbb{R}^n , its *variance* is defined as

$$\text{Var}[\mathbf{X}] \doteq \text{Cov}[\mathbf{X}, \mathbf{X}] \quad (1.43)$$

$$= \mathbb{E} \left[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top \right] \quad (1.44)$$

$$= \mathbb{E} [\mathbf{X}\mathbf{X}^\top] - \mathbb{E}[\mathbf{X}] \cdot \mathbb{E}[\mathbf{X}]^\top \quad (1.45)$$

$$= \begin{bmatrix} \text{Cov}[X_1, X_1] & \cdots & \text{Cov}[X_1, X_n] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_n, X_1] & \cdots & \text{Cov}[X_n, X_n] \end{bmatrix}. \quad (1.46)$$

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.

Exercise 1.28: 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, $\mathbf{x}^\top \Sigma \mathbf{x} \geq 0$ for any $\mathbf{x} \in \mathbb{R}^n$.

▷ *Solution*

Two useful properties of variance are the following:

- It follows from eq. (1.38) that for any linear map $A \in \mathbb{R}^{m \times n}$ and vector $\mathbf{b} \in \mathbb{R}^m$,

$$\text{Var}[A\mathbf{X} + \mathbf{b}] = A\text{Var}[\mathbf{X}]A^\top. \quad (1.47)$$

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

- It follows from the definition of variance (1.44) that for any two random vectors \mathbf{X} and \mathbf{Y} ,

$$\text{Var}[\mathbf{X} + \mathbf{Y}] = \text{Var}[\mathbf{X}] + \text{Var}[\mathbf{Y}] + 2\text{Cov}[\mathbf{X}, \mathbf{Y}]. \quad (1.48)$$

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

$$\text{Var}[\mathbf{X} \mid \mathbf{Y}] \doteq \mathbb{E} \left[(\mathbf{X} - \mathbb{E}[\mathbf{X} \mid \mathbf{Y}])(\mathbf{X} - \mathbb{E}[\mathbf{X} \mid \mathbf{Y}])^\top \mid \mathbf{Y} \right]. \quad (1.49)$$

Intuitively, the conditional variance is the remaining variance when we use $\mathbb{E}[\mathbf{X} \mid \mathbf{Y}]$ to predict \mathbf{X} rather than if we used $\mathbb{E}[\mathbf{X}]$.

Theorem 1.29 (Law of total variance, LOTV).

$$\text{Var}[\mathbf{X}] = \mathbb{E}[\text{Var}[\mathbf{X} \mid \mathbf{Y}]] + \text{Var}[\mathbb{E}[\mathbf{X} \mid \mathbf{Y}]]. \quad (1.50)$$

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 2.3, we will see that both terms have a meaningful characterization in the context of Bayesian learning.

Proof. To simplify the notation, we present only a proof for the univariate setting.¹⁵

$$\begin{aligned} \text{Var}[X] &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 \\ &= \mathbb{E}[\mathbb{E}[X^2 \mid Y]] - \mathbb{E}[\mathbb{E}[X \mid Y]]^2 \\ &= \mathbb{E}[\text{Var}[X \mid Y] + \mathbb{E}[X \mid Y]^2] - \mathbb{E}[\mathbb{E}[X \mid Y]]^2 \\ &= \mathbb{E}[\text{Var}[X \mid Y]] + \left(\mathbb{E}[\mathbb{E}[X \mid Y]^2] - \mathbb{E}[\mathbb{E}[X \mid Y]]^2 \right) \\ &= \mathbb{E}[\text{Var}[X \mid Y]] + \text{Var}[\mathbb{E}[X \mid Y]]. \quad \square \end{aligned}$$

¹⁵ In the *univariate* setting (as opposed to the *multivariate* setting) we consider a single random variable.

by the tower rule (1.33)

by the definition of variance (1.45)

by the definition of variance (1.45)

1.1.12 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

$$P_Y(y) = \mathbb{P}(Y \leq y) = \mathbb{P}(g(X) \leq y) = \mathbb{P}(X \leq g^{-1}(y)). \quad (1.51)$$

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 .¹⁶

There is an analogous change of variables formula for the multivariate setting.

¹⁶ 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).

Fact 1.30 (Change of variables formula). *Let \mathbf{X} be a random vector in \mathbb{R}^n with density p_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 \mathbf{g} as follows:

$$p_{\mathbf{Y}}(\mathbf{y}) = p_{\mathbf{X}}(\mathbf{g}^{-1}(\mathbf{y})) \cdot \left| \det(\mathbf{D}\mathbf{g}^{-1}(\mathbf{y})) \right| \quad (1.52)$$

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

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

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

$$\mathbf{g}_{\#}p_{\mathbf{X}} \doteq p_{\mathbf{Y}} \quad \text{where } \mathbf{Y} = \mathbf{g}(\mathbf{X}). \quad (1.55)$$

1.2 Supervised Learning

In *supervised learning*, 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, $\mathcal{D} \doteq \{(x_i, y_i)\}_{i=1}^n$, where $x_i \in \mathcal{X}$ are *inputs* and the $y_i \in \mathcal{Y}$ are *labels*, and we want to find a function \hat{f} that best-approximates f . It is common to consider a class of functions \mathcal{F} (called the *function class*) that serve as candidates for \hat{f} , where each function is described by some parameters θ .

It is often assumed that the observations are noisy, that is,

$$y_i = f(x_i) + \epsilon_i(x_i), \quad (1.56)$$

where $\epsilon_i(x_i)$ is some independent zero-mean noise, for example Gaussian. Often, the dependence of the noise on x_i is omitted for brevity. When the noise distribution may depend on x_i , the noise is said to be *heteroscedastic* and otherwise the noise is called *homoscedastic*. We will describe this model in greater detail when introducing Bayesian linear regression in chapter 2.

Typically, we differentiate between the task of *regression* where $\mathcal{Y} \doteq \mathbb{R}^k$ (often, we learn a scalar label, so $k = 1$), 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 we defined $\Delta^{\mathcal{C}}$

Jacobian Given a vector-valued function,

$$\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad \mathbf{x} \mapsto \begin{bmatrix} g_1(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{bmatrix},$$

where $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$, the *Jacobian* of \mathbf{g} at $\mathbf{x} \in \mathbb{R}^n$ is

$$\mathbf{D}\mathbf{g}(\mathbf{x}) \doteq \begin{bmatrix} \frac{\partial g_1(\mathbf{x})}{\partial x(1)} & \cdots & \frac{\partial g_1(\mathbf{x})}{\partial x(n)} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_m(\mathbf{x})}{\partial x(1)} & \cdots & \frac{\partial g_m(\mathbf{x})}{\partial x(n)} \end{bmatrix}. \quad (1.53)$$

Observe that for a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$\mathbf{D}f(\mathbf{x}) = \nabla f(\mathbf{x})^{\top}. \quad (1.54)$$

as the set of all probability distributions over the set of classes \mathcal{C} . Recall from remark 1.10 that $\Delta^{\mathcal{C}}$ is an $(m - 1)$ -dimensional convex polytope in the space of probabilities $[0, 1]^m$.

1.2.1 Population Risk and Empirical Risk

We denote the underlying (and unknown) joint distribution over input-label pairs (x, y) by \mathcal{P} . The notion of “approximation error” is typically captured by a loss function $\ell(\hat{y}; y) \in \mathbb{R}$ which is small when the prediction \hat{y} is “close” to the true label y and large otherwise. As our objective is to best-approximate the mappings $(x, y) \sim \mathcal{P}$, we aim to minimize

$$\mathbb{E}_{(x, y) \sim \mathcal{P}} [\ell(\hat{f}(x); y)]. \quad (1.57)$$

This quantity is also called the *population risk*. However, the underlying distribution \mathcal{P} is unknown to us. All that we can work with is the training data for which we assume $\mathcal{D} \stackrel{\text{iid}}{\sim} \mathcal{P}$. It is therefore natural to consider minimizing

$$\frac{1}{n} \sum_{i=1}^n \ell(\hat{f}(x_i); y_i), \quad \mathcal{D} = \{(x_i, y_i)\}_{i=1}^n, \quad (1.58)$$

which is known as the *empirical risk*.

However, selecting \hat{f} by minimizing the empirical risk can be problematic. The reason is that in this case the model \hat{f} and the empirical risk depend on the same data \mathcal{D} , implying that the empirical risk will not be an unbiased estimator of the population risk.¹⁷ This can result in a model which fits the training data too closely (called *overfitting*), and which is failing to generalize to unseen data. We will discuss some solutions to this problem in section 4.4 when covering model selection in the context of Gaussian processes.

¹⁷ In section 1.4.1, we discuss estimator bias in detail.

1.3 Bayesian Learning and Inference

Bayesian learning is the process of updating a Bayesian prior $\mathbb{P}(A)$ to a Bayesian posterior $\mathbb{P}(A | B)$ upon observing B . Thereby, we typically replace A with the parameterization of a model θ , and B with labeled training data $\{(x_i, y_i)\}_{i=1}^n$.

At this point, it is helpful to differentiate between learning and inference. By *learning* (or estimation) we refer to the aforementioned process of learning a model from data. In contrast, by *inference* (or prediction) we refer to the process of using our learned model to predict labels at new inputs.¹⁸

At the core of Bayesian learning is Bayes’ rule.

¹⁸ In the literature, these terms are sometimes used differently, where our learning is referred to as “inference” and our inference is referred to as “prediction”.

Theorem 1.31 (Bayes' rule). *Given random vectors \mathbf{X} in \mathbb{R}^n and \mathbf{Y} in \mathbb{R}^m , we have for any $\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m$,*

$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}. \quad (1.59)$$

Proof. Bayes' rule is a direct consequence of the definition of conditional densities (1.13) and the product rule (1.14). \square

Remark 1.32: Notation of prior and posterior

We generally represent the prior distribution by p and the posterior distribution by $p(\cdot | \mathbf{y})$. Note that sometimes, to highlight their distinction, $p(\cdot)$ is used to refer to the prior.

It is useful to consider the meaning of each term separately:

- $p(\mathbf{x} | \mathbf{y})$ (*posterior*) is the updated belief about \mathbf{x} after observing \mathbf{y} ,
- $p(\mathbf{x})$ (*prior*) is the initial belief about \mathbf{x} ,
- $p(\mathbf{y} | \mathbf{x})$ (*(conditional) likelihood*) describes how likely the observations \mathbf{y} are under a given value \mathbf{x} ,
- $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x})$ (*joint likelihood* or *generative model*) combines prior and likelihood,
- $p(\mathbf{y})$ (*marginal likelihood*) describes how likely the observations \mathbf{y} are across all values of \mathbf{x} .

The marginal likelihood can be computed using the law of total probability (1.15),

$$p(\mathbf{y}) = \int_{\mathbf{X}(\Omega)} p(\mathbf{y} | \mathbf{x})p(\mathbf{x}) d\mathbf{x}, \quad (1.60)$$

or the sum rule. Note, however, that as the marginal likelihood does not depend on \mathbf{x} ,

$$p(\mathbf{x} | \mathbf{y}) \propto p(\mathbf{y} | \mathbf{x})p(\mathbf{x}), \quad (1.61)$$

so we do not need to compute it to optimize for \mathbf{x} (e.g., find the gradient with respect to \mathbf{x}). This is a helpful fact, which we will make use of quite often.

Exercise 1.33: Bayes' rule

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.

▷ *Solution*

Remark 1.34: Improper priors

Not always is it required that the prior $p(x)$ is a valid distribution (i.e., integrates to 1). Consider for example, the “uniform prior” $p(x) \propto \mathbb{1}\{x \in I\}$ where $I \subseteq \mathbb{R}$ is an infinitely large interval. 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.

1.3.1 Using Bayes’ rule for supervised learning

Returning to the task of learning a model from training data, we interpret $p(\theta)$ as the degree of our belief that the model parameterized by θ “describes the data best”. The likelihood describes how likely the training data is under a particular model. Often, we assume that our data was sampled independently, allowing us to write the likelihood as a product:

$$p(y_{1:n} \mid x_{1:n}, \theta) = \prod_{i=1}^n p(y_i \mid x_i, \theta). \quad (1.62)$$

The posterior represents our belief about the best model after seeing the training data. Using Bayes’ rule, we can write it as¹⁹

$$p(\theta \mid x_{1:n}, y_{1:n}) = \frac{1}{Z} p(\theta) \prod_{i=1}^n p(y_i \mid x_i, \theta) \quad \text{where} \quad (1.63)$$

$$Z \doteq \int p(\theta) \prod_{i=1}^n p(y_i \mid x_i, \theta) d\theta \quad (1.64)$$

and again assuming that the training data was sampled independently. We often refer to Z as the *normalizing constant*.

Finally, we can use our learned model for inference (predictions) at a new input x^* by conditioning on θ ,

$$\begin{aligned} p(y^* \mid x^*, x_{1:n}, y_{1:n}) &= \int p(y^*, \theta \mid x^*, x_{1:n}, y_{1:n}) d\theta \\ &= \int p(y^* \mid x^*, \theta) p(\theta \mid x_{1:n}, y_{1:n}) d\theta. \end{aligned} \quad (1.65)$$

Here, the distribution over models $p(\theta \mid x_{1:n}, y_{1:n})$ is called the *posterior* and the distribution over predictions $p(y^* \mid x^*, x_{1:n}, y_{1:n})$ is called the *predictive posterior*.

¹⁹ We generally assume that

$$p(\theta \mid x_{1:n}) = p(\theta).$$

For our purposes, you can think of the inputs $x_{1:n}$ as fixed deterministic parameters.

by the sum rule (1.10)

by the product rule (1.14) and $y^* \perp x_{1:n}, y_{1:n} \mid \theta$

1.3.2 Conjugate Priors

In general, the integrals in the normalizing constant Z (1.64) and the predictive distribution (1.65) cannot be expressed in closed-form.

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 apply the same learning mechanism recursively.

Under some conditions the Gaussian is self-conjugate (cf. section 2.2). That is, if we have a Gaussian prior and a Gaussian likelihood, then our posterior will also be a Gaussian.

Example 1.35: Bernoulli and beta 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

$$\text{Bin}(k; n, \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k} \quad (1.66)$$

and the PDF of the beta distribution,

$$\text{Beta}(\theta; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}, \quad (1.67)$$

where Γ denotes the gamma function. We assume $\theta \sim \text{Beta}(\alpha, \beta)$ and $k | \theta \sim \text{Bin}(n, \theta)$. Let $n_H = k$ be the number of heads and $n_T = n - k$ the number of tails in the binomial trial k . Now,

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

using Bayes' rule (1.59)

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

1.4 Parameter Estimation

A very common — albeit non-Bayesian — approach to learning is to reduce the posterior distribution $p(\theta | x_{1:n}, y_{1:n})$ to a point estimate of θ . We will now take a brief look at this approach, but first define the concept of an estimator more formally.

1.4.1 Estimators

Suppose we are given a collection of independent samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ from some random vector \mathbf{X} . Often, the exact distribution of \mathbf{X} is unknown to us, but we still want to “estimate” some property of this distribution, for example its mean. We denote the property that we aim to estimate from our sample by $\boldsymbol{\theta}$. For example, if our goal is estimating the mean of \mathbf{X} , then $\boldsymbol{\theta} \doteq \mathbb{E}[\mathbf{X}]$.

An *estimator* for a parameter $\boldsymbol{\theta}$ is a random vector \mathbf{U} that is a function of n sample variables $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(n)}$ whose distribution is identical to the distribution of \mathbf{X} . Any concrete sample $\mathbf{x}^{(i)} \sim \mathbf{X}^{(i)}$ yields a concrete estimate $\mathbf{u} \sim \mathbf{U}$ of $\boldsymbol{\theta}$. How these samples are composed is not restricted, but their composition should lead to a “good estimate” of $\boldsymbol{\theta}$. Mainly, there are two measures of goodness of an estimator: its expectation and its variance.

Clearly, we want $\mathbb{E}[\mathbf{U}] = \boldsymbol{\theta}$. Estimators that satisfy this property are called *unbiased*. The *bias*, $\mathbb{E}[\mathbf{U} - \boldsymbol{\theta}]$, of an unbiased estimator is $\mathbf{0}$.

Example 1.36: Estimating expectations

Perhaps the most common estimator is the *sample mean* $\bar{\mathbf{X}}_n$, which is an unbiased estimator for $\mathbb{E}[\mathbf{X}]$. It is defined as

$$\bar{\mathbf{X}}_n \doteq \frac{1}{n} \sum_{i=1}^n \mathbf{X}^{(i)} \quad (1.68)$$

where $\mathbf{X}^{(i)} \stackrel{\text{iid}}{\sim} \mathbf{X}$ are independent samples of the random vector \mathbf{X} .

Using a sample mean to estimate an expectation is often also called *Monte Carlo sampling*, and the resulting estimate is referred to as a *Monte Carlo average*. In section 1.4.2, we will see that the variance of the sample mean is small for reasonably large n .

Example 1.37: Estimating variances

Similarly to the sample mean, the *sample variance* (or *sample covariance matrix*) of a random vector \mathbf{X} ,

$$\mathbf{S}_n^2 \doteq \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}^{(i)} - \bar{\mathbf{X}}_n)(\mathbf{X}^{(i)} - \bar{\mathbf{X}}_n)^\top \quad (1.69)$$

$$= \frac{n}{n-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}^{(i)} \mathbf{X}^{(i)\top} - \bar{\mathbf{X}}_n \bar{\mathbf{X}}_n^\top \right) \quad (1.70)$$

where $\mathbf{X}^{(i)} \stackrel{\text{iid}}{\sim} \mathbf{X}$ are independent samples of \mathbf{X} , is an unbiased

estimator of $\text{Var}[\mathbf{X}]$.

Intuitively, the normalizing factor is $1/n-1$ because S_n^2 depends on the sample mean \bar{X}_n , which is obtained using the same samples. For this reason, S_n^2 has $n-1$ degrees of freedom.²⁰

Exercise 1.38: Sample variance

Let X be a zero-mean random variable. Confirm that the sample variance of X is indeed unbiased.

▷ *Solution*

The second desirable property of an estimator U is that its variance is small.²¹ A common measure for the variance of an estimator of θ is the *mean squared error*,

$$\text{MSE}(U) \doteq \mathbb{E}_U[(U - \theta)^2]. \quad (1.71)$$

We say that an estimator is *consistent* if its mean squared error converges to zero as $n \rightarrow \infty$.

Remark 1.39: The bias-variance tradeoff

Often, bias and variance of an estimator have to be traded as it is difficult to achieve both a small bias and a small variance. This is known as the *bias-variance tradeoff*.

There exists a natural characterization of the mean squared error in terms of bias and variance that highlights this tradeoff. Namely,

$$\text{MSE}(U) = \mathbb{E}_U[(U - \theta)^2] = \text{Var}[U] + \mathbb{E}_U[U - \theta]^2, \quad (1.72)$$

the mean squared error of an estimator can be written as the sum of its variance and its squared bias.

We further say that an estimator U is *sharply concentrated* around θ if for any $\epsilon > 0$,

$$\mathbb{P}(|U - \theta| > \epsilon) \leq \exp(-\Omega(\epsilon)) \doteq \delta, \quad (1.73)$$

where $\Omega(\epsilon)$ denotes the class of functions that grow at least linearly in ϵ .²² Thus, if an estimator is sharply concentrated, its absolute error is bounded by an exponentially quickly decaying error probability δ as ϵ grows.

²⁰ That is, any sample $\mathbf{X}^{(i)}$ can be recovered using all other samples and the sample mean.

²¹ The variance of estimators U is typically studied component wise.

using (1.45) and $\text{Var}_U[U - \theta] = \text{Var}[U]$

²² That is, $h \in \Omega(\epsilon)$ iff $\lim_{\epsilon \rightarrow \infty} \frac{h(\epsilon)}{\epsilon} > 0$. With slight abuse of notation, we force h to be positive (so as to ensure that the argument to the exponential function is negative) whereas in the traditional definition of Landau symbols, h is only required to grow linearly in absolute value.

Remark 1.40: Heavy tails

It is often said that a sharply concentrated estimator U has *small tails*, where “tails” refer to the “ends” of a PDF. Let us examine the difference between a *light-tailed* and a *heavy-tailed* distribution more closely.

Definition 1.41. A distribution P_X is said to have a *heavy (right) tail* iff its *tail distribution function*

$$\bar{P}_X(x) \doteq 1 - P_X(x) = \mathbb{P}(X > x) \quad (1.74)$$

decays slower than that of the exponential distribution, that is,

$$\limsup_{x \rightarrow \infty} \frac{\bar{P}_X(x)}{e^{-\lambda x}} = \infty \quad (1.75)$$

for all $\lambda > 0$. When $\limsup_{x \rightarrow \infty} \frac{\bar{P}_X(x)}{\bar{P}_Y(x)} > 1$, the (right) tail of X is said to be *heavier* than the (right) tail of Y .

It is immediate from the definitions that the distribution of an unbiased estimator is light-tailed if and only if the estimator is sharply concentrated, so both notions are equivalent.

A Gaussian $X \sim \mathcal{N}(0, 1)$ is light-tailed since its tail distribution function is bounded by

$$\bar{P}_X(x) = \int_x^\infty \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \leq \int_x^\infty \frac{t}{x} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt = \frac{e^{-x^2/2}}{\sqrt{2\pi}x}. \quad (1.76)$$

Another light-tailed distribution is the *Laplace distribution*

$$\text{Laplace}(x; \mu, h) \propto \exp\left(-\frac{|x - \mu|}{h}\right) \quad (1.77)$$

with mean μ and so-called length scale parameter h . It places more probability mass directly at the mean at the expense of heavier tails. An example of a heavy-tailed distribution is the *log-normal distribution*. A random variable Y is logarithmically normal distributed with parameters μ and σ^2 iff $\log Y \sim \mathcal{N}(\mu, \sigma^2)$.

For more details and examples of heavy-tailed distributions, refer to “The Fundamentals of Heavy Tails” (Nair et al., 2022).

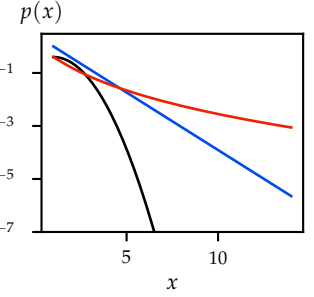


Figure 1.5: Shown are the right tails of the PDFs of a **Gaussian** with mean 1 and variance 1, a **exponential distribution** with mean 1 and parameter $\lambda = 1$, and a **log-normal distribution** with mean 1 and variance 1 on a log-scale.

using $\frac{t}{x} \geq 1$

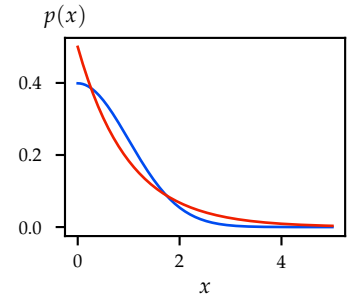


Figure 1.6: Shown are the right tails of the PDFs of a **Laplace distribution** with mean 0 and length scale 1 and a **Gaussian** with mean 0 and variance 1.

1.4.2 Mean Estimation and Concentration Inequalities

We have seen that we desire two properties in estimators: namely, (1) that they are unbiased; and (2) that their variance is small.²³ For

²³ That is, they are consistent and, ideally, their variance converges quickly.

estimating expectations with the sample mean \bar{X}_n which is also known as a *Monte Carlo approximation*, we will see that both properties follow from standard results in probability theory.

It is immediate from the definition of the sample mean (1.68) that it is an unbiased estimator for $\mathbb{E}[X]$. We will now see that the sample mean is also consistent, and even better, converges with exponentially decaying error probability.

First, let us recall the common notions of convergence of a sequence of random variables X_n .

Definition 1.42 (Convergence of random variables). Let $\{X_n\}_{n \in \mathbb{N}}$ be a sequence of random variables and X another random variable. We say that,

1. X_n converges to X *almost surely* (also called *convergence with probability 1*) if

$$\mathbb{P}\left(\left\{\omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)\right\}\right) = 1, \quad (1.78)$$

and we write $X_n \xrightarrow{\text{a.s.}} X$ as $n \rightarrow \infty$.

2. X_n converges to X *in probability* if for any $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \epsilon) = 0, \quad (1.79)$$

and we write $X_n \xrightarrow{\mathbb{P}} X$ as $n \rightarrow \infty$.

3. X_n converges to X *in distribution* if for all points $x \in X(\Omega)$ at which P_X is continuous,

$$\lim_{n \rightarrow \infty} P_{X_n}(x) = P_X(x), \quad (1.80)$$

and we write $X_n \xrightarrow{\mathcal{D}} X$ as $n \rightarrow \infty$.

It can be shown that as $n \rightarrow \infty$,

$$X_n \xrightarrow{\text{a.s.}} X \implies X_n \xrightarrow{\mathbb{P}} X \implies X_n \xrightarrow{\mathcal{D}} X. \quad (1.81)$$

In the following, we will use some classical concentration inequalities.

Exercise 1.43: Simple concentration inequalities

1. Prove *Markov's inequality* which says that if X is a non-negative random variable, then for any $\epsilon > 0$,

$$\mathbb{P}(X \geq \epsilon) \leq \frac{\mathbb{E}X}{\epsilon}. \quad (1.82)$$

2. Prove *Chebyshev's inequality* which says that if X is a random

variable with finite and non-zero variance, then for any $\epsilon > 0$,

$$\mathbb{P}(|X - \mathbb{E}X| \geq \epsilon) \leq \frac{\text{Var}X}{\epsilon^2}. \quad (1.83)$$

▷ *Solution*

The law of large numbers is a classical result in statistics and directly implies that the sample mean is consistent.

Fact 1.44 (Strong law of large numbers, SLLN). *Given the random variable $X : \Omega \rightarrow \mathbb{R}$ with independent samples $x_i \sim X$ and finite variance. Let f be a “nice” function and $Y \doteq f(X)$. We have,*

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n f(x_i) \xrightarrow{\text{a.s.}} \mathbb{E}[f(X)], \quad (1.84)$$

as $n \rightarrow \infty$.

Exercise 1.45: Weak law of large numbers

Prove this weaker version of the SLLN which is known as the *weak law of large numbers* (WLLN).

Let $X : \Omega \rightarrow \mathbb{R}$ be a random variable with finite variance. Then,

$$\bar{X}_n \xrightarrow{\mathbb{P}} \mathbb{E}X. \quad (1.85)$$

▷ *Solution*

Under more restrictive assumptions it is even possible to show that the sample mean \bar{X}_n is sharply concentrated around the mean $\mathbb{E}X$, which is a much stronger property than consistency. One such assumption is provided by the notion of a sub-Gaussian random variable.

Definition 1.46 (Sub-Gaussian random variable). A random variable $X : \Omega \rightarrow \mathbb{R}$ is called σ -sub-Gaussian for $\sigma > 0$ if for all $\lambda \in \mathbb{R}$,

$$\mathbb{E} \left[e^{\lambda(X - \mathbb{E}X)} \right] \leq \exp \left(\frac{\sigma^2 \lambda^2}{2} \right). \quad (1.86)$$

Intuitively, a random variable is sub-Gaussian if its tail probabilities decay at least as fast as those of a Gaussian. The class of sub-Gaussian random variables is therefore a subclass of the light-tailed distributions discussed in remark 1.40.

Example 1.47: Examples of sub-Gaussian random variables

- If a random variable X is σ -sub-Gaussian, then so is $-X$.
- If $X \sim \mathcal{N}(\mu, \sigma^2)$ then a simple calculation shows that

$$\varphi_X(\lambda) \doteq \mathbb{E}[e^{\lambda X}] = \exp\left(\mu\lambda + \frac{\sigma^2\lambda^2}{2}\right) \quad \text{for any } \lambda \in \mathbb{R} \quad (1.87)$$

which is called the “moment-generating function” of the normal distribution, and implying that X is σ -sub-Gaussian.

- If the random variable $X : \Omega \rightarrow \mathbb{R}$ satisfies $a \leq X \leq b$ with probability 1 then X is $\frac{b-a}{2}$ -sub-Gaussian. This is also known as *Hoeffding’s lemma*.

Theorem 1.48 (Hoeffding’s inequality). *Let $X : \Omega \rightarrow \mathbb{R}$ be a σ -sub-Gaussian random variable. Then,*

$$\mathbb{P}(|\bar{X}_n - \mathbb{E}X| \geq \epsilon) \leq 2 \exp\left(-\frac{n\epsilon^2}{2\sigma^2}\right). \quad (1.88)$$

In words, the absolute error of the expectation estimate is bounded by an exponentially quickly decaying error probability δ . Solving for n , we obtain that for

$$n \geq \frac{2\sigma^2}{\epsilon^2} \log \frac{2}{\delta} \quad (1.89)$$

the probability that the absolute error is greater than ϵ is at most δ .

Proof. Let $S_n \doteq n\bar{X}_n = X_1 + \dots + X_n$. We have for any $\lambda, \epsilon > 0$ that

$$\begin{aligned} \mathbb{P}(\bar{X}_n - \mathbb{E}X \geq \epsilon) &= \mathbb{P}(S_n - \mathbb{E}S_n \geq n\epsilon) \\ &= \mathbb{P}\left(e^{\lambda(S_n - \mathbb{E}S_n)} \geq e^{n\epsilon\lambda}\right) && \text{using that } z \mapsto e^{\lambda z} \text{ is increasing} \\ &\leq e^{-n\epsilon\lambda} \mathbb{E}[e^{\lambda(S_n - \mathbb{E}S_n)}] && \text{using Markov’s inequality (1.82)} \\ &= e^{-n\epsilon\lambda} \prod_{i=1}^n \mathbb{E}[e^{\lambda(X_i - \mathbb{E}X)}] && \text{using independence of the } X_i \\ &\leq e^{-n\epsilon\lambda} \prod_{i=1}^n e^{\sigma^2\lambda^2/2} && \text{using the characterizing property of a } \sigma\text{-sub-Gaussian random variable (1.86)} \\ &= \exp\left(-n\epsilon\lambda + \frac{n\sigma^2\lambda^2}{2}\right). \end{aligned}$$

Minimizing the expression with respect to λ , we set $\lambda = \epsilon/\sigma^2$, and obtain

$$\mathbb{P}(\bar{X}_n - \mathbb{E}X \geq \epsilon) \leq \min_{\lambda > 0} \left\{ \exp\left(-n\epsilon\lambda + \frac{n\sigma^2\lambda^2}{2}\right) \right\} = \exp\left(-\frac{n\epsilon^2}{2\sigma^2}\right).$$

The theorem then follows from

$$\mathbb{P}(|\bar{X}_n - \mathbb{E}X| \geq \epsilon) = \mathbb{P}(\bar{X}_n - \mathbb{E}X \geq \epsilon) + \mathbb{P}(\bar{X}_n - \mathbb{E}X \leq -\epsilon)$$

and noting that the second term can be bounded analogously to the first term by considering the random variables $-X_1, \dots, -X_n$. \square

The law of large numbers and Hoeffding's inequality show that we can estimate $\mathbb{E}[f(X)]$ very precisely with “few” samples using a sample mean. Crucially, we require that the samples x_i are from the same distribution as the random variable X and that the samples are independent.

1.4.3 Maximum Likelihood Estimation

Let us return to finding point estimates of our model. Intuitively, the “best model” should have the property that the training data is as likely as possible under this model when compared to all other models. This is precisely the *maximum likelihood estimate* (or MLE):

$$\hat{\theta}_{\text{MLE}} \doteq \arg \max_{\theta} p(y_{1:n} \mid x_{1:n}, \theta) \quad (1.90)$$

$$= \arg \max_{\theta} \prod_{i=1}^n p(y_i \mid x_i, \theta)$$

using the independence of the training data (1.62)

$$= \arg \max_{\theta} \sum_{i=1}^n \log p(y_i \mid x_i, \theta). \quad (1.91)$$

using that log is a monotonic function

Taking the logarithm of the likelihood is a standard technique. The resulting $\log p(y_{1:n} \mid x_{1:n}, \theta)$ is called *log-likelihood*, and its negative is known as the *negative log-likelihood* $\ell_{\text{nll}}(\theta; \mathcal{D})$,

$$= \arg \min_{\theta} \ell_{\text{nll}}(\theta; \mathcal{D}). \quad (1.92)$$

Example 1.49: MLE of Bernoulli random variables

Let us consider an i.i.d. sample $x_{1:n}$ of $\text{Bern}(p)$ distributed random variables. We want to estimate the parameter p using a MLE. We have,

$$\begin{aligned} \hat{p}_{\text{MLE}} &= \arg \max_p \mathbb{P}(x_{1:n} \mid p) \\ &= \arg \max_p \sum_{i=1}^n \log \mathbb{P}(x_i \mid p) \\ &= \arg \max_p \sum_{i=1}^n \log p^{x_i} (1-p)^{1-x_i} \end{aligned}$$

using the Bernoulli PMF, see example 1.9

$$= \arg \max_p \sum_{i=1}^n x_i \log p + (1 - x_i) \log(1 - p).$$

Computing the first derivative with respect to p , we see that the objective is maximized by

$$= \frac{1}{n} \sum_{i=1}^n x_i. \quad (1.93)$$

Thus, the maximum likelihood estimate for the parameter p of $X \sim \text{Bern}(p)$ coincides with the sample mean \bar{X} .

1.4.4 Maximum A Posteriori Estimation

Often, the MLE overfits to the training data. The danger of overfitting can be reduced by taking a “more Bayesian” approach and maximizing over the entire posterior distribution instead of only maximizing the likelihood. This results in the *maximum a posteriori estimate* (or MAP estimate):

$$\hat{\theta}_{\text{MAP}} \doteq \arg \max_{\theta} p(\theta \mid x_{1:n}, y_{1:n}) \quad (1.94)$$

$$= \arg \max_{\theta} p(y_{1:n} \mid x_{1:n}, \theta) \cdot p(\theta) \quad (1.95)$$

$$= \arg \max_{\theta} \log p(\theta) + \sum_{i=1}^n \log p(y_i \mid x_i, \theta). \quad (1.96)$$

$$= \arg \min_{\theta} \underbrace{-\log p(\theta)}_{\text{regularization}} + \underbrace{\ell_{\text{nl}}(\theta; \mathcal{D})}_{\text{quality of fit}}. \quad (1.97)$$

Here, the *log-prior* $\log p(\theta)$ acts as a regularizer. For example,

- if $p(\theta) = \mathcal{N}(\theta; \mathbf{0}, \sigma_p^2 I)$ then $-\log p(\theta) = \frac{\sigma_p^2}{2} \|\theta\|_2^2$,
- if $p(\theta) = \text{Laplace}(\theta; \mathbf{0}, h)$ then $-\log p(\theta) = \frac{h}{2} \|\theta\|_1^2$, and
- for a uniform prior (i.e., a prior that is independent of θ), MAP estimation reduces to likelihood maximization.

Note that in case the posterior is Gaussian, the MAP estimate (i.e., mode of the posterior distribution) corresponds to the mean.

1.5 Optimization

Finding parameter estimates is one of the many examples where we seek to minimize some function ℓ .²⁴ The field of optimization has a rich history, which we will not explore in much detail here. What will be important for us is that given that the function to be optimized (called the *objective* or *loss*) fulfills certain smoothness properties, optimization is a well-understood problem and can often be solved exactly

by Bayes’ rule (1.59)

using that log is a monotonic function

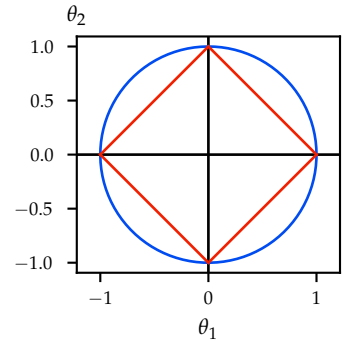


Figure 1.7: Level sets of **Gaussian** and **Laplace** regularization. It can be seen that Laplace regularization is more effective in encouraging sparse solutions (that is, solutions where many components are set to exactly 0).

²⁴ W.l.o.g. we assume that we want to minimize ℓ . If we wanted to maximize the objective, we can simply minimize its negation.

(e.g., when the objective is convex) or “approximately” when the objective is non-convex. In fact, we will see that it is often advantageous to frame problems as optimization problems when suitable because the machinery to solve these problems is so extensive.

Readings

For a more thorough reminder of optimization methods, read chapter 7 of “Mathematics for machine learning” (Deisenroth et al., 2020).

1.5.1 Stationary Points

In this section, we derive some basic facts about unconstrained optimization problems. Given some function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, we want to find

$$\min_{x \in \mathbb{R}^n} f(x). \quad (1.98)$$

We say that a point $x^* \in \mathbb{R}^n$ is a (global) *optimum* of f if $f(x^*) \leq f(x)$ for any $x \in \mathbb{R}^n$.

Consider the more general problem of minimizing f over some subset $S \subseteq \mathbb{R}^n$, that is, to minimize the function $f_S : S \rightarrow \mathbb{R}, x \mapsto f(x)$. If there exists some open subset $S \subseteq \mathbb{R}^n$ including x^* such that x^* is optimal with respect to the function f_S , then x^* is called a *local optimum* of f .

Remark 1.50: Differentiability

We will generally assume that f is continuously (Fréchet) differentiable on \mathbb{R}^n . That is, at any point $x \in \mathbb{R}^n$, there exists $\nabla f(x)$ such that for any $\delta \in \mathbb{R}^n$,

$$f(x + \delta) = f(x) + \nabla f(x)^\top \delta + o(\|\delta\|_2), \quad (1.99)$$

where $\lim_{\delta \rightarrow 0} \frac{o(\|\delta\|_2)}{\|\delta\|_2} = 0$, and $\nabla f(x)$ is continuous on \mathbb{R}^n .

Equation (1.99) is also called a *first-order expansion* of f at x .

Definition 1.51 (Stationary point). Given a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, a point $x \in \mathbb{R}^n$ where $\nabla f(x) = \mathbf{0}$ is called a *stationary point* of f .

Being a stationary point is not sufficient for optimality. Take for example the point $x = 0$ of $f(x) = x^3$. Such a point that is stationary but not (locally) optimal is called a *saddle point*.

Theorem 1.52 (First-order optimality condition). *If $x \in \mathbb{R}^n$ is a local extremum of a differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, then $\nabla f(x) = \mathbf{0}$.*

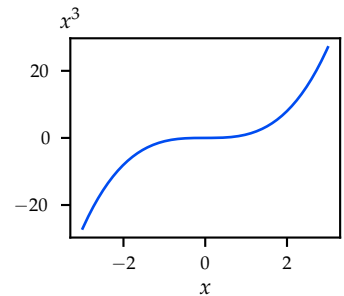


Figure 1.8: Example of a saddle point at $x = 0$.

Proof. Assume x is a local minimum of f . Then, for all $d \in \mathbb{R}^n$ and for all small enough $\lambda \in \mathbb{R}$, we have $f(x) \leq f(x + \lambda d)$, so

$$\begin{aligned} 0 &\leq f(x + \lambda d) - f(x) \\ &= \lambda \nabla f(x)^\top d + o(\lambda \|d\|_2). \end{aligned}$$

Dividing by λ and taking the limit $\lambda \rightarrow 0$, we obtain

$$0 \leq \nabla f(x)^\top d + \lim_{\lambda \rightarrow 0} \frac{o(\lambda \|d\|_2)}{\lambda} = \nabla f(x)^\top d.$$

Take $d \doteq -\nabla f(x)$. Then, $0 \leq -\|\nabla f(x)\|_2^2$, so $\nabla f(x) = 0$. \square

1.5.2 Convexity

Convex functions are a subclass of functions where finding global minima is substantially easier than for general functions.

Definition 1.53 (Convex function). A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is *convex* iff

$$\forall x, y \in \mathbb{R}^n : \forall \theta \in [0, 1] : f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y). \quad (1.100)$$

That is, any line between two points on f lies “above” f . If the inequality of eq. (1.100) is strict, we say that f is *strictly convex*.

If the function f is convex, we say that the function $-f$ is *concave*. The above is also known as the *0th-order characterization of convexity*.

Theorem 1.54 (First-order characterization of convexity). Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable, then f is convex iff

$$\forall x, y \in \mathbb{R}^n : f(y) \geq f(x) + \nabla f(x)^\top (y - x). \quad (1.101)$$

Observe that the right-hand side of the inequality is an affine function with slope $\nabla f(x)$ based at $f(x)$.

Proof. In the following, we will make use of directional derivatives.

Exercise 1.55: Directional derivatives

Given a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ that is differentiable at $x \in \mathbb{R}^n$, the *directional derivative* of f at x in the direction $d \in \mathbb{R}^n$ is

$$Df(x)[d] \doteq \lim_{\lambda \rightarrow 0} \frac{f(x + \lambda d) - f(x)}{\lambda}. \quad (1.102)$$

Show that

$$Df(x)[d] = \nabla f(x)^\top d. \quad (1.103)$$

using a first-order expansion of f around x

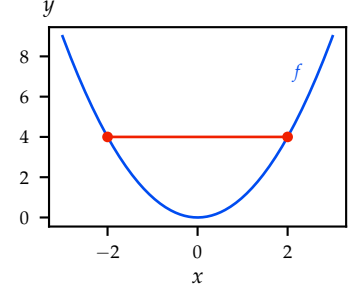


Figure 1.9: Example of a convex function. Any line between two points on f lies “above” f .

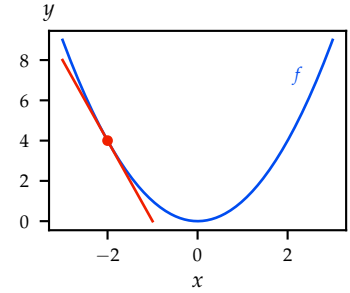


Figure 1.10: The first-order characterization characterizes convexity in terms of affine lower bounds. Shown is an affine lower bound based at $x = -2$.

Hint: Consider a first-order expansion of f at \mathbf{x} in direction $\lambda \mathbf{d}$.

▷ *Solution*

- “ \Rightarrow ”: Fix any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. As f is convex,

$$f((1 - \theta)\mathbf{x} + \theta\mathbf{y}) \leq (1 - \theta)f(\mathbf{x}) + \theta f(\mathbf{y}),$$

for all $\theta \in [0, 1]$. We can rearrange to

$$f(\underbrace{(1 - \theta)\mathbf{x} + \theta\mathbf{y}}_{\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})}) - f(\mathbf{x}) \leq \theta(f(\mathbf{y}) - f(\mathbf{x})).$$

Dividing by θ yields,

$$\frac{f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\theta} \leq f(\mathbf{y}) - f(\mathbf{x}).$$

Taking the limit $\theta \rightarrow 0$ on both sides gives the directional derivative at \mathbf{x} in direction $\mathbf{y} - \mathbf{x}$,

$$\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) = Df(\mathbf{x})[\mathbf{y} - \mathbf{x}] \leq f(\mathbf{y}) - f(\mathbf{x}).$$

- “ \Leftarrow ”: Fix any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and let $\mathbf{z} \doteq \theta\mathbf{y} + (1 - \theta)\mathbf{x}$. We have,

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}), \quad \text{and} \\ f(\mathbf{x}) &\geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}). \end{aligned}$$

We also have $\mathbf{y} - \mathbf{z} = (1 - \theta)(\mathbf{y} - \mathbf{x})$ and $\mathbf{x} - \mathbf{z} = \theta(\mathbf{x} - \mathbf{y})$. Hence,

$$\begin{aligned} \theta f(\mathbf{y}) + (1 - \theta)f(\mathbf{x}) &\geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top \underbrace{(\theta(\mathbf{y} - \mathbf{z}) + (1 - \theta)(\mathbf{x} - \mathbf{z}))}_0 \\ &= f(\theta\mathbf{y} + (1 - \theta)\mathbf{x}). \end{aligned} \quad \square$$

Theorem 1.56. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex and differentiable function. Then, if $\mathbf{x}^* \in \mathbb{R}^n$ is a stationary point of f , \mathbf{x}^* is a global minimum of f .*

Proof. By the first-order characterization of convexity, we have for any $\mathbf{y} \in \mathbb{R}^n$,

$$f(\mathbf{y}) \geq f(\mathbf{x}^*) + \underbrace{\nabla f(\mathbf{x}^*)^\top}_{0} (\mathbf{y} - \mathbf{x}^*) = f(\mathbf{x}^*). \quad \square$$

Generally, the main difficulty in solving convex optimization problems lies therefore in finding stationary points (or points that are sufficiently close to stationary points).

Remark 1.57: Second-order characterization of convexity

We say that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously (Fréchet) differentiable iff for any point $\mathbf{x} \in \mathbb{R}^n$, there exists $\nabla f(\mathbf{x})$ and $\mathbf{H}_f(\mathbf{x})$ such that for any $\delta \in \mathbb{R}^n$,

$$f(\mathbf{x} + \delta) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \delta + \frac{1}{2} \delta^\top \mathbf{H}_f(\mathbf{x}) \delta + o(\|\delta\|_2^2), \quad (1.104)$$

where $\lim_{\delta \rightarrow 0} \frac{o(\|\delta\|_2^2)}{\|\delta\|_2^2} = 0$, and $\nabla f(\mathbf{x})$ and $\mathbf{H}_f(\mathbf{x})$ are continuous on \mathbb{R}^n . Equation (1.104) is also called a *second-order expansion* of f at \mathbf{x} .

Twice continuously differentiable functions admit a natural characterization of convexity in terms of the Hessian.

Fact 1.58 (Second-order characterization of convexity). *Consider a twice continuously differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Then, f is convex iff $\mathbf{H}_f(\mathbf{x})$ is positive semi-definite for all $\mathbf{x} \in \mathbb{R}^n$.*

It follows that f is concave iff $\mathbf{H}_f(\mathbf{x})$ is negative semi-definite for all $\mathbf{x} \in \mathbb{R}^n$.

1.5.3 Stochastic Gradient Descent

In this course, we primarily employ so-called first-order methods, which rely on (estimates of) the gradient of the objective function to determine a direction of local improvement. The main idea behind these methods is to repeatedly take a step in the opposite direction of the gradient scaled by a learning rate η_t , which may depend on the current iteration t .

We will often want to minimize a stochastic optimization objective

$$L(\boldsymbol{\theta}) \doteq \mathbb{E}_{\mathbf{x} \sim p}[\ell(\boldsymbol{\theta}; \mathbf{x})] \quad (1.107)$$

where ℓ and its gradient $\nabla \ell$ are known.

Based on our discussion in previous subsections, it is a natural first step to look for stationary points of L , that is, the roots of ∇L .

Fact 1.59 (Robbins-Monro (RM) algorithm). *Let $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be an unknown function of which we want to find a root and suppose that we have access to unbiased noisy observations $\mathbf{G}(\boldsymbol{\theta})$ of $g(\boldsymbol{\theta})$. The scheme*

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t - \eta_t \mathbf{g}^{(t)}(\boldsymbol{\theta}_t), \quad (1.108)$$

where $\mathbf{g}^{(t)}(\boldsymbol{\theta}_t) \sim \mathbf{G}(\boldsymbol{\theta}_t)$ are independent and unbiased estimates of

Hessian The Hessian of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at a point $\mathbf{x} \in \mathbb{R}^n$ is

$$\mathbf{H}_f(\mathbf{x}) \doteq \mathbf{H}f(\mathbf{x})$$

$$\doteq \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x(1) \partial x(1)} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x(1) \partial x(n)} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x(n) \partial x(1)} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x(n) \partial x(n)} \end{bmatrix} \quad (1.105)$$

$$= (\mathbf{D}\nabla f(\mathbf{x}))^\top \quad (1.106)$$

$$\in \mathbb{R}^{n \times n}.$$

Thus, a Hessian captures the curvature of f . If the Hessian of f is positive definite at a point \mathbf{x} , then f is “curved up around \mathbf{x} ”.

Hessians are symmetric when the second partial derivatives are continuous, due to Schwartz’s theorem.

$g(\theta_t)$, is known as the Robbins-Monro algorithm.

It can be shown that if the sequence of learning rates $\{\eta_t\}_{t=0}^{\infty}$ is chosen such that the Robbins-Monro conditions,²⁵

$$\eta_t \geq 0 \quad \forall t, \quad \sum_{t=0}^{\infty} \eta_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} \eta_t^2 < \infty, \quad (1.109)$$

and additional regularity assumptions are satisfied,²⁶ then we have that

$$g(\theta_t) \xrightarrow{\text{a.s.}} \mathbf{0} \quad \text{as } t \rightarrow \infty. \quad (1.110)$$

That is, the RM-algorithm converges to a root almost surely.²⁷

Using Robbins-Monro to find a root of ∇L is known as *stochastic gradient descent*. In particular, when ℓ is convex and the RM-conditions are satisfied, Robbins-Monro converges to a stationary point (and hence, a global minimum) of L almost surely.

Moreover, it can be shown that for general ℓ and satisfied RM-conditions, stochastic gradient descent converges to a local minimum almost surely. Intuitively, the randomness in the gradient estimates allows the algorithm to “jump past” saddle points.

A commonly used strategy to obtain unbiased gradient estimates is to take the sample mean of the gradient with respect to some set of samples B (also called a *batch*) as is shown in alg. 1.60.

Algorithm 1.60: Stochastic gradient descent, SGD

```

1 initialize  $\theta$ 
2 repeat
3   draw mini-batch  $B \doteq \{x^{(1)}, \dots, x^{(m)}\}, x^{(i)} \sim p$ 
4    $\theta \leftarrow \theta - \eta_t \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \ell(\theta; x^{(i)})$ 
5 until converged
```

Example 1.61: Minimizing training loss

A common application of SGD in the context of machine learning is the following: $p \doteq \text{Unif}(\{x_1, \dots, x_n\})$ is a uniform distribution over the training inputs, yielding the objective

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; x_i), \quad (1.111)$$

where $\ell(\theta; x_i)$ is the loss of a model parameterized by θ for train-

²⁵ $\eta_t = 1/t$ is an example of a learning rate satisfying the RM-conditions.

²⁶ For more details, refer to “Online learning and stochastic approximations” (Bottou, 1998).

²⁷ See eq. (1.78) for a definition of almost sure convergence of a sequence of random variables.

ing input x_i and m is a fixed *batch size*. Here, computing the gradients exactly (i.e., choosing $m = n$) is expensive when the size of the training set is large.

A commonly used alternative to sampling each batch independently, is to split the training set into equally sized batches and perform a gradient descent step with respect to each of them. Typically, the optimum is not found after a single pass through the training data, so the same procedure is repeated multiple times. One such iteration is called an *epoch*.

Remark 1.62: Regularization via weight decay

A common technique to improve the “stability” of minima found through gradient-based methods is to “regularize” the loss by adding an explicit bias favoring “simpler” solutions. That is, given a loss ℓ measuring the quality of fit, we consider the regularized loss

$$\ell'(\theta; x) \doteq \underbrace{\ell(\theta; x)}_{\text{quality of fit}} + \underbrace{r(\theta)}_{\text{regularization}} \quad (1.112)$$

where $r(\theta)$ is large for “complex” and small for “simple” choices of θ , respectively. A common choice is $r(\theta) = \frac{\lambda}{2} \|\theta\|_2^2$ for some $\lambda > 0$, which is known as ℓ_2 -regularization.

Recall from section 1.4.4 that in the context of likelihood maximization, this choice of r corresponds to imposing the Gaussian prior $\mathcal{N}(\mathbf{0}, \lambda I)$, and finding the MAP estimate. Imposing, for example, a Laplace prior leads to ℓ_1 -regularization.

Using ℓ_2 -regularization, we obtain for the gradient,

$$\nabla_{\theta} \ell'(\theta; x) = \nabla_{\theta} \ell(\theta; x) + \lambda \theta. \quad (1.113)$$

Thus, a gradient descent step changes to

$$\theta \leftarrow (1 - \lambda \eta_t) \theta - \eta_t \nabla_{\theta} \ell(\theta; x). \quad (1.114)$$

That is, in each step, θ decays towards zero at the rate $\lambda \eta_t$. This regularization method is also called *weight decay*.

Remark 1.63: Adaptive learning rates and momentum

Commonly, in optimization, constant learning rates are used to accelerate mixing. The performance can be further improved by taking an adaptive gradient step with respect to the geometry of the cost function, which is done by commonly used algorithms such as *Adagrad* and *Adam*. The methods employed by these algorithms are known as *adaptive learning rates* and *momentum*.

For more details on momentum read section 7.1.2 of “Mathematics for machine learning” (Deisenroth et al., 2020) and for an overview of the aforementioned optimization algorithms refer to “An overview of gradient descent optimization algorithms” (Ruder, 2016).

We will have a brief look at other approaches to stochastic optimization in section 6.3.

1.6 Multivariate Normal Distribution

Using general distributions for learning and inference is computationally very expensive when the number of dimensions is large — even in the discrete setting. 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. fig. 1.11). One strategy to get around this computational blowup is to restrict the class of distributions.

It turns out that Gaussians have many extremely useful properties: they have a compact representation and — as we will see in later chapters — they allow for closed-form learning and inference.

In eq. (1.6), 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

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \doteq \frac{1}{\sqrt{\det(2\pi\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (1.115)$$

where $\boldsymbol{\mu} \in \mathbb{R}^n$ is the mean, $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$ the covariance matrix, and $\boldsymbol{\Lambda} \doteq \boldsymbol{\Sigma}^{-1}$ the so-called *precision matrix*. \mathbf{X} is also called a *Gaussian random vector* (GRV). $\mathcal{N}(\mathbf{0}, \mathbf{I})$ is the multivariate *standard normal distribution*. In this definition, we assume that the covariance matrix $\boldsymbol{\Sigma}$ is invertible,

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
\dots	\dots	\dots	\dots	
1	\dots	1	1	0.0003

Figure 1.11: A table representing a joint distribution of n binary random variables. The number of parameters is $2^n - 1$.

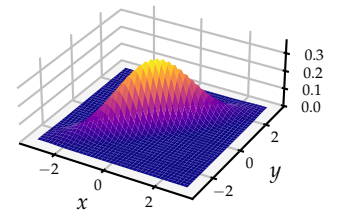
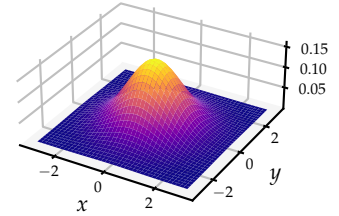


Figure 1.12: Shown are two-dimensional Gaussians with mean $\mathbf{0}$ and covariance matrices

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

respectively.

i.e., does not have the eigenvalue 0.

Remark 1.64: Zero eigenvalues of covariance matrices

It can be shown that a covariance matrix has a zero eigenvalue iff there exists a deterministic affine²⁸ relationship between some of the variables in the joint distribution.

As we have already seen that a covariance matrix does not have negative eigenvalues,²⁹ this ensures that Σ and Λ are positive definite.³⁰

We call a Gaussian *isotropic* iff its covariance matrix is of the form $\Sigma = \sigma^2 I$ for some $\sigma^2 \in \mathbb{R}$. In this case, the sublevel sets of the PDF are perfect spheres.

Note that a Gaussian can be represented using only $\mathcal{O}(n^2)$ parameters. In the case of a diagonal covariance matrix (which as we will see corresponds to independent univariate Gaussians) we just need $\mathcal{O}(n)$ parameters.

Exercise 1.65: 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

$$\mathcal{N}(x; \mu, \Sigma) \propto \mathcal{N}(x; \mu_1, \Sigma_1) \cdot \mathcal{N}(x; \mu_2, \Sigma_2) \quad (1.116)$$

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.

▷ *Solution*

1.6.1 Properties of Gaussians

In this section, we study some of the special properties of Gaussians.

Theorem 1.66 (Marginal and conditional distribution). *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,*

$$\mathbf{X}_A \sim \mathcal{N}(\mu_A, \Sigma_{AA}), \quad (1.117)$$

and that for any such conditional distribution,

$$\mathbf{X}_A \mid \mathbf{X}_B = \mathbf{x}_B \sim \mathcal{N}(\mu_{A|B}, \Sigma_{A|B}) \quad \text{where} \quad (1.118)$$

$$\mu_{A|B} \doteq \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (\mathbf{x}_B - \mu_B), \quad (1.119)$$

²⁸ An *affine transformation* is a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m, x \mapsto Ax + b$ where $A \in \mathbb{R}^{m \times n}$ is an invertible linear map and $b \in \mathbb{R}^m$ is called a translation.

²⁹ see exercise 1.28

³⁰ The inverse of a positive definite matrix is also positive definite.

By μ_A we denote $[\mu_{i_1}, \dots, \mu_{i_k}]$ where $A = \{i_1, \dots, i_k\}$. Σ_{AA} is defined analogously.

Here, μ_A characterizes the prior belief and $\Sigma_{AB} \Sigma_{BB}^{-1} (\mathbf{x}_B - \mu_B)$ represents “how different” \mathbf{x}_B is from what is expected.

$$\Sigma_{A|B} \doteq \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}. \quad (1.120)$$

Observe, that the variance can only shrink! Moreover, how much the variance is reduced depends purely on *where* the observations are made (e.g., the choice of B) but not on *what* the observations are. Also observe that $\mu_{A|B}$ depends affinely on μ_B .

Exercise 1.67: Marginal / conditional distribution of a Gaussian

Prove theorem 1.66. 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).

▷ *Solution*

The following theorem states that the independence of Gaussians is characterized by their correlation.

Theorem 1.68 (Independence of Gaussians). *Two jointly Gaussian random vectors, \mathbf{X} and \mathbf{Y} , are independent if and only if \mathbf{X} and \mathbf{Y} are uncorrelated.*

Proof. Recall that independence of any random vectors \mathbf{X} and \mathbf{Y} implies that they are uncorrelated,³¹ the converse implication is a special property of Gaussian random vectors.

³¹ This is because the expectation of their product factorizes, see eq. (1.26).

Consider the joint Gaussian random vector $\mathbf{Z} \doteq [\mathbf{X} \ \mathbf{Y}] \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ and assume that $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_X, \Sigma_X)$ and $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}_Y, \Sigma_Y)$ are uncorrelated.

Then, Σ can be expressed as

$$\Sigma = \begin{bmatrix} \Sigma_X & \mathbf{0} \\ \mathbf{0} & \Sigma_Y \end{bmatrix},$$

implying that the PDF of \mathbf{Z} factorizes,

$$\begin{aligned} \mathcal{N}([x \ y]; \mu, \Sigma) &= \frac{1}{\sqrt{\det(2\pi\Sigma_X) \cdot \det(2\pi\Sigma_Y)}} \exp \left(-\frac{1}{2}(x - \mu_X)^\top \Sigma_X^{-1}(x - \mu_X) \right. \\ &\quad \left. -\frac{1}{2}(y - \mu_Y)^\top \Sigma_Y^{-1}(y - \mu_Y) \right) \\ &= \mathcal{N}(x; \mu_X, \Sigma_X) \cdot \mathcal{N}(y; \mu_Y, \Sigma_Y). \quad \square \end{aligned}$$

In particular, this implies that random variables that are jointly Gaussian are mutually independent iff the covariance matrix of the corresponding multivariate normal distribution is diagonal.

Exercise 1.69: Closedness properties of Gaussians

Next, we derive two important closedness properties of Gaussians.

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

$$\varphi_{\mathbf{X}}(\mathbf{t}) \doteq \mathbb{E} \exp(\mathbf{t}^\top \mathbf{X}), \quad \text{for all } \mathbf{t} \in \mathbb{R}^n. \quad (1.121)$$

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

$$\varphi_{\mathbf{X}}(\mathbf{t}) = \exp \left(\mathbf{t}^\top \mu + \frac{1}{2} \mathbf{t}^\top \Sigma \mathbf{t} \right). \quad (1.122)$$

This generalizes the MGF of the univariate Gaussian from eq. (1.87).

Prove the following facts.

1. **Fact 1.70** (Closedness under affine transformations). *Given an n -dimensional Gaussian $\mathbf{X} \sim \mathcal{N}(\mu, \Sigma)$, and $A \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$,*

$$A\mathbf{X} + \mathbf{b} \sim \mathcal{N}(A\mu + \mathbf{b}, A\Sigma A^\top). \quad (1.123)$$

2. **Fact 1.71** (Additivity). *Given two independent Gaussian random vectors $\mathbf{X} \sim \mathcal{N}(\mu, \Sigma)$ and $\mathbf{X}' \sim \mathcal{N}(\mu', \Sigma')$ in \mathbb{R}^n ,*

$$\mathbf{X} + \mathbf{X}' \sim \mathcal{N}(\mu + \mu', \Sigma + \Sigma'). \quad (1.124)$$

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

▷ *Solution*

Observe that eq. (1.123) implies that the random vector $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where $\boldsymbol{\Sigma}$ is positive definite is equivalently characterized as

$$\mathbf{X} = \boldsymbol{\Sigma}^{1/2} \mathbf{Y} + \boldsymbol{\mu}. \quad (1.125)$$

where $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\boldsymbol{\Sigma}^{1/2}$ is the square root of $\boldsymbol{\Sigma}$.³² Importantly, this implies together with theorem 1.66 and fact 1.71 that:

³² More details on the square root of a symmetric and positive semi-definite matrix can be found in section 1.6.3.

- Any marginal of a GRV is a GRV.
- Any affine transformation of a GRV is a GRV.

Exercise 1.72: 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.

▷ *Solution*

1.6.2 Conditional Linear Gaussians

It is often useful to consider a Gaussian random variable X in terms of another Gaussian random variable Y . Suppose $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$ and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ are jointly Gaussian. Then, using the closed-form expression of Gaussian conditionals (1.118), their conditional distribution is given as

$$X \mid Y = y \sim \mathcal{N}(\mu_{X|Y}, \sigma_{X|Y}^2) \quad \text{where} \quad (1.126)$$

$$\mu_{X|Y} \doteq \mu_X + \sigma_{XY} \sigma_Y^{-2} (y - \mu_Y), \quad (1.127)$$

$$\sigma_{X|Y}^2 \doteq \sigma_X^2 - \sigma_{XY} \sigma_Y^{-2} \sigma_{YX}. \quad (1.128)$$

Here, σ_{XY} refers to the covariance of X and Y . Again, note that $\mu_{X|Y}$ depends affinely on μ_X .

Observe that this admits an equivalent characterization of X as an affine function of Y with added independent Gaussian noise. Formally, we define

$$X \doteq aY + b + \epsilon \quad \text{where} \quad (1.129)$$

$$a \doteq \sigma_{XY} \sigma_Y^{-2}, \quad (1.130)$$

$$b \doteq \mu_X - \sigma_{XY} \sigma_Y^{-2} \mu_Y, \quad (1.131)$$

$$\epsilon \sim \mathcal{N}(0, \sigma_{X|Y}^2). \quad (1.132)$$

It directly follows from the closedness of Gaussians under linear transformations (1.123) that the characterization of X via (1.129) is equivalent to $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$, and hence, *any* Gaussian X can be modeled as

a *conditional linear Gaussian*, i.e., a linear function of another Gaussian Y with additional independent Gaussian noise.

1.6.3 Quadratic Forms

Definition 1.73 (Quadratic form). Given a symmetric matrix $A \in \mathbb{R}^{n \times n}$, the *quadratic form* induced by A is defined as

$$f : \mathbb{R}^n \rightarrow \mathbb{R}, \quad x \mapsto x^\top A x = \sum_{i=1}^n \sum_{j=1}^n A(i, j) \cdot x(i) \cdot x(j). \quad (1.133)$$

We call A a *positive definite* matrix if all eigenvalues of A are positive. Equivalently, we have $f(x) > 0$ for all $x \in \mathbb{R}^n \setminus \{0\}$ and $f(0) = 0$. Similarly, A is called *positive semi-definite* if all eigenvalues of A are non-negative, or equivalently, if $f(x) \geq 0$ for all $x \in \mathbb{R}^n$. In particular, if A is positive definite then $\sqrt{f(x)}$ is a norm (called the *Mahalanobis norm* induced by A), and is often denoted by $\|x\|_A$.

If A is positive definite, then the sublevel sets of its induced quadratic form f are convex ellipsoids. Not coincidentally, the same is true for the sublevel sets of the PDF of a normal distribution $\mathcal{N}(\mu, \Sigma)$, which we have seen an example of in fig. 1.12. Hereby, the axes of the ellipsoid and their corresponding squared lengths are the eigenvectors and eigenvalues of Σ , respectively.

Remark 1.74: Correspondence of quadratic forms and Gaussians

Observe that the PDF of a zero-mean multivariate Gaussian $\mathcal{N}(0, \Sigma)$ is an exponentiated and appropriately scaled quadratic form induced by the positive definite precision matrix Σ^{-1} . The constant factor is chosen such that the resulting function is a valid probability density function, that is, sums to one.

One important property of positive definiteness of Σ is that Σ can be decomposed into the product of a lower-triangular matrix with its transpose. This is known as the *Cholesky decomposition*.

Fact 1.75 (Cholesky decomposition, symmetric matrix-form). *For any symmetric and positive (semi-)definite matrix $A \in \mathbb{R}^{n \times n}$, there is a decomposition of the form*

$$A = \mathcal{L} \mathcal{L}^\top \quad (1.134)$$

where $\mathcal{L} \in \mathbb{R}^{n \times n}$ is lower triangular and positive (semi-)definite.

We will not prove this fact, but it is not hard to see that a decomposition exists (it takes more work to show that \mathcal{L} is lower-triangular).

Let A be a symmetric and positive (semi-)definite matrix. By the spectral theorem of symmetric matrices, $A = V\Lambda V^\top$, where Λ is a diagonal matrix of eigenvalues and V is an orthonormal matrix of corresponding eigenvectors. Consider the matrix

$$A^{1/2} \doteq V\Lambda^{1/2}V^\top \quad (1.135)$$

where $\Lambda^{1/2} \doteq \text{diag}\{\sqrt{\Lambda(i,i)}\}$, also called the *square root* of A . Then,

$$\begin{aligned} A^{1/2}A^{1/2} &= V\Lambda^{1/2}V^\top V\Lambda^{1/2}V^\top \\ &= V\Lambda^{1/2}\Lambda^{1/2}V^\top \\ &= V\Lambda V^\top = A. \end{aligned} \quad (1.136)$$

It is immediately clear from the definition that $A^{1/2}$ is also symmetric and positive (semi-)definite.

Quadratic forms of positive semi-definite matrices are a generalization of the Euclidean norm, as

$$\|x\|_A^2 = x^\top Ax = x^\top A^{1/2}A^{1/2}x = (A^{1/2}x)^\top A^{1/2}x = \|A^{1/2}x\|_2^2, \quad (1.137)$$

and in particular,

$$\log \mathcal{N}(x; \mu, \Sigma) = -\frac{1}{2} \|x - \mu\|_{\Sigma^{-1}}^2 + \text{const} \quad (1.138)$$

$$= -\frac{1}{2} [x^\top \Sigma^{-1}x - 2\mu^\top \Sigma^{-1}x] + \text{const}, \quad (1.139)$$

$$\log \mathcal{N}(x; \mathbf{0}, I) = -\frac{1}{2} \|x\|_2^2 + \text{const}. \quad (1.140)$$

PART I

Probabilistic Machine Learning

Bayesian Linear Regression

2.1 Linear Regression

We consider the supervised learning setting that we have introduced in section 1.2. Given a set of (x, y) pairs, linear regression aims to find a linear model that fits the data optimally.¹ Given the linear model

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

for $\mathbf{x} \in \mathbb{R}^d$ and labeled training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, we want to find optimal weights \mathbf{w} .

Remark 2.1: Implicit use of intercepts

In general, a linear function

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

also has an intercept term w_0 . Observe, however, that if we define $\mathbf{x}' \doteq (\mathbf{x}, 1)$ and $\mathbf{w}' \doteq (\mathbf{w}, w_0)$, then we have $\mathbf{w}'^\top \mathbf{x}' = \mathbf{w}^\top \mathbf{x} + w_0$, implying that w.l.o.g. it suffices to study linear functions without the intercept term w_0 .

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 $\mathbf{f} \doteq [f_1 \cdots f_n]^\top$. It requires little calculation to verify that $\mathbf{f} = \mathbf{X}\mathbf{w}$.

¹ As we have discussed in section 1.2, regression models can also be used for classification. One example of a linear regression model for classification is Bayesian logistic regression, which we will see in example 5.2.

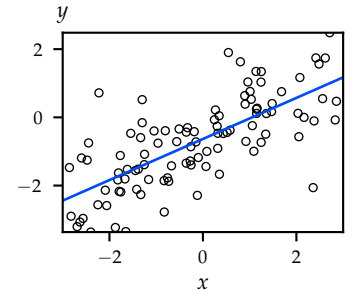


Figure 2.1: Example of linear regression with the least squares estimator (shown in blue).

There are many ways of estimating w from data, the most common being the *least squares estimator*,

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

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. This reduces the potential of overfitting to the training data.²

² Ridge regression is more robust than standard linear regression in the presence of multicollinearity. *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.

Exercise 2.2: Closed-form linear regression

1. Show 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 losses are strictly convex functions).

2. For a $n \times m$ matrix A and vector $x \in \mathbb{R}^m$, we call $\Pi_A x$ the orthogonal projection of x onto $\text{span}\{A\} = \{Ax' \mid x' \in \mathbb{R}^m\}$. In particular, an orthogonal projection satisfies $x - \Pi_A x \perp Ax'$ for all $x' \in \mathbb{R}^m$.

Show that \hat{w}_{ls} from eq. (2.4) is such that $X\hat{w}_{\text{ls}}$ is the unique closest point to y on $\text{span}\{X\}$, i.e., it satisfies $X\hat{w}_{\text{ls}} = \Pi_X y$.

▷ *Solution*

Readings

For a more thorough introduction to linear regression, read

- chapter 4 of “Introduction to Machine Learning” (Krause and Yang, 2022) or
- chapter 9 of “Mathematics for machine learning” (Deisenroth et al., 2020).

2.1.1 Maximum Likelihood Estimation

Let us assume that the observations

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

are made under the influence of i.i.d. homoscedastic Gaussian noise $\epsilon_i \sim \mathcal{N}(0, \sigma_n^2)$.³ Recall from section 1.6.2 that this is equivalently characterized by the Gaussian likelihood,

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

Computing the maximum likelihood estimate for the weights,

$$\begin{aligned} \hat{\mathbf{w}}_{\text{MLE}} &= \arg \max_{\mathbf{w} \in \mathbb{R}^d} \log p(y_{1:n} \mid \mathbf{x}_{1:n}, \mathbf{w}) \\ &= \arg \max_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \log p(y_i \mid x_i, \mathbf{w}) \\ &= \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \mathbf{w}^\top x_i)^2, \end{aligned}$$

we observe that $\hat{\mathbf{w}}_{\text{MLE}} = \hat{\mathbf{w}}_{\text{ls}}$.

Example 2.3: Variance of least squares

It is useful to understand the variance of the weights $\hat{\mathbf{w}}_{\text{ls}}$ in terms of the data (X, \mathbf{y}) because this allows us to determine the variance of a prediction,

$$\begin{aligned} \text{Var}[y^* \mid \mathbf{x}^*, X] &= \text{Var}[\hat{\mathbf{w}}_{\text{ls}}^\top \mathbf{x}^* \mid \mathbf{x}^*, X] \\ &= \mathbf{x}^{*\top} \text{Var}[\hat{\mathbf{w}}_{\text{ls}} \mid X] \mathbf{x}^*. \end{aligned} \quad (2.8)$$

The variance of the weights of the least squares estimator is

$$\begin{aligned} \text{Var}[\hat{\mathbf{w}}_{\text{ls}} \mid X] &= \text{Var}[(X^\top X)^{-1} X^\top \mathbf{y} \mid X] \\ &= (X^\top X)^{-1} X^\top \text{Var}[\mathbf{y} \mid X] ((X^\top X)^{-1} X^\top)^\top \\ &= (X^\top X)^{-1} X^\top \text{Var}[\mathbf{y} \mid X] X (X^\top X)^{-1}. \end{aligned}$$

Due to the Gaussian likelihood (2.7), $\text{Var}[\mathbf{y} \mid X] = \sigma_n^2 \mathbf{I}$, so

$$= \sigma_n^2 (X^\top X)^{-1}. \quad (2.9)$$

Exercise 2.4: Variance of least squares around training data

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

▷ *Solution*

³ ϵ_i is often called *additive white Gaussian noise*. In practice, the noise variance σ_n^2 is typically unknown and also has to be estimated. It is a straightforward exercise to check that the maximum likelihood estimate of σ_n^2 given fixed weights \mathbf{w} is given by $\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^\top x_i)^2$.

using independence of the training data

plugging in the Gaussian likelihood

using that y^* is a linear function in $\hat{\mathbf{w}}_{\text{ls}}$ evaluated at \mathbf{x}^*

using eq. (1.47)

using eq. (2.4)

using eq. (1.47)

using $(A^\top)^{-1} = (A^{-1})^\top$

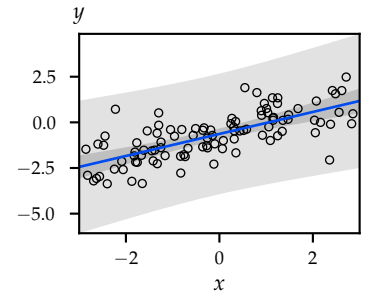


Figure 2.2: The same linear regression model as in fig. 2.1 assuming $\sigma_n = 1$. The estimated variance of the weights (i.e., epistemic uncertainty) is shown in dark gray, the aleatoric uncertainty in light gray.

2.2 Weight-space View

There are two views on Bayesian linear regression, the weight-space view and the function-space view. We begin by discussing the common view in terms of a weight parameterization. In section 2.6, we introduce the interpretation in terms of functions, which will then naturally lead us to Gaussian processes.

The most immediate and natural Bayesian interpretation of linear regression is to simply impose a prior on the weights w . We will use a Gaussian prior,

$$w \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I}), \quad (2.10)$$

and the Gaussian likelihood from eq. (2.7).⁴

Remark 2.5: Maximum entropy principle

The choice of using a Gaussian prior may seem somewhat arbitrary at first sight (except for the nice analytical properties of Gaussians that we have seen in section 1.6). The following intuitive principle gives a more fundamental argument in favor of Gaussian priors.

The *entropy* of a distribution p is defined as

$$H[p] \doteq \mathbb{E}_{x \sim p}[-\log p(x)]. \quad (2.11)$$

Intuitively, if the entropy of p is high, then the distribution p exhibits high “uncertainty”. We give a thorough introduction to entropy in section 5.4.

When choosing a prior, the *maximum entropy principle* states that one should choose the distribution maximizing entropy subject to all the constraints that are to be imposed on the prior. Intuitively, this encodes the least information necessary in the prior. In philosophy, this principle is also known as *Occam’s razor*. It turns out that \mathcal{N} has the maximum entropy among all distributions on \mathbb{R}^d with known mean and covariance. You will prove this in the univariate setting in exercise 5.20.

Thus, the maximum entropy principle can serve as an argument for why it might be reasonable to assume a Gaussian prior.

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

$$\begin{aligned} \log p(w \mid x_{1:n}, y_{1:n}) \\ = \log p(w) + \log p(y_{1:n} \mid x_{1:n}, w) + \text{const} \end{aligned}$$

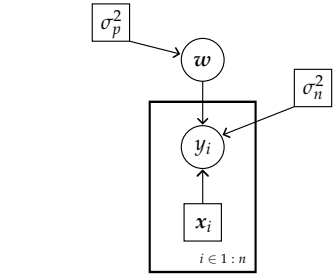


Figure 2.3: Directed graphical model of Bayesian linear regression in plate notation.

⁴ The choice of hyperparameters such as the prior variance σ_p^2 is discussed in greater detail in section 4.4.

by Bayes’ rule (1.59)

$$\begin{aligned}
&= \log p(\mathbf{w}) + \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \mathbf{w}) + \text{const} && \text{using independence of the samples} \\
&= -\frac{1}{2} \left[\sigma_p^{-2} \|\mathbf{w}\|_2^2 + \sigma_n^{-2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 \right] + \text{const} && \text{using the Gaussian prior and likelihood} \\
&= -\frac{1}{2} \left[\sigma_p^{-2} \|\mathbf{w}\|_2^2 + \sigma_n^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \right] + \text{const} && \text{using } \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \\
&= -\frac{1}{2} \left[\sigma_p^{-2} \mathbf{w}^\top \mathbf{w} + \sigma_n^{-2} (\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y}) \right] + \text{const} \\
&= -\frac{1}{2} \left[\mathbf{w}^\top (\sigma_n^{-2} \mathbf{X}^\top \mathbf{X} + \sigma_p^{-2} \mathbf{I}) \mathbf{w} - 2\sigma_n^{-2} \mathbf{y}^\top \mathbf{X} \mathbf{w} \right] + \text{const}. \quad (2.12)
\end{aligned}$$

Observe that the log-posterior is a quadratic form in \mathbf{w} , so the posterior distribution is Gaussian (cf. eq. (1.139)):

$$\mathbf{w} | \mathbf{x}_{1:n}, y_{1:n} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}). \quad (2.13)$$

where we can read off the mean and covariance to be

$$\boldsymbol{\Sigma} \doteq (\sigma_n^{-2} \mathbf{X}^\top \mathbf{X} + \sigma_p^{-2} \mathbf{I})^{-1}, \quad (2.14a)$$

$$\boldsymbol{\mu} \doteq \sigma_n^{-2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y}. \quad (2.14b)$$

This also shows that Gaussians with known variance and linear likelihood are self-conjugate, a property that we had hinted at in section 1.3.2. It can be shown more generally that Gaussians with known variance are self-conjugate to any Gaussian likelihood (Murphy, 2007).

For general distributions the posterior can typically not be expressed in closed-form — this is a very special property of Gaussians!

2.2.1 Maximum A Posteriori Estimation

Computing the MAP estimate for the weights,

$$\begin{aligned}
\hat{\mathbf{w}}_{\text{MAP}} &= \arg \max_{\mathbf{w}} \log p(y_{1:n} | \mathbf{x}_{1:n}, \mathbf{w}) + \log p(\mathbf{w}) \\
&= \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{\sigma_n^2}{\sigma_p^2} \|\mathbf{w}\|_2^2, \quad (2.15)
\end{aligned}$$

using that the likelihood and prior are Gaussian

we observe that this is identical to ridge regression with weight decay $\lambda \doteq \sigma_n^2/\sigma_p^2$, and hence, $\hat{\mathbf{w}}_{\text{MAP}} = \hat{\mathbf{w}}_{\text{ridge}}$. Note that this is simply the MLE loss with an additional ℓ_2 regularization term (originating from the prior) that encourages keeping weights small.

Also, recall that the MAP estimate corresponds to the mode of the posterior distribution, which in the case of a Gaussian is simply its mean $\boldsymbol{\mu}$. As to be expected, $\boldsymbol{\mu}$ coincides with the analytical solution to ridge regression from eq. (2.5).

Example 2.6: 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.

A commonly used alternative to ridge regression is the *least absolute shrinkage and selection operator* (or *lasso*), which uses ℓ_1 regularization,

$$\hat{\mathbf{w}}_{\text{lasso}} \doteq \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1. \quad (2.16)$$

It turns out that lasso can also be viewed as Bayesian learning, using a Laplace prior $\mathbf{w} \sim \text{Laplace}(\mathbf{0}, \ell)$ with length scale ℓ instead of a Gaussian prior. Recall that for the PDF of the Laplace distribution we have,

$$\text{Laplace}(\mathbf{x}; \boldsymbol{\mu}, \ell) \propto \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|_1}{\ell}\right). \quad (2.17)$$

Computing the MAP estimate for the weights yields,

$$\begin{aligned} \hat{\mathbf{w}}_{\text{MAP}} &= \arg \max_{\mathbf{w}} \log p(y_{1:n} | \mathbf{x}_{1:n}, \mathbf{w}) + \log p(\mathbf{w}) \\ &= \arg \min_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \frac{\sigma_n^2}{\ell} \|\mathbf{w}\|_1. \end{aligned} \quad (2.18)$$

using that the likelihood is Gaussian and the prior is Laplacian

This coincides with the optimization objective of lasso with weight decay $\lambda \doteq \sigma_n^2/\ell$.

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

$$y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n} \sim \mathcal{N}(f^*, \sigma_n^2). \quad (2.19)$$

Note that using point estimates like the MAP estimate does not quantify uncertainty in the weights. The MAP estimate simply collapses all mass of the posterior around its mode. This is especially harmful when we are highly unsure about the best model (e.g., because we have observed insufficient data).

2.2.2 Bayesian Inference

Rather than selecting a single weight vector $\hat{\mathbf{w}}$ to make predictions, we can use the full posterior distribution. Doing so is also called “Bayesian

inference”, and in the setting of linear regression known as *Bayesian linear regression*.

For a test point x^* , we let $f^* \doteq w^\top x^*$. Using the closedness of Gaussians under linear transformations (1.123),

$$f^* \mid x^*, x_{1:n}, y_{1:n} \sim \mathcal{N}(\mu^\top x^*, x^{*\top} \Sigma x^*). \quad (2.20)$$

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

$$y^* \mid x^*, x_{1:n}, y_{1:n} \sim \mathcal{N}(\mu^\top x^*, x^{*\top} \Sigma x^* + \sigma_n^2). \quad (2.21)$$

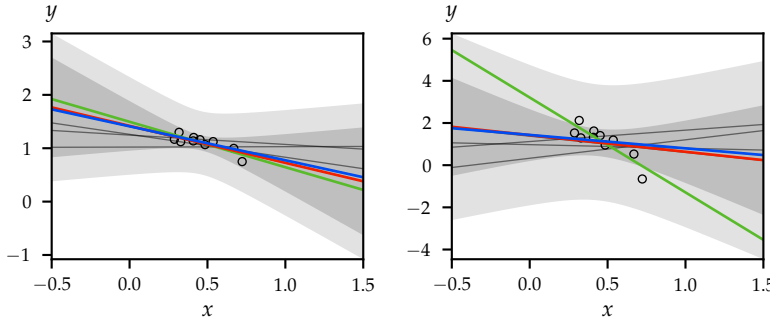


Figure 2.4: Comparison of **linear regression (MLE)**, **ridge regression (MAP estimate)**, and **Bayesian linear regression** when the data is generated according to

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

The **true mean** is shown in blue, the MLE in green, 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. Plausible realizations are shown in black. On the left, $\sigma_n = 0.15$. On the right, $\sigma_n = 0.7$.

Exercise 2.7: Bayesian linear regression

Suppose you are given the following observations

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

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

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

▷ *Solution*

2.3 Aleatoric and Epistemic Uncertainty

The predictive posterior distribution from eq. (2.21) highlights a decomposition of uncertainty wherein $\mathbf{x}^{\star\top} \boldsymbol{\Sigma} \mathbf{x}^{\star}$ 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”).

A natural Bayesian 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 non-linear 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 inputs.

Epistemic and aleatoric uncertainty can be formally defined in terms of the law of total variance (1.50),

$$\text{Var}[y^* \mid \mathbf{x}^*] = \mathbb{E}_{\boldsymbol{\theta}}[\text{Var}_{y^*}[y^* \mid \mathbf{x}^*, \boldsymbol{\theta}]] + \text{Var}_{\boldsymbol{\theta}}[\mathbb{E}_{y^*}[y^* \mid \mathbf{x}^*, \boldsymbol{\theta}]]. \quad (2.22)$$

Here, the mean variability of predictions y^* averaged across all models $\boldsymbol{\theta}$ is the *aleatoric uncertainty*. In contrast, the variability of the mean prediction y^* under each model $\boldsymbol{\theta}$ is the *epistemic uncertainty*.

Exercise 2.8: Aleatoric and epistemic uncertainty of BLR

Prove for Bayesian linear regression that $\mathbf{x}^{\star\top} \boldsymbol{\Sigma} \mathbf{x}^{\star}$ is the epistemic uncertainty and σ_n^2 the aleatoric uncertainty in y^* under the decomposition of eq. (2.22).

▷ *Solution*

2.4 Recursive Bayesian Updates

A more general feature of supervised Bayesian learning, which often leads to efficient algorithms, is its recursive structure. Suppose we are

given a prior $p(\boldsymbol{\theta})$ and observations $y_{1:n}$. We define

$$p^{(t)}(\boldsymbol{\theta}) \doteq p(\boldsymbol{\theta} \mid y_{1:t}) \quad (2.23)$$

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

$$\begin{aligned} p^{(t+1)}(\boldsymbol{\theta}) &= p(\boldsymbol{\theta} \mid y_{1:t+1}) \\ &\propto p(\boldsymbol{\theta} \mid y_{1:t}) \cdot p(y_{t+1} \mid \boldsymbol{\theta}, y_{1:t}) && \text{using Bayes' rule (1.59)} \\ &= p^{(t)}(\boldsymbol{\theta}) \cdot p(y_{t+1} \mid \boldsymbol{\theta}). && \text{using } y_{t+1} \perp y_{1:t} \mid \boldsymbol{\theta}, \text{ see fig. 2.3} \end{aligned} \quad (2.24)$$

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

Thus, as data arrives online (i.e., in “real-time”), we can obtain the new posterior and use it to replace our prior. In general, the posterior distribution can be complicated. For Bayesian linear regression with a Gaussian prior and likelihood, however, the posterior distribution is also a Gaussian,

$$p^{(t)}(\boldsymbol{w}) = \mathcal{N}(\boldsymbol{w}; \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}), \quad (2.25)$$

which can be stored efficiently using only $\mathcal{O}(d^2)$ parameters.

Exercise 2.9: Online Bayesian linear regression

1. Can you design an algorithm that updates the posterior (as opposed to recalculating it from scratch using eq. (2.13)) in a smarter way? The requirement is that the memory should not grow as $\mathcal{O}(t)$.
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 $\mathcal{O}(d^2)$?

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

▷ *Solution*

This 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 3.2, we show that online Bayesian linear regression is, in fact, an example of a Kalman filter.

2.5 Non-linear Regression

We can use linear regression not only to learn linear functions. The trick is to apply a non-linear 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 2.10: 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

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

Note that the feature dimension e is $\sum_{i=0}^m \binom{d+i-1}{i} = \Theta(d^m)$.⁵ 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.

2.6 Function-space View

Let us 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 .

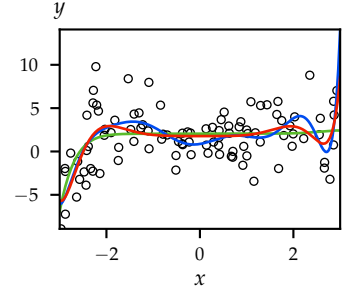


Figure 2.5: Applying linear regression with a feature space of polynomials of degree 10. The **least squares estimate** is shown in blue, **ridge regression** in red, and **lasso** in green.

⁵ 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 “ $\circ \mid \mid \circ \mid$ ” where \circ denotes a selected spot and \mid 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 $\binom{d+i-1}{i}$.

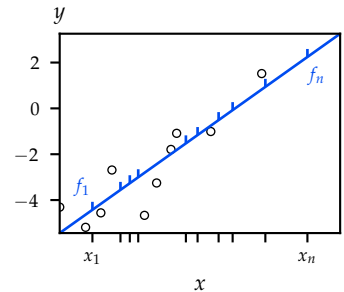


Figure 2.6: An illustration of the function-space view. The model is described by the points (x_i, f_i) .

You may say, that nothing has changed — and you would be right: that is precisely the point. However, instead of considering a prior over the weights $w \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$ 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 (1.123), we obtain the equivalent prior

$$f | X \sim \mathcal{N}(\Phi \mathbb{E}[w], \underbrace{\Phi \text{Var}[w] \Phi^\top}_K) = \mathcal{N}(\mathbf{0}, \sigma_p^2 \underbrace{\Phi \Phi^\top}_K) \quad (2.26)$$

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 $K(i, j) = \sigma_p^2 \cdot \phi(x_i)^\top \phi(x_j)$.

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*

$$k(x, x') \doteq \sigma_p^2 \cdot \phi(x)^\top \phi(x') \quad (2.27)$$

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

$$K = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix} \quad (2.28)$$

Observe that by definition of the kernel matrix in eq. (2.26), the kernel matrix is a covariance matrix and the kernel function measures the covariance of the function values $f(x)$ and $f(x')$ given inputs x and x' . That is,

$$k(x, x') = \text{Cov}[f(x), f(x')]. \quad (2.29)$$

Moreover, note that we have reformulated⁶ 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), and encodes our prior beliefs. This is known as the *kernel trick*.

⁶ we often say “kernelized”

2.6.1 Learning and Inference

We have already kernelized the Bayesian linear regression prior. The posterior distribution $f | X, \mathbf{y}$ is again Gaussian due to the closedness properties of Gaussians, analogously to our derivation of the prior kernel matrix in eq. (2.26).

It remains to show that we can also rely on the kernel trick during inference. Given the test point \mathbf{x}^* , we define

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

We immediately obtain $\tilde{\mathbf{f}} = \tilde{\Phi} \mathbf{w}$. Analogously to our analysis of inference from the weight-space view, we add the label noise to obtain the estimate $\tilde{\mathbf{y}} = \tilde{\mathbf{f}} + \tilde{\epsilon}$ where $\tilde{\epsilon} \doteq [\epsilon_1 \cdots \epsilon_n \epsilon^*]^\top \sim \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I})$ is the independent label noise. Applying the same reasoning as we did for the prior, we obtain

$$\tilde{\mathbf{f}} \mid \mathbf{X}, \mathbf{x}^* \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{K}}) \quad (2.30)$$

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

$$\tilde{\mathbf{y}} \mid \mathbf{X}, \mathbf{x}^* \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{K}} + \sigma_n^2 \mathbf{I}). \quad (2.31)$$

Finally, we can conclude from the closedness of Gaussian random vectors under conditional distributions (1.118) that the predictive posterior $y^* \mid \mathbf{x}^*, \mathbf{X}, \mathbf{y}$ follows again a normal distribution. We will do a full derivation of the posterior and predictive posterior in section 4.1.

2.6.2 Efficient Polynomial Regression

But how does the kernel trick address our concerns about efficiency raised in section 2.5? 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 $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$.
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 2.11. *For the polynomial feature transformation ϕ up to degree m from example 2.10, it can be shown that*

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^m \quad (2.32)$$

up to constant factors.

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

3

Kalman Filters

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. Kalman filters are an example of 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*. We will discuss this process more broadly in the next section. A Kalman filter is simply a Bayes filter using a Gaussian distribution over the states and conditional linear Gaussians to describe the evolution of states and observations.

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*,

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

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

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, $\boldsymbol{\epsilon}$ 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 section 1.6.2, their joint distribution (over all variables) is also Gaussian. This means that predicting the future states

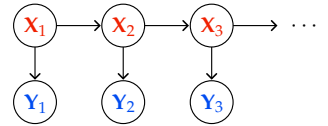


Figure 3.1: Directed graphical model of a Kalman filter with hidden states \mathbf{X}_t and observables \mathbf{Y}_t .

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 2.4).

From the directed graphical model of a Kalman filter shown in fig. 3.1, 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,

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

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 (1.14)

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

Example 3.2: 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 2.4. Observe that this algorithm performs Bayesian filtering. That is, it keeps track of the weights given noisy labels y_i . Moreover, we have used a Gaussian prior and likelihood. This is precisely the setting of a Kalman filter!

More concretely, our motion model is given as

- $\mathbf{x}_t = \mathbf{w}^{(t)}$;
- $\mathbf{F} = \mathbf{I}$; and
- $\boldsymbol{\epsilon} = \mathbf{0}$.

Our (time-dependent) sensor model is given as

- $\mathbf{H}_t = \text{diag}\{\mathbf{x}_t\}$; and
- $\boldsymbol{\eta}_t = \boldsymbol{\epsilon}_t$.

Using the sensor model (3.3), this yields d -dimensional labels,

$$\mathbf{y}_t = \text{diag}\{\mathbf{x}_t\} \mathbf{w}^{(t)} + \boldsymbol{\epsilon}_t \in \mathbb{R}^d. \quad (3.8)$$

Observe that the label y_t can simply be recovered by summing over all elements of the vector,

$$y_t = \sum_{i=1}^d \mathbf{y}_t(i). \quad (3.9)$$

3.1 Bayesian Filtering

As we have already mentioned, *Bayesian filtering* is the process of keeping track of an agent's state using noisy observations. Bayesian filtering is described by the following recursive scheme with the two phases, *conditioning* (also called *update*) and *prediction*.

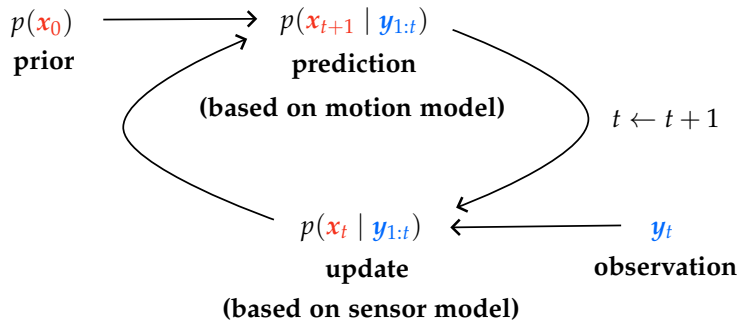


Figure 3.2: Schematic view of Bayesian filtering.

Algorithm 3.3: Bayesian filtering

- 1 start with a prior over initial states $p(\mathbf{x}_0)$
 - 2 **for** $t = 1$ **to** ∞ **do**
 - 3 assume we have $p(\mathbf{x}_t | \mathbf{y}_{1:t-1})$
 - 4 compute $p(\mathbf{x}_t | \mathbf{y}_{1:t})$ using the new observation \mathbf{y}_t (conditioning)
 - 5 compute $p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t})$ (prediction)
-

Let us consider the conditioning step first:

$$\begin{aligned} p(\mathbf{x}_t | \mathbf{y}_{1:t}) &= \frac{1}{Z} p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{y}_{1:t-1}) \\ &= \frac{1}{Z} p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) p(\mathbf{y}_t | \mathbf{x}_t). \end{aligned} \quad (3.10)$$

using Bayes' rule (1.59)

using the conditional independence structure (3.6)

For the prediction step, we obtain,

$$\begin{aligned}
 p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:t}) &= \int p(\mathbf{x}_{t+1}, \mathbf{x}_t \mid \mathbf{y}_{1:t}) d\mathbf{x}_t && \text{using the sum rule (1.10)} \\
 &= \int p(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{y}_{1:t}) p(\mathbf{x}_t \mid \mathbf{y}_{1:t}) d\mathbf{x}_t && \text{using the product rule (1.14)} \\
 &= \int p(\mathbf{x}_{t+1} \mid \mathbf{x}_t) p(\mathbf{x}_t \mid \mathbf{y}_{1:t}) d\mathbf{x}_t. && \text{using the conditional independence structure (3.4)}
 \end{aligned} \tag{3.11}$$

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

Remark 3.4: 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 $\mathbf{X}_k \mid \mathbf{y}_{1:t}$ where $t > k$. Note that if $k = t$, then Bayesian smoothing coincides with Bayesian filtering.

Analogously to eq. (3.10),

$$p(\mathbf{x}_k \mid \mathbf{y}_{1:t}) \propto p(\mathbf{x}_k \mid \mathbf{y}_{1:k}) p(\mathbf{y}_{k+1:t} \mid \mathbf{x}_k). \tag{3.12}$$

If we assume a Gaussian prior and conditional Gaussian transition and dynamics models (this is called *Kalman smoothing*), then by the closedness properties of Gaussians, $\mathbf{X}_k \mid \mathbf{y}_{1:t}$ is a Gaussian. Indeed, all terms of eq. (3.12) are Gaussian PDFs and as seen in eq. (1.116), 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,

$$\begin{aligned}
 p(\mathbf{y}_{k+1:t} \mid \mathbf{x}_k) &= \int p(\mathbf{y}_{k+1:t} \mid \mathbf{x}_k, \mathbf{x}_{k+1}) p(\mathbf{x}_{k+1} \mid \mathbf{x}_k) d\mathbf{x}_{k+1} && \text{using the sum rule (1.10) and product rule (1.14)} \\
 &= \int p(\mathbf{y}_{k+1:t} \mid \mathbf{x}_{k+1}) p(\mathbf{x}_{k+1} \mid \mathbf{x}_k) d\mathbf{x}_{k+1} && \text{using the conditional independence structure (3.5)} \\
 &= \int p(\mathbf{y}_{k+1} \mid \mathbf{x}_{k+1}) p(\mathbf{y}_{k+2:t} \mid \mathbf{x}_{k+1}) p(\mathbf{x}_{k+1} \mid \mathbf{x}_k) d\mathbf{x}_{k+1} && \text{using the conditional independence structure (3.6)}
 \end{aligned} \tag{3.13}$$

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 eq. (3.12) can be computed efficiently: one using a (forward) Kalman filter; the other using a “backward” Kalman filter. More concretely, in time $\mathcal{O}(t)$, 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 Update

Before introducing the general Kalman update, let us consider a simple example.

Example 3.5: Random walk in 1d

We use the simple motion and sensor models,³

$$X_{t+1} | x_t \sim \mathcal{N}(x_t, \sigma_x^2), \quad (3.14a)$$

$$Y_t | x_t \sim \mathcal{N}(x_t, \sigma_y^2). \quad (3.14b)$$

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

$$\mu_{t+1} \doteq \frac{\sigma_y^2 \mu_t + (\sigma_t^2 + \sigma_x^2) y_{t+1}}{\sigma_t^2 + \sigma_x^2 + \sigma_y^2}, \quad \sigma_{t+1}^2 \doteq \frac{(\sigma_t^2 + \sigma_x^2) \sigma_y^2}{\sigma_t^2 + \sigma_x^2 + \sigma_y^2}. \quad (3.15)$$

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

$$\lambda \doteq \frac{\sigma_t^2 + \sigma_x^2}{\sigma_t^2 + \sigma_x^2 + \sigma_y^2} = 1 - \frac{\sigma_y^2}{\sigma_t^2 + \sigma_x^2 + \sigma_y^2} \in [0, 1]. \quad (3.16)$$

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

$$\mu_{t+1} = (1 - \lambda) \mu_t + \lambda y_{t+1} \quad (3.17)$$

$$= \mu_t + \lambda (y_{t+1} - \mu_t). \quad (3.18)$$

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,

$$\sigma_{t+1}^2 = \lambda \sigma_y^2 = (1 - \lambda) (\sigma_t^2 + \sigma_x^2). \quad (3.19)$$

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

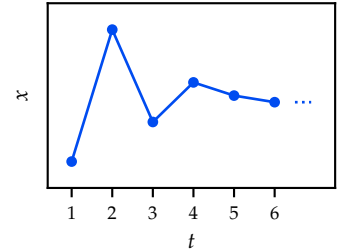


Figure 3.3: Hidden states during a random walk in one dimension.

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

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

Exercise 3.6: Kalman update

Derive the predictive distribution $X_{t+1} \mid y_{1:t+1}$ (3.15) of the Kalman filter described in the above example using your knowledge about multivariate Gaussians from section 1.6.

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

▷ *Solution*

The general formulas for the *Kalman update* follow the same logic. Given the prior belief $\mathbf{X}_t \mid \mathbf{y}_{1:t} \sim \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$, we have

$$\mathbf{X}_{t+1} \mid \mathbf{y}_{1:t+1} \sim \mathcal{N}(\boldsymbol{\mu}_{t+1}, \boldsymbol{\Sigma}_{t+1}) \quad \text{where} \quad (3.20)$$

$$\boldsymbol{\mu}_{t+1} \doteq \mathbf{F}\boldsymbol{\mu}_t + \mathbf{K}_{t+1}(\mathbf{y}_{t+1} - \mathbf{H}\mathbf{F}\boldsymbol{\mu}_t), \quad (3.21)$$

$$\boldsymbol{\Sigma}_{t+1} \doteq (\mathbf{I} - \mathbf{K}_{t+1}\mathbf{H})(\mathbf{F}\boldsymbol{\Sigma}_t\mathbf{F}^\top + \boldsymbol{\Sigma}_x). \quad (3.22)$$

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

$$\mathbf{K}_{t+1} \doteq (\mathbf{F}\boldsymbol{\Sigma}_t\mathbf{F}^\top + \boldsymbol{\Sigma}_x)\mathbf{H}^\top(\mathbf{H}(\mathbf{F}\boldsymbol{\Sigma}_t\mathbf{F}^\top + \boldsymbol{\Sigma}_x)\mathbf{H}^\top + \boldsymbol{\Sigma}_y)^{-1}. \quad (3.23)$$

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.

3.3 Predicting

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

$$\mathbf{X}_{t+1} \mid \mathbf{y}_{1:t} \sim \mathcal{N}(\hat{\boldsymbol{\mu}}_{t+1}, \hat{\boldsymbol{\Sigma}}_{t+1}), \quad (3.24)$$

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

For the mean,

$$\begin{aligned}
 \hat{\mu}_{t+1} &= \mathbb{E}[x_{t+1} \mid \mathbf{y}_{1:t}] \\
 &= \mathbb{E}[F\mathbf{x}_t + \boldsymbol{\epsilon}_t \mid \mathbf{y}_{1:t}] && \text{using the motion model (3.2)} \\
 &= F\mathbb{E}[\mathbf{x}_t \mid \mathbf{y}_{1:t}] && \text{using linearity of expectation (1.24) and } \mathbb{E}[\boldsymbol{\epsilon}_t] = \mathbf{0} \\
 &= F\boldsymbol{\mu}_t. && \text{using the mean of the Kalman update}
 \end{aligned} \tag{3.25}$$

For the covariance matrix,

$$\begin{aligned}
 \hat{\Sigma}_{t+1} &= \mathbb{E}\left[(x_{t+1} - \hat{\mu}_{t+1})(x_{t+1} - \hat{\mu}_{t+1})^\top \mid \mathbf{y}_{1:t}\right] && \text{using the definition of the covariance matrix (1.46)} \\
 &= F\mathbb{E}\left[(x_t - \boldsymbol{\mu}_t)(x_t - \boldsymbol{\mu}_t)^\top \mid \mathbf{y}_{1:t}\right]F^\top + \mathbb{E}[\boldsymbol{\epsilon}_t\boldsymbol{\epsilon}_t^\top] && \text{using (3.25), the motion model (3.2) and that } \boldsymbol{\epsilon}_t \text{ is independent of the observations} \\
 &= F\Sigma_t F^\top + \Sigma. &&
 \end{aligned} \tag{3.26}$$

Exercise 3.7: 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?

▷ *Solution*

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).

4

Gaussian Processes

We obtain Gaussian processes by extending the function-space view of Bayesian linear regression which we introduced in section 2.6 to infinite domains.¹ We are still concerned with the problem of estimating the value of a function $f : \mathcal{X} \rightarrow \mathbb{R}$ at arbitrary points $x^* \in \mathcal{X}$ given training data $\{x_i, y_i\}_{i=1}^n$, where the labels are assumed to be corrupted by homoscedastic Gaussian noise,

$$y_i = f(x_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma_n^2).$$

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.

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.

We use a set \mathcal{X} to index the collection of random variables. 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 $A \doteq \{x_1, \dots, x_m\} \subseteq \mathcal{X}$, we have

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

where

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

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

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

¹ In section 2.6, our domain was given by the set of weights (weight-space view) or the set of (noisy) observations (function-space view).

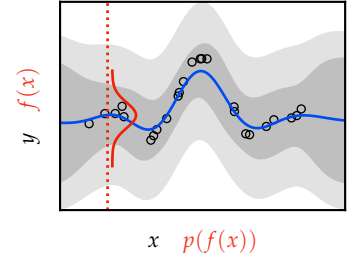


Figure 4.1: A Gaussian process can be interpreted as an infinite-dimensional Gaussian over functions. At any location x in the domain, this yields a distribution over values $f(x)$ shown in red. The blue line corresponds to the MAP estimate (i.e., mean function of the Gaussian process), the dark gray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

The random variable f_x represents the value of the function $f(x)$ at location $x \in \mathcal{X}$, we thus write $f(x) \doteq f_x$. Intuitively, a Gaussian process can be interpreted as a normal distribution over functions — and is therefore often called an infinite-dimensional Gaussian.

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

Readings

Chapter 2 of “Gaussian processes for machine learning” (Williams and Rasmussen, 2006) provides additional background.

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

4.1 Learning and Inference

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

$$\begin{bmatrix} y \\ f^* \end{bmatrix} \mid x^*, x_{1:n} \sim \mathcal{N}(\tilde{\mu}, \tilde{K}), \quad \text{where} \quad (4.4)$$

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

Deriving the conditional distribution using (1.118), we obtain that the Gaussian process posterior is given by

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

$$\mu'(x) \doteq \mu(x) + k_{x, A}^\top (K_{AA} + \sigma_n^2 \mathbf{I})^{-1} (y_A - \mu_A), \quad (4.7)$$

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

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 in section 2.6.1 in the context of Bayesian linear regression, but did not make the predictive posterior explicit. Using eq. (4.6) and our assumption of homoscedastic noise, the predictive posterior at the test point x^* is

$$y^* \mid x^*, x_{1:n}, y_{1:n} \sim \mathcal{N}(\mu^*, k^* + \sigma_n^2). \quad (4.9)$$

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 $\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$. We have already seen in eq. (1.125) that

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

where $\mathbf{K}^{1/2}$ is the square root of \mathbf{K} and $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is standard Gaussian noise.

2. For the second approach, recall the product rule (1.14),

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

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:

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

This general approach is known as *forward sampling*.

We will discuss approximate sampling methods in section 4.5.

4.3 Kernel Functions

We have previously introduced kernels in section 2.6 in the context of a function-space view of Bayesian linear regression. There we have already seen that kernel functions define the “shape” of the functions that are realized from a Gaussian distribution over functions, namely a Gaussian process. In this section, we generalize the notion of kernels to encompass function classes beyond linear and polynomial functions.

Definition 4.2 (Kernel function). A *kernel function* k is defined as

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

for locations \mathbf{x} and \mathbf{x}' such that

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

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

We say that a kernel function is *positive definite* if K_{AA} is positive 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. Indeed, it can be shown that a function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ satisfying the above two conditions corresponds to a (not necessarily finite dimensional) feature representation in some feature space.⁴ The corresponding feature space is a reproducing kernel Hilbert space (RKHS) which we will study in greater detail in section 4.3.4.

Intuitively, the kernel function evaluated at locations x and x' describes how $f(x)$ and $f(x')$ are related. If x and x' are “close”, then $f(x)$ and $f(x')$ are usually taken to be positively correlated, encoding a “smooth” function.

4.3.1 Common Kernels

We will now define some commonly used kernels. Often an additional factor σ^2 (*output scale*) is added, which we assume 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 2.5 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 2.6) and comparing the derived kernel function (2.27) with the definition of the linear kernel (4.13).

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

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

where ℓ is its *length scale*. The larger the length scale ℓ , the smoother the resulting functions.⁵

⁴ That is, there exists a feature transformation $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^e$ where $e \in \mathbb{R} \cup \{\infty\}$ such that $k(x, x') = \phi(x)^\top \phi(x')$. This is known as *Mercer's theorem*, see theorem 4.2 of “Gaussian processes for machine learning” (Williams and Rasmussen, 2006).

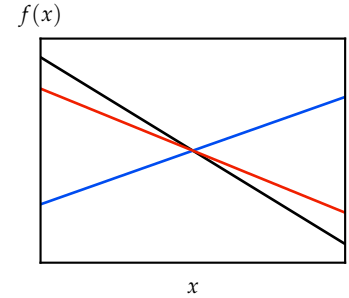


Figure 4.2: Functions sampled according to a Gaussian process with a linear kernel and $\phi = \text{id}$.

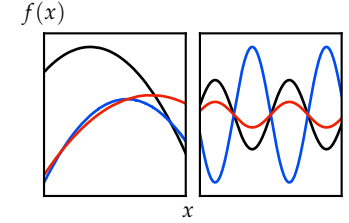


Figure 4.3: Functions sampled according to a Gaussian process with a linear kernel and $\phi(x) = [1, x, x^2]$ (left) and $\phi(x) = \sin(x)$ (right).

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

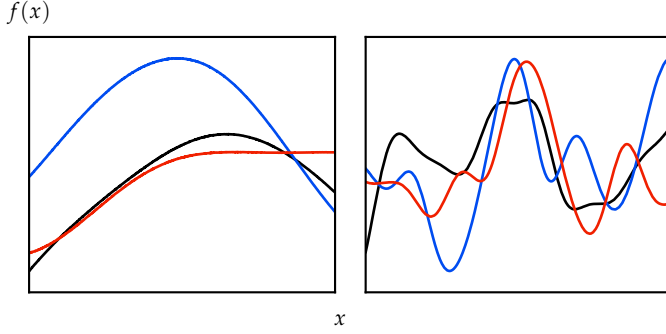


Figure 4.4: Functions sampled according to a Gaussian process with a Gaussian kernel and length scales $\ell = 5$ (left) and $\ell = 1$ (right).

Exercise 4.4: Feature space of Gaussian kernel

1. Show that the univariate Gaussian kernel with length scale $\ell = 1$ 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?

▷ *Solution*

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

$$k(\mathbf{x}, \mathbf{x}'; \ell) \doteq \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_1}{\ell}\right). \quad (4.15)$$

As can be seen in fig. 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. It is defined as

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

where Γ is the Gamma function, K_ν the modified Bessel function of the second kind, and ℓ a length scale parameter. The result-

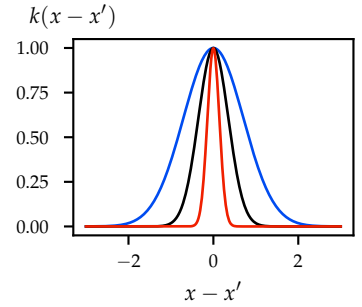


Figure 4.5: Gaussian kernel with length scales $\ell = 1$, $\ell = 0.5$, and $\ell = 0.2$.

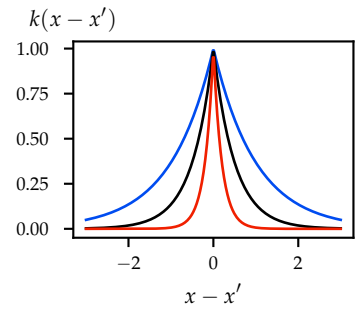


Figure 4.6: Laplace kernel with length scales $\ell = 1$, $\ell = 0.5$, and $\ell = 0.2$.

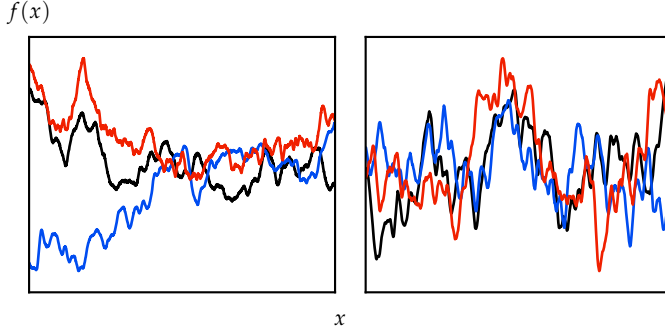


Figure 4.7: Functions sampled according to a Gaussian process with a Laplace kernel and length scales $\ell = 10\,000$ (left) and $\ell = 10$ (right).

ing functions are $\lceil \nu \rceil - 1$ times mean square differentiable.⁶ 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. 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.

Exercise 4.5: A Kalman filter as a Gaussian process

Next we will show that the Kalman filter from example 3.5 can be seen as a Gaussian process. To this end, we define

$$f : \mathbb{N}_0 \rightarrow \mathbb{R}, \quad t \mapsto X_t. \quad (4.21)$$

Assuming that $X_0 \sim \mathcal{N}(0, \sigma_0^2)$ and $X_{t+1} \doteq X_t + \epsilon_t$ with independent noise $\epsilon_t \sim \mathcal{N}(0, \sigma_x^2)$, show that

$$f \sim \mathcal{GP}(0, k_{\text{KF}}) \quad \text{where} \quad (4.22)$$

$$k_{\text{KF}}(t, t') \doteq \sigma_0^2 + \sigma_x^2 \min\{t, t'\}. \quad (4.23)$$

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).

▷ *Solution*

4.3.2 Kernel Composition

Given 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:

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

⁶Mean square continuity and mean square differentiability are a generalization of continuity and differentiation to random processes, where the limit is generalized to a limit in the “mean square sense”.

A sequence of random variables $\{X_n\}_{n \in \mathbb{N}}$ is said to converge to the random variable X in *mean square* if

$$\lim_{n \rightarrow \infty} \mathbb{E}[(X_n - X)^2] = 0. \quad (4.17)$$

Using Markov’s inequality (1.82) it can be seen that convergence in mean square implies convergence in probability. Moreover, convergence in mean square implies

$$\lim_{n \rightarrow \infty} \mathbb{E}X_n = \mathbb{E}X \quad \text{and} \quad (4.18)$$

$$\lim_{n \rightarrow \infty} \mathbb{E}X_n^2 = \mathbb{E}X^2. \quad (4.19)$$

Whereas a deterministic function $f(x)$ is said to be continuous at a point x^* if $\lim_{x \rightarrow x^*} f(x) = f(x^*)$, a random process $f(x)$ is *mean square continuous* at x^* if

$$\lim_{x \rightarrow x^*} \mathbb{E}[(f(x) - f(x^*))^2] = 0. \quad (4.20)$$

It can be shown that a random process is mean square continuous at x^* iff its kernel function $k(x, x')$ is continuous at $x = x' = x^*$.

Similarly, a random process $f(x)$ is *mean square differentiable* at a point x in direction i if $(f(x + \delta e_i) - f(x))/\delta$ converges in mean square as $\delta \rightarrow 0$ where e_i is the unit vector in direction i . This notion can be extended to higher-order derivatives.

For the precise notions of mean square continuity and mean square differentiability in the context of Gaussian processes, refer to section 4.1.1 of “Gaussian processes for machine learning” (Williams and Rasmussen, 2006).

cients 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

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

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 $\mathcal{GP}(\mu, k')$.

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” and
- 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 4.6 (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 $\tilde{k}(\|\mathbf{x} - \mathbf{x}'\|_2) = k(\mathbf{x}, \mathbf{x}')$.

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

Example 4.7: 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.

4.3.4 Reproducing Kernel Hilbert Spaces

Recall that Gaussian processes keep track of a posterior distribution $f \mid \mathbf{x}_{1:n}, y_{1:n}$ over functions. It can be shown that the corresponding MAP estimate \hat{f} corresponds to a regularized optimization problem in a suitable space of functions known as the *reproducing kernel Hilbert space*. 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 4.8 (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

$$\mathcal{H}_k(\mathcal{X}) \doteq \left\{ f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) : n \in \mathbb{N}, \mathbf{x}_i \in \mathcal{X}, \alpha_i \in \mathbb{R} \right\} \quad (4.24)$$

with inner product

$$\langle f, g \rangle_k \doteq \sum_{i=1}^n \sum_{j=1}^{n'} \alpha_i \alpha'_j k(\mathbf{x}_i, \mathbf{x}'_j), \quad (4.25)$$

where $g(\mathbf{x}) = \sum_{j=1}^{n'} \alpha'_j k(\mathbf{x}, \mathbf{x}'_j)$ and norm $\|f\|_k = \sqrt{\langle f, f \rangle_k}$.

It is straightforward to check that for any $\mathbf{x} \in \mathcal{X}$, $k(\mathbf{x}, \cdot) \in \mathcal{H}_k(\mathcal{X})$. 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*. That is, evaluations of RKHS functions f are inner products in $\mathcal{H}_k(\mathcal{X})$ parameterized by the “feature map” $k(\mathbf{x}, \cdot)$.

Exercise 4.9: 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

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq \|f\|_k \|k(\mathbf{x}, \cdot) - k(\mathbf{y}, \cdot)\|_k.$$

▷ *Solution*

Foundational in characterizing the solution to regularized optimization problems within RKHSs is the *representer theorem* (Schölkopf et al., 2001).

Fact 4.10 (Representer theorem). *Let k be a kernel and let $\lambda > 0$. For $f \in \mathcal{H}_k(\mathcal{X})$ and training data \mathcal{D} comprising n observations, let $\ell(f, \mathcal{D}) \in \mathbb{R} \cup \{\infty\}$ denote any loss function. Then, any minimizer of*

$$\hat{f} \doteq \arg \min_{f \in \mathcal{H}_k(\mathcal{X})} \ell(f, \mathcal{D}) + \lambda \|f\|_k^2 \quad (4.26)$$

admits a representation of the form

$$\hat{f}(\mathbf{x}) = \hat{\boldsymbol{\alpha}}^\top \mathbf{k}_{\mathbf{x},A} \quad \text{for some } \hat{\boldsymbol{\alpha}} \in \mathbb{R}^n. \quad (4.27)$$

Exercise 4.11: MAP estimate of Gaussian processes

We will now show that the MAP estimate of GP regression corresponds to the solution of the following regularized linear regression problem in the RKHS,

$$\hat{f} \doteq \arg \min_{f \in \mathcal{H}_k(\mathcal{X})} -\log p(y_{1:n} | \mathbf{x}_{1:n}, f) + \frac{1}{2} \|f\|_k^2. \quad (4.28)$$

Here, the first term corresponds to the likelihood, measuring the “quality of fit”. The regularization term encodes smoothness assumptions on f as specified by the reproducing kernel Hilbert space of the used kernel function.

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 eq. (4.28) is equivalent to

$$\hat{\boldsymbol{\alpha}} \doteq \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{K}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_{\mathbf{K}}^2 \quad (4.29)$$

for some $\lambda > 0$ which is also known as *kernel ridge regression*. Determine λ .

2. Show that eq. (4.29) with the λ determined in (1) is equivalent to the MAP estimate of GP regression.

Hint: Recall from eq. (4.6) that the MAP estimate at a point \mathbf{x}^ is $\mathbb{E}[f^* | \mathbf{x}^*, \mathbf{X}, \mathbf{y}] = \mathbf{k}_{\mathbf{x}^*, A}^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$.*

▷ *Solution*

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 Bayesian 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 $\mathcal{D}^{\text{train}} \doteq \{(\mathbf{x}_i^{\text{train}}, y_i^{\text{train}})\}_{i=1}^n$ and validation sets $\mathcal{D}^{\text{val}} \doteq \{(\mathbf{x}_i^{\text{val}}, y_i^{\text{val}})\}_{i=1}^m$. 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),

$$\hat{f}_j \doteq \arg \max_f p(f | \mathbf{x}_{1:n}^{\text{train}}, y_{1:n}^{\text{train}}). \quad (4.30)$$

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

$$\hat{\theta} \doteq \arg \max_{\theta_j} p(y_{1:m}^{\text{val}} | \mathbf{x}_{1:m}^{\text{val}}, \hat{f}_j). \quad (4.31)$$

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

Remark 4.12: Approximating population risk

Why is it useful to separate the data into a training and a validation set? Recall from section 1.2.1 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

$$\frac{1}{m} \sum_{i=1}^m \ell(y_i^{\text{val}} | \mathbf{x}_i^{\text{val}}, \hat{f}_j) \approx \mathbb{E}_{(x,y) \sim \mathcal{P}} [\ell(y | \mathbf{x}, \hat{f}_j)], \quad (4.32)$$

using Monte Carlo sampling.⁸ In words, for reasonably large m , minimizing the empirical risk as we do in eq. (4.31) 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

$$\hat{\theta}_{\text{MLE}} \doteq \arg \max_{\theta} p(y_{1:n} | \mathbf{x}_{1:n}, \theta) \quad (4.33)$$

$$\begin{aligned} &= \arg \max_{\theta} \int p(y_{1:n}, f | \mathbf{x}_{1:n}, \theta) df \\ &= \arg \max_{\theta} \int p(y_{1:n} | \mathbf{x}_{1:n}, f, \theta) p(f | \theta) df. \end{aligned} \quad (4.34)$$

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

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 (1.88) can be used to gauge how large m should be.

using the definition of marginal likelihood in Bayes’ rule (1.59)

by conditioning on f using the sum rule (1.10)

using the product rule (1.14)

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 .

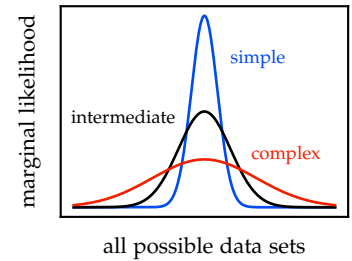


Figure 4.8: A schematic illustration of the marginal likelihood of a simple, intermediate, and complex model across all possible data sets.

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 eq. (4.3) that

$$y_{1:n} \mid x_{1:n}, \theta \sim \mathcal{N}(\mathbf{0}, K_{f,\theta} + \sigma_n^2 \mathbf{I}) \quad (4.35)$$

where $K_{f,\theta}$ denotes the kernel matrix at the inputs $x_{1:n}$ depending on the kernel function parameterized by θ . We write $K_{y,\theta} \doteq K_{f,\theta} + \sigma_n^2 \mathbf{I}$. Continuing from eq. (4.33), we obtain

$$\begin{aligned} \hat{\theta}_{\text{MLE}} &= \arg \max_{\theta} \mathcal{N}(\mathbf{y}; \mathbf{0}, K_{y,\theta}) \\ &= \arg \min_{\theta} \frac{1}{2} \mathbf{y}^\top K_{y,\theta}^{-1} \mathbf{y} + \frac{1}{2} \log \det(K_{y,\theta}) + \frac{n}{2} \log 2\pi \end{aligned} \quad (4.36)$$

taking the negative logarithm

$$= \arg \min_{\theta} \frac{1}{2} \mathbf{y}^\top K_{y,\theta}^{-1} \mathbf{y} + \frac{1}{2} \log \det(K_{y,\theta}) \quad (4.37)$$

the last term is independent of θ

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

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,

$$\frac{\partial}{\partial \theta_j} \log p(y_{1:n} \mid x_{1:n}, \theta) = \frac{1}{2} \text{tr} \left((\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - K_{y,\theta}^{-1}) \frac{\partial K_{y,\theta}}{\partial \theta_j} \right) \quad (4.38)$$

where $\boldsymbol{\alpha} \doteq K_{y,\theta}^{-1} \mathbf{y}$ and $\text{tr}(\mathbf{M})$ is the trace of a matrix \mathbf{M} . This optimization problem is, in general, non-convex. Figure 4.9 gives an example of two local optima according to empirical Bayes.

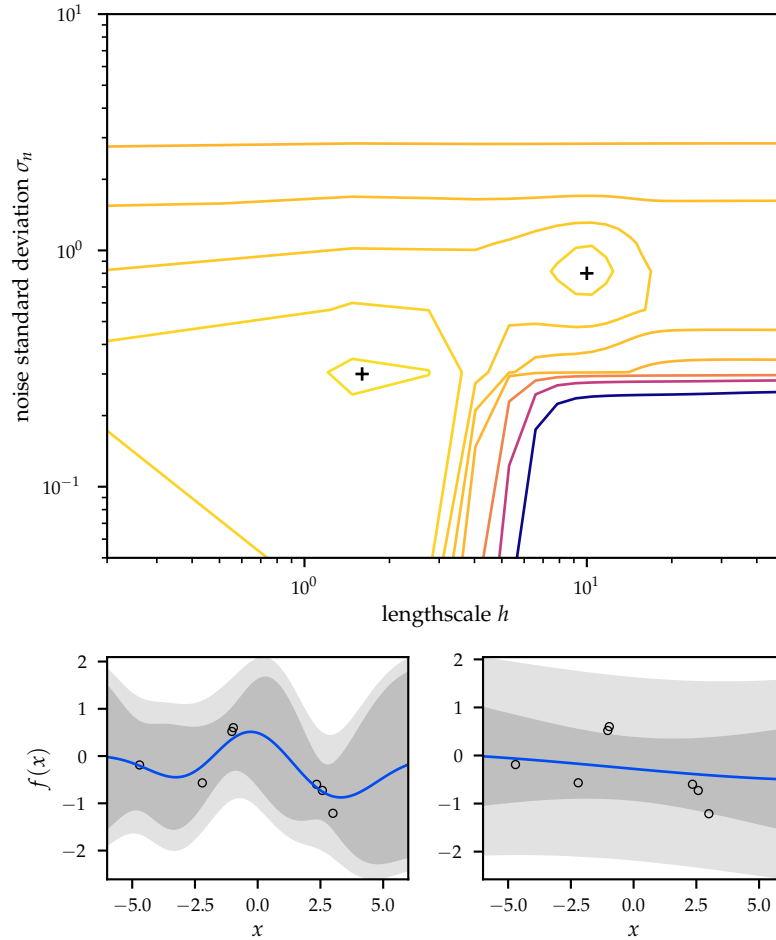


Figure 4.9: The top plot shows contour lines of an empirical Bayes with two local optima. The bottom two plots show the Gaussian processes corresponding to the two optimal models. The left model with smaller lengthscale is chosen within a more flexible class of models, while the right model explains more observations through noise. Adapted from figure 5.5 of “Gaussian processes for machine learning” (Williams and Rasmussen, 2006).

Exercise 4.13: Gradient of the marginal likelihood

In this exercise, we derive eq. (4.38).

Recall that we were considering a dataset (X, y) of noise-perturbed evaluations $y_i = f(x_i) + \epsilon_i$ where $\epsilon_i \sim \mathcal{N}(0, \sigma_n^2)$ 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 $p(y | X, \theta)$.

1. Derive eq. (4.38).

Hint: You can use the following identities:

(a) for any invertible matrix M ,

$$\frac{\partial}{\partial \theta_j} M^{-1} = -M^{-1} \frac{\partial M}{\partial \theta_j} M^{-1} \quad \text{and} \quad (4.39)$$

(b) for any symmetric positive definite matrix M ,

$$\frac{\partial}{\partial \theta_j} \log \det(M) = \text{tr} \left(M^{-1} \frac{\partial M}{\partial \theta_j} \right). \quad (4.40)$$

2. Assume now that the covariance function for the noisy targets (i.e., including the noise contribution) can be expressed as $k_{y,\theta}(x, x') = \theta_0 \tilde{k}(x, x')$ where \tilde{k} is a valid kernel independent of θ_0 .⁹ Show that $\frac{\partial}{\partial \theta_0} \log p(y | X, \theta) = 0$ 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 y by a scalar s ?

▷ *Solution*

Taking a step back, observe that taking a Bayesian 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 θ . Continuing on our Bayesian adventure, we could place a prior $p(\theta)$ on them too.¹⁰ We could use it to obtain the MAP estimate (still a point estimate!) which adds an additional regularization term

$$\hat{\theta}_{\text{MAP}} \doteq \arg \max_{\theta} p(\theta | x_{1:n}, y_{1:n}) \quad (4.41)$$

$$= \arg \min_{\theta} -\log p(\theta) - \log p(y_{1:n} | x_{1:n}, \theta). \quad (4.42)$$

An alternative approach is to consider the full posterior distribution over parameters θ . The resulting predictive distribution is, however,

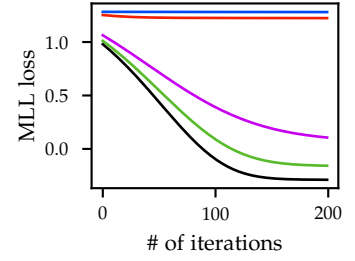


Figure 4.10: An example of model selection by maximizing the log likelihood (without hyperpriors) using a **linear**, **quadratic**, **Laplace**, **Matérn** ($\nu = 3/2$), and **Gaussian** kernel, respectively. They are used to learn the function

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

using SGD with learning rate 0.1.

⁹ That is, $K_{y,\theta}(i, j) = k_{y,\theta}(x_i, x_j)$.

¹⁰ Such a prior is called *hyperprior*.

using Bayes' rule (1.59) and then taking the negative logarithm

intractable,

$$p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}) = \int \int p(y^* | \mathbf{x}^*, f) \cdot p(f | \mathbf{x}_{1:n}, y_{1:n}, \boldsymbol{\theta}) \cdot p(\boldsymbol{\theta}) df d\boldsymbol{\theta}. \quad (4.43)$$

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 Bayesian 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.

Readings

Chapter 5 of “Gaussian processes for machine learning” (Williams and Rasmussen, 2006).

4.5 Approximations

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

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 \mathbf{x} , a very simple approximation is to only condition on those samples \mathbf{x}' that are “close” (where $|k(\mathbf{x}, \mathbf{x}')| \geq \tau$ 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 $\boldsymbol{\phi} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ that approximates the kernel,

$$k(\mathbf{x}, \mathbf{x}') \approx \boldsymbol{\phi}(\mathbf{x})^\top \boldsymbol{\phi}(\mathbf{x}'). \quad (4.44)$$

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

One example of this approach are *random Fourier features*, which we will discuss in the following.

Remark 4.14: 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 $x \in \mathbb{R}$,

$$e^{ix} = \cos x + i \sin x \quad (4.45)$$

where i is the imaginary unit of complex numbers. The formula is illustrated in fig. 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 amplitude of the frequencies ξ ,

$$\hat{f}(\xi) \doteq \int_{\mathbb{R}^d} f(x) e^{-i2\pi\xi^\top x} dx. \quad (4.46)$$

\hat{f} is called the *Fourier transform* of f . f is called the *inverse Fourier transform* of \hat{f} , and can be computed using

$$f(x) = \int_{\mathbb{R}^d} \hat{f}(\xi) e^{i2\pi\xi^\top x} d\xi. \quad (4.47)$$

It is common to write $\omega \doteq 2\pi\xi$. See fig. 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,

$$k(x - x') = \int_{\mathbb{R}^d} p(\omega) e^{i\omega^\top (x - x')} d\omega. \quad (4.48)$$

Fact 4.15 (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.*

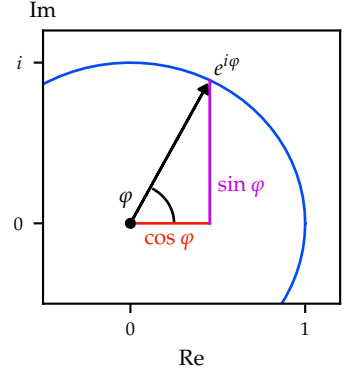


Figure 4.11: Illustration of Euler's formula. It can be seen that $e^{i\varphi}$ corresponds to a (counter-clockwise) rotation on the unit circle as φ varies from 0 to 2π .

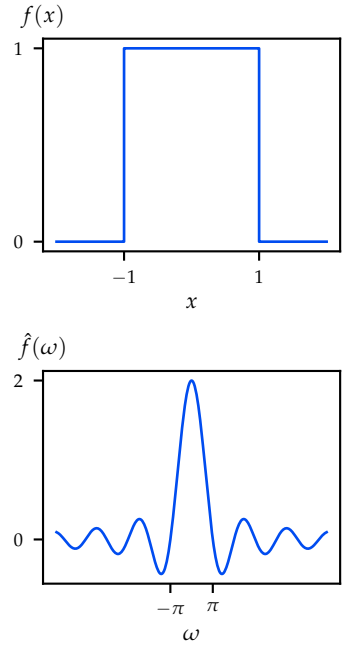


Figure 4.12: The Fourier transform of a rectangular pulse,

$$f(x) \doteq \begin{cases} 1 & x \in [-1, 1] \\ 0 & \text{otherwise,} \end{cases}$$

is given by

$$\begin{aligned} \hat{f}(\omega) &= \int_{-1}^1 e^{-i\omega x} dx = \frac{1}{i\omega} (e^{i\omega} - e^{-i\omega}) \\ &= \frac{2 \sin(\omega)}{\omega}. \end{aligned}$$

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 4.16: Eigenvalue spectrum of stationary kernels

It turns out that $e^{i\omega^\top(x-x')}$ are the eigenfunctions of a stationary kernel with respect to Lebesgue measure. Intuitively, $p(\omega)$ corresponds therefore to the average magnitude of the eigenfunction with frequency ω . The decay of eigenvalues provides information about the smoothness of a kernel: a “rougher” process has more power at high frequencies, and therefore their eigenvalues decay more slowly.

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 4.17: Spectral density of the Gaussian kernel

The Gaussian kernel with length scale ℓ has the spectral density

$$\begin{aligned} p(\omega) &= \int_{\mathbb{R}^d} k(x - x'; \ell) e^{-i\omega^\top(x-x')} d(x - x') \\ &= \int_{\mathbb{R}^d} \exp\left(-\frac{\|x\|_2^2}{2\ell^2} - i\omega^\top x\right) dx \\ &= (2\ell^2\pi)^{d/2} \exp\left(-\ell^2 \frac{\|\omega\|_2^2}{2}\right). \end{aligned} \quad (4.49)$$

using the definition of the Fourier transform (4.46)

using the definition of the Gaussian kernel (4.14)

The key idea is now to interpret the kernel as an expectation,

$$\begin{aligned} k(x - x') &= \int_{\mathbb{R}^d} p(\omega) e^{i\omega^\top(x-x')} d\omega \\ &= \mathbb{E}_{\omega \sim p} \left[e^{i\omega^\top(x-x')} \right] \\ &= \mathbb{E}_{\omega \sim p} \left[\cos(\omega^\top x - \omega^\top x') + i \sin(\omega^\top x - \omega^\top x') \right]. \end{aligned}$$

from eq. (4.48)

by the definition of expectation (1.23b)

using Euler's formula (4.45)

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,

$$\begin{aligned} &= \mathbb{E}_{\omega \sim p} [\cos(\omega^\top x - \omega^\top x')] \\ &= \mathbb{E}_{\omega \sim p} \mathbb{E}_{b \sim \text{Unif}([0, 2\pi])} [\cos((\omega^\top x + b) - (\omega^\top x' + b))] \\ &= \mathbb{E}_{\omega \sim p} \mathbb{E}_{b \sim \text{Unif}([0, 2\pi])} [\cos(\omega^\top x + b) \cos(\omega^\top x' + b) \\ &\quad + \sin(\omega^\top x + b) \sin(\omega^\top x' + b)] \end{aligned}$$

expanding with $b - b$

using the angle subtraction identity, $\cos(\alpha - \beta) = \cos \alpha \cos \beta + \sin \alpha \sin \beta$

$$\begin{aligned}
&= \mathbb{E}_{\omega \sim p} \mathbb{E}_{b \sim \text{Unif}([0, 2\pi])} \left[2 \cos(\omega^\top x + b) \cos(\omega^\top x' + b) \right] \\
&= \mathbb{E}_{\omega \sim p, b \sim \text{Unif}([0, 2\pi])} [z_{\omega, b}(x) \cdot z_{\omega, b}(x')] \quad (4.50)
\end{aligned}$$

where $z_{\omega, b}(x) \doteq \sqrt{2} \cos(\omega^\top x + b)$,

$$\approx \frac{1}{m} \sum_{i=1}^m z_{\omega^{(i)}, b^{(i)}}(x) \cdot z_{\omega^{(i)}, b^{(i)}}(x') \quad (4.51)$$

for independent samples $\omega^{(i)} \stackrel{\text{iid}}{\sim} p$ and $b^{(i)} \stackrel{\text{iid}}{\sim} \text{Unif}([0, 2\pi])$,

$$= z(x)^\top z(x') \quad (4.52)$$

where

$$z(x) \doteq \frac{1}{\sqrt{m}} [z_{\omega^{(1)}, b^{(1)}}(x), \dots, z_{\omega^{(m)}, b^{(m)}}(x)]^\top \quad (4.53)$$

is the (randomized) feature map of random Fourier features.

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 4.18 (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*

$$\begin{aligned}
&\mathbb{P} \left(\sup_{x, x' \in \mathcal{M}} |z(x)^\top z(x') - k(x - x')| \geq \epsilon \right) \\
&\leq 2^8 \left(\frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon} \right)^2 \exp \left(-\frac{m\epsilon^2}{8(d+2)} \right) \quad (4.54)
\end{aligned}$$

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 .

Note that the error probability decays exponentially fast in the dimension of the Fourier feature space.

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 1.36

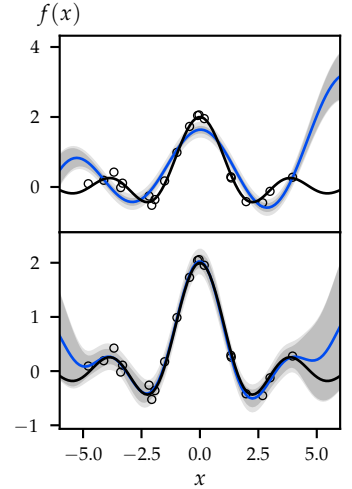


Figure 4.13: Example of random Fourier features with where the number of features m is 5 (top) and 10 (bottom), respectively. The noise-free true function is shown in black and the mean of the Gaussian process is shown in blue.

Exercise 4.19: Uniform convergence of Fourier features

In this exercise, we will prove the above theorem.

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 $T \leq (4 \text{diam}(\mathcal{M})/r)^d$. 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 $\Delta^* = \arg \max_{\Delta \in \mathcal{M}_\Delta} \|\nabla f(\Delta)\|_2$.

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 4.18.

Hint: You may use that

$$(a) \quad \alpha r^{-d} + \beta r^2 = 2\beta^{\frac{d}{d+2}} \alpha^{\frac{2}{d+2}} \text{ for } r = (\alpha/\beta)^{\frac{1}{d+2}} \text{ and}$$

$$(b) \quad \frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon} \geq 1.$$

5. Show that for the Gaussian kernel (4.14), $\sigma_p^2 = \frac{d}{\ell^2}$.

Hint: First show $\sigma_p^2 = -\text{tr}(\mathbf{H}_\Delta k(\mathbf{0}))$.

▷ *Solution*

Remark 4.20: Solving GPs using Kalman filters and smoothers

It turns out that for kernel functions whose spectral density can be represented as the rational function

$$p(\omega) = \frac{\text{constant}}{\text{polynomial in } \omega^2}, \quad (4.55)$$

Kalman filters and smoothers can be used for learning and sampling. In the following, we give a brief summary of this method. For more details, read “Kalman filtering and smoothing solutions

to temporal Gaussian process regression models” (Hartikainen and Särkkä, 2010).

If the spectral density is of this form it can be shown that the random process f can be represented as the output of a linear time invariant stochastic differential equation.¹¹ Such Gaussian processes are also known as *Gaussian Markov Processes*. It can be shown that the Matérn kernel satisfies this property (and hence, GPs with a Matérn kernel can be solved *exactly* using this method); the spectral density of the Gaussian kernel (see eq. (4.49)) can be approximated using a second-order Taylor expansion.

The continuous linear time invariant model can be transformed into a Kalman filter where the motion and sensor models depend on the spectral density of the kernel function (Hartikainen and Särkkä, 2010).

Then, computing the marginal posterior $f \mid x_{1:n}, y_{1:n}$ is simply a Kalman smoothing problem as seen in remark 3.4, and hence, can be solved in linear time in n .

This approach also lends itself to marginal likelihood maximization, as the terms of the optimization objective are computed as by-products of the Kalman filter algorithm.

4.5.3 Data Sampling

Another natural approach is to only consider a (random) subset of the training data during learning. The naïve approach is to subsample uniformly at random. Not very surprisingly, we can do much better.

One subsampling method is the so-called *inducing points method* (Quinonero-Candela and Rasmussen, 2005). The idea is to summarize the data around so-called inducing points.¹² 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,

$$p(f^*, f) = \int_{\mathbb{R}^k} p(f^*, f, u) du = \int_{\mathbb{R}^k} p(f^*, f \mid u) p(u) du, \quad (4.56)$$

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 $u \sim \mathcal{N}(\mathbf{0}, K_{UU})$.

¹¹ This can also be thought of as a first-order continuous (vector) Markov process

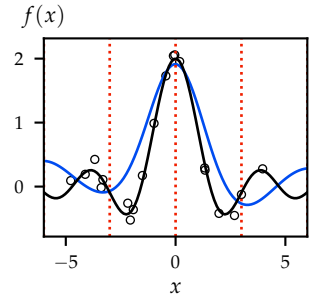


Figure 4.14: Inducing points u are shown as vertical dotted red lines. The noise-free true function is shown in black and the mean of the Gaussian process is shown in blue. Observe that the true function is approximated “well” around the inducing points.

¹² The inducing points can be treated as hyperparameters.

using the sum rule (1.10) and product rule (1.14)

The key idea is to approximate the joint prior, assuming that f^* and f are conditionally independent given \mathbf{u} ,

$$p(f^*, f) \approx \int_{\mathbb{R}^k} p(f^* | \mathbf{u}) p(f | \mathbf{u}) p(\mathbf{u}) d\mathbf{u}. \quad (4.57)$$

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 (1.118),

$$p(f | \mathbf{u}) \sim \mathcal{N}(f; K_{AU} K_{UU}^{-1} \mathbf{u}, K_{AA} - Q_{AA}), \quad (4.58a)$$

$$p(f^* | \mathbf{u}) \sim \mathcal{N}(f^*; K_{\star U} K_{UU}^{-1} \mathbf{u}, K_{\star\star} - Q_{\star\star}) \quad (4.58b)$$

where $Q_{ab} \doteq K_{aU} K_{UU}^{-1} K_{Ub}$. Intuitively, K_{AA} represents the prior covariance and Q_{AA} represents the covariance “explained” by the inducing points.¹³

Computing the full covariance matrix is expensive. In the following, we mention two approximations to the covariance of the training conditional (and testing conditional).

¹³ For more details, refer to section 2 of “A unifying view of sparse approximate Gaussian process regression” (Quinero-Candela and Rasmussen, 2005).

Example 4.21: Subset of regressors

The *subset of regressors* (SoR) approximation is defined as

$$q_{\text{SoR}}(f | \mathbf{u}) \doteq \mathcal{N}(f; K_{AU} K_{UU}^{-1} \mathbf{u}, \mathbf{0}), \quad (4.59a)$$

$$q_{\text{SoR}}(f^* | \mathbf{u}) \doteq \mathcal{N}(f^*; K_{\star U} K_{UU}^{-1} \mathbf{u}, \mathbf{0}). \quad (4.59b)$$

Comparing to eq. (4.58), SoR simply forgets about all variance and covariance.

Exercise 4.22: Subset of regressors

1. Using an SoR approximation, prove the following:

$$q_{\text{SoR}}(f, f^*) = \mathcal{N}\left(\begin{bmatrix} f \\ f^* \end{bmatrix}; \mathbf{0}, \begin{bmatrix} Q_{AA} & Q_{A\star} \\ Q_{\star A} & Q_{\star\star} \end{bmatrix}\right) \quad (4.60)$$

$$q_{\text{SoR}}(f^* | \mathbf{y}) = \mathcal{N}(f^*; Q_{\star A} \tilde{Q}_{AA}^{-1} \mathbf{y}, Q_{\star\star} - Q_{\star A} \tilde{Q}_{AA}^{-1} Q_{A\star}) \quad (4.61)$$

where $\tilde{Q}_{ab} \doteq Q_{ab} + \sigma_n^2$.

2. Derive that the resulting model is a degenerate Gaussian process with covariance function

$$k_{\text{SoR}}(\mathbf{x}, \mathbf{x}') \doteq \mathbf{k}_{\mathbf{x}, U}^\top K_{UU}^{-1} \mathbf{k}_{\mathbf{x}', U}. \quad (4.62)$$

▷ *Solution*

Example 4.23: Fully independent training conditional

The *fully independent training conditional* (FITC) approximation is defined as

$$q_{\text{FITC}}(f | \mathbf{u}) \doteq \mathcal{N}(f; \mathbf{K}_{AU} \mathbf{K}_{UU}^{-1} \mathbf{u}, \text{diag}\{\mathbf{K}_{AA} - \mathbf{Q}_{AA}\}), \quad (4.63a)$$

$$q_{\text{FITC}}(f^* | \mathbf{u}) \doteq \mathcal{N}(f^*; \mathbf{K}_{*U} \mathbf{K}_{UU}^{-1} \mathbf{u}, \text{diag}\{\mathbf{K}_{**} - \mathbf{Q}_{**}\}). \quad (4.63b)$$

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.

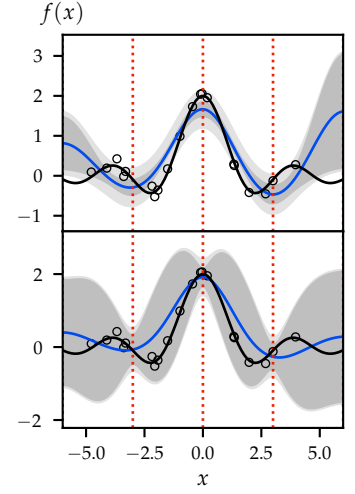


Figure 4.15: Comparison of SoR (top) and FITC (bottom). The inducing points \mathbf{u} are shown as vertical dotted red lines. The noise-free true function is shown in black and the mean of the Gaussian process is shown in blue.

Variational Inference

We have seen how to do learning and 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 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 learning and inference (where the fundamental difficulty lies in solving high-dimensional integrals) to an optimization problem. We have already seen in section 1.5 that optimizing (stochastic) objectives is a well-understood problem with efficient algorithms that perform well in practice.

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 approximation) of the posterior distribution around its mode. Let

$$\psi(\theta) \doteq \log p(\theta \mid x_{1:n}, y_{1:n}) \quad (5.2)$$

denote the log-posterior. Then, using a second-order Taylor approximation (1.104) around the mode $\hat{\theta}$ of ψ (i.e., the MAP estimate), we obtain

$$\begin{aligned} \psi(\theta) &\approx \hat{\psi}(\theta) \doteq \psi(\hat{\theta}) + (\theta - \hat{\theta})^\top \nabla \psi(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^\top H_\psi(\hat{\theta})(\theta - \hat{\theta}) \\ &= \psi(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^\top H_\psi(\hat{\theta})(\theta - \hat{\theta}). \end{aligned} \quad (5.3) \quad \text{using } \nabla \psi(\hat{\theta}) = 0$$

Compare this expression to the log of the PDF of a Gaussian,

$$\log \mathcal{N}(\theta; \hat{\theta}, \Lambda^{-1}) = -\frac{1}{2}(\theta - \hat{\theta})^\top \Lambda(\theta - \hat{\theta}) + \text{const}. \quad (5.4)$$

Thus,

$$\begin{aligned} \hat{\psi}(\theta) &= \psi(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^\top H_\psi(\hat{\theta})(\theta - \hat{\theta}) \\ &= \log \mathcal{N}(\theta; \hat{\theta}, -H_\psi(\hat{\theta})^{-1}) + \text{const}, \end{aligned} \quad (5.5)$$

where we note that $\psi(\hat{\theta})$ is a constant. We therefore define

$$\Lambda \doteq -H_\psi(\hat{\theta}) = -H_\theta \log p(\theta \mid x_{1:n}, y_{1:n})|_{\theta=\hat{\theta}}. \quad (5.6)$$

Observe that Λ is symmetric. Moreover, ψ is a concave function and hence $H_\psi(\theta)$ is negative semi-definite,¹ implying that Λ is positive semi-definite. Using that the inverse of a positive semi-definite matrix is also positive semi-definite, Λ^{-1} is a valid covariance matrix. The *Laplace approximation* q of p is then given by

$$q(\theta) \doteq \mathcal{N}(\theta; \hat{\theta}, \Lambda^{-1}) \propto \exp(\hat{\psi}(\theta)). \quad (5.7)$$

Intuitively, 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 fig. 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 Bayesian inference.

¹ see remark 1.57

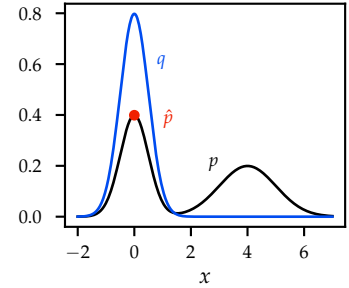


Figure 5.1: The Laplace approximation q greedily selects the mode of the true posterior distribution p and matches the curvature around the mode \hat{p} . As shown here, the Laplace approximation can be extremely overconfident when p is not approximately Gaussian.

Example 5.1: Laplace approximation of a Gaussian

Consider approximating the Gaussian density $p(\theta) = \mathcal{N}(\theta; \mu, \Sigma)$ using a Laplace approximation.

We know that the mode of p is μ , which we can verify by computing the gradient,

$$\nabla_{\theta} \log p(\theta) = -\frac{1}{2}(2\Sigma^{-1}\theta - 2\Sigma^{-1}\mu) \stackrel{!}{=} 0 \iff \theta = \mu. \quad (5.8)$$

For the Hessian of $\log p(\theta)$, we get

$$H_{\theta} \log p(\theta) = (D_{\theta}(\Sigma^{-1}\mu - \Sigma^{-1}\theta))^{\top} = -(\Sigma^{-1})^{\top} = -\Sigma^{-1}. \quad (5.9)$$

Thus, the Laplace approximation of a Gaussian $p(\theta)$ is exact (because the second-order Taylor approximation of $\log p(\theta)$ is exact).

5.1.1 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*,

$$\sigma(z) \doteq \frac{1}{1 + \exp(-z)} \in (0, 1), \quad z = w^{\top}x, \quad (5.10)$$

is used to obtain the class probabilities. *Bayesian logistic regression* corresponds to Bayesian linear regression with a Bernoulli likelihood,

$$y \mid x, w \sim \text{Bern}(\sigma(w^{\top}x)), \quad (5.11)$$

where $y \in \{-1, 1\}$ is the binary class label.² Observe that given a data point (x, y) , the probability of a correct classification is

$$p(y \mid x, w) = \begin{cases} \sigma(w^{\top}x) & \text{if } y = 1 \\ 1 - \sigma(w^{\top}x) & \text{if } y = -1 \end{cases} = \sigma(yw^{\top}x), \quad (5.12)$$

as the logistic function σ is symmetric around 0. Also, recall that Bayesian linear regression used the prior

$$p(w) = \mathcal{N}(w; \mathbf{0}, \sigma_p^2 \mathbf{I}) \propto \exp\left(-\frac{1}{2\sigma_p^2} \|w\|_2^2\right).$$

Let us first find the posterior mode, that is, the MAP estimate of the weights:

$$\hat{w} = \arg \max_w p(w \mid x_{1:n}, y_{1:n})$$

using $(A^{-1})^{\top} = (A^{\top})^{-1}$ and symmetry of Σ

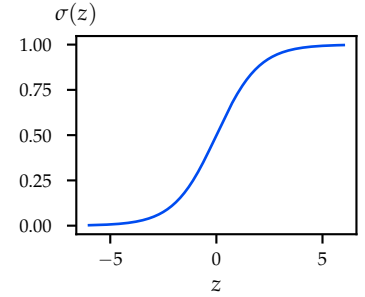


Figure 5.2: The logistic function squashes the linear function $w^{\top}x$ onto the interval $(0, 1)$.

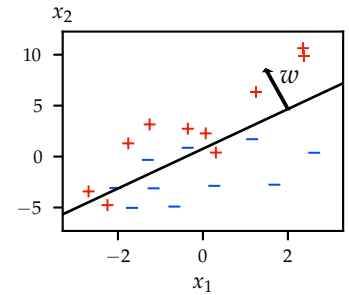


Figure 5.3: Logistic regression classifies data into two classes with a linear decision boundary.

² The same approach extends to Gaussian processes where it is known as *Gaussian process classification*. See exercise 5.5 and “Scalable variational Gaussian process classification” (Hensman et al., 2015).

$$\begin{aligned}
&= \arg \max_w p(\mathbf{w}) p(y_{1:n} \mid \mathbf{x}_{1:n}, \mathbf{w}) && \text{using Bayes' rule (1.59)} \\
&= \arg \max_w \log p(\mathbf{w}) + \log p(y_{1:n} \mid \mathbf{x}_{1:n}, \mathbf{w}) && \text{taking the logarithm} \\
&= \arg \max_w -\frac{1}{2\sigma_p^2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \log \sigma(y_i \mathbf{w}^\top \mathbf{x}_i) && \text{using independence of the observations} \\
&&& \text{and eq. (5.12)} \\
&= \arg \min_w \frac{1}{2\sigma_p^2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)). && \text{using the definition of } \sigma \text{ (5.10)} \tag{5.13}
\end{aligned}$$

Note that for $\lambda = 1/2\sigma_p^2$, the above optimization is equivalent to standard (regularized) logistic regression where

$$\ell_{\log}(\mathbf{w}^\top \mathbf{x}; y) \doteq \log(1 + \exp(-y \mathbf{w}^\top \mathbf{x})) \tag{5.14}$$

is called *logistic loss*. You show in exercise 5.3 that the gradient of the logistic loss is given by

$$\nabla_{\mathbf{w}} \ell_{\log}(\mathbf{w}^\top \mathbf{x}; y) = -y \mathbf{x} \cdot \sigma(-y \mathbf{w}^\top \mathbf{x}). \tag{5.15}$$

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,

$$\mathbf{w} \leftarrow \mathbf{w}(1 - 2\lambda\eta_t) + \eta_t y \mathbf{x} \sigma(-y \mathbf{w}^\top \mathbf{x}), \tag{5.16}$$

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

$$\pi_i \doteq \mathbb{P}(y_i = 1 \mid \mathbf{x}_i, \hat{\mathbf{w}}) = \sigma(\hat{\mathbf{w}}^\top \mathbf{x}_i) \tag{5.17}$$

the probability of \mathbf{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

$$\begin{aligned}
\Lambda &= -\mathbf{H}_{\mathbf{w}} \log p(\mathbf{w} \mid \mathbf{x}_{1:n}, y_{1:n})|_{\mathbf{w}=\hat{\mathbf{w}}} \\
&= -\mathbf{H}_{\mathbf{w}} \log p(y_{1:n} \mid \mathbf{x}_{1:n}, \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} - \mathbf{H}_{\mathbf{w}} \log p(\mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} \\
&= \sum_{i=1}^n \mathbf{H}_{\mathbf{w}} \ell_{\log}(\mathbf{w}^\top \mathbf{x}_i; y_i) \Big|_{\mathbf{w}=\hat{\mathbf{w}}} + \sigma_p^{-2} \mathbf{I}
\end{aligned}$$

using the definition of the logistic loss (5.14)

$$\begin{aligned}
 &= \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \pi_i (1 - \pi_i) + \sigma_p^{-2} \mathbf{I} \\
 &= \mathbf{X}^\top \text{diag}_{i \in [n]} \{\pi_i (1 - \pi_i)\} \mathbf{X} + \sigma_p^{-2} \mathbf{I}.
 \end{aligned} \tag{5.18}$$

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, $\mathbf{\Lambda}$ does not depend on the normalization constant of the posterior distribution which is hard to compute.

Overall, we have that $\mathcal{N}(\hat{\mathbf{w}}, \mathbf{\Lambda}^{-1})$ is the Laplace approximation of $p(\mathbf{w} \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$.

using the Hessian of the logistic loss
(5.19)

Exercise 5.3: Logistic loss

1. Derive the gradient of ℓ_{\log} as given in eq. (5.15).
2. Show that

$$\mathbf{H}_{\mathbf{w}} \ell_{\log}(\mathbf{w}^\top \mathbf{x}; y) = \mathbf{x} \mathbf{x}^\top \cdot \sigma(\mathbf{w}^\top \mathbf{x}) \cdot (1 - \sigma(\mathbf{w}^\top \mathbf{x})). \tag{5.19}$$

Hint: Begin by deriving the first derivative of the logistic function, and use the chain rule of multivariate calculus,

$$\underbrace{D_{\mathbf{x}}(f \circ g)}_{\mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}} = \underbrace{(D_{g(\mathbf{x})} f \circ g)}_{\mathbb{R}^n \rightarrow \mathbb{R}^{m \times k} \cdot \mathbb{R}^{k \times n}} \cdot D_{\mathbf{x}} g \tag{5.20}$$

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} ?

▷ *Solution*

5.2 Inference with a Variational Posterior

How can we perform inference using our variational approximation? We simply approximate the (intractable) true posterior with our variational posterior,

$$\begin{aligned}
 p(y^* \mid \mathbf{x}^*, \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) &= \int p(y^* \mid \mathbf{x}^*, \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) d\boldsymbol{\theta} \\
 &\approx \int p(y^* \mid \mathbf{x}^*, \boldsymbol{\theta}) q_{\lambda}(\boldsymbol{\theta}) d\boldsymbol{\theta}
 \end{aligned} \tag{5.21}$$

using the sum rule (1.10)

At first sight, it may seem as though we did not gain much by our approximation. However, observe that interpreting the integral from a function-space view, the final prediction y^* is conditionally independent from the model parameters $\boldsymbol{\theta}$ given the “latent value” f^* :³

$$= \int \int p(y^* \mid f^*) p(f^* \mid \mathbf{x}^*, \boldsymbol{\theta}) q_{\lambda}(\boldsymbol{\theta}) d\boldsymbol{\theta} df^*$$

³ This is discussed in greater detail in section 2.6.

once more, using the sum rule (1.10)

$$= \int p(y^* | f^*) \int p(f^* | \mathbf{x}^*, \boldsymbol{\theta}) q_\lambda(\boldsymbol{\theta}) d\boldsymbol{\theta} df^* \quad (5.22) \quad \text{rearranging terms}$$

Generally, $f^* = g(\boldsymbol{\theta}; \mathbf{x}^*)$ where g is a deterministic function representing our “model” parameterized by $\boldsymbol{\theta}$. Then, $p(f^* | \mathbf{x}^*, \boldsymbol{\theta}) = \delta_{g(\boldsymbol{\theta}; \mathbf{x}^*)}(f^*)$ is a point density (cf. example 1.18) and

$$(g(\cdot; \mathbf{x}^*)_{\#} q_\lambda)(f^*) = \int p(f^* | \mathbf{x}^*, \boldsymbol{\theta}) q_\lambda(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (5.23)$$

using sum rule (1.10) and product rule (1.14) to condition on $\boldsymbol{\theta}$

is the pushforward of $q_\lambda(\boldsymbol{\theta})$ under perturbation $g(\cdot; \mathbf{x}^*)$. We discuss such a “change of variables” in greater detail in section 1.1.12. Alternatively, the integral in eq. (5.23) can be viewed as an expectation over the variational posterior q_λ and approximated using Monte Carlo sampling.

Example 5.4: Inference for Bayesian logistic regression

In the case of Bayesian logistic regression with a Laplace approximation, the posterior over weights $q_\lambda(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \hat{\mathbf{w}}, \boldsymbol{\Lambda}^{-1})$ is a Gaussian and $f^* = g(\mathbf{w}; \mathbf{x}^*) = \mathbf{w}^\top \mathbf{x}^*$ given fixed weights \mathbf{w} . We therefore have due to the closedness properties of Gaussians (1.123),

$$g(\cdot; \mathbf{x}^*)_{\#} q_\lambda = \mathcal{N}(\hat{\mathbf{w}}^\top \mathbf{x}^*, \mathbf{x}^{*\top} \boldsymbol{\Lambda}^{-1} \mathbf{x}^*). \quad (5.24)$$

Crucially, this is a one-dimensional Gaussian!

As we have seen in eq. (5.12), for Bayesian logistic regression, the prediction y^* depends deterministically on the predicted “latent value” f^* : $p(y^* | f^*) = \sigma(y^* f^*)$. Combining equations (5.22), (5.23), and (5.24), we obtain

$$p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}) \approx \int \sigma(y^* f^*) \cdot \mathcal{N}(f^*; \hat{\mathbf{w}}^\top \mathbf{x}^*, \mathbf{x}^{*\top} \boldsymbol{\Lambda}^{-1} \mathbf{x}^*) df^*. \quad (5.25)$$

We have replaced the (high-dimensional) integral over the model parameters $\boldsymbol{\theta}$ 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.

Exercise 5.5: 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 ,

$$f \sim \mathcal{GP}(0, k), \quad y \mid x, f \sim \text{Bern}(\sigma(f(x))) \quad (5.26)$$

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 \mathbf{x}^* , we first compute the distribution of the latent variable f^* ,

$$p(f^* \mid \mathbf{x}_{1:n}, y_{1:n}, \mathbf{x}^*) = \int p(f^* \mid \mathbf{x}_{1:n}, \mathbf{x}^*, f) p(f \mid \mathbf{x}_{1:n}, y_{1:n}) df \quad (5.27)$$

using sum rule (1.10) and product rule (1.14), and $f^* \perp y_{1:n} \mid f$

where $p(f \mid \mathbf{x}_{1:n}, y_{1:n})$ is the posterior over the latent variables.

1. Assuming that we can efficiently compute $p(f^* \mid \mathbf{x}_{1:n}, y_{1:n}, \mathbf{x}^*)$ (approximately), describe how we can find the predictive posterior $p(y^* = +1 \mid \mathbf{x}_{1:n}, y_{1:n}, \mathbf{x}^*)$.
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 (5.27) is analytically intractable. A common technique is to approximate the latent posterior $p(f \mid \mathbf{x}_{1:n}, y_{1:n})$ with a Gaussian using a Laplace approximation $q \doteq \mathcal{N}(\hat{f}, \Lambda^{-1})$. 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 $p(y_i \mid f_i) = \sigma(y_i f_i)$, where we use the symmetry of the probit likelihood around 0.
 - (b) Assume that $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$ is the linear kernel ($\sigma_p = 1$) and that σ is the logistic function (5.10). Show for this setting that the matrix Λ derived in (a) is equivalent to the precision matrix Λ' of the Laplace approximation of

Bayesian logistic regression (5.18).⁴ 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.⁵ 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,

$$f \sim \mathcal{GP}(0, k), \quad y = \mathbb{1}\left\{ \underbrace{f(\mathbf{x}) + \epsilon}_{\text{GP regression}} \geq 0 \right\}, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2). \quad (5.28)$$

Show that the model from eq. (5.26) 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 eq. (5.28).

▷ *Solution*

⁴ This should not be surprising since — as already mentioned — Gaussian process classification is a generalization of Bayesian logistic regression.

⁵ Using the Laplace-approximated latent posterior, $[f^* \ f]$ are jointly Gaussian. Thus, it directly follows from theorem 1.66 that the marginal distribution over f^* is also a Gaussian.

5.3 Blueprint of Variational Inference

General Bayesian inference poses the challenge of approximating the optimal Bayesian 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 trade-off 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 Bayesian 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 benefit is that we can reduce the (generally intractable) problem of high-dimensional integration to the (often much more tractable) problem of optimization.

Definition 5.6 (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 5.7: Family of independent Gaussians

A straightforward example for a variational family is the family of independent Gaussians,

$$\mathcal{Q} \doteq \left\{ q(\theta) = \mathcal{N}(\theta; \mu, \text{diag}_{i \in [d]} \{\sigma_i^2\}) \right\}, \quad (5.29)$$

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!

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:

$$q^* \doteq \arg \min_{q \in \mathcal{Q}} \text{KL}(q\|p) = \arg \min_{\lambda \in \Lambda} \text{KL}(q_\lambda\|p). \quad (5.30)$$

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 eq. (5.30) can be solved efficiently.

5.4 Information Theory

One of our main objectives throughout this course 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 .

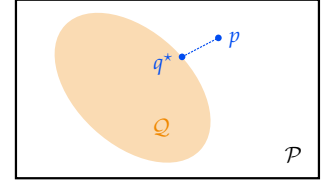


Figure 5.4: An illustration of variational inference in the space of distributions \mathcal{P} . The variational distribution $q^* \in \mathcal{Q}$ is the optimal approximation of the true posterior p .

5.4.1 Surprise

The *surprise* about an event with probability u is defined as

$$S[u] \doteq -\log u. \quad (5.31)$$

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 5.8 (Axiomatic characterization of surprise). *The axioms*

1. $S[1] = 0$
2. $S[u] > S[v] \implies u < v$ (anti-monotonicity)
3. S is continuous,
4. $S[uv] = S[u] + S[v]$ for independent events,

characterize S up to a positive constant factor.

Proof. Cauchy's functional equation, $f(x+y) = f(x) + f(y)$, 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

$$f(e^x e^y) = f(e^{x+y}) = g(x+y) = g(x) + g(y) = f(e^x) + f(e^y).$$

It therefore follows from the third and fourth axioms of surprise that $S[u] = c \log u$ for any $c \in \mathbb{R}$. The second axiom of surprise implies that $c < 0$, and thus, $S[u] = -c' \log u$ for any $c' > 0$. Finally, observe that the first axiom of surprise is satisfied for any such c' . \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 fig. 5.6. Information theory is the study of uncertainty in terms of surprise.

Throughout this course 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 already encountered the "surprise space" was in the context of likelihood maximization (cf. section 1.4.3) where we used that the log-transform linearizes products of probabilities. We will see later in section 6.3.2 that in many cases the surprise $S[p(x)]$ can also be interpreted as a "cost" or "energy" associated with the state x .

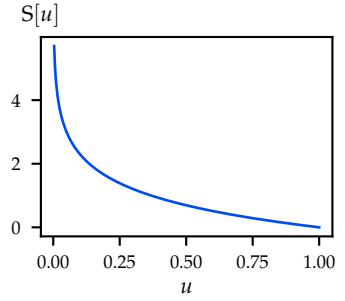


Figure 5.5: Surprise $S[u]$ associated with an event of probability u .

certain events are not surprising
we are more surprised by unlikely events
no jumps in surprise for infinitesimal changes of probability

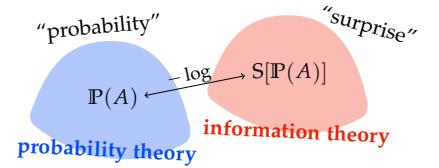


Figure 5.6: Illustration of the probability space and the corresponding "surprise space".

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,

$$H[p] \doteq \mathbb{E}_{x \sim p}[S[p(x)]] = \mathbb{E}_{x \sim p}[-\log p(x)]. \quad (5.32)$$

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 $p(x) \leq 1$ ($\forall x$).⁶ For discrete distributions it is common to use the logarithm with base 2 rather than the natural logarithm.⁷ Thus, entropy is given by

$$H[p] = -\sum_x p(x) \log_2 p(x) \quad (\text{if } p \text{ is discrete}), \quad (5.33a)$$

$$H[p] = -\int p(x) \log p(x) dx \quad (\text{if } p \text{ is continuous}). \quad (5.33b)$$

From the fourth axiom of surprise, we immediately obtain the following property of entropy.

Corollary 5.9. *For a multivariate distribution p that can be factorized into $p(x_{1:d}) = \prod_{i=1}^d p_i(x_i)$, we have that*

$$H[p] = \sum_{i=1}^d H[p_i]. \quad (5.34)$$

Example 5.10: 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 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

⁶ The entropy of a continuous distribution can be negative. For example,

$$H[\text{Unif}([a, b])] = -\int \frac{1}{b-a} \log \frac{1}{b-a} dx = \log(b-a)$$

which is negative if $b-a < 1$.

⁷ Recall that $\log_2 x = \frac{\log x}{\log 2}$, that is, logarithms to a different base only differ by a constant factor.

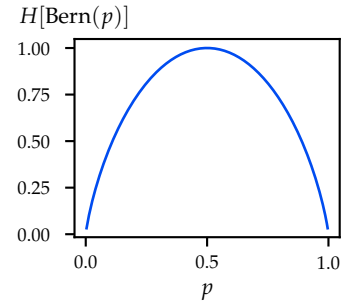


Figure 5.7: Entropy $H[\text{Bern}(p)]$ of a Bernoulli experiment with success probability p .

sample $x \sim p$, or in other words, the average “information” carried by a sample x .

Example 5.11: 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 (5.33b), we obtain,

$$\begin{aligned} H[\mathcal{N}(\mu, \sigma^2)] &= - \int \frac{1}{Z} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \\ &\quad \cdot \log\left(\frac{1}{Z} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)\right) dx \\ &= \log Z \underbrace{\int \frac{1}{Z} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) dx}_1 \\ &\quad + \int \frac{1}{Z} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \frac{(x - \mu)^2}{2\sigma^2} dx \\ &= \log Z + \frac{1}{2\sigma^2} \mathbb{E}[(x - \mu)^2] \\ &= \log(\sigma\sqrt{2\pi}) + \frac{1}{2} \\ &= \log(\sigma\sqrt{2\pi e}). \end{aligned} \quad (5.35)$$

In general, the entropy of a Gaussian is

$$H[\mathcal{N}(\mu, \Sigma)] = \frac{1}{2} \log \det(2\pi e \Sigma) = \frac{1}{2} \log((2\pi e)^d \det(\Sigma)). \quad (5.36)$$

Observe that the surprise $S[u]$ is convex in u . Jensen’s inequality is a useful tool when working with expectations of convex functions such as entropy.

Fact 5.12 (Jensen’s Inequality). *Given a random variable $X : \Omega \rightarrow \mathbb{R}$ and a convex function $g : \mathbb{R} \rightarrow \mathbb{R}$, we have*

$$g(\mathbb{E}[X]) \leq \mathbb{E}[g(X)]. \quad (5.37)$$

For a concave function h , $-h$ is convex, and hence

$$h(\mathbb{E}[X]) \geq \mathbb{E}[h(X)]. \quad (5.38)$$

using LOTUS (1.30)

using $\mathbb{E}[(x - \mu)^2] = \text{Var}[x] = \sigma^2$ (1.44)

using $\log \sqrt{e} = 1/2$

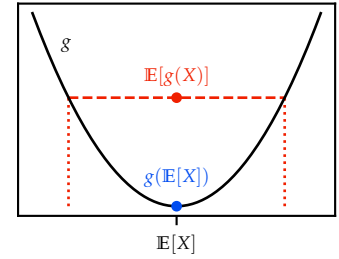


Figure 5.8: An illustration of Jensen’s inequality. Due to the convexity of g , we have that g evaluated at $\mathbb{E}[X]$ will always be below the average of evaluations of g .

Exercise 5.13: Jensen's inequality

1. Prove the finite form of Jensen's inequality.

Theorem 5.14 (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,*

$$f(\theta_1 \mathbf{x}_1 + \dots + \theta_k \mathbf{x}_k) \leq \theta_1 f(\mathbf{x}_1) + \dots + \theta_k f(\mathbf{x}_k). \quad (5.39)$$

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 , $H[p] \leq \log_2 n$.

▷ *Solution*

Readings

To read more about entropy, refer to chapter 2 of “Elements of information theory” (Cover and Thomas, 2006). In chapter 8, we introduce the related concepts of conditional entropy and mutual information.

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 5.15 (Cross-entropy). The *cross-entropy* of a distribution q relative to the distribution p is

$$H[p||q] \doteq \mathbb{E}_{x \sim p}[S[q(x)]] = \mathbb{E}_{x \sim p}[-\log q(x)]. \quad (5.40)$$

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 ,

$$H[p||q] = H[p] + KL(p||q) \geq H[p]. \quad (5.41)$$

$KL(p||q) \geq 0$ is shown in exercise 5.19

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.

Remark 5.16: Minimizing cross-entropy

Minimizing cross-entropy is a commonly used objective for *density estimation* problems, where we want to learn some true distribution p using a parameterized distribution q_θ . As $H[p]$ is fixed, minimizing cross-entropy corresponds to minimizing the KL-divergence between p and q_θ .

For 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} ,

$$\begin{aligned}\ell_{\text{bce}}(\hat{y}; y) &\doteq H[\text{Bern}(y) \parallel \text{Bern}(\hat{y})] \\ &= - \sum_{x \in \{0, 1\}} \text{Bern}(x; y) \log \text{Bern}(x; \hat{y}) \\ &= -y \log \hat{y} - (1 - y) \log(1 - \hat{y}).\end{aligned}\tag{5.42}$$

This loss function is also known as the *binary cross-entropy loss*.

Exercise 5.17: Binary cross-entropy loss

Show that the logistic loss (5.14) is equivalent to the binary cross-entropy loss with $\hat{y} = \sigma(\hat{f})$. That is,

$$\ell_{\log}(\hat{f}; y) = \ell_{\text{bce}}(\hat{y}; y).\tag{5.43}$$

▷ *Solution*

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 5.18 (Kullback-Leibler divergence, KL-divergence). Given two distributions p and q , the *Kullback-Leibler divergence* (or *relative entropy*) of q with respect to p ,

$$\text{KL}(p \parallel q) \doteq H[p \parallel q] - H[p]\tag{5.44}$$

$$= \mathbb{E}_{\theta \sim p}[S[q(\theta)] - S[p(\theta)]]\tag{5.45}$$

$$= \mathbb{E}_{\theta \sim p} \left[\log \frac{p(\theta)}{q(\theta)} \right],\tag{5.46}$$

measures how different q is from a reference distribution p .

In words, $\text{KL}(p \parallel q)$ measures the *additional* expected surprise when observing samples from p that is due to assuming the (wrong) distribu-

tion q and which not inherent in the distribution p already.⁸ Intuitively, the KL-divergence measures the information loss when approximating p with q .

⁸The KL-divergence only captures the additional expected surprise as the surprise inherent in p (as measured by $H[p]$) is subtracted.

The KL-divergence has the following properties:

- $\text{KL}(p\|q) \geq 0$ for any distributions p and q ,
- $\text{KL}(p\|q) = 0$ if and only if $p = q$ almost surely, and
- there exist distributions p and q such that $\text{KL}(p\|q) \neq \text{KL}(q\|p)$.

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.

Exercise 5.19: 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

$$f(\theta_1 x_1 + \dots + \theta_k x_k) \leq \theta_1 f(x_1) + \dots + \theta_k f(x_k) \quad (5.47)$$

iff $x_1 = \dots = x_k$. This is a slight adaptation of Jensen's inequality in finite-form, which you proved in exercise 5.13.

▷ *Solution*

Exercise 5.20: 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 $H[g] \geq H[f]$.

▷ *Solution*

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 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 $p(\boldsymbol{\theta}_{1:n}) > q(\boldsymbol{\theta}_{1:n})$ and vice versa. Assuming that the samples are independent and rewriting the inequality slightly, we choose p if

$$\prod_{i=1}^n \frac{p(\boldsymbol{\theta}_i)}{q(\boldsymbol{\theta}_i)} > 1, \quad \text{or equivalently if} \quad \sum_{i=1}^n \log \frac{p(\boldsymbol{\theta}_i)}{q(\boldsymbol{\theta}_i)} > 0.$$

Assume that $\boldsymbol{\theta}_i \sim p$. Then, using the law of large numbers (1.84),

$$\frac{1}{n} \sum_{i=1}^n \log \frac{p(\boldsymbol{\theta}_i)}{q(\boldsymbol{\theta}_i)} \xrightarrow{\text{a.s.}} \mathbb{E}_{\boldsymbol{\theta} \sim p} \left[\log \frac{p(\boldsymbol{\theta})}{q(\boldsymbol{\theta})} \right] = \text{KL}(p \| q) \quad (5.48)$$

as $n \rightarrow \infty$. Recall that assuming $p \neq q$, we have that $\text{KL}(p \| q) > 0$ with probability 1. In particular, this shows that as $n \rightarrow \infty$, we correctly choose p with probability 1.

Example 5.21: 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

$$\begin{aligned} \text{KL}(\text{Bern}(p) \| \text{Bern}(q)) &= \sum_{x \in \{0,1\}} \text{Bern}(x; p) \log \frac{\text{Bern}(x; p)}{\text{Bern}(x; q)} \\ &= p \log \frac{p}{q} + (1-p) \log \frac{(1-p)}{(1-q)}. \end{aligned} \quad (5.49)$$

Exercise 5.22: 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 . Show that the KL-divergence of p and q is given by

$$\begin{aligned} \text{KL}(p \| q) &= \frac{1}{2} \left(\text{tr}(\boldsymbol{\Sigma}_q^{-1} \boldsymbol{\Sigma}_p) + (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)^\top \boldsymbol{\Sigma}_q^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q) \right. \\ &\quad \left. - d + \log \frac{\det(\boldsymbol{\Sigma}_q)}{\det(\boldsymbol{\Sigma}_p)} \right). \end{aligned} \quad (5.50)$$

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 $\mathbf{A} \in \mathbb{R}^{d \times d}$,

$$\mathbb{E} \left[(\mathbf{X} - \mathbf{m})^\top \mathbf{A} (\mathbf{X} - \mathbf{m}) \right] = (\boldsymbol{\mu} - \mathbf{m})^\top \mathbf{A} (\boldsymbol{\mu} - \mathbf{m}) + \text{tr}(\mathbf{A} \boldsymbol{\Sigma}) \quad (5.51)$$

▷ *Solution*

Example 5.23: KL-divergence of independent Gaussians

For independent Gaussians with unit variance $p \doteq \mathcal{N}(\mu_p, I)$ and $q \doteq \mathcal{N}(\mu_q, I)$, the expression simplifies to the squared Euclidean distance,

$$\text{KL}(p\|q) = \frac{1}{2} \|\mu_q - \mu_p\|_2^2. \quad (5.52)$$

If we approximate independent Gaussians with variances σ_i^2 ,

$$p \doteq \mathcal{N}(\mu, \text{diag}\{\sigma_1^2, \dots, \sigma_d^2\}),$$

by a standard normal distribution, $q \doteq \mathcal{N}(\mathbf{0}, I)$, the expression simplifies to

$$\text{KL}(p\|q) = \frac{1}{2} \sum_{i=1}^d (\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2). \quad (5.53)$$

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,

$$q_1^* \doteq \arg \min_{q \in \mathcal{Q}} \text{KL}(p\|q) \quad q_2^* \doteq \arg \min_{q \in \mathcal{Q}} \text{KL}(q\|p)$$

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 prediction. 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 (5.30) 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

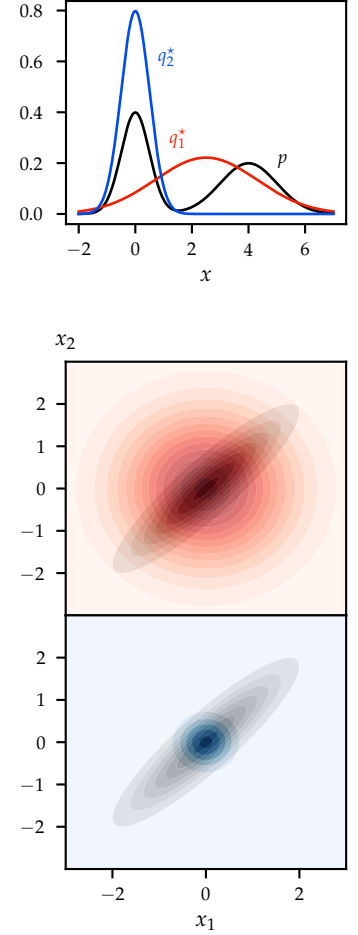


Figure 5.9: Comparison of the **forward** KL-divergence q_1^* and the **reverse** KL-divergence q_2^* when used to approximate the **true posterior** p . The first plot shows the PDFs in a one-dimensional feature space where p is a mixture of two univariate Gaussians. The second plot shows contour lines 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^* .

$p(\cdot \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n})$ is in the variational family \mathcal{Q} , then

$$\arg \min_{q \in \mathcal{Q}} \text{KL}(q \| p(\cdot \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n})) \stackrel{\text{a.s.}}{=} p(\cdot \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) \quad (5.54)$$

as $\min_{q \in \mathcal{Q}} \text{KL}(q \| p(\cdot \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:n})) \stackrel{\text{a.s.}}{=} 0$, so minimizing reverse-KL still recovers the true posterior almost surely.

Remark 5.24: 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 (5.50), we obtain for the reverse KL-divergence,

$$\begin{aligned} \text{KL}(q \| p) &= \frac{1}{2} \left(\text{tr}(\text{diag}\{\sigma_i^{-2}\}) + \boldsymbol{\mu}^\top \text{diag}\{\sigma_i^{-2}\} \boldsymbol{\mu} - d \right. \\ &\quad \left. + \log \det(\text{diag}\{\sigma_i^2\}) \right) \\ &= \frac{1}{2} \sum_{i=1}^d \left(\sigma_i^{-2} + \frac{\mu_i^2}{\sigma_i^2} - 1 + \log \sigma_i^2 \right). \end{aligned} \quad (5.55)$$

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 eq. (5.53). 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 .

Exercise 5.25: 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 $q(x, y) = p(x)p(y)$?

▷ *Solution*

5.4.6 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 relate minimizing *forward-KL* to two other well-known inference algorithms.

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 $p(x)$ is a generative model, and we aim to learn the parameterized model likelihood q_λ .⁹

Lemma 5.26 (Forward KL-divergence as MLE). *Given some generative model $p(x)$ and a likelihood $q_\lambda(x) = q(x | \lambda)$ (that we use to approximate the true data distribution), we have*

$$\arg \min_{\lambda \in \Lambda} \text{KL}(p\|q_\lambda) \stackrel{\text{a.s.}}{=} \arg \max_{\lambda \in \Lambda} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \log q(x^{(i)} | \lambda), \quad (5.56)$$

where $x^{(i)} \stackrel{\text{iid}}{\sim} p$ are independent samples from the true data distribution.

Proof.

$$\begin{aligned} \text{KL}(p\|q_\lambda) &= H[p\|q_\lambda] - H[p] \\ &= \mathbb{E}_{(x,y) \sim p} [-\log q(x | \lambda)] + \text{const} \\ &\stackrel{\text{a.s.}}{=} - \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \log q(x^{(i)} | \lambda) + \text{const} \end{aligned}$$

where $x^{(i)} \stackrel{\text{iid}}{\sim} p$ are independent samples. □

⁹ Note that here, we aim to learn model parameters λ for estimating the probability of x , whereas in the setting of variational Bayesian inference, we want to learn parameters λ of a distribution over θ and θ parameterizes a distribution over x .

using the definition of KL-divergence (5.44)

dropping $H[p]$ and using the definition of cross-entropy (5.40)

using Monte Carlo sampling, i.e., the law of large numbers (1.84)

This interpretation is not immediately useful for Bayesian inference (i.e., in the setting where p is a posterior distribution over model parameters θ) as a maximum likelihood estimate requires samples from p which we cannot obtain in this case.¹⁰

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 \mathbf{a} and \mathbf{B} of the first and second moment of p ,¹¹ and if q_λ is a Gaussian with parameters $\lambda = \{\mu, \Sigma\}$, then moment matching chooses λ as the solution to

$$\begin{aligned}\mathbb{E}_p[\theta] &\approx \mathbf{a} \stackrel{!}{=} \mu = \mathbb{E}_{q_\lambda}[\theta] \\ \mathbb{E}_p[\theta\theta^\top] &\approx \mathbf{B} \stackrel{!}{=} \Sigma + \mu\mu^\top = \mathbb{E}_{q_\lambda}[\theta\theta^\top].\end{aligned}$$

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

$$\mathcal{N}(\theta; \mu, \Sigma) = \frac{1}{Z(\lambda)} \exp(\lambda^\top s(\theta)) \quad \text{where} \quad (5.57)$$

$$\lambda \doteq \begin{bmatrix} \Sigma^{-1}\mu \\ \text{vec}[\Sigma^{-1}] \end{bmatrix} \quad (5.58)$$

$$s(\theta) \doteq \begin{bmatrix} \theta \\ \text{vec}[-\frac{1}{2}\theta\theta^\top] \end{bmatrix} \quad (5.59)$$

and $Z(\lambda) \doteq \int \exp(\lambda^\top s(\theta)) d\theta$, and we will confirm this in just a moment.¹² The family of distributions with densities of the form (5.57) — with an additional scaling constant $h(\theta)$ which is often 1 — is called 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 eqs. (5.58) and (5.59), consider

$$\begin{aligned}\mathcal{N}(\theta; \mu, \Sigma) &\propto \exp\left(-\frac{1}{2}(\theta - \mu)^\top \Sigma^{-1}(\theta - \mu)\right) \\ &\propto \exp\left(\text{tr}\left(-\frac{1}{2}\theta^\top \Sigma^{-1}\theta\right) + \theta^\top \Sigma^{-1}\mu\right)\end{aligned}$$

¹⁰ It is possible to obtain “approximate” samples using Markov chain Monte Carlo (MCMC) methods which we discuss in chapter 6.

¹¹ 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 (1.45)

¹² 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 $n \cdot m$.

expanding the inner product and using that $\text{tr}(x) = x$ for all $x \in \mathbb{R}$

$$\begin{aligned}
 &= \exp\left(\text{tr}\left(-\frac{1}{2}\boldsymbol{\theta}\boldsymbol{\theta}^\top\boldsymbol{\Sigma}^{-1}\right) + \boldsymbol{\theta}^\top\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}\right) \\
 &= \exp\left(\text{vec}\left[-\frac{1}{2}\boldsymbol{\theta}\boldsymbol{\theta}^\top\right]^\top\text{vec}\left[\boldsymbol{\Sigma}^{-1}\right] + \boldsymbol{\theta}^\top\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}\right).
 \end{aligned}$$

using that the trace is invariant under cyclic permutations

using $\text{tr}(\mathbf{AB}) = \text{vec}[\mathbf{A}]^\top \text{vec}[\mathbf{B}]$ for any $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$

This allows us to express the forward KL-divergence as

$$\begin{aligned}
 \text{KL}(p\|q_\lambda) &= \int p(\boldsymbol{\theta}) \log \frac{p(\boldsymbol{\theta})}{q_\lambda(\boldsymbol{\theta})} d\boldsymbol{\theta} \\
 &= - \int p(\boldsymbol{\theta}) \cdot \boldsymbol{\lambda}^\top \mathbf{s}(\boldsymbol{\theta}) d\boldsymbol{\theta} + \log Z(\boldsymbol{\lambda}) + \text{const.}
 \end{aligned}$$

using that $\int p(\boldsymbol{\theta}) \log p(\boldsymbol{\theta}) d\boldsymbol{\theta}$ is constant

Differentiating with respect to the natural parameters $\boldsymbol{\lambda}$ gives

$$\begin{aligned}
 \nabla_{\boldsymbol{\lambda}} \text{KL}(p\|q_\lambda) &= - \int p(\boldsymbol{\theta}) \mathbf{s}(\boldsymbol{\theta}) d\boldsymbol{\theta} + \frac{1}{Z(\boldsymbol{\lambda})} \int \mathbf{s}(\boldsymbol{\theta}) \exp(\boldsymbol{\lambda}^\top \mathbf{s}(\boldsymbol{\theta})) d\boldsymbol{\theta} \\
 &= -\mathbb{E}_{\boldsymbol{\theta} \sim p}[\mathbf{s}(\boldsymbol{\theta})] + \mathbb{E}_{\boldsymbol{\theta} \sim q_\lambda}[\mathbf{s}(\boldsymbol{\theta})].
 \end{aligned}$$

Hence, for any minimizer of $\text{KL}(p\|q_\lambda)$, we have that the sufficient statistics under p and q_λ match:

$$\mathbb{E}_p[\mathbf{s}(\boldsymbol{\theta})] = \mathbb{E}_{q_\lambda}[\mathbf{s}(\boldsymbol{\theta})]. \quad (5.60)$$

Therefore, in the Gaussian case,

$$\mathbb{E}_p[\boldsymbol{\theta}] = \mathbb{E}_{q_\lambda}[\boldsymbol{\theta}] \quad \text{and} \quad \mathbb{E}_p\left[-\frac{1}{2}\boldsymbol{\theta}\boldsymbol{\theta}^\top\right] = \mathbb{E}_{q_\lambda}\left[-\frac{1}{2}\boldsymbol{\theta}\boldsymbol{\theta}^\top\right],$$

implying that

$$\mathbb{E}_p[\boldsymbol{\theta}] = \boldsymbol{\mu} \quad \text{and} \quad \text{Var}_p[\boldsymbol{\theta}] = \mathbb{E}_p[\boldsymbol{\theta}\boldsymbol{\theta}^\top] - \mathbb{E}_p[\boldsymbol{\theta}] \cdot \mathbb{E}_p[\boldsymbol{\theta}]^\top = \boldsymbol{\Sigma} \quad (5.61)$$

using eq. (1.45)

where $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ 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 .

5.5 Evidence Lower Bound

To complete the blueprint of variational inference from section 5.3, it remains to show that the reverse KL-divergence can be minimized efficiently. We have

$$\begin{aligned}
 \text{KL}(q\|p(\cdot \mid \mathbf{x}_{1:n}, y_{1:n})) &= \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta} \mid \mathbf{x}_{1:n}, y_{1:n})} \right] \\
 &= \mathbb{E}_{\boldsymbol{\theta} \sim q} \left[\log \frac{p(y_{1:n} \mid \mathbf{x}_{1:n}) q(\boldsymbol{\theta})}{p(y_{1:n}, \boldsymbol{\theta} \mid \mathbf{x}_{1:n})} \right] \\
 &= \log p(y_{1:n} \mid \mathbf{x}_{1:n}) \\
 &\quad - \underbrace{\mathbb{E}_{\boldsymbol{\theta} \sim q} [\log p(y_{1:n}, \boldsymbol{\theta} \mid \mathbf{x}_{1:n})]}_{-L(q, p; \mathcal{D})} - H[q]
 \end{aligned}$$

using the definition of the KL-divergence (5.44)

using the definition of conditional probability (1.11)

using linearity of expectation (1.24)

where $L(q, p; \mathcal{D})$ is called the *evidence lower bound* (ELBO) given the data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$. This gives the relationship

$$L(q, p; \mathcal{D}) = \underbrace{\log p(y_{1:n} | x_{1:n})}_{\text{const}} - \text{KL}(q \| p(\cdot | x_{1:n}, y_{1:n})). \quad (5.62)$$

Thus, maximizing the ELBO coincides with minimizing reverse-KL.

The ELBO can also be expressed in various other forms:

$$L(q, p; \mathcal{D}) \doteq \mathbb{E}_{\theta \sim q}[\log p(y_{1:n}, \theta | x_{1:n})] + H[q] \quad (5.63a)$$

$$= \mathbb{E}_{\theta \sim q}[\log p(y_{1:n}, \theta | x_{1:n}) - \log q(\theta)] \quad (5.63b)$$

$$= \mathbb{E}_{\theta \sim q}[\log p(y_{1:n} | x_{1:n}, \theta) + \log p(\theta) - \log q(\theta)] \quad (5.63c)$$

$$= \mathbb{E}_{\theta \sim q}[\log p(y_{1:n} | x_{1:n}, \theta)] - \text{KL}(q \| p(\cdot)). \quad (5.63d)$$

using the definition of entropy (5.32)

using the product rule (1.14)

using the definition of KL-divergence (5.44)

where we denote by $p(\cdot)$ the prior distribution.

Interpreting eq. (5.63d), maximizing the ELBO selects a variational distribution q that is close to the prior distribution $p(\cdot)$ while also maximizing the likelihood of the data $p(y_{1:n} | x_{1:n}, \theta)$. Similarly, interpreting eq. (5.63a), maximizing the ELBO selects q that has large joint likelihood $p(y_{1:n}, \theta | x_{1:n})$ (i.e., the data is likely and q is close to the prior) and large entropy $H[q]$. 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 (cf. remark 2.5).¹³

Recalling that KL-divergence is non-negative, it follows from eq. (5.62) that the evidence lower bound is a (uniform¹⁴) lower bound to the evidence $\log p(y_{1:n} | x_{1:n})$:

$$\log p(y_{1:n} | x_{1:n}) \geq L(q, p; \mathcal{D}). \quad (5.64)$$

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.

Remark 5.27: Free energy and minimizing surprise

Observe that the evidence can also be interpreted as the negative

¹³ The following example illustrates why maximizing entropy is reasonable. Consider a criminal trial with three suspects. The collected evidence shows that suspect 3 can not have committed the crime, however it does not yield any information about suspects 1 and 2. Clearly, any distribution respecting the data must assign zero probability of having committed the crime to suspect 3. However, any distribution interpolating between $(1, 0, 0)$ and $(0, 1, 0)$ respects the data. In this case, the maximum entropy principle suggests that the desired distribution is $(\frac{1}{2}, \frac{1}{2}, 0)$, and indeed, any alternative distribution seems unreasonable.

¹⁴ That is, the bound holds for any variational distribution q (with full support).

surprise about the observations under the prior distribution $p(\theta)$,

$$S[p(y_{1:n} | x_{1:n})] \leq \underbrace{\mathbb{E}_{\theta \sim q}[S[p(y_{1:n}, \theta | x_{1:n})]]}_{\text{called energy}} - H[q] \quad (5.65)$$

$$= -L(q, p; \mathcal{D}) \quad (5.66)$$

for any variational distribution q . Here, $-L(q, p; \mathcal{D})$ is also called the (*variational*) *free energy* with respect to q . Free energy can also be characterized as

$$-L(q, p; \mathcal{D}) = \underbrace{\mathbb{E}_{\theta \sim q}[S[p(y_{1:n} | x_{1:n}, \theta)]]}_{\text{"inaccuracy"}} + \underbrace{\text{KL}(q \| p(\cdot))}_{\text{"complexity"}} \quad (5.67)$$

$$= \underbrace{S[p(y_{1:n} | x_{1:n})]}_{\text{"extrinsic" value}} + \underbrace{\text{KL}(q \| p(\cdot | x_{1:n}, y_{1:n}))}_{\text{"epistemic" value}}. \quad (5.68)$$

The “extrinsic” value is independent of the variational distribution q and can be thought of as a fixed “problem cost”, whereas the “epistemic” value can be interpreted as the approximation error or “solution cost”. In particular, the minimum of free energy across all variational distributions q (when q may be any distribution) coincides with the surprise about the observations almost surely.

Maximizing the ELBO can therefore be interpreted as minimizing a proxy to the surprise about the observations (averaged across all models θ based on the prior). This is in contrast to maximum likelihood estimation, which picks a *single* model θ that minimizes the surprise about the observations.

There exists the hypothesis that “physical systems” pursue paths of least surprise, i.e., their dynamics minimize the surprise in their observations.¹⁵ Intuitively, to minimize surprise, there is a natural tradeoff between 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, and 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.

As we have alluded to previously when introducing the ELBO, the two terms constituting free energy model this natural trade-off. Therein, the *entropy* (which is maximized) encourages q to have uncertainty, i.e., to “explore”. In contrast, the *energy* (which is minimized) encourages q to fit the observed data closely, i.e., to “exploit”. This highlights a deep relationship of free energy

by negating eq. (5.64)

the two terms are commonly referred to as “inaccuracy” and “complexity”, mirroring our previous interpretation of eq. (5.63d)

using the derivation of eq. (5.62)

¹⁵ Refer to the *free energy principle* which was originally introduced by the neuroscientist Karl Friston (Friston, 2010).

to decision-making under uncertainty (which is called *planning*) where the *exploration-exploitation dilemma* arises.¹⁶ For this reason, this idea also plays a central role in approaches to reinforcement learning such as entropy regularization.¹⁷

Example 5.28: ELBO for Bayesian logistic regression

Recall that Bayesian logistic regression uses the prior distribution $w \sim \mathcal{N}(\mathbf{0}, I)$.¹⁸

Suppose we use the variational family \mathcal{Q} of all diagonal Gaussians from example 5.7. We have already seen in eq. (5.53) that for a prior $p \sim \mathcal{N}(\mathbf{0}, I)$ 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,

$$\begin{aligned} \mathbb{E}_{w \sim q_\lambda} [\log p(y_{1:n} | \mathbf{x}_{1:n}, w)] &= \mathbb{E}_{w \sim q_\lambda} \left[\sum_{i=1}^n \log p(y_i | x_i, w) \right] \\ &= \mathbb{E}_{w \sim q_\lambda} \left[- \sum_{i=1}^n \ell_{\log}(w; x_i, y_i) \right] \end{aligned} \quad (5.69)$$

¹⁶ We discuss this in great detail in the second half of this course. See section 9.1 for an introduction to the “exploration-exploitation dilemma”.

¹⁷ We explore this connection in section 12.5.

¹⁸ We omit the scaling factor σ_p^2 here for simplicity.

using independence of the data

substituting the logistic loss (5.14)

Exercise 5.29: 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,

$$\mathcal{Q} \doteq \{q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}. \quad (5.70)$$

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:

$$\tilde{q} = \arg \min_{q \in \mathcal{Q}} \text{KL}(q \| p(\cdot | \mathcal{D})).$$

2. Show that $\tilde{q} = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ satisfies

$$\begin{aligned} \mathbf{0} &= \mathbb{E}_{\boldsymbol{\theta} \sim \tilde{q}} [\nabla_{\boldsymbol{\theta}} \log p(y_{1:n}, \boldsymbol{\theta} \mid \mathbf{x}_{1:n})] \\ \boldsymbol{\Sigma}^{-1} &= -\mathbb{E}_{\boldsymbol{\theta} \sim \tilde{q}} [\mathbf{H}_{\boldsymbol{\theta}} \log p(y_{1:n}, \boldsymbol{\theta} \mid \mathbf{x}_{1:n})]. \end{aligned} \quad (5.71)$$

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 $p = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$,

$$\begin{aligned} \nabla_{\boldsymbol{\mu}} \mathbb{E}_{\mathbf{x} \sim p} [f(\mathbf{x})] &= \mathbb{E}_{\mathbf{x} \sim p} [\nabla_{\mathbf{x}} f(\mathbf{x})] \\ \nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{\mathbf{x} \sim p} [f(\mathbf{x})] &= \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p} [\mathbf{H}_{\mathbf{x}} f(\mathbf{x})]. \end{aligned} \quad (5.72)$$

Recall the Laplace approximation of $p(\boldsymbol{\theta} \mid \mathcal{D}) \propto p(y_{1:n}, \boldsymbol{\theta} \mid \mathbf{x}_{1:n})$,

$$\begin{aligned} \mathbf{0} &= \nabla_{\boldsymbol{\theta}} \log p(y_{1:n}, \boldsymbol{\theta} \mid \mathbf{x}_{1:n}) \\ \boldsymbol{\Sigma}^{-1} &= -\mathbf{H}_{\boldsymbol{\theta}} \log p(y_{1:n}, \boldsymbol{\theta} \mid \mathbf{x}_{1:n})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \end{aligned} \quad (5.73)$$

where $\hat{\boldsymbol{\theta}}$ is the MAP estimate. Crucially, the Laplace approximation is fitted *locally* at $\hat{\boldsymbol{\theta}}$. Comparing eq. (5.73) to eq. (5.71), we see that Gaussian variational inference satisfies the conditions of the Laplace approximation *on average*. For more details, refer to “The variational Gaussian approximation revisited” (Opper and Archambeau, 2009).

▷ *Solution*

5.5.1 Gradient of Evidence Lower Bound

The next problem is, of course, solving the optimization problem of maximizing the ELBO. In stochastic gradient descent, we already know a tool which we can use. However, SGD requires unbiased gradient estimates of the loss, which is given as $\ell(\boldsymbol{\lambda}; \mathcal{D}) \doteq -L(q_{\boldsymbol{\lambda}}, p; \mathcal{D})$. Thus, we need to obtain gradient estimates of

$$\nabla_{\boldsymbol{\lambda}} L(q_{\boldsymbol{\lambda}}, p; \mathcal{D}) = \nabla_{\boldsymbol{\lambda}} \mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\lambda}}} [\log p(y_{1:n} \mid \mathbf{x}_{1:n}, \boldsymbol{\theta})] - \nabla_{\boldsymbol{\lambda}} \text{KL}(q_{\boldsymbol{\lambda}} \parallel p(\cdot)). \quad (5.74)$$

using the definition of the ELBO (5.63d)

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 eq. (5.50).

Exercise 5.30: 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 $\boldsymbol{\mu} \doteq [\mu_1 \cdots \mu_d]$ and $\boldsymbol{\lambda} \doteq [\mu_1 \cdots \mu_d \sigma_1 \cdots \sigma_d]$. Show that the gradient of $\text{KL}(q_\lambda \| p(\cdot))$ with respect to $\boldsymbol{\lambda}$ is given by

$$\nabla_{\boldsymbol{\mu}} \text{KL}(q_\lambda \| p(\cdot)) = \sigma_p^{-2} \boldsymbol{\mu}, \quad \text{and} \quad (5.75a)$$

$$\nabla_{[\sigma_1 \cdots \sigma_d]} \text{KL}(q_\lambda \| p(\cdot)) = \left[\frac{\sigma_1}{\sigma_p^2} - \frac{1}{\sigma_1} \quad \cdots \quad \frac{\sigma_d}{\sigma_p^2} - \frac{1}{\sigma_d} \right]. \quad (5.75b)$$

▷ *Solution*

Obtaining the gradient of the evidence lower bound is more difficult. This is because the expectation integrates over the measure q_λ , which depends on the variational parameters $\boldsymbol{\lambda}$. Thus, we cannot move the gradient operator inside the expectation as we have done previously in remark 1.22. 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”:

$$\begin{aligned} \nabla_{\boldsymbol{\lambda}} \mathbb{E}_{\boldsymbol{\theta} \sim q_\lambda} [\log p(y_{1:n} | \mathbf{x}_{1:n}, \boldsymbol{\theta})] \\ = \mathbb{E}_{\boldsymbol{\theta} \sim q_\lambda} [\log p(y_{1:n} | \mathbf{x}_{1:n}, \boldsymbol{\theta}) \underbrace{\nabla_{\boldsymbol{\lambda}} \log q_\lambda(\boldsymbol{\theta})}_{\text{score function}}], \end{aligned} \quad (5.76)$$

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 5.31 (Reparameterization trick). *Given a random variable $\boldsymbol{\epsilon} \sim \phi$ (which is independent of $\boldsymbol{\lambda}$) and given a differentiable and invertible function $\mathbf{g} : \mathbb{R}^d \rightarrow \mathbb{R}^d$. We let $\boldsymbol{\theta} \doteq \mathbf{g}(\boldsymbol{\epsilon}; \boldsymbol{\lambda})$. Then,*

$$q_\lambda(\boldsymbol{\theta}) = \phi(\boldsymbol{\epsilon}) \cdot |\det(\mathbf{D}_{\boldsymbol{\epsilon}} \mathbf{g}(\boldsymbol{\epsilon}; \boldsymbol{\lambda}))|^{-1}, \quad (5.77)$$

$$\mathbb{E}_{\boldsymbol{\theta} \sim q_\lambda} [\mathbf{f}(\boldsymbol{\theta})] = \mathbb{E}_{\boldsymbol{\epsilon} \sim \phi} [\mathbf{f}(\mathbf{g}(\boldsymbol{\epsilon}; \boldsymbol{\lambda}))] \quad (5.78)$$

for a “nice” function $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^e$.

Proof. By the change of variables formula (1.52) and using $\boldsymbol{\epsilon} = \mathbf{g}^{-1}(\boldsymbol{\theta}; \boldsymbol{\lambda})$,

$$\begin{aligned} q_\lambda(\boldsymbol{\theta}) &= \phi(\boldsymbol{\epsilon}) \cdot \left| \det(\mathbf{D}_{\boldsymbol{\theta}} \mathbf{g}^{-1}(\boldsymbol{\theta}; \boldsymbol{\lambda})) \right| \\ &= \phi(\boldsymbol{\epsilon}) \cdot \left| \det((\mathbf{D}_{\boldsymbol{\epsilon}} \mathbf{g}(\boldsymbol{\epsilon}; \boldsymbol{\lambda}))^{-1}) \right| \\ &= \phi(\boldsymbol{\epsilon}) \cdot |\det(\mathbf{D}_{\boldsymbol{\epsilon}} \mathbf{g}(\boldsymbol{\epsilon}; \boldsymbol{\lambda}))|^{-1}. \end{aligned}$$

by the inverse function theorem,
 $\mathbf{D} \mathbf{g}^{-1}(\mathbf{y}) = \mathbf{D} \mathbf{g}(\mathbf{x})^{-1}$
 using $\det(\mathbf{A}^{-1}) = \det(\mathbf{A})^{-1}$

Equation (5.78) is a direct consequence of the law of the unconscious statistician (1.30). \square

In other words, the reparameterization trick requires finding a function $g(\cdot; \lambda)$ and density ϕ such that $q_\lambda = g(\cdot; \lambda)_\# \phi$ is the pushforward of ϕ under perturbation g .

Applying the reparameterization trick, we can use our analysis of eq. (1.28) to swap the order of gradient and expectation,

$$\nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [f(\theta)] = \mathbb{E}_{\epsilon \sim \phi} [\nabla_\lambda f(g(\epsilon; \lambda))]. \quad (5.79)$$

If we can find a function g such that $\theta \sim q_\lambda$ for some sampling distribution ϕ independent of λ , we can use the reparameterization trick to simplify our gradient. We call a distribution q_λ which admits reparameterization *reparameterizable*.

Example 5.32: 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 eq. (1.123) that a Gaussian random vector $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ following a standard normal distribution can be transformed to follow the Gaussian distribution q_λ by using the linear transformation,

$$\theta = g(\epsilon; \lambda) \doteq \Sigma^{1/2} \epsilon + \mu. \quad (5.80)$$

In particular, we have

$$\epsilon = g^{-1}(\theta; \lambda) = \Sigma^{-1/2}(\theta - \mu) \quad \text{and} \quad (5.81)$$

$$\phi(\epsilon) = q_\lambda(\theta) \cdot \left| \det(\Sigma^{1/2}) \right|. \quad (5.82)$$

by solving eq. (5.80) for ϵ

using the reparameterization trick (i.e., the change of variables formula) (5.77)

Exercise 5.33: Reparameterizable distributions

1. Let $X \sim \text{Unif}([a, b])$ for any $a \leq b$. That is,

$$p_X(x) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{otherwise.} \end{cases} \quad (5.83)$$

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. Let $X \sim \text{Cauchy}(m, \tau)$ be a random variable that follows a *Cauchy distribution* with location m and scale τ . The CDF of a Cauchy distribution is defined as

$$P_X(x) \doteq \frac{1}{\pi} \arctan\left(\frac{x-m}{\tau}\right) + \frac{1}{2}. \quad (5.84)$$

Show that $\text{Cauchy}(0, 1)$ can be reparameterized in terms of $\text{Unif}([0, 1])$.

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.

▷ *Solution*

In the following, we write $\mathbf{C} \doteq \Sigma^{1/2}$. Let us now derive the gradient estimate for the evidence lower bound assuming the Gaussian variational approximation from example 5.32. This approach extends to any reparameterizable distribution.

$$\begin{aligned} \nabla_{\lambda} \mathbb{E}_{\theta \sim q_{\lambda}} [\log p(y_{1:n} \mid \mathbf{x}_{1:n}, \theta)] \\ = \nabla_{\mathbf{C}, \mu} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\log p(y_{1:n} \mid \mathbf{x}_{1:n}, \theta) \Big|_{\theta = \mathbf{C}\epsilon + \mu} \right] \end{aligned} \quad (5.85)$$

using the reparameterization trick (5.78)

$$\begin{aligned} = n \cdot \nabla_{\mathbf{C}, \mu} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\frac{1}{n} \sum_{i=1}^n \log p(y_i \mid \mathbf{x}_i, \theta) \Big|_{\theta = \mathbf{C}\epsilon + \mu} \right] \\ = n \cdot \nabla_{\mathbf{C}, \mu} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{i \sim \text{Unif}([n])} \left[\log p(y_i \mid \mathbf{x}_i, \theta) \Big|_{\theta = \mathbf{C}\epsilon + \mu} \right] \end{aligned}$$

using independence of the data and extending with n/n

interpreting the sum as an expectation

$$= n \cdot \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{i \sim \text{Unif}([n])} \left[\nabla_{\mathbf{C}, \mu} \log p(y_i \mid \mathbf{x}_i, \theta) \Big|_{\theta = \mathbf{C}\epsilon + \mu} \right] \quad (5.86)$$

using eq. (1.28)

$$\approx n \cdot \frac{1}{m} \sum_{j=1}^m \nabla_{\mathbf{C}, \mu} \log p(y_{i(j)} \mid \mathbf{x}_{i(j)}, \theta) \Big|_{\theta = \mathbf{C}\epsilon^{(j)} + \mu} \quad (5.87)$$

using Monte Carlo sampling

where $\epsilon^{(j)} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$ 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*. The only requirement is that we can obtain unbiased gradient estimates

from the evidence lower bound (and the likelihood). 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.

Optional Readings

- Ranganath, Gerrish, and Blei (2014).
Black box variational inference.
- Duvenaud and Adams (2015).
Black-box stochastic variational inference in five lines of python.
- Titsias and Lázaro-Gredilla (2014).
Doubly stochastic variational Bayes for non-conjugate inference.
- Hensman, Matthews, and Ghahramani (2015).
Scalable variational Gaussian process classification.
- Mohamed, Rosca, Figurnov, and Mnih (2020).
Monte Carlo Gradient Estimation in Machine Learning.

Markov Chain Monte Carlo Methods

Variational inference approximates the entire posterior distribution. However, note that the key challenge in Bayesian inference is not learning the posterior distribution, but using the posterior distribution for predictions,

$$p(y^* | x^*, x_{1:n}, y_{1:n}) = \int p(y^* | x^*, \theta) p(\theta | x_{1:n}, y_{1:n}) d\theta. \quad (6.1)$$

This integral can be interpreted as an expectation over the posterior distribution,

$$= \mathbb{E}_{\theta \sim p(\cdot | x_{1:n}, y_{1:n})} [p(y^* | x^*, \theta)]. \quad (6.2)$$

Observe that the likelihood $f(\theta) \doteq p(y^* | x^*, \theta)$ 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,

$$\approx \frac{1}{m} \sum_{i=1}^m f(\theta^{(i)}) \quad (6.3)$$

for independent $\theta^{(i)} \stackrel{\text{iid}}{\sim} p(\cdot | x_{1:n}, y_{1:n})$. The law of large numbers (1.84) and Hoeffding's inequality (1.88) imply that this estimator is consistent and sharply concentrated.¹

¹ For more details, refer to section 1.4.2.

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.

Even if we were able to approximate Z , this generally does not yield an efficient sampling method. For example, the inverse transform sampling discussed in section 1.1.5 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 begin with, let us revisit the fundamental theory behind Markov chains.

Definition 6.1 (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*,

$$X_{t+1} \perp X_{0:t-1} \mid X_t, \quad (6.6)$$

is satisfied.

Intuitively, the Markov property states that future behavior is independent of past states given the present state.

Remark 6.2: 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 6.29).

We restrict our attention to *time-homogeneous* Markov chains,³ which can be characterized by a *transition function*,

$$p(x' \mid x) \doteq \mathbb{P}(X_{t+1} = x' \mid X_t = x). \quad (6.7)$$

² A *stochastic process* is a sequence of random variables.

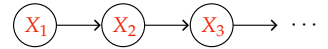


Figure 6.1: Directed graphical model of a Markov chain. The random variable X_{t+1} is conditionally independent of the random variables $X_{0:t-1}$ given X_t .

³ That is, the transition probabilities do not change over time.

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 $X_0 \sim q_0$.

Exercise 6.3: Markov chain update

Prove that one iteration of the Markov chain can be expressed as

$$q_{t+1} = q_t \mathbf{P}. \quad (6.9)$$

▷ *Solution*

It is implied directly that we can write the state of the Markov chain at time $t + k$ as

$$q_{t+k} = q_t \mathbf{P}^k. \quad (6.10)$$

Exercise 6.4: k -step transitions

Prove that 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. We denote this entry by $p^{(k)}(x' | x)$.

▷ *Solution*

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 6.5 (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 eq. (6.9) that equivalently, π is stationary w.r.t. a transition matrix P iff

$$\pi = \pi P. \quad (6.12)$$

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 6.6: 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 fig. 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 6.7: When does a Markov chain converge?

Even if a Markov chain has a unique stationary distribution, it must not converge to it. Consider example (3) in fig. 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 $q_t = (0, 1)$ and $q_{2t} = (1, 0)$ 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 6.8 (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 6.9: 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

$$P' \doteq \frac{1}{2}P + \frac{1}{2}I. \quad (6.17)$$

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.

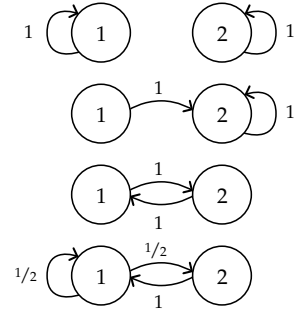


Figure 6.2: Transition graphs of Markov chains: (1) is not ergodic as its transition diagram is not strongly connected; (2) is not ergodic for the same reason; (3) is irreducible but periodic and therefore not ergodic; (4) is ergodic with stationary distribution $\pi(1) = 2/3, \pi(2) = 1/3$.

Take now π to be a stationary distribution of P . We have that π is also a stationary distribution of P' as

$$\pi P' = \frac{1}{2} \pi P + \frac{1}{2} \pi I = \frac{1}{2} \pi + \frac{1}{2} \pi = \pi. \quad (6.18)$$

using (6.12)

Fact 6.10 (Fundamental theorem of ergodic Markov chains). *An ergodic Markov chain has a unique stationary distribution π (with full support) and*

$$\lim_{t \rightarrow \infty} q_t = \pi \quad (6.19)$$

irrespectively of the initial distribution q_0 .

Proof. Refer to theorem 4.9 of “Markov chains and mixing times” (Levin and Peres, 2017) for a proof. \square

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 6.11: 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 6.12 (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 (5.44), 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 6.13 (Mixing time). For a Markov chain with stationary distribution π , its *mixing time* with respect to the total variation distance 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

$$\tau_{\text{TV}}(\epsilon) \in \mathcal{O}(\text{poly}(n, \log(1/\epsilon))). \quad (6.24)$$

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.

Exercise 6.14: 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”?

▷ *Solution*

Readings

For an in-depth treatment of Markov chains refer to chapter 1 of “Markov chains and mixing times” (Levin and Peres, 2017).

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 6.15 (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 6.16. *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 eq. (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 6.17 (Ergodic theorem). *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 $f : S \rightarrow \mathbb{R}$,*

$$\frac{1}{n} \sum_{i=1}^n f(x_i) \xrightarrow{\text{a.s.}} \sum_{x \in S} \pi(x) f(x) = \mathbb{E}_{x \sim \pi}[f(x)] \quad (6.28)$$

as $n \rightarrow \infty$ where $x_i \sim X_i \mid x_{i-1}$.

Proof. See appendix C of “Markov chains and mixing times” (Levin and Peres, 2017) for a proof. \square

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 chain does not yet approximate the posterior well. Typically, the first t_0 samples are therefore discarded,

$$\mathbb{E}[f(X)] \approx \frac{1}{T - t_0} \sum_{t=t_0+1}^T f(X_t). \quad (6.29)$$

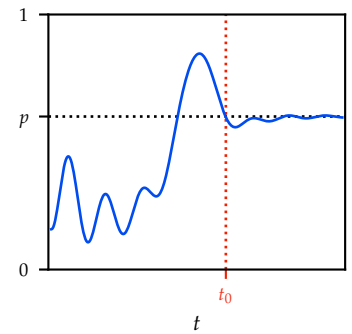


Figure 6.3: Illustration of the “burn-in” time t_0 of a Markov chain approximating the posterior $p(y^* = 1 \mid X, y)$ of Bayesian logistic regression. The true posterior p is shown in gray. The distribution of the Markov chain at time t is shown in red.

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* $r(\mathbf{x}' | \mathbf{x})$ which, given we are in state \mathbf{x} , proposes a new state \mathbf{x}' . Metropolis and Hastings showed that using the *acceptance distribution* $\text{Bern}(\alpha(\mathbf{x}' | \mathbf{x}))$ where

$$\alpha(\mathbf{x}' | \mathbf{x}) \doteq \min \left\{ 1, \frac{p(\mathbf{x}')r(\mathbf{x} | \mathbf{x}')}{p(\mathbf{x})r(\mathbf{x}' | \mathbf{x})} \right\} \quad (6.30)$$

$$= \min \left\{ 1, \frac{q(\mathbf{x}')r(\mathbf{x} | \mathbf{x}')}{q(\mathbf{x})r(\mathbf{x}' | \mathbf{x})} \right\} \quad (6.31)$$

similarly to the detailed balance equation, the normalizing constant Z cancels

to decide whether to follow the proposal yields a Markov chain with stationary distribution $p(\mathbf{x}) = \frac{1}{Z}q(\mathbf{x})$.

Algorithm 6.18: Metropolis-Hastings algorithm

```

1 initialize  $\mathbf{x} \in \mathbb{R}^n$ 
2 for  $t = 1$  to  $T$  do
3   sample  $\mathbf{x}' \sim r(\mathbf{x}' | \mathbf{x})$ 
4   sample  $u \sim \text{Unif}([0, 1])$ 
5   if  $u \leq \alpha(\mathbf{x}' | \mathbf{x})$  then update  $\mathbf{x} \leftarrow \mathbf{x}'$ 
6   else update  $\mathbf{x} \leftarrow \mathbf{x}$ 

```

Theorem 6.19 (Metropolis-Hastings theorem). *Given an arbitrary proposal distribution r , the stationary distribution of the Markov chain simulated by the Metropolis-Hastings algorithm is $p(\mathbf{x}) = \frac{1}{Z}q(\mathbf{x})$.*

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,

$$\begin{aligned} p(x) \cdot p(x' | x) &= \frac{1}{Z} q(x) p(x' | x) && \text{using the definition of the distribution } p \\ &= \frac{1}{Z} q(x) r(x' | x) \alpha(x' | x) && \text{using the transition probabilities of the Markov chain} \\ &= \frac{1}{Z} q(x) r(x' | x) \frac{q(x') r(x | x')}{q(x) r(x' | x)} && \text{using the definition of the acceptance distribution } \alpha \\ &= \frac{1}{Z} q(x') r(x | x') \\ &= \frac{1}{Z} q(x') r(x | x') \alpha(x | x') && \text{using the definition of the acceptance distribution } \alpha \\ &= \frac{1}{Z} q(x') p(x | x') && \text{using the transition probabilities of the Markov chain} \\ &= p(x') \cdot p(x | x'). && \text{using the definition of the distribution } p \quad \square \end{aligned}$$

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 $p(\cdot | x)$ have full support. For example, if the proposal distribution $r(\cdot | x)$ has full support, the full support of $p(\cdot | x)$ follows immediately from eq. (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 alg. 6.20.

Algorithm 6.20: 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
       $p(x_i \mid \mathbf{x}_{-i})$ 

```

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 $p(\cdot \mid \mathbf{x}_{-i})$ have full support.

Theorem 6.21 (Gibbs sampling as Metropolis-Hastings). *Gibbs sampling is a Metropolis-Hastings algorithm. For any fixed $i \in [n]$, it has proposal distribution*

$$r_i(\mathbf{x}' \mid \mathbf{x}) \doteq \begin{cases} p(x'_i \mid \mathbf{x}'_{-i}) & \mathbf{x}' \text{ differs from } \mathbf{x} \text{ only in entry } i \\ 0 & \text{otherwise} \end{cases} \quad (6.33)$$

and acceptance distribution $\alpha_i(\mathbf{x}' \mid \mathbf{x}) \doteq 1$.

Proof. We show that $\alpha_i(\mathbf{x}' \mid \mathbf{x}) = 1$ follows from the definition of an acceptance distribution in Metropolis-Hastings (6.31) and the choice of proposal distribution (6.33).

By (6.31),

$$\alpha_i(\mathbf{x}' \mid \mathbf{x}) = \min \left\{ 1, \frac{p(\mathbf{x}') r_i(\mathbf{x} \mid \mathbf{x}')}{p(\mathbf{x}) r_i(\mathbf{x}' \mid \mathbf{x})} \right\}$$

Note that $p(\mathbf{x}) = p(x_i, \mathbf{x}_{-i}) = p(x_i \mid \mathbf{x}_{-i}) p(\mathbf{x}_{-i})$ using the product rule (1.14). Therefore,

$$\begin{aligned} &= \min \left\{ 1, \frac{p(x'_i \mid \mathbf{x}'_{-i}) p(\mathbf{x}'_{-i}) r_i(\mathbf{x} \mid \mathbf{x}')}{p(x_i \mid \mathbf{x}_{-i}) p(\mathbf{x}_{-i}) r_i(\mathbf{x}' \mid \mathbf{x})} \right\} \\ &= \min \left\{ 1, \frac{p(x'_i \mid \mathbf{x}'_{-i}) p(\mathbf{x}'_{-i}) p(x_i \mid \mathbf{x}_{-i})}{p(x_i \mid \mathbf{x}_{-i}) p(\mathbf{x}_{-i}) p(x'_i \mid \mathbf{x}'_{-i})} \right\} && \text{using the proposal distribution (6.33)} \\ &= \min \left\{ 1, \frac{p(\mathbf{x}'_{-i})}{p(\mathbf{x}_{-i})} \right\} \\ &= 1. && \square \quad \text{using that } p(\mathbf{x}'_{-i}) = p(\mathbf{x}_{-i}) \end{aligned}$$

If the index i is chosen uniformly at random as in alg. 6.20, 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 6.21 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 6.22 (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

$$p(x_i | \mathbf{x}_{-i}) = \frac{p(x_i, \mathbf{x}_{-i})}{\sum_{x_i} p(x_i, \mathbf{x}_{-i})} = \frac{q(x_i, \mathbf{x}_{-i})}{\sum_{x_i} q(x_i, \mathbf{x}_{-i})}. \quad (6.34)$$

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.

using the definition of condition probability (1.11) and the sum rule (1.10)

Exercise 6.23: 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 (1.67).

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 priors

$$\mu \sim \mathcal{N}(\mu_0, \lambda_0^{-1}) \quad \text{and} \quad \lambda \sim \text{Gamma}(\alpha, \beta).$$

We would like to sample from the posterior $p(\mu, \lambda | x_{1:n})$. Show that

$$\mu | \lambda, x_{1:n} \sim \mathcal{N}(m_\lambda, l_\lambda^{-1}) \quad \text{and} \quad \lambda | \mu, x_{1:n} \sim \text{Gamma}(a_\mu, b_\mu),$$

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. The *Pareto distribution* was originally used to model the distribution of wealth in a society, but is also used to model many

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 .

other phenomena such as the size of cities, the frequency of words, and the returns on stocks. Formally, the Pareto distribution is defined by the following PDF,

$$\text{Pareto}(x; \alpha, c) = \frac{\alpha c^\alpha}{x^{\alpha+1}} \mathbb{1}\{x \geq c\}, \quad x \in \mathbb{R} \quad (6.35)$$

where the *tail index* $\alpha > 0$ models the “weight” of the right tail of the distribution (larger α means lighter tail) and $c > 0$ corresponds to a cutoff threshold. The distribution is supported on $[c, \infty)$, and as $\alpha \rightarrow \infty$ it approaches a point density at c .

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 essentially corresponds to a uniform prior (i.e., “no 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.

▷ *Solution*

Readings

For more on Gibbs sampling read chapter 11.3 of “Pattern recognition and machine learning” (Bishop and Nasrabadi, 2006).

6.3 Sampling as Optimization

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 6.24 (Gibbs distribution). Formally, a *Gibbs distribution* (also called a *Boltzmann distribution*) is a continuous distribution p whose PDF is of the form

$$p(x) = \frac{1}{Z} \exp(-f(x)). \quad (6.36)$$

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ 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 (5.57) with sufficient statistic $-f(x)$.

⁶ This can easily be seen as $\exp(\cdot) > 0$.

Exercise 6.25: Maximum entropy property of Gibbs distribution

1. Let X be a random variable supported on the finite set $\mathcal{T} \subset \mathbb{R}$.⁷ 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.

2. What happens for $T \rightarrow \{0, \infty\}$?

▷ *Solution*

⁷ The same result can be shown to hold for arbitrary compact subsets.

Observe that the posterior distribution can always be interpreted as a Gibbs distribution as long as prior and likelihood have full support,

$$\begin{aligned} p(\theta \mid x_{1:n}, y_{1:n}) &= \frac{1}{Z} p(\theta) p(y_{1:n} \mid x_{1:n}, \theta) \\ &= \frac{1}{Z} \exp(-[-\log p(\theta) - \log p(y_{1:n} \mid x_{1:n}, \theta)]). \end{aligned} \quad (6.37)$$

using Bayes' rule (1.59)

Thus, defining the energy function

$$f(\theta) \doteq -\log p(\theta) - \log p(y_{1:n} \mid x_{1:n}, \theta) \quad (6.38)$$

$$= -\log p(\theta) - \sum_{i=1}^n \log p(y_i \mid x_i, \theta), \quad (6.39)$$

yields

$$p(\theta \mid x_{1:n}, y_{1:n}) = \frac{1}{Z} \exp(-f(\theta)). \quad (6.40)$$

Note that f coincides with the loss function used for MAP estimation (1.96). For a uniform prior, the regularization term vanishes and the energy reduces to the negative log-likelihood $\ell_{\text{NLL}}(\theta; \mathcal{D})$ (i.e., the loss function of maximum likelihood estimation (1.91)).

Exercise 6.26: Energy function of Bayesian logistic regression

Recall from eq. (5.13) that the energy function of Bayesian logistic regression is

$$f(w) = \lambda \|w\|_2^2 + \sum_{i=1}^n \log(1 + \exp(-y_i w^\top x_i)), \quad (6.41)$$

which coincided with the standard optimization objective of (reg-

ularized) logistic regression.

Show that the posterior distribution of Bayesian logistic regression is log-concave.

▷ *Solution*

Using that the posterior is a Gibbs distribution, we can rewrite the acceptance distribution of Metropolis-Hastings,

$$\alpha(\mathbf{x}' | \mathbf{x}) = \min \left\{ 1, \frac{r(\mathbf{x} | \mathbf{x}')}{r(\mathbf{x}' | \mathbf{x})} \exp(f(\mathbf{x}) - f(\mathbf{x}')) \right\}. \quad (6.42)$$

this is obtained by substituting the PDF of a Gibbs distribution for the posterior

Example 6.27: Metropolis-Hastings with Gaussian proposals

Let us consider Metropolis-Hastings with the Gaussian proposal distribution,

$$r(\mathbf{x}' | \mathbf{x}) \doteq \mathcal{N}(\mathbf{x}'; \mathbf{x}, \tau \mathbf{I}). \quad (6.43)$$

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

$$\alpha(\mathbf{x}' | \mathbf{x}) = \min \{ 1, \exp(f(\mathbf{x}) - f(\mathbf{x}')) \}. \quad (6.44)$$

Intuitively, when a state with lower energy is proposed, that is $f(\mathbf{x}') \leq f(\mathbf{x})$, 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 $f(\mathbf{x}) - f(\mathbf{x}')$. 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 $p(\mathbf{x}) \propto \exp(-f(\mathbf{x}))$. As we have seen in exercise 6.25, this Gibbs distribution is the associated maximum entropy distribution. Observe that

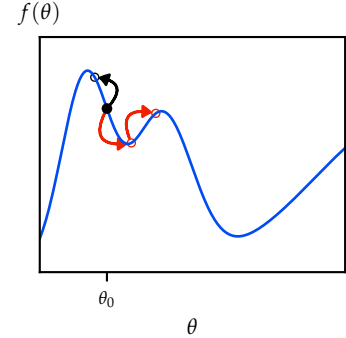


Figure 6.4: Metropolis-Hastings and Langevin dynamics minimize the energy function $f(\theta)$ shown in blue. Suppose we start at the black dot θ_0 , then the black and red arrows denote possible subsequent samples. Metropolis-Hastings uses an “uninformed” search direction, whereas Langevin dynamics uses the gradient of $f(\theta)$ to make “more promising” proposals. The random proposals help get past local optima.

the surprise about \mathbf{x} under this distribution is given by

$$S[p(\mathbf{x})] = f(\mathbf{x}) + \log Z. \quad (6.45)$$

That is, up to a constant shift, the energy of \mathbf{x} *coincides* with the surprise about \mathbf{x} . Energies are therefore sufficient for comparing the “likelihood” of points, and they do not require normalization.⁸

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 $\ell(\mathbf{x})$ can be thought of as an energy function with an associated maximum entropy distribution $p(\mathbf{x}) \propto \exp(-\ell(\mathbf{x}))$.

⁸ Intuitively, an energy can be used to compare the “likelihood” of two points \mathbf{x} and \mathbf{x}' whereas the probability \mathbf{x} makes a statement about the “likelihood” of \mathbf{x} relative to *all* other points.

Exercise 6.28: 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 $x'_i \sim p(\cdot \mid \mathbf{x}_{-i})$.

Show that

$$\mathbb{E}_{x'_i \sim p(\cdot \mid \mathbf{x}_{-i})} [f(\mathbf{x}')] \leq f(\mathbf{x}) - S[p(x_i \mid \mathbf{x}_{-i})] + H[p(\cdot \mid \mathbf{x}_{-i})]. \quad (6.46)$$

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., $S[p(x_i \mid \mathbf{x}_{-i})] \geq H[p(\cdot \mid \mathbf{x}_{-i})]$.

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.

▷ *Solution*

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 6.27 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.

A natural idea is to shift the proposal distribution perpendicularly to the gradient of the energy function. This yields the following proposal distribution,

$$r(\mathbf{x}' | \mathbf{x}) = \mathcal{N}(\mathbf{x}'; \mathbf{x} - \eta_t \nabla f(\mathbf{x}), 2\eta_t \mathbf{I}). \quad (6.47)$$

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 $\alpha(\mathbf{x}' | \mathbf{x}) \rightarrow 1$ 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 eq. (6.47) is known as the *unadjusted Langevin algorithm* (ULA).

Observe that if the stationary distribution is the posterior distribution

$$p(\boldsymbol{\theta} | \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) = \frac{1}{Z} \exp \left(\underbrace{\log p(\boldsymbol{\theta}) + \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta})}_{-f(\boldsymbol{\theta})} \right)$$

with energy f as we discussed in section 6.3.1, then the proposal $\boldsymbol{\theta}'$ of MALA/LMC can be equivalently formulated as

$$\begin{aligned} \boldsymbol{\theta}' &\leftarrow \boldsymbol{\theta} - \eta_t \nabla f(\boldsymbol{\theta}) + \epsilon \\ &= \boldsymbol{\theta} + \eta_t \left(\nabla \log p(\boldsymbol{\theta}) + \sum_{i=1}^n \nabla \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) \right) + \epsilon \end{aligned} \quad (6.48)$$

where $\epsilon \sim \mathcal{N}(\mathbf{0}, 2\eta_t \mathbf{I})$.

Remark 6.29: 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 6.30 (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,⁹ 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 stochastic differential equation (SDE)

$$d\boldsymbol{\theta}_t = \underbrace{-\nabla f(\boldsymbol{\theta}_t) dt}_{\text{drift}} + \underbrace{\sqrt{2} d\mathbf{W}_t}_{\text{noise}}, \quad (6.49)$$

Such a stochastic process is also called a *diffusion (process)* and eq. (6.49) 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 eq. (6.49) 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 $\boldsymbol{\theta}_k \approx \boldsymbol{\theta}(\tau_k)$ 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

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \nabla f(\boldsymbol{\theta}_k) \Delta t_k + \sqrt{2} \Delta \mathbf{W}_k \quad (6.50)$$

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 eq. (6.48).

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).

Exercise 6.31: Mixing time of Langevin dynamics

In this exercise, we will show that for certain Gibbs distributions $p(\boldsymbol{\theta}) \propto \exp(-f(\boldsymbol{\theta}))$, 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*

$$dx_t = -\nabla f(x_t) dt. \quad (6.51)$$

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 x_t to a fixed point (also called an *equilibrium*).

⁹ 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 $s < t$.

Let us assume that f is α -strongly convex for some $\alpha > 0$, that is,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\alpha}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n. \quad (6.52)$$

In words, f is lower bounded by a quadratic function with curvature α . Moreover, assume w.l.o.g. that f minimized at $f(\mathbf{0}) = 0$.¹⁰

1. Show that f satisfies the *Polyak-Łojasiewicz (PL) inequality*, i.e.,

$$f(\mathbf{x}) \leq \frac{1}{2\alpha} \|\nabla f(\mathbf{x})\|_2^2 \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (6.53)$$

2. Prove $\frac{d}{dt}f(\mathbf{x}_t) \leq -2\alpha f(\mathbf{x}_t)$.

Thus, $\mathbf{0}$ is the fixed point of eq. (6.51) 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

$$g(t) \leq g(0) \exp\left(\int_0^t \beta(s) ds\right) \quad \forall t \in [0, T]. \quad (6.54)$$

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 $p(\boldsymbol{\theta}) \propto \exp(-f(\boldsymbol{\theta}))$. We will use that the evolution of $\{\boldsymbol{\theta}_t\}_{t \geq 0}$ following the Langevin dynamics (6.49) is equivalently characterized by their densities $\{q_t\}_{t \geq 0}$ following the *Fokker-Planck equation*

$$\frac{\partial q_t}{\partial t} = \nabla \cdot (q_t \nabla f) + \Delta q_t. \quad (6.55)$$

Here, $\nabla \cdot$ and Δ are the divergence and Laplacian operators, respectively.¹¹ 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).¹²

4. Show that $\Delta q_t = \nabla \cdot (q_t \nabla \log q_t)$, implying that the Fokker-Planck equation simplifies to

$$\frac{\partial q_t}{\partial t} = \nabla \cdot \left(q_t \nabla \log \frac{q_t}{p} \right). \quad (6.56)$$

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 the similarity of the integrand of $\text{KL}(q_t \| p)$, $q_t \log \frac{q_t}{p}$, to eq. (6.56). We will therefore use the KL-divergence as Lyapunov function.

¹⁰ This can always be achieved by shifting the coordinate system and subtracting a constant from f .

¹¹ For ease of notation, we omit the explicit dependence of q_t , p , and f on $\boldsymbol{\theta}$.

¹² 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 \mathbf{x} , \mathbf{F} points towards \mathbf{x} , then the divergence at \mathbf{x} is negative as the volume shrinks. If \mathbf{F} points away from \mathbf{x} , then the divergence at \mathbf{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 $\varphi(\mathbf{x})$ is smaller than the average value of φ in a small neighborhood of \mathbf{x} , then the Laplacian at \mathbf{x} is positive.

Regarding the Fokker-Planck equation (6.55), 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.

5. Prove $\frac{d}{dt}\text{KL}(q_t\|p) = -J(q_t\|p)$. Here,

$$J(q_t\|p) \doteq \mathbb{E}_{\theta \sim q_t} \left[\left\| \nabla \log \frac{q_t(\theta)}{p(\theta)} \right\|_2^2 \right] \quad (6.57)$$

denotes the *relative Fisher information* of p with respect to q_t .
Hint: For any distribution q on \mathbb{R}^n ,

$$\int_{\mathbb{R}^n} (\nabla \cdot q\mathbf{F}) \varphi \, dx = - \int_{\mathbb{R}^n} q \nabla \varphi \cdot \mathbf{F} \, dx \quad (6.58)$$

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 :

$$\text{KL}(q\|p) \leq \frac{1}{2\alpha} J(q\|p). \quad (6.59)$$

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., $\tau_{\text{TV}}(\epsilon) \in \mathcal{O}(\text{poly}(n, \log(1/\epsilon)))$.

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).

▷ *Solution*

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 alg. 6.32 (Welling and Teh, 2011). Observe that SGLD (6.60) differs from MALA/LMC (6.48) only by using a sample-based approximation of the gradient.

Algorithm 6.32: 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  $\epsilon \sim \mathcal{N}(\mathbf{0}, 2\eta_t \mathbf{I})$ 
5    $\theta \leftarrow \theta + \eta_t \left( \nabla \log p(\theta) + \frac{n}{m} \sum_{j=1}^m \nabla \log p(y_{i_j} \mid x_{i_j}, \theta) \right) + \epsilon$  // (6.60)

```

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.

Readings

For more details on SGLD, read the original paper, “Bayesian learning via stochastic gradient Langevin dynamics” (Welling and Teh, 2011).

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

$$p(\mathbf{x}) \propto \exp(-f(\mathbf{x}))$$

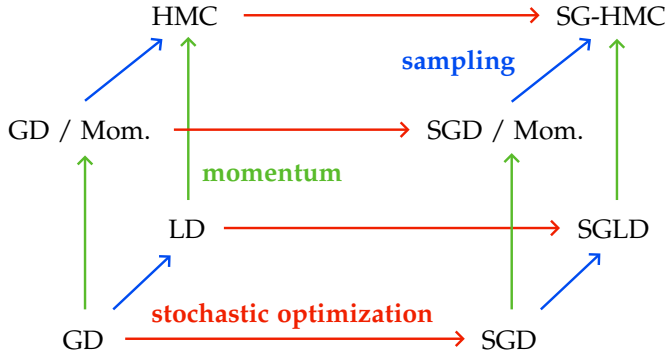


Figure 6.5: A commutative diagram of optimization algorithms. Langevin dynamics (LD) is the non-stochastic variant of SGLD.

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”.

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,

$$p(x, y) \propto \exp\left(-\frac{1}{2m} \|y\|_2^2 - f(x)\right). \quad (6.61)$$

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,¹³

$$\frac{dx}{dt} = \nabla_y H, \quad \frac{dy}{dt} = -\nabla_x H, \quad (6.62)$$

¹³ That is, HMC follows the trajectory of these dynamics for some time.

reaching some new point (x', y') and *projecting* back to the state space by selecting x' as the new sample. This is illustrated in fig. 6.6. In the next iteration, we resample the momentum $y' \sim p(\cdot \mid x')$ and repeat the procedure.

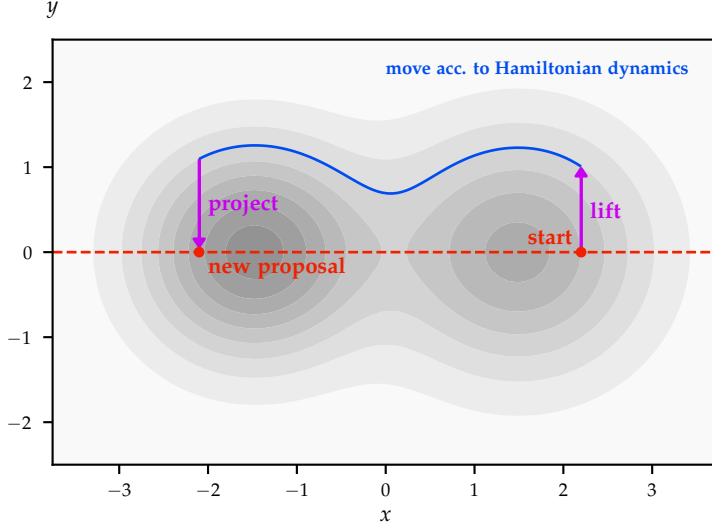


Figure 6.6: Illustration of Hamiltonian Monte Carlo. Shown is the contour plot of a distribution p , which is a mixture of two Gaussians, in the (x, y) -space.

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.

In an implementation of this algorithm, one has to solve eq. (6.62) numerically rather than exactly. Typically, this is done using the *Leapfrog method*, which for a step size τ computes

$$y(t + \tau/2) = y(t) - \frac{\tau}{2} \nabla_x f(x(t)) \quad (6.63a)$$

$$x(t + \tau) = x(t) + \frac{\tau}{m} y(t + \tau/2) \quad (6.63b)$$

$$y(t + \tau) = y(t + \tau/2) - \frac{\tau}{2} \nabla_x f(x(t + \tau)). \quad (6.63c)$$

Then, one repeats this procedure L times to arrive at a point (x', 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

$$\alpha((x', y') \mid (x, y)) \doteq \min\{1, \exp(H(x', y') - H(x, y))\}. \quad (6.64)$$

It follows that $p(x, y)$ is the stationary distribution of the Markov chain underlying HMC. Due to the independence of x and y , this also implies that the projection to x yields a Markov chain with stationary distribution $p(x)$.

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 eq. (6.62) 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.¹⁴ 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 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 $(\mathbf{x}', -\mathbf{y}')$.¹⁵ 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 $p(\mathbf{x}', -\mathbf{y}') = p(\mathbf{x}', \mathbf{y}')$,¹⁶ so we can safely disregard the direction of time when computing the acceptance probability in eq. (6.64).

¹⁴ The 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.

¹⁵ More formally, the proposal distribution is the Dirac delta at $(\mathbf{x}', -\mathbf{y}')$.

¹⁶ We use here that $p(\mathbf{y} | \mathbf{x})$ was chosen to be symmetric around zero.

Exercise 6.33: 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 (6.47) can be seen as a special case of HMC if only one Leapfrog step is used ($L = 1$) and $m = 1$.

▷ *Solution*

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. 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.

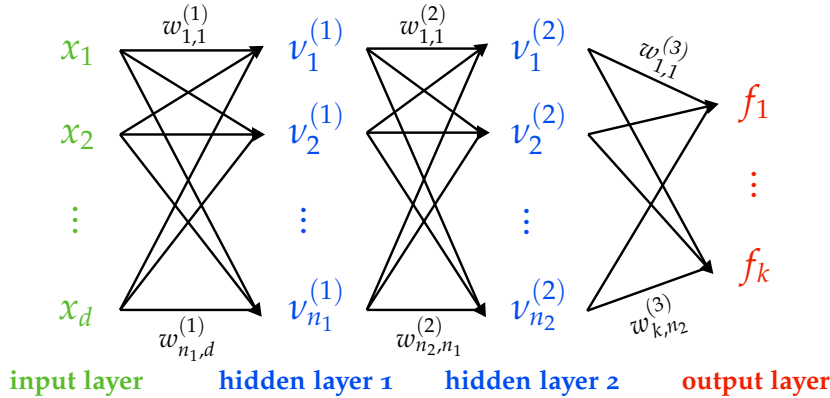
Bayesian Deep Learning

7.1 Artificial Neural Networks

In practice, it is often seen that models perform better when labels may nonlinearly depend on the inputs. One widely used family of nonlinear functions are *artificial “deep” neural networks*,¹

$$f : \mathbb{R}^d \rightarrow \mathbb{R}^k, \quad f(x; \theta) \doteq \varphi(W_L \varphi(W_{L-1}(\cdots \varphi(W_1 x)))) \quad (7.1)$$

where $\theta \doteq [W_1, \dots, W_L]$ 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.



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. This type of neural network is also called a *multilayer perceptron*.

A neural network can be visualized by a *computation graph*. An example for such a computation graph is given in fig. 7.1. The columns of the computation graph are commonly called *layers*, whereby the 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 $\mathbf{f} \doteq [f_1, \dots, f_k]$. The *activations* of an individual (hidden) layer l of the neural network are described by

$$\mathbf{v}^{(l)} \doteq \boldsymbol{\varphi}(\mathbf{W}_l \mathbf{v}^{(l-1)}) \quad (7.2)$$

where $\mathbf{v}^{(0)} = \mathbf{x}$. The activation of the i -th node is $v_i^{(l)} = \mathbf{v}^{(l)}(i)$.

7.1.1 Activation Functions

The non-linearity 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 (5.10) 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.

Remark 7.1: Universal approximation

It is important that the activation function is non-linear because otherwise, any composition of layers would still represent a linear function. Non-linear activation functions allow the network to represent arbitrary functions. The *universal approximation theorem* states that any artificial neural network with just a single hidden layer (with arbitrary width) and activation function φ where φ is not a polynomial can approximate any continuous function to an arbitrary degree of 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, $\mathbf{f} = [f_1, \dots, f_c]$, and normalize them into a probability distribution. Often, the *softmax function* is used for normalization,

$$\sigma_i(\mathbf{f}) \doteq \frac{\exp(f_i)}{\sum_{j=1}^c \exp(f_j)} \quad (7.5)$$

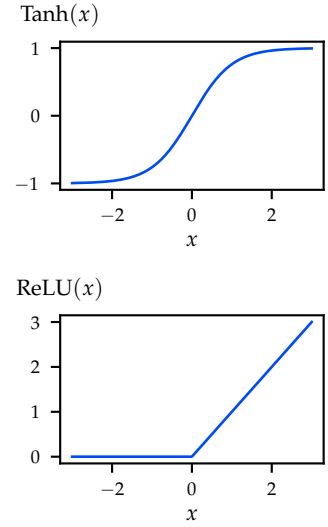


Figure 7.2: The Tanh and ReLU activation functions, respectively.

where $\sigma_i(\mathbf{f})$ corresponds to the probability mass of class i . Note that the softmax corresponds to a Gibbs distribution with energies $-\mathbf{f}$. In particular, exercise 6.25 implies that the softmax is the maximum entropy distribution over c classes satisfying $\mathbb{E}_{i \sim \sigma(\mathbf{f})} [f_i] \in \mathbb{R}$ for any logits \mathbf{f} .

Exercise 7.2: 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 (5.10) for the univariate model $f \doteq f_1 - f_0$. That is, $\sigma_1(\mathbf{f}) = \sigma(f)$ and $\sigma_0(\mathbf{f}) = 1 - \sigma(f)$.

Thus, the softmax function is a generalization of the logistic function to more than two classes.

▷ *Solution*

7.1.3 Maximum Likelihood Estimation

We will study neural networks under the lens of supervised learning (cf. section 1.2) where we are provided some independently-sampled (noisy) data $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ generated according to an unknown process $(\mathbf{x}, \mathbf{y}) \sim p$, 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 $\ell(\theta; \mathbf{y})$ which quantifies the approximation error of the network outputs $\mathbf{f}(\mathbf{x}; \theta)$. 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,

$$\ell_{\text{mse}}(\theta; \mathcal{D}) \doteq \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i; \theta) - y_i)^2. \quad (7.6)$$

As we have already seen in section 2.1.1 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 $\mathbf{y} \in \{0, 1\}^c$ 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 $q_\theta(\cdot | \mathbf{x})$ the resulting probability distribution over classes with PMF $[\sigma_1(\mathbf{f}(\mathbf{x}; \theta)), \dots, \sigma_c(\mathbf{f}(\mathbf{x}; \theta))]$, and aim to find θ such that $q_\theta(\mathbf{y} | \mathbf{x}) \approx p(\mathbf{y} | \mathbf{x})$. In this context, it is common to minimize the cross-entropy,

$$\mathbb{H}[p || q_\theta] = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p} [-\log q_\theta(\mathbf{y} | \mathbf{x})]$$

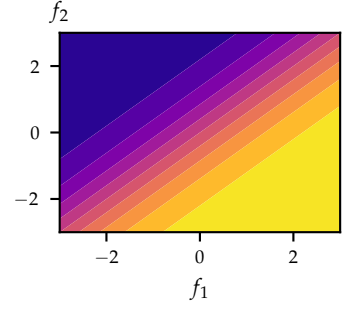


Figure 7.3: Softmax $\sigma_1(f_1, f_2)$ for a binary classification problem. Blue denotes a small probability and yellow denotes a large probability of belonging to class 1, respectively.

³ That is, exactly one component of \mathbf{y} is 1 and all others are 0, indicating to which class the given example belongs.

using the definition of cross-entropy (5.40)

$$\approx -\underbrace{\frac{1}{n} \sum_{i=1}^n \log q_{\theta}(\mathbf{y}_i \mid \mathbf{x}_i)}_{\doteq \ell_{\text{ce}}(\theta; \mathcal{D})} \quad (7.7) \quad \text{using Monte Carlo sampling}$$

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

$$\propto -\sum_{i=1}^n \log q_{\theta}(\mathbf{y}_i \mid \mathbf{x}_i) = \ell_{\text{nll}}(\theta; \mathcal{D}) \quad (7.8)$$

Recall that $\ell_{\text{nll}}(\theta; \mathcal{D})$ 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 exercise 5.3 that for a two-class classification problem the cross-entropy loss is equivalent to the logistic loss.

⁴ We have previously seen this in section 5.4.6. Note that minimizing the cross-entropy, $H[p \parallel q_{\theta}]$, is equivalent to minimizing forward-KL, $\text{KL}(p \parallel q_{\theta})$.

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 $\theta = \mathbf{W}_{1:L}$ 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)*.

Readings

For a thorough introduction to backpropagation, refer to section 4.2 of “Computational Intelligence Lab” (Hofmann, 2022).

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 adapt our Bayesian approach to neural networks? We use the same strategy which we already used for Bayesian logistic regression, we impose a Gaussian prior on the weights,

$$\theta \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I}). \quad (7.9)$$

Similarly, we can use a Gaussian likelihood to describe how well the data is described by the model f ,

$$y \mid \mathbf{x}, \theta \sim \mathcal{N}(f(\mathbf{x}; \theta), \sigma_n^2). \quad (7.10)$$

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 1.6 and many times since.

7.2.1 Maximum A Posteriori Estimation

Before studying Bayesian inference, let us first consider MAP estimation in the context of neural networks.

The MAP estimate of the weights is obtained by

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} \log p(\theta) + \sum_{i=1}^n \log p(y_i \mid \mathbf{x}_i, \theta). \quad (7.11)$$

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,

$$\log p(y_i \mid \mathbf{x}_i, \theta) = -\frac{(y_i - f(\mathbf{x}_i; \theta))^2}{2\sigma_n^2} + \text{const.} \quad (7.12)$$

Obtaining the MAP estimate instead, simply corresponds to adding an ℓ_2 -regularization term to the squared error loss,

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} \frac{1}{2\sigma_p^2} \|\theta\|_2^2 + \frac{1}{2\sigma_n^2} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \theta))^2. \quad (7.13)$$

As we have already seen in remark 1.62 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,

$$\theta \leftarrow \theta(1 - \lambda \eta_t) + \eta_t \sum_{i=1}^n \nabla \log p(y_i \mid \mathbf{x}_i, \theta) \quad (7.14)$$

where $\lambda \doteq 1/\sigma_p^2$. The gradients of the likelihood can be obtained using automatic differentiation.

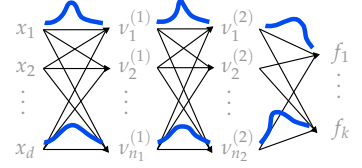


Figure 7.4: Bayesian neural networks model a distribution over the weights of a neural network.

using that for an isotropic Gaussian prior, $\log p(\theta) = -\frac{1}{2\sigma_p^2} \|\theta\|_2^2 + \text{const}$

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* as the noise is assumed to be uniform across the domain. In many settings, however, the noise is inherently *heteroscedastic*. That is, the noise varies depending on the input and which “region” of the domain the input is from. This behavior is visualized in fig. 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

$$y \mid x, \theta \sim \mathcal{N}(\mu(x; \theta), \sigma^2(x; \theta)) \quad \text{where} \quad (7.15a)$$

$$\mu(x; \theta) \doteq f_1(x; \theta), \quad (7.15b)$$

$$\sigma^2(x; \theta) \doteq \exp(f_2(x; \theta)). \quad (7.15c)$$

Using this model, the likelihood term from eq. (7.11) is

$$\begin{aligned} \log p(y_i \mid x_i, \theta) &= \log \mathcal{N}(y_i; \mu(x_i; \theta), \sigma^2(x_i; \theta)) \\ &= \log \frac{1}{\sqrt{2\pi\sigma^2(x_i; \theta)}} - \frac{(y_i - \mu(x_i; \theta))^2}{2\sigma^2(x_i; \theta)} \\ &= \underbrace{\log \frac{1}{\sqrt{2\pi}}}_{\text{const}} - \frac{1}{2} \left[\log \sigma^2(x_i; \theta) + \frac{(y_i - \mu(x_i; \theta))^2}{\sigma^2(x_i; \theta)} \right]. \end{aligned} \quad (7.16)$$

Hence, the model can either explain the label y_i by an accurate model $\mu(x_i; \theta)$ or by a large variance $\sigma^2(x_i; \theta)$, 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 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.

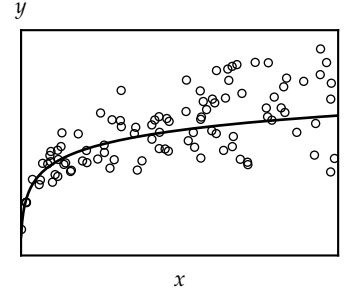


Figure 7.5: Illustration of data with variable (heteroscedastic) noise. The noise increases as the inputs increase in magnitude. The noise-free function is shown in black.

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 5.7.⁷ Recall the fundamental objective of variational inference (5.30),

$$\begin{aligned} & \arg \min_{q \in \mathcal{Q}} \text{KL}(q \| p(\cdot | \mathbf{x}_{1:n}, \mathbf{y}_{1:n})) \\ &= \arg \max_{q \in \mathcal{Q}} L(q, p; \mathcal{D}) \\ &= \arg \max_{q \in \mathcal{Q}} \mathbb{E}_{\theta \sim q} [\log p(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}, \theta)] - \text{KL}(q \| p(\cdot)). \end{aligned}$$

using eq. (5.62)

using eq. (5.63d)

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 (5.85),

$$\mathbb{E}_{\theta \sim q} [\log p(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}, \theta)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\log p(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}, \theta) \Big|_{\theta = \Sigma^{1/2} \epsilon + \mu} \right].$$

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,

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) &= \int p(y^* | \mathbf{x}^*, \theta) p(\theta | \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) d\theta \\ &= \mathbb{E}_{\theta \sim p(\cdot | \mathbf{x}_{1:n}, \mathbf{y}_{1:n})} [p(y^* | \mathbf{x}^*, \theta)] \\ &\approx \mathbb{E}_{\theta \sim q_\lambda} [p(y^* | \mathbf{x}^*, \theta)] \\ &\approx \frac{1}{m} \sum_{i=1}^m p(y^* | \mathbf{x}^*, \theta^{(i)}) \end{aligned} \tag{7.17}$$

using the sum rule (1.10) and product rule (1.14)

interpreting the integral as an expectation over the posterior

approximating the posterior with the variational posterior q_λ

using Monte Carlo sampling

for independent samples $\theta^{(i)} \stackrel{\text{iid}}{\sim} q_\lambda$,

$$= \frac{1}{m} \sum_{i=1}^m \mathcal{N}(y^*; \mu(\mathbf{x}^*; \theta^{(i)}), \sigma^2(\mathbf{x}^*; \theta^{(i)})). \tag{7.18}$$

using the neural network

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,

$$\mathbb{E}[y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, \mathbf{y}_{1:n}] \approx \frac{1}{m} \sum_{i=1}^m \mu(\mathbf{x}^*; \theta^{(i)}) \doteq \bar{\mu}(\mathbf{x}^*), \tag{7.19}$$

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.

⁸ see eq. (5.50)

and the variance of our predictions,

$$\begin{aligned}\text{Var}[y^* \mid \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}] &= \mathbb{E}_{\theta}[\text{Var}_{y^*}[y^* \mid \mathbf{x}^*, \theta]] + \text{Var}_{\theta}[\mathbb{E}_{y^*}[y^* \mid \mathbf{x}^*, \theta]] \\ &= \mathbb{E}_{\theta}[\sigma^2(\mathbf{x}^*; \theta)] + \text{Var}_{\theta}[\mu(\mathbf{x}^*; \theta)].\end{aligned}$$

using the law of total variance (1.50)

using the likelihood (7.15a)

Recall from eq. (2.22) 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 $\theta^{(i)}$,

$$\begin{aligned}\text{Var}[y^* \mid \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}] &\approx \frac{1}{m} \sum_{i=1}^m \sigma^2(\mathbf{x}^*; \theta^{(i)}) \\ &\quad + \frac{1}{m-1} \sum_{i=1}^m (\mu(\mathbf{x}^*; \theta^{(i)}) - \bar{\mu}(\mathbf{x}^*))^2\end{aligned}\tag{7.20}$$

using a sample mean (1.68) and sample variance (1.69).

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

$$p(y^* \mid \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}) \approx \frac{1}{T} \sum_{t=1}^T p(y^* \mid \mathbf{x}^*, \theta^{(t)}).$$

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*.

⁹ That is, combine the individual samples of weights $\theta^{(i)}$.

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,

$$\theta \sim \mathcal{N}(\mu, \Sigma), \quad \text{where} \tag{7.21a}$$

$$\mu \doteq \frac{1}{T} \sum_{i=1}^T \theta^{(i)}, \tag{7.21b}$$

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 (1.68)

$$\Sigma \doteq \frac{1}{T-1} \sum_{i=1}^T (\theta^{(i)} - \mu)(\theta^{(i)} - \mu)^\top, \quad (7.21c)$$

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 7.3: 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 eq. (7.21)), is known as the *stochastic weight averaging-Gaussian* (SWAG) method (Izmailov et al., 2018).

using a sample variance (1.69)

using the characterization of sample variance in terms of estimators of the first and second moment (1.70)

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 fig. 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

$$q(\theta \mid \lambda) \doteq \prod_{j=1}^d q_j(\theta_j \mid \lambda_j) \quad (7.24)$$

where d is the number of weights of the neural network and

$$q_j(\theta_j \mid \lambda_j) \doteq p\delta_0(\theta_j) + (1-p)\delta_{\lambda_j}(\theta_j). \quad (7.25)$$

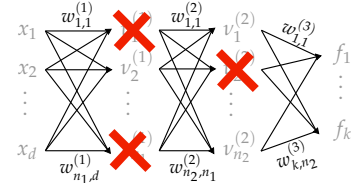


Figure 7.6: Illustration of dropout regularization. Some vertices of the computation graph are randomly omitted. In contrast, dropconnect regularization randomly omits edges of the computation graph.

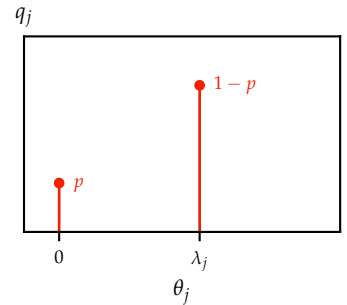


Figure 7.7: Interpretation of dropconnect regularization as variational inference. The only coordinates where the variational posterior q_j has positive probability are 0 and λ_j .

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 \mathbf{z} with entries $z(i) \sim \text{Bern}(p)$, yielding $\mathbf{z} \odot \lambda$ 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 eq. (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,

$$\begin{aligned} & \arg \max_{\lambda \in \Lambda} L(q_\lambda, p; \mathcal{D}) \\ &= \arg \max_{\lambda \in \Lambda} \mathbb{E}_{\theta \sim q_\lambda} [\log p(y_{1:n} | \mathbf{x}_{1:n}, \theta)] - \text{KL}(q \| p(\cdot)) \\ &\approx \arg \min_{\lambda \in \Lambda} -\frac{1}{m} \sum_{i=1}^m \log p(y_{1:n} | \mathbf{x}_{1:n}, \theta) \Big|_{\theta = \mathbf{z}^{(i)} \odot \lambda + \epsilon^{(i)}} + \frac{p}{2} \|\lambda\|_2^2 \end{aligned} \quad (7.27)$$

¹¹ see example 1.18

¹² $\mathbf{A} \odot \mathbf{B}$ denotes the Hadamard (element-wise) product.

using eq. (5.63d)

using Monte Carlo sampling

where we reparameterize $\theta \sim q_\lambda$ by $\theta = \mathbf{z} \odot \lambda + \epsilon$ with $\mathbf{z}^{(i)} \sim \text{Bern}(p)$ and $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. 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,

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{x}_{1:n}, \mathbf{y}_{1:n}) &\approx \mathbb{E}_{\theta \sim q_\lambda} [p(y^* | \mathbf{x}^*, \theta)] \\ &\approx \frac{1}{m} \sum_{i=1}^m p(y^* | \mathbf{x}^*, \theta^{(i)}) \end{aligned} \quad (7.28)$$

using Monte Carlo sampling

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 eq. (7.17). In words, we average the predictions of m neural networks for each of which we randomly “drop out” weights.

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

$$\begin{aligned} p(y^* \mid x^*, x_{1:n}, y_{1:n}) &= \mathbb{E}_{\theta \sim p(\cdot \mid x_{1:n}, y_{1:n})} [p(y^* \mid x^*, \theta)] \\ &\approx \frac{1}{m} \sum_{i=1}^m p(y^* \mid x^*, \theta^{(i)}). \end{aligned} \quad (7.29) \quad \text{using bootstrapping}$$

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 sufficient diversity between models (Lakshminarayanan et al., 2017).

Note that eq. (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 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.

7.4 Calibration

A key challenge of Bayesian deep learning (and also other Bayesian 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.

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 $\mathbf{x}_{1:m}^{\text{val}}$ of size m given the training set $\mathbf{x}_{1:n}^{\text{train}}$ 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,

$$\begin{aligned}
 & \log p(y_{1:m}^{\text{val}} | \mathbf{x}_{1:m}^{\text{val}}, \mathbf{x}_{1:n}^{\text{train}}, y_{1:n}^{\text{train}}) \\
 &= \log \int p(y_{1:m}^{\text{val}} | \mathbf{x}_{1:m}^{\text{val}}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{x}_{1:n}^{\text{train}}, y_{1:n}^{\text{train}}) d\boldsymbol{\theta} && \text{using the sum rule (1.10) and product rule (1.14)} \\
 &\approx \log \int p(y_{1:m}^{\text{val}} | \mathbf{x}_{1:m}^{\text{val}}, \boldsymbol{\theta}) q_{\lambda}(\boldsymbol{\theta}) d\boldsymbol{\theta} && \text{approximating with the variational posterior} \\
 &= \log \int \prod_{i=1}^m p(y_i^{\text{val}} | \mathbf{x}_i^{\text{val}}, \boldsymbol{\theta}) q_{\lambda}(\boldsymbol{\theta}) d\boldsymbol{\theta} && \text{using the independence of the data}
 \end{aligned} \tag{7.30}$$

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,

$$\begin{aligned}
 &= \log \mathbb{E}_{\boldsymbol{\theta} \sim q_{\lambda}} \left[\prod_{i=1}^m p(y_i^{\text{val}} | \mathbf{x}_i^{\text{val}}, \boldsymbol{\theta}) \right] && \text{interpreting the integral as an expectation over the variational posterior} \\
 &\geq \mathbb{E}_{\boldsymbol{\theta} \sim q_{\lambda}} \left[\sum_{i=1}^m \log p(y_i^{\text{val}} | \mathbf{x}_i^{\text{val}}, \boldsymbol{\theta}) \right] && \text{using Jensen's inequality (5.37)}
 \end{aligned} \tag{7.31}$$

$$\approx \frac{1}{k} \sum_{j=1}^k \sum_{i=1}^m \log p(y_i^{\text{val}} | \mathbf{x}_i^{\text{val}}, \boldsymbol{\theta}^{(j)}) \tag{7.32}$$

using Monte Carlo sampling

for independent samples $\boldsymbol{\theta}^{(j)} \stackrel{\text{iid}}{\sim} q_{\lambda}$.

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).¹³

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, $\mathbb{P}(Y_i = 1 | \mathbf{x}_i)$. 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 define B_m as the set of samples falling into bin m and let

$$\text{freq}(B_m) \doteq \frac{1}{|B_m|} \sum_{i \in B_m} \mathbb{1}\{Y_i = 1\} \tag{7.33}$$

be the proportion of samples in bin m that belong to class 1 and let

$$\text{conf}(B_m) \doteq \frac{1}{|B_m|} \sum_{i \in B_m} \mathbb{P}(Y_i = 1 | \mathbf{x}_i) \tag{7.34}$$

¹³ Reliability diagrams generalize beyond this restricted example.

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,

$$\ell_{\text{ECE}} \doteq \sum_{m=1}^M \frac{|B_m|}{n} |\text{freq}(B_m) - \text{conf}(B_m)| \quad (7.35)$$

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,

$$\ell_{\text{MCE}} \doteq \max_{m \in [M]} |\text{freq}(B_m) - \text{conf}(B_m)|. \quad (7.36)$$

7.4.3 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,

$$\min_{M, f, a} \sum_{m=1}^M \sum_{i=1}^n \mathbb{1}\{a_m \leq \mathbb{P}(Y_i = 1 \mid \mathbf{x}_i) < a_{m+1}\} (f_m - y_i)^2 \quad (7.37a)$$

$$\text{subject to } 0 = a_1 \leq \dots \leq a_{M+1} = 1, \quad (7.37b)$$

$$f_1 \leq \dots \leq f_M \quad (7.37c)$$

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.

3. *Platt scaling* adjusts the logits z_i of the output layer to

$$q_i \doteq \sigma(az_i + b) \quad (7.38)$$

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

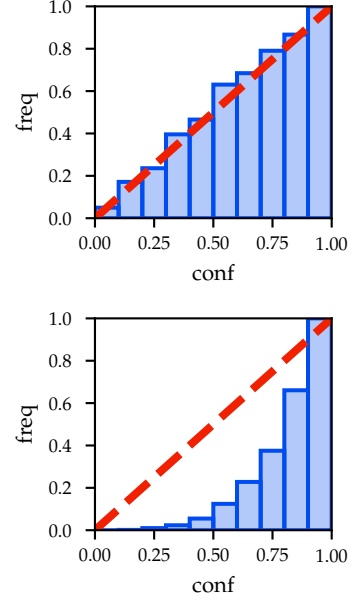


Figure 7.8: Examples of reliability diagrams with ten bins. A perfectly calibrated model approximates the diagonal dashed red line. The first reliability diagram shows a well-calibrated model. In contrast, the second reliability diagram shows an overconfident model.

$T \in \mathbb{R}$,

$$q_i \doteq \sigma\left(\frac{z_i}{T}\right). \quad (7.39)$$

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 exercise 6.25, 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?.

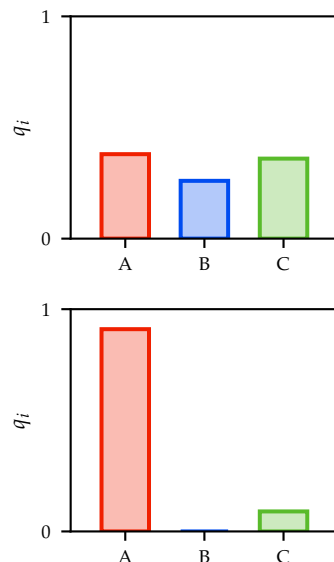


Figure 7.9: Illustration of temperature scaling for a classifier with three classes. On the top, we have a prediction with a high temperature, yielding a very uncertain prediction (in favor of class A). Below, we have a prediction with a low temperature, yielding a prediction that is strongly in favor of class A . Note that the ranking ($A \succ C \succ B$) is preserved.

PART II

Sequential Decision-Making

8

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*. In section 8.5, we briefly discuss other ways of quantifying usefulness of 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 $x \sim \mathbf{X}$. 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., $\mathbf{Y} = y$ for a random vector \mathbf{Y}),

$$H[\mathbf{X} \mid \mathbf{Y} = y] \doteq \mathbb{E}_{x \sim p(x|y)}[-\log p(x \mid y)]. \quad (8.1)$$

Instead of averaging over the surprise of samples from the distribution $p(x)$ (like the entropy $H[\mathbf{X}]$), this quantity simply averages over the surprise of samples from the conditional distribution $p(x \mid y)$.

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

$$H[\mathbf{X} \mid \mathbf{Y}] \doteq \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y})}[H[\mathbf{X} \mid \mathbf{Y} = \mathbf{y}]] \quad (8.2)$$

$$= \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y})}[-\log p(\mathbf{x} \mid \mathbf{y})]. \quad (8.3)$$

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 $H[\mathbf{X} \mid \mathbf{Y} = \mathbf{y}]$ and the conditional entropy $H[\mathbf{X} \mid \mathbf{Y}]$. The former simply corresponds to a Bayesian update of our uncertainty in \mathbf{X} after we have observed the realization $\mathbf{y} \sim \mathbf{Y}$. 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} ,

$$H[\mathbf{X}, \mathbf{Y}] \doteq \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y})}[-\log p(\mathbf{x}, \mathbf{y})], \quad (8.4)$$

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 (1.14) 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,

$$I(\mathbf{X}; \mathbf{Y}) = I(\mathbf{Y}; \mathbf{X}). \quad (8.11)$$

Exercise 8.4: Mutual information and KL divergence

Show that by expanding the definition of mutual information,

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y}) &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p} \left[\log \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} \right] \\ &= \text{KL}(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x})p(\mathbf{y})) \end{aligned} \quad (8.12)$$

$$= \mathbb{E}_{\mathbf{y} \sim p} [\text{KL}(p(\mathbf{x} | \mathbf{y}) \| p(\mathbf{x}))], \quad (8.13)$$

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} .

▷ *Solution*

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. exercise 5.19), 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

$$0 \leq I(\mathbf{X}; \mathbf{Y}) = H[\mathbf{X}] - H[\mathbf{X} | \mathbf{Y}]. \quad (8.14)$$

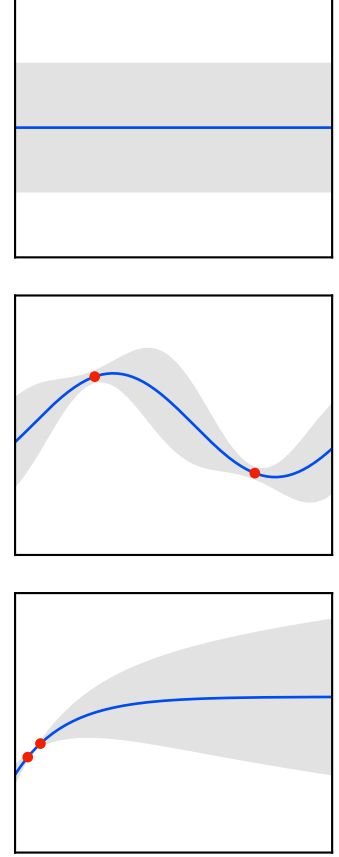


Figure 8.1: Information gain. The first graph shows the prior. The second graph shows a selection of samples with a large information gain (large reduction in uncertainty). The third graph shows a selection of samples with a small information gain (small reduction in uncertainty).

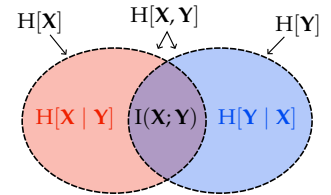


Figure 8.2: Relationship between mutual information and entropy, expressed as a Venn diagram.

Example 8.5: Mutual information of Gaussians

Given the Gaussian random vector $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and the noisy observation $\mathbf{Y} = \mathbf{X} + \epsilon$ where $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I})$, we want to find the information gain of \mathbf{X} when observing \mathbf{Y} . Using our definitions from this chapter, we obtain

$$\begin{aligned}
 I(\mathbf{X}; \mathbf{Y}) &= I(\mathbf{Y}; \mathbf{X}) \\
 &= H[\mathbf{Y}] - H[\mathbf{Y} | \mathbf{X}] \\
 &= H[\mathbf{Y}] - H[\epsilon] \\
 &= \frac{1}{2} \log \left((2\pi e)^d \det(\boldsymbol{\Sigma} + \sigma_n^2 \mathbf{I}) \right) - \frac{1}{2} \log \left((2\pi e)^d \det(\sigma_n^2 \mathbf{I}) \right) \\
 &= \frac{1}{2} \log \frac{\det(\boldsymbol{\Sigma} + \sigma_n^2 \mathbf{I})}{\det(\sigma_n^2 \mathbf{I})} \\
 &= \frac{1}{2} \log \det(\mathbf{I} + \sigma_n^{-2} \boldsymbol{\Sigma}). \tag{8.15}
 \end{aligned}$$

Intuitively, the larger the noise σ_n^2 in relation to the covariance of \mathbf{X} , the smaller the information gain.

using symmetry (8.11)

by mutual information (8.9)

given \mathbf{X} , the only randomness in \mathbf{Y} originates from ϵ

using the entropy of Gaussians (5.36)

8.2.1 Data Processing Inequality and Sufficient Statistics

Definition 8.6 (Conditional mutual information). The *conditional mutual information* of \mathbf{X} and \mathbf{Y} given \mathbf{Z} is defined as

$$I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) \doteq H[\mathbf{X} | \mathbf{Z}] - H[\mathbf{X} | \mathbf{Y}, \mathbf{Z}]. \tag{8.16}$$

$$= H[\mathbf{X}, \mathbf{Z}] + H[\mathbf{Y}, \mathbf{Z}] - H[\mathbf{Z}] - H[\mathbf{X}, \mathbf{Y}, \mathbf{Z}] \tag{8.17}$$

$$= I(\mathbf{X}; \mathbf{Y}, \mathbf{Z}) - I(\mathbf{X}; \mathbf{Z}). \tag{8.18}$$

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:

$$I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) = I(\mathbf{Y}; \mathbf{X} | \mathbf{Z}). \tag{8.19}$$

Exercise 8.7: Non-monotonicity of cond. mutual information

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, we will see in this exercise that 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} .

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 $\mathbf{X} \perp \mathbf{Z} \mid \mathbf{Y}$ then $I(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z}) \leq I(\mathbf{X}; \mathbf{Y})$.
3. Note that the condition $\mathbf{X} \perp \mathbf{Z} \mid \mathbf{Y}$ 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

$$I(\mathbf{X}; \mathbf{Z}) \leq I(\mathbf{X}; \mathbf{Y}). \quad (8.20)$$

which is also known as the *data processing inequality*, and which says that processing cannot increase the information contained in a signal.

▷ *Solution*

The data processing inequality 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 (8.20), $I(\lambda; s(\mathbf{X})) \leq I(\lambda; \mathbf{X})$. 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 .

8.2.2 Synergy and Redundancy

To understand the behavior of mutual information under conditioning, it is helpful to consider the *interaction information*

$$I(\mathbf{X}; \mathbf{Y}; \mathbf{Z}) \doteq I(\mathbf{X}; \mathbf{Y}) - I(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z}). \quad (8.21)$$

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} .

Exercise 8.8: 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 $I(Y; X_1; X_2)$.

▷ *Solution*

Readings

Read sections 2.2 through 2.6 of “Elements of information theory” (Cover and Thomas, 2006) for additional background on entropy, mutual information and the KL-divergence.

8.2.3 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 $x \in \mathcal{X}$,

$$y_x \doteq f_x + \epsilon_x, \quad (8.22)$$

$f_x \doteq f(x)$, and ϵ_x is some zero-mean Gaussian noise. For a set of observations $S = \{x_1, \dots, x_n\}$, we can write $\mathbf{y}_S = \mathbf{f}_S + \boldsymbol{\epsilon}$ where

$$\mathbf{y}_S \doteq \begin{bmatrix} y_{x_1} \\ \vdots \\ y_{x_n} \end{bmatrix}, \quad \mathbf{f}_S \doteq \begin{bmatrix} f_{x_1} \\ \vdots \\ f_{x_n} \end{bmatrix}, \quad \text{and} \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I}).$$

Note that both \mathbf{y}_S and \mathbf{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,

$$I(S) \doteq I(\mathbf{f}_S; \mathbf{y}_S) = H[\mathbf{f}_S] - H[\mathbf{f}_S \mid \mathbf{y}_S]. \quad (8.23)$$

Here, $H[\mathbf{f}_S]$ corresponds to the uncertainty about \mathbf{f}_S before obtaining the observations \mathbf{y}_S and $H[\mathbf{f}_S \mid \mathbf{y}_S]$ corresponds to the uncertainty about \mathbf{f}_S after obtaining the observations \mathbf{y}_S .

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 8.9 (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

$$\Delta_F(x \mid A) \doteq F(A \cup \{x\}) - F(A). \quad (8.24)$$

Intuitively, the marginal gain describes how much “adding” the additional x to A increases the value of F .

Exercise 8.10: Marginal gain of maximizing mutual information

Show that in the context of maximizing mutual information, the marginal gain is

$$\Delta_I(x \mid A) = I(f_x; y_x \mid \mathbf{y}_A) \quad (8.25)$$

$$= H[y_x \mid \mathbf{y}_A] - H[\epsilon_x]. \quad (8.26)$$

▷ *Solution*

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 $H[\epsilon_x]$. Altogether, the marginal gain represents the reduction in uncertainty by observing $\{x\}$.

Definition 8.11 (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

$$F(A \cup \{x\}) - F(A) \geq F(B \cup \{x\}) - F(B). \quad (8.27)$$

Equivalently, using our definition of marginal gain, we have that F is submodular iff for any $x \in \mathcal{X}$ and any $A \subseteq B \subseteq \mathcal{X}$,

$$\Delta_F(x \mid A) \geq \Delta_F(x \mid B). \quad (8.28)$$

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 8.12 (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

$$F(A) \leq F(B). \quad (8.29)$$

If F is also submodular, then F is called *monotone submodular*.

Theorem 8.13. *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,

$$I \text{ is submodular} \iff \Delta_I(x \mid A) \geq \Delta_I(x \mid B)$$

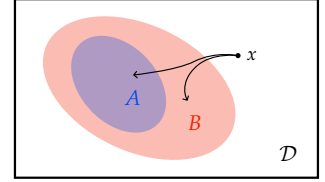


Figure 8.3: Monotone submodularity. The effect of “adding” x to the smaller set A is larger than the effect of adding x to the larger set B .

by submodularity in terms of marginal gain (8.28)

$$\begin{aligned}
&\iff H[y_x | y_A] - H[\epsilon_x] \geq H[y_x | y_B] - H[\epsilon_x] && \text{using eq. (8.26)} \\
&\iff H[y_x | y_A] \geq H[y_x | y_B]. && H[\epsilon_x] \text{ cancels}
\end{aligned}$$

Due to the “information never hurts” principle (8.8) of entropy and as $A \subseteq B$, this is always true. Moreover,

$$\begin{aligned}
I \text{ is monotone} &\iff I(A) \leq I(B) && \text{by the definition of monotonicity (8.29)} \\
&\iff I(f_A; y_A) \leq I(f_B; y_B) && \text{using the definition of } I \text{ (8.23)} \\
&\iff I(f_B; y_A) \leq I(f_B; y_B) && \text{using } I(f_B; y_A) = I(f_A; y_A) \text{ as } y_A \perp f_B | f_A \\
&\iff H[f_B] - H[f_B | y_A] \leq H[f_B] - H[f_B | y_B] && \text{using the definition of MI (8.9)} \\
&\iff H[f_B | y_A] \geq H[f_B | y_B], && H[f_B] \text{ cancels}
\end{aligned}$$

which is also satisfied due to the “information never hurts” principle (8.8). \square

Exercise 8.14: Submodularity means no synergy

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. To this end, show that for all $A \subseteq B$,

$$I(f_x; y_x; y_{B \setminus A} | y_A) \geq 0. \quad (8.30)$$

▷ *Solution*

8.4 Maximizing Mutual Information

As we cannot efficiently pick a set $S \subseteq \mathcal{X}$ to maximize mutual information, 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 this greedy approach provides a good approximation.

Theorem 8.15 (Greedy submodular function maximization). *If F is monotone submodular, then greedily maximizing F provides a $(1 - 1/e)$ -approximation,²*

² $1 - 1/e \approx 0.632$

$$F(S_n) \geq \left(1 - \frac{1}{e}\right) \max_{\substack{S \subseteq \mathcal{X} \\ |S|=n}} F(S). \quad (8.31)$$

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 (8.29) of F . We write,

$\{x_1^*, \dots, x_n^*\} \doteq S^*$. We have,

$$F(S^*) \leq F(S^* \cup S_t)$$

using monotonicity (8.29)

$$= F(S_t) + \sum_{i=1}^n \Delta_F(x_i^* \mid S_t \cup \{x_1^*, \dots, x_{i-1}^*\})$$

using the definition of marginal gain (8.24)

$$\leq F(S_t) + \sum_{x^* \in S^*} \Delta_F(x^* \mid S_t)$$

using submodularity (8.28)

$$\leq F(S_t) + n(F(S_{t+1}) - F(S_t)).$$

using that $S_{t+1} = S_t \cup \{x\}$ is chosen such that $\Delta_F(x \mid S_t)$ is maximized (8.32)

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

Readings

The original proof of greedy maximization for submodular functions was given by “An analysis of approximations for maximizing submodular set functions” (Nemhauser et al., 1978).

For more background on maximizing submodular functions, refer to “Submodular function maximization” (Krause and Golovin, 2014).

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,

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} \Delta_I(x \mid S_t) \quad (8.32)$$

$$= \arg \max_{x \in \mathcal{X}} I(f_x; y_x \mid y_{S_t}). \quad (8.33)$$

using eq. (8.25)

Note that f_x and y_x are univariate random variables. Thus, using our formula for the mutual information of conditional linear Gaussians (8.15), we can simplify to,

$$= \arg \max_{x \in \mathcal{X}} \frac{1}{2} \log \left(1 + \frac{\sigma_t^2(x)}{\sigma_n^2} \right) \quad (8.34)$$

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),

$$= \arg \max_{x \in \mathcal{X}} \sigma_t^2(x). \quad (8.35)$$

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. See fig. 8.4, for an example.

Looking at eq. (8.34) suggests a natural fix. Instead of only considering the epistemic uncertainty $\sigma_t^2(x)$, we can also consider the aleatoric uncertainty $\sigma_n^2(x)$ by modeling heteroscedastic noise, yielding

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} \frac{1}{2} \log \left(1 + \frac{\sigma_t^2(x)}{\sigma_n^2(x)} \right) = \arg \max_{x \in \mathcal{X}} \frac{\sigma_t^2(x)}{\sigma_n^2(x)}. \quad (8.36)$$

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 ,

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} H[y_x \mid x_{1:t}, y_{1:t}]. \quad (8.37)$$

The entropy of a categorical distribution is simply a finite sum of weighted surprise terms.

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

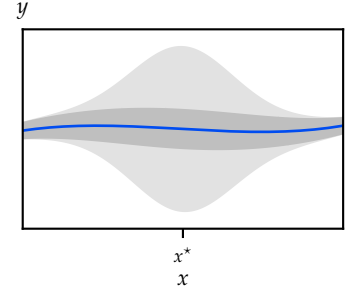


Figure 8.4: Uncertainty sampling with heteroscedastic noise. The epistemic uncertainty of the model is shown in a dark gray. The aleatoric uncertainty of the data is shown in a light gray. Uncertainty sampling would repeatedly pick points around x^* as they maximize the epistemic uncertainty, even though the aleatoric uncertainty at x^* is much larger than at the boundary.

³ see section 1.2

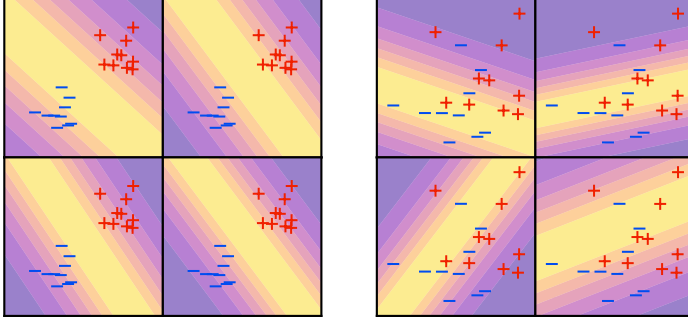


Figure 8.5: Uncertainty sampling in classification. The area with high uncertainty (as measured by entropy) is highlighted in yellow. The shown figures display each a sequence of model updates, each after one new observation. In the left figure, the classes are well-separated and uncertainty is dominated 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.

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,

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} I(\theta; y_x \mid x_{1:t}, y_{1:t}) \quad (8.38)$$

$$= \arg \max_{x \in \mathcal{X}} I(y_x; \theta \mid x_{1:t}, y_{1:t})$$

$$= \arg \max_{x \in \mathcal{X}} H[y_x \mid x_{1:t}, y_{1:t}] - H[y_x \mid \theta, x_{1:t}, y_{1:t}]$$

$$= \arg \max_{x \in \mathcal{X}} H[y_x \mid x_{1:t}, y_{1:t}] - \mathbb{E}_{\theta \mid x_{1:t}, y_{1:t}} [H[y_x \mid \theta, x_{1:t}, y_{1:t}]] \quad (8.39)$$

$$= \arg \max_{x \in \mathcal{X}} \underbrace{H[y_x \mid x_{1:t}, y_{1:t}]}_{\text{entropy of predictive posterior}} - \mathbb{E}_{\theta \mid x_{1:t}, y_{1:t}} \underbrace{H[y_x \mid \theta]}_{\text{entropy of likelihood}}. \quad (8.40)$$

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 (5.32) and assuming $y_x \perp x_{1:t}, y_{1:t} \mid \theta$

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 eq. (8.37). The second term mirrors our description of [aleatoric uncertainty](#) from section 2.3. 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, eq. (8.39) 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 $p(y \mid x, \theta)$ 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 Beyond Mutual Information: Experimental Design

The problem of identifying which experiments to conduct in order to maximize learning is studied extensively in the field of *experimental design*. 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).

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*. In the following, we give an overview of the most common design criteria. We continue to focus on the setting of Gaussian likelihood and prior, and for simplicity we assume that the domain \mathcal{X} is finite, so our objective is to “learn” $f \doteq f_{\mathcal{X}}$.⁴

- A *d-optimal design* selects S to minimize the determinant of the posterior covariance matrix which in the Gaussian case corresponds to minimizing posterior entropy (MacKay, 1992):

$$\arg \min_S \det(\text{Var}[f \mid y_S]) = \arg \min_S H[f \mid y_S] = \arg \max_S I(f; y_S). \quad (8.41)$$

⁴ If the domain is not finite, we can alternatively learn the parameters θ of a parametric model f akin to BALD (cf. section 8.4.3).

using the entropy of a Gaussian (5.36) and the definition of mutual information (8.9)

This is the criterion that we have been studying so far,⁵ and as seen in section 8.4.1, in a regression problem it coincides with uncertainty sampling when the noise is homoscedastic. Geometrically, this criterion can be interpreted as selecting the design which leads to the largest reduction of volume of the uncertainty ellipsoid.

- A common alternative is *a-optimal design* which selects S to minimize the posterior average variance, i.e., the trace of the posterior covariance matrix:

$$\arg \min_S \text{tr}(\text{Var}[\mathbf{f} \mid \mathbf{y}_S]). \quad (8.42)$$

This criterion may be more useful when the goal is to obtain a good estimate of the posterior means (i.e., make marginal predictions), and one does not care about the posterior interaction between random variables.

- Yet another alternative is *e-optimal design* which selects S to minimize the maximum eigenvalue of the posterior covariance matrix:

$$\arg \min_S \lambda_{\max}(\text{Var}[\mathbf{f} \mid \mathbf{y}_S]). \quad (8.43)$$

Geometrically, this can be interpreted as minimizing the diameter of the uncertainty ellipsoid.

Observe that these design criteria are all grounded in different “scalarizations” of the posterior covariance matrix.

8.6 Beyond Global Learning

Active learning and experimental design are typically concerned with reducing the uncertainty about f *globally* across the domain. However, in many practical applications, we are only interested in a *local* region of the domain.

A common example is determining where f is maximized. The case where we *additionally* seek to only evaluate a limited number of suboptimal points while searching for the maximum will be discussed thoroughly in the next chapter, but for now, we only want to find the location of the maximum. In view of *multi-armed bandits* (which we will introduce in section 9.2.1), this is also known as the *best-arm identification* problem when the domain is finite.

One viable approach is to use standard design criteria that reduce uncertainty globally, however, they attempt to reduce uncertainty in areas of the domain that are already known to be suboptimal. By restricting the “points of interest”, one can typically achieve a significantly better sample complexity.

⁵ Note that since $f_{\mathcal{X} \setminus S} \perp \mathbf{y}_S \mid f_S$, it follows from eq. (8.18) that

$$I(\mathbf{f}; \mathbf{y}_S) = I(f_S; \mathbf{y}_S) + \underbrace{I(f_{\mathcal{X} \setminus S}; \mathbf{y}_S \mid f_S)}_0.$$

8.6.1 Entropy Search

One canonical family of methods is inspired by the d-optimality criterion. Let $\mathbf{x}^* \doteq \arg \max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ denote the location of the maximum and select the observation that maximizes the information gain about the location of the maximum,

$$\mathbf{x}_{t+1} \doteq \arg \max_{\mathbf{x} \in \mathcal{X}} I(\mathbf{x}^*; y_{\mathbf{x}} \mid \mathbf{x}_{1:t}, y_{1:t}). \quad (8.44)$$

Using symmetry of mutual information (8.11), we can write this in two different ways:

$$= \arg \min_{\mathbf{x} \in \mathcal{X}} \mathbb{E}_{y_{\mathbf{x}} \mid \mathbf{x}_{1:t}, y_{1:t}} [H[\mathbf{x}^* \mid \mathbf{x}_{1:t}, y_{1:t}, y_{\mathbf{x}}]] \quad (8.45)$$

$$= \arg \max_{\mathbf{x} \in \mathcal{X}} H[y_{\mathbf{x}} \mid \mathbf{x}_{1:t}, y_{1:t}] - \mathbb{E}_{\mathbf{x}^* \mid \mathbf{x}_{1:t}, y_{1:t}} [H[y_{\mathbf{x}} \mid \mathbf{x}_{1:t}, y_{1:t}, \mathbf{x}^*]] \quad (8.46)$$

where to simplify the expressions, we used the definition of mutual information (8.9) and the definition of conditional entropy (8.2). The algorithm in eq. (8.45) is known as *entropy search* (ES) (Hennig and Schuler, 2012), and perhaps the most natural as it selects \mathbf{x}_{t+1} to minimize the posterior entropy of \mathbf{x}^* . However, approximating the entropy of the high-dimensional posterior distribution over \mathbf{x}^* is difficult.

More amenable to approximations is the formulation from eq. (8.46) which is known as *predictive entropy search* (PES) (Hernández-Lobato et al., 2014). Observe that the derivation of PES is analogous to BALD (8.39), only that here we condition on the location of the maximum \mathbf{x}^* rather than on the model parameters θ . The first term of eq. (8.46) is the entropy of the predictive distribution of $y_{\mathbf{x}}$ which we know to be

$$H[y_{\mathbf{x}} \mid \mathbf{x}_{1:t}, y_{1:t}] = \frac{1}{2} \log \left(2\pi e (\sigma_f^2(\mathbf{x}) + \sigma_n^2) \right). \quad \text{using eq. (5.35)}$$

The expectation of the second term can be approximated using Monte carlo sampling by drawing m functions from the posterior distribution and taking the maximum of each function as a sample of \mathbf{x}^* . The main computational difficulty lies in approximating the entropy of the “predictive” distribution $y_{\mathbf{x}} \mid \mathbf{x}_{1:t}, y_{1:t}, \mathbf{x}^*$, but unlike ES, this distribution is one-dimensional.

The dependence of ES and PES on the posterior over \mathbf{x}^* is problematic as approximations are typically expensive to compute. An alternative approach is to measure information gain with respect to the maximum value $f^* \doteq \max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ rather than the location of the maximum \mathbf{x}^* :

$$\mathbf{x}_{t+1} \doteq \arg \max_{\mathbf{x} \in \mathcal{X}} I(f^*; y_{\mathbf{x}} \mid \mathbf{x}_{1:t}, y_{1:t}) \quad (8.47)$$

$$= \arg \max_{x \in \mathcal{X}} H[y_x \mid \mathbf{x}_{1:t}, y_{1:t}] - \mathbb{E}_{f^* \mid \mathbf{x}_{1:t}, y_{1:t}} [H[y_x \mid \mathbf{x}_{1:t}, y_{1:t}, f^*]]. \quad (8.48)$$

using the definition of mutual information (8.9) and the definition of conditional entropy (8.2)

The *max-value entropy search* (MES) approximation ignores the observation noise in $y_x \mid \mathbf{x}_{1:t}, y_{1:t}, f^*$ and approximates $f_x \mid \mathbf{x}_{1:t}, y_{1:t}, f^*$ by a truncated normal, for which the entropy can be approximated well in closed form (Wang and Jegelka, 2017). Since MES ignores the observation noise, it may significantly overestimate the information gain when the signal-to-noise ratio is low. In contrast, *output-space predictive entropy search* (OPES) approximates $f_x \mid \mathbf{x}_{1:t}, y_{1:t}, f^*$ by a Gaussian, which gives $H[y_x \mid \mathbf{x}_{1:t}, y_{1:t}, f^*] \approx \frac{1}{2} \log(2\pi e(\sigma_*^2(\mathbf{x}) + \sigma_n^2))$ where $\sigma_*^2(\mathbf{x})$ denotes the variance of the Gaussian approximation (Hoffman and Ghahramani, 2015).

It may not be immediately obvious why the maximum function value yields sufficient information to find the location of the maximum. In practice, “optimal value identification” often works quite well. Intuitively, determining the optimal value requires sufficiently exploring all approximately optimal points.

Optional Readings

More details on entropy search methods can be found in sections 8.8 and 8.9 of “Bayesian optimization” (Garnett, 2023).

For an example of a local criterion akin to a-optimality, see “Truncated variance reduction: A unified approach to bayesian optimization and level-set estimation” (Bogunovic et al., 2016).

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 course. 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^* ,

$$y_t = f^*(x_t) + \epsilon_t. \quad (9.2)$$

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

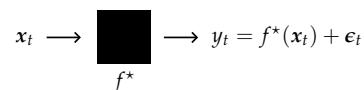


Figure 9.1: Illustration of Bayesian optimization. We pass an input x_t into the unknown function f^* to obtain noisy observations y_t .

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^* .

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 between 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 un-*

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 .

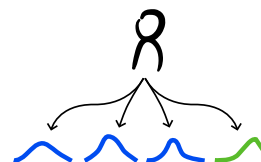


Figure 9.2: Illustration of a multi-armed bandit with four arms, each with a different reward distribution. The agent tries to identify the arm with the most beneficial reward distribution shown in green.

certainty”, 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.

Readings

Refer to

- “Regret analysis of stochastic and nonstochastic multi-armed bandit problems” (Bubeck et al., 2012),
- “Bandit algorithms” (Lattimore and Szepesvári, 2020), or
- chapter 2 of “Reinforcement learning: An introduction” (Sutton and Barto, 2018)

for a comprehensive overview of multi-armed bandits.

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.

Exercise 9.2: Convergence to the static optimum

Show that any algorithm achieves sublinear regret if and only if it converges to the static optimum, that is,

$$\lim_{t \rightarrow \infty} f^*(x_t) = \max_x f^*(x). \quad (9.6)$$

Hint: Use that if a sequence a_n converges to a as $n \rightarrow \infty$, then we have for the sequence

$$b_n \doteq \frac{1}{n} \sum_{i=1}^n a_i \quad (9.7)$$

that $\lim_{n \rightarrow \infty} b_n = a$. This is also known as the Cesàro mean.

▷ *Solution*

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 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 neural networks. The various acquisition functions F are used in the same way as is illustrated in alg. 9.3.

One possible acquisition function is uncertainty sampling (8.34), 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.

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.

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*,

$$\text{cost}(\text{ALG}) \leq \alpha \cdot \text{cost}(\text{OPT}),$$

where OPT corresponds to the dynamic optimal choice (in hindsight).

Algorithm 9.3: Bayesian optimization (with GPs)

```

1 initialize  $f \sim \mathcal{GP}(\mu_0, k_0)$ 
2 for  $t = 1$  to  $T$  do
3   choose  $\mathbf{x}_t = \arg \max_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}; \mu_{t-1}, k_{t-1})$ 
4   observe  $y_t = f(\mathbf{x}_t) + \epsilon_t$ 
5   perform a Bayesian update to obtain  $\mu_t$  and  $k_t$ 

```

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 fig. 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 very 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),

$$\mathbf{x}_{t+1} \doteq \arg \max_{\mathbf{x} \in \mathcal{X}} \mu_t(\mathbf{x}) + \beta_{t+1} \sigma_t(\mathbf{x}), \quad (9.8)$$

where $\sigma_t(\mathbf{x}) \doteq \sqrt{k_t(\mathbf{x}, \mathbf{x})}$ is the standard deviation at \mathbf{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.

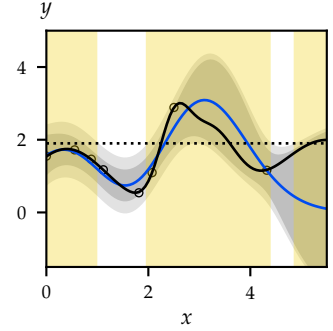


Figure 9.3: Optimism in Bayesian optimization. The **unknown function** is shown in black, our **model** in blue with gray confidence bounds. The dotted black line denotes the maximum lower bound. We can therefore focus our exploration to the yellow regions where the upper confidence bound is higher than the maximum lower bound.

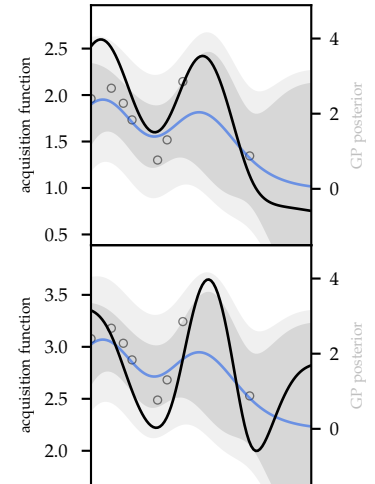


Figure 9.4: Plot of the UCB acquisition function for $\beta = 0.25$ and $\beta = 1$, respectively.

⁸ Due to the monotonicity of $(\cdot)^2$, it does not matter whether we optimize the variance or standard deviation at \mathbf{x} .

As can be seen in fig. 9.4, the UCB acquisition function is generally non-convex. For selecting the next point, we can use approximate 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 fig. 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 $1 - \delta$,

$$\forall t \geq 1, \forall \mathbf{x} \in \mathcal{X} : f^*(\mathbf{x}) \in \mathcal{C}_t(\mathbf{x}) \doteq [\mu_{t-1}(\mathbf{x}) \pm \beta_t(\delta) \cdot \sigma_{t-1}(\mathbf{x})]. \quad (9.9)$$

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 an 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*

$$\beta_t(\delta) = \mathcal{O}\left(\sqrt{\log(|\mathcal{X}|t/\delta)}\right) \quad (9.10)$$

satisfies $\mathbb{P}(\forall t \geq 1, \mathbf{x} \in \mathcal{X} : f^(\mathbf{x}) \in \mathcal{C}_t(\mathbf{x})) \geq 1 - \delta$.*

Exercise 9.5: Bayesian confidence intervals

In this exercise, we derive the above theorem.

1. For fixed $t \geq 1$ and $\mathbf{x} \in \mathcal{X}$, prove

$$\mathbb{P}(f^*(\mathbf{x}) \notin \mathcal{C}_t(\mathbf{x}) \mid \mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1}) \leq e^{-\beta_t^2/2}. \quad (9.11)$$

Hint: Bound $\mathbb{P}(Z > c)$ for $Z \sim \mathcal{N}(0, 1)$ and $c > 0$.

2. Prove theorem 9.4.

▷ *Solution*

Under the assumption of well-calibrated confidence intervals, we can bound the regret of UCB.

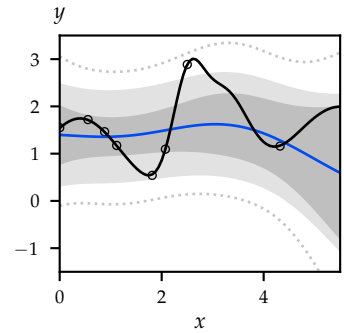


Figure 9.5: Re-scaling the confidence bounds. The dotted gray lines represent updated confidence bounds.

Theorem 9.6 (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

$$R_T = \mathcal{O}\left(\beta_T(\delta)\sqrt{\gamma_T T}\right) \quad (9.12)$$

where

$$\gamma_T \doteq \max_{\substack{S \subseteq \mathcal{X} \\ |S|=T}} \mathcal{I}(f_S; \mathbf{y}_S) \quad (9.13)$$

is the maximum information gain after T rounds.

Exercise 9.7: Regret of GP-UCB

To develop some intuition, we will derive the above theorem.

1. Show that if eq. (9.9) 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 $\mathbf{f}_T \doteq \mathbf{f}_{S_T}$ and $\mathbf{y}_T \doteq \mathbf{y}_{S_T}$. Prove

$$\mathcal{I}(\mathbf{f}_T; \mathbf{y}_T) = \frac{1}{2} \sum_{t=1}^T \log \left(1 + \frac{\sigma_{t-1}^2(\mathbf{x}_t)}{\sigma_n^2} \right). \quad (9.14)$$

3. Combine (1) and (2) to show theorem 9.6. 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 $C \doteq M / \log(1 + M)$.

▷ *Solution*

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 9.8 (Information gain of common kernels). *Due to submodularity, we have the following bounds on the information gain of common kernels:*

- linear kernel

$$\gamma_T = \mathcal{O}(d \log T), \quad (9.15)$$

- Gaussian kernel

$$\gamma_T = \mathcal{O}\left((\log T)^{d+1}\right), \quad (9.16)$$

- Matérn kernel for $\nu > \frac{1}{2}$

$$\gamma_T = \mathcal{O}\left(T^{\frac{d}{2\nu+d}} (\log T)^{\frac{2\nu}{2\nu+d}}\right). \quad (9.17)$$

Proof. Refer to theorem 5 of “Gaussian process optimization in the bandit setting: No regret and experimental design” (Srinivas et al., 2010) and remark 2 of “On information gain and regret bounds in gaussian process bandits” (Vakili et al., 2021). \square

The information gain of common kernels is illustrated in fig. 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 9.6 does not yield sublinear regret. However, we can see from theorem 9.8 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*.

Exercise 9.9: 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 $\gamma_T = \mathcal{O}(d \log T)$.
2. Deduce from (1) and theorem 9.6 that $\lim_{T \rightarrow \infty} R_T/T = 0$.

▷ *Solution*

As mentioned, the size of the confidence intervals can also be analyzed under a frequentist assumption on f^* .

Theorem 9.10 (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*

$$\beta_t(\delta) = \|f^*\|_k + \sigma_n \sqrt{2(\gamma_t + \log(1/\delta))} \quad (9.18)$$

satisfies $\mathbb{P}(\forall t \geq 1, \mathbf{x} \in \mathcal{X} : f^(\mathbf{x}) \in \mathcal{C}_t(\mathbf{x})) \geq 1 - \delta$.*

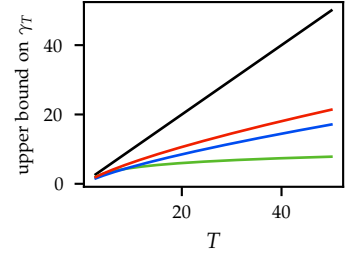


Figure 9.6: Information gain of **independent**, **linear**, **Gaussian**, and **Matérn** ($\nu \approx 0.5$) kernels with $d = 2$ (up to constant factors). The kernels with sublinear information gain have strong diminishing returns (due to their strong dependence between “close” points). In contrast, 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 the norm of the underlying reproducing kernel Hilbert space $\mathcal{H}_k(\mathcal{X})$. We remark that theorem 9.10 holds also under the looser assumption that observations are perturbed by σ_n -sub-Gaussian noise (cf. eq. (1.86)) instead of Gaussian noise.

Remark 9.11: Bayesian vs frequentist assumption

Theorems 9.4 and 9.10 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).

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

$$I_t(x) \doteq (f(x) - \hat{f}_t)_+ \quad (9.19)$$

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,

$$x_{t+1} \doteq \arg \max_{x \in \mathcal{X}} \mathbb{P}(I_t(x) > 0 \mid x_{1:t}, y_{1:t}) \quad (9.20)$$

$$= \arg \max_{x \in \mathcal{X}} \mathbb{P}(f(x) > \hat{f}_t \mid x_{1:t}, y_{1:t}) \quad (9.21)$$

$$= \arg \max_{x \in \mathcal{X}} \Phi \left(\frac{\mu_t(x) - \hat{f}_t}{\sigma_t(x)} \right) \quad (9.22)$$

using linear transformations of Gaussians (1.123)

where Φ denotes the CDF of the standard normal distribution and we use that $f(x) \mid x_{1:t}, y_{1:t} \sim \mathcal{N}(\mu_t(x), \sigma_t^2(x))$. 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),

$$\mathbf{x}_{t+1} \doteq \arg \max_{\mathbf{x} \in \mathcal{X}} \mathbb{E}[I_t(\mathbf{x}) \mid \mathbf{x}_{1:t}, y_{1:t}]. \quad (9.23)$$

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).

Exercise 9.12: Closed-form expected improvement

Let us denote the acquisition function of EI from eq. (9.23) by $\text{EI}_t(\mathbf{x})$. In this exercise, we derive a closed-form expression.

1. Show that

$$\text{EI}_t(\mathbf{x}) = \int_{(\hat{f}_t - \mu_t(\mathbf{x}))/\sigma_t(\mathbf{x})}^{+\infty} (\mu_t(\mathbf{x}) + \sigma_t(\mathbf{x})\epsilon - \hat{f}_t) \cdot \phi(\epsilon) d\epsilon \quad (9.24)$$

where ϕ is the PDF of the univariate standard normal distribution (1.6).

Hint: Reparameterize the posterior distribution

$$f(\mathbf{x}) \mid \mathbf{x}_{1:t}, y_{1:t} \sim \mathcal{N}(\mu_t(\mathbf{x}), \sigma_t^2(\mathbf{x}))$$

using a standard normal distribution.

2. Using the above expression, show that

$$\text{EI}_t(\mathbf{x}) = (\mu_t(\mathbf{x}) - \hat{f}_t) \Phi\left(\frac{\mu_t(\mathbf{x}) - \hat{f}_t}{\sigma_t(\mathbf{x})}\right) + \sigma_t(\mathbf{x}) \phi\left(\frac{\mu_t(\mathbf{x}) - \hat{f}_t}{\sigma_t(\mathbf{x})}\right) \quad (9.25)$$

where $\Phi(u) \doteq \int_{-\infty}^u \phi(\epsilon) d\epsilon$ denotes the CDF of the standard normal distribution.

Note that the first term of eq. (9.25) 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 $z_t(\mathbf{x}) \doteq (\mu_t(\mathbf{x}) - \hat{f}_t)/\sigma_t(\mathbf{x})$.

▷ *Solution*

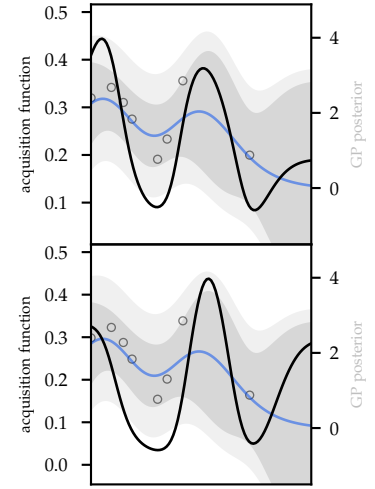


Figure 9.7: Plot of the PI and EI acquisition functions, respectively.

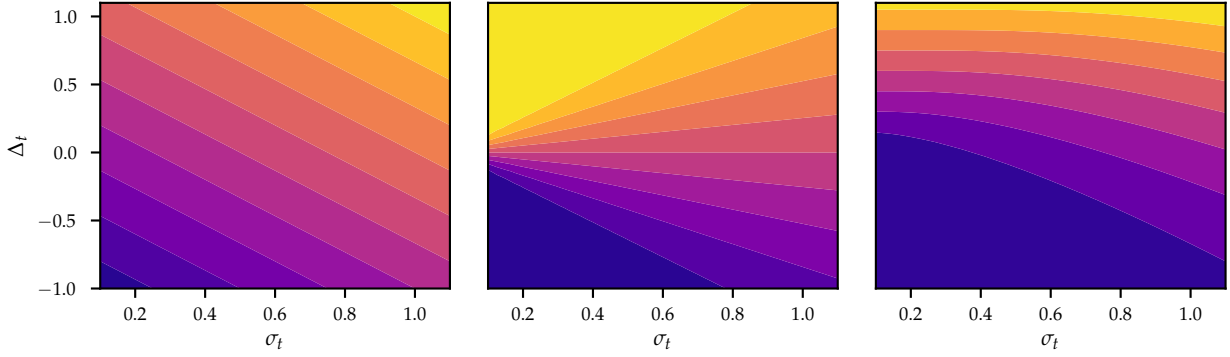


Figure 9.8: Contour lines of acquisition functions for varying $\Delta_t = \mu_t(x) - \hat{f}_t$ and σ_t . The first graph shows contour lines of UCB with $\beta_t = 0.75$, the second of PI, and the third of EI.

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),

$$\pi(x \mid x_{1:t}, y_{1:t}) \doteq \mathbb{P}_{f \mid x_{1:t}, y_{1:t}} \left(f(x) = \max_{x'} f(x') \right) \quad (9.26)$$

$$x_{t+1} \sim \pi(\cdot \mid x_{1:t}, y_{1:t}). \quad (9.27)$$

This approach 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,

$$\pi(x \mid x_{1:t}, y_{1:t}) = \mathbb{E}_{f \mid x_{1:t}, y_{1:t}} [\mathbb{1}\{f(x) = \max_{x'} f(x')\}], \quad (9.28)$$

which we can approximate using Monte Carlo sampling (typically using a single sample),

$$\approx \mathbb{1}\{\tilde{f}_{t+1}(x) = \max_{x'} \tilde{f}_{t+1}(x')\} \quad (9.29)$$

where $\tilde{f}_{t+1} \sim p(\cdot \mid x_{1:t}, y_{1:t})$ 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 $\tilde{f}_{t+1} \sim p(\cdot \mid x_{1:t}, y_{1:t})$ from our posterior distri-

bution. Then, we simply maximize \tilde{f}_{t+1} ,

$$\mathbf{x}_{t+1} \doteq \arg \max_{\mathbf{x} \in \mathcal{X}} \tilde{f}_{t+1}(\mathbf{x}). \quad (9.30)$$

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

For their analysis of Thompson sampling, Russo and Van Roy (2016) introduce the *regret-information ratio*

$$\Psi_t(\mathbf{x}) \doteq \frac{\Delta(\mathbf{x})^2}{I_t(\mathbf{x})} \quad (9.31)$$

where $\Delta(\mathbf{x}) \doteq \max_{\mathbf{x}'} f^*(\mathbf{x}') - f^*(\mathbf{x})$ is the instantaneous regret if \mathbf{x} were to be observed (so that $r_t = \Delta(\mathbf{x}_t)$) and $I_t(\mathbf{x})$ is some function capturing the “information gain” associated with observing \mathbf{x} in iteration $t + 1$. They 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 regret-information ratio.

Theorem 9.13 (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*

$$R_T \leq \sqrt{\gamma_T \bar{\Psi}_T T}. \quad (9.32)$$

Proof. By eq. (9.31), $r_t = \Delta(\mathbf{x}_t) = \sqrt{\Psi_{t-1}(\mathbf{x}_t) \cdot I_{t-1}(\mathbf{x}_t)}$. Hence,

$$\begin{aligned} R_T &= \sum_{t=1}^T r_t \\ &= \sum_{t=1}^T \sqrt{\Psi_{t-1}(\mathbf{x}_t) \cdot I_{t-1}(\mathbf{x}_t)} \\ &\leq \sqrt{\sum_{t=1}^T \Psi_{t-1}(\mathbf{x}_t) \cdot \sum_{t=1}^T I_{t-1}(\mathbf{x}_t)} && \text{using the Cauchy-Schwarz inequality} \\ &\leq \sqrt{\gamma_T \bar{\Psi}_T T}. && \square \text{ using the assumptions on } I_t(\cdot) \text{ and } \Psi_t(\cdot) \end{aligned}$$

Example 9.14: How to measure “information gain”?

One possibility of measuring the “information gain” is

$$I_t(\mathbf{x}) \doteq I(f_{\mathbf{x}}; y_{\mathbf{x}} \mid \mathbf{x}_{1:t}, y_{1:t}) \quad (9.33)$$

which — as you may recall — is precisely the marginal gain of the utility $I(S) = I(f_S; y_S)$ we were studying in chapter 8. In this case,

$$\begin{aligned} \sum_{t=1}^T I_{t-1}(\mathbf{x}_t) &= \sum_{t=1}^T I(f_{\mathbf{x}_t}; y_{\mathbf{x}_t} \mid \mathbf{x}_{1:t-1}, y_{1:t-1}) \\ &= \sum_{t=1}^T I(f_{\mathbf{x}_{1:T}}; y_{\mathbf{x}_t} \mid \mathbf{x}_{1:t-1}, y_{1:t-1}) \\ &= I(f_{\mathbf{x}_{1:T}}; y_{\mathbf{x}_{1:T}}) \\ &\leq \gamma_T. \end{aligned}$$

using $I(\mathbf{X}, \mathbf{Z}; \mathbf{Y}) \geq I(\mathbf{X}; \mathbf{Y})$ which follows from eqs. (8.14) and (8.18) and is called *monotonicity of MI*

by repeated application of eq. (8.18), also called the *chain rule of MI*

by definition of γ_T (9.13)

The regret bound from theorem 9.13 suggests an algorithm which in each iteration chooses the point which minimizes the regret-information ratio (9.31). 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^* ,

$$\hat{\Delta}_t(\mathbf{x}) \doteq \max_{\mathbf{x}' \in \mathcal{X}} u_t(\mathbf{x}') - l_t(\mathbf{x}). \quad (9.34)$$

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. eq. (9.9)) so that the confidence interval is well-calibrated and $\Delta(\mathbf{x}) \leq \hat{\Delta}_t(\mathbf{x})$ with high probability. The resulting algorithm

$$\mathbf{x}_{t+1} \doteq \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \hat{\Psi}_t(\mathbf{x}) \doteq \frac{\hat{\Delta}_t(\mathbf{x})^2}{I_t(\mathbf{x})} \right\} \quad (9.35)$$

is known as *information-directed sampling* (IDS).

Theorem 9.15 (Regret of IDS, lemma 8 of Kirschner and Krause (2018)). *Let $\beta_t(\delta)$ be chosen “correctly” for a fixed $\delta \in (0, 1)$. Then, if the measure of information gain is $I_t(\mathbf{x}) = I(f_{\mathbf{x}}; y_{\mathbf{x}} \mid \mathbf{x}_{1:t}, y_{1:t})$, with probability at least $1 - \delta$, IDS has cumulative regret*

$$R_T = \mathcal{O}\left(\beta_T(\delta) \sqrt{\gamma_T T}\right). \quad (9.36)$$

Exercise 9.16: Regret of IDS

We derive the above theorem.

1. Prove that for all $t \geq 1$, $\hat{\Delta}_t(x_t^{\text{UCB}}) \leq 2\beta_{t+1}\sigma_t(x_t^{\text{UCB}})$ where we denote by $x_t^{\text{UCB}} \doteq \arg \max_{x \in \mathcal{X}} u_t(x)$ the point maximizing the upper confidence bound after iteration t .
2. Using (1) and the assumption on I_t , bound $\hat{\Psi}_t(x_t^{\text{UCB}})$.
Hint: You may find the hint of exercise 9.7 (3) useful.
3. Complete the proof of theorem 9.15.

▷ Solution

Regret bounds such as theorem 9.15 can be derived also for different measures of information gain. For example, the argument of exercise 9.16 also goes through for the “greedy” measure

$$I_t(x) \doteq I(f_{x_t^{\text{UCB}}}; y_x \mid x_{1:t}, y_{1:t}) \quad (9.37)$$

which focuses exclusively on reducing the uncertainty at x_t^{UCB} rather than globally. We compare the two measures of information gain in fig. 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 (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*.

9.4 Model Selection

Selecting a model of f^* is much harder than in the i.i.d. data setting of supervised learning. There are mainly the two following 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

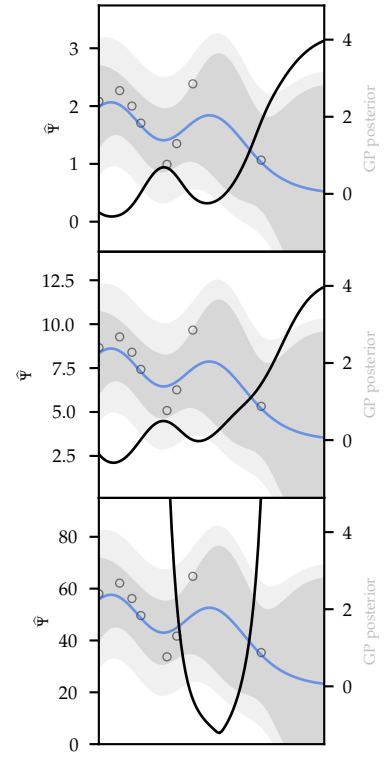


Figure 9.9: Plot of the surrogate regret-information ratio $\hat{\Psi}$: IDS selects its minimizer. The first two plots use the “global” information gain measure from example 9.14 with $\beta = 0.25$ and $\beta = 0.5$, respectively. The third plot uses the “greedy” information gain measure from eq. (9.37) and $\beta = 1$.

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.

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.

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

$$p(x' \mid x, a) \doteq \mathbb{P}(X_{t+1} = x' \mid X_t = x, A_t = a) \quad (10.1)$$

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.

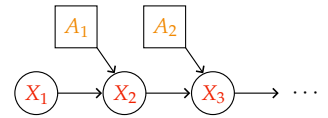


Figure 10.1: Directed graphical model of a Markov decision process with hidden states X_t and actions A_t .

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 $t > 0$,

$$\pi(a \mid x) \doteq \mathbb{P}(A_t = a \mid X_t = x). \quad (10.3)$$

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 12.4.5 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,

$$p^\pi(x' \mid x) \doteq \mathbb{P}(X_{t+1}^\pi = x' \mid X_t^\pi = x) = \sum_{a \in A} \pi(a \mid x) p(x' \mid x, a). \quad (10.4)$$

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

$$\begin{array}{lll} \text{(instantaneous)} & \text{(finite-horizon)} & \text{(mean payoff)} \\ G_t \doteq R_t & , \quad G_t \doteq \sum_{m=0}^{T-1} R_{t+m}, & \text{and } G_t \doteq \liminf_{T \rightarrow \infty} \frac{1}{T} \sum_{m=0}^T R_{t+m}. \end{array}$$

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.

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 definition of the discounted payoff (10.5); corresponds to one step in the induced Markov chain

10.1 Bellman Expectation Equation

Let us now see how we can compute the value function,

$$\begin{aligned}
 v^\pi(x) &= \mathbb{E}_\pi[G_0 \mid X_0 = x] \\
 &= \mathbb{E}_\pi \left[\sum_{m=0}^{\infty} \gamma^m R_m \mid X_0 = x \right] \\
 &= \mathbb{E}_\pi \left[\gamma^0 R_0 \mid X_0 = x \right] + \gamma \mathbb{E}_\pi \left[\sum_{m=0}^{\infty} \gamma^m R_{m+1} \mid X_0 = x \right] \\
 &= r(x, \pi(x)) + \gamma \mathbb{E}_{x'} \left[\mathbb{E}_\pi \left[\sum_{m=0}^{\infty} \gamma^m R_{m+1} \mid X_1 = x' \right] \mid X_0 = x \right] \\
 &= r(x, \pi(x)) + \gamma \sum_{x' \in X} p(x' \mid x, \pi(x)) \mathbb{E}_\pi \left[\sum_{m=0}^{\infty} \gamma^m R_{m+1} \mid X_1 = x' \right] \\
 &= r(x, \pi(x)) + \gamma \sum_{x' \in X} p(x' \mid x, \pi(x)) \mathbb{E}_\pi \left[\sum_{m=0}^{\infty} \gamma^m R_m \mid X_0 = x' \right] \\
 &= r(x, \pi(x)) + \gamma \sum_{x' \in X} p(x' \mid x, \pi(x)) \mathbb{E}_\pi[G_0 \mid X_0 = x'] \\
 &= r(x, \pi(x)) + \gamma \sum_{x' \in X} p(x' \mid x, \pi(x)) \cdot v^\pi(x'). \tag{10.10} \\
 &= r(x, \pi(x)) + \gamma \mathbb{E}_{x' \mid x, \pi(x)}[v^\pi(x')]. \tag{10.11}
 \end{aligned}$$

using the definition of the value function (10.7)

using the definition of the discounted payoff (10.5)

using linearity of expectation (1.24)

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,

$$v^\pi(x) = \sum_{a \in A} \pi(a \mid x) \left(r(x, a) + \gamma \sum_{x' \in X} p(x' \mid x, a) v^\pi(x') \right) \tag{10.12}$$

$$= \mathbb{E}_{a \sim \pi(x)} \left[r(x, a) + \gamma \mathbb{E}_{x' \mid x, a} [v^\pi(x')] \right] \tag{10.13}$$

$$= \mathbb{E}_{a \sim \pi(x)} [q^\pi(x, a)]. \tag{10.14}$$

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,

$$v^\pi(x) = q^\pi(x, \pi(x)). \quad (10.17)$$

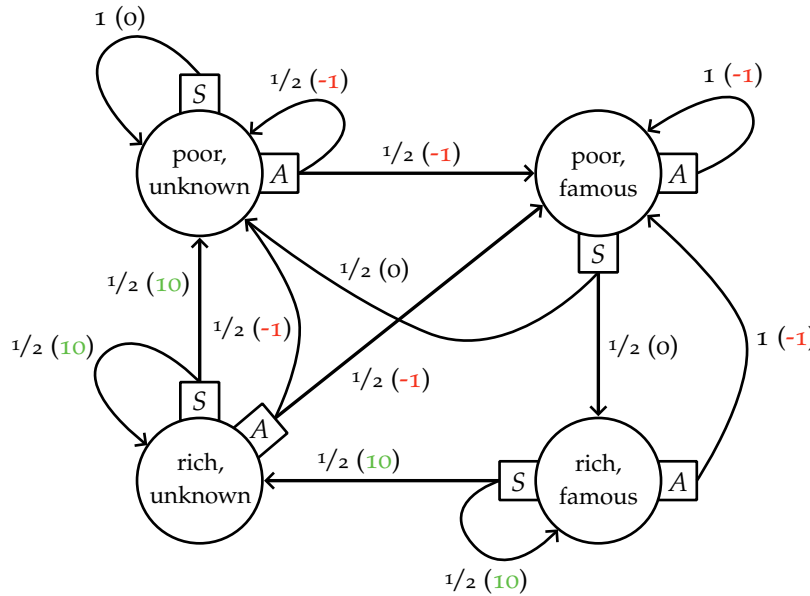


Figure 10.2: Suppose you are building a company. The shown MDP models “how to become rich and famous”. Here, the action S is short for *saving* and the action A is short for *advertising*.

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.

Exercise 10.8: Value functions

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 .

▷ *Solution*

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,

$$\begin{aligned} \mathbf{v}^\pi &\doteq \begin{bmatrix} v^\pi(1) \\ \vdots \\ v^\pi(n) \end{bmatrix}, \quad \mathbf{r}^\pi \doteq \begin{bmatrix} r(1, \pi(1)) \\ \vdots \\ r(n, \pi(n)) \end{bmatrix}, \quad \text{and} \\ \mathbf{P}^\pi &\doteq \begin{bmatrix} p(1 | 1, \pi(1)) & \cdots & p(n | 1, \pi(1)) \\ \vdots & \ddots & \vdots \\ p(1 | n, \pi(n)) & \cdots & p(n | n, \pi(n)) \end{bmatrix} \end{aligned} \quad (10.18)$$

and a little bit of linear algebra, the Bellman expectation equation (10.10) is equivalent to

$$\mathbf{v}^\pi = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{v}^\pi \quad (10.19)$$

$$\iff (\mathbf{I} - \gamma \mathbf{P}^\pi) \mathbf{v}^\pi = \mathbf{r}^\pi$$

$$\iff \mathbf{v}^\pi = (\mathbf{I} - \gamma \mathbf{P}^\pi)^{-1} \mathbf{r}^\pi. \quad (10.20)$$

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 \mathbf{v}^π , we can use that it is the unique fixed-point of the affine mapping $\mathbf{B}^\pi : \mathbb{R}^n \rightarrow \mathbb{R}^n$,

$$\mathbf{B}^\pi \mathbf{v} \doteq \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{v}. \quad (10.21)$$

Using this fact (which we will prove in just a moment), we can use fixed-point iteration of \mathbf{B}^π .

Algorithm 10.9: Fixed-point iteration

```

1 initialize  $\mathbf{v}^\pi$  (e.g., as  $\mathbf{0}$ )
2 for  $t = 1$  to  $T$  do
3    $\mathbf{v}^\pi \leftarrow \mathbf{B}^\pi \mathbf{v}^\pi = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{v}^\pi$ 
```

Fixed-point iteration has computational advantages, for example, for sparse transitions.

Theorem 10.10. \mathbf{v}^π is the unique fixed-point of \mathbf{B}^π .

Proof. It is immediate from Bellman's expectation equation (10.19) and the definition of \mathbf{B}^π (10.21) that \mathbf{v}^π is a fixed-point of \mathbf{B}^π . To prove uniqueness, we will show that \mathbf{B}^π is a contraction.

Remark 10.11: 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

$$\begin{aligned} \|B^\pi v - B^\pi v'\|_\infty &= \|r^\pi + \gamma P^\pi v - r^\pi - \gamma P^\pi v'\|_\infty \\ &= \gamma \|P^\pi(v - v')\|_\infty \\ &\leq \gamma \max_{x \in X} \sum_{x' \in X} p(x' | x, \pi(x)) \cdot |v(x') - v'(x')| \\ &\leq \gamma \|v - v'\|_\infty. \end{aligned} \quad (10.24)$$

Thus, by eq. (10.22), B^π is a contraction and by Banach's fixed-point theorem v^π is its unique fixed-point. \square

Let v_t^π be the value function estimate after t iterations. Then, we have for the convergence of fixed-point iteration,

$$\begin{aligned} \|v_t^\pi - v^\pi\|_\infty &= \|B^\pi v_{t-1}^\pi - B^\pi v^\pi\|_\infty \\ &\leq \gamma \|v_{t-1}^\pi - v^\pi\|_\infty \\ &= \gamma^t \|v_0^\pi - v^\pi\|_\infty. \end{aligned} \quad (10.25)$$

This shows that fixed-point iteration converges to v^π exponentially quickly.

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

$$\pi \geq \pi' \iff v^\pi(x) \geq v^{\pi'}(x) \quad (\forall x \in X). \quad (10.27)$$

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 $|\sum_i a_i| \leq \sum_i |a_i|$ using $\sum_{x' \in X} p(x' | x, \pi(x)) = 1$ and $|v(x') - v'(x')| \leq \|v - v'\|_\infty$

using the update rule of fixed-point iteration and $B^\pi v^\pi = v^\pi$

using (10.24)

by induction

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 10.12 (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 ,

$$\pi_v(x) \doteq \arg \max_{a \in A} r(x, a) + \gamma \sum_{x' \in X} p(x' | x, a) \cdot v(x'). \quad (10.30)$$

Exercise 10.13: Greedy policies

Show that if q and v arise from the same policy, that is, q is defined in terms of v as per eq. (10.9), then

$$\pi_v \equiv \pi_q. \quad (10.31)$$

This implies that we can use v and q interchangeably.

▷ *Solution*

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 fig. 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 10.14 (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,

$$\pi^*(x) = \arg \max_{a \in A} q^*(x, a). \quad (10.32)$$

Proof. It is a direct consequence of eq. (10.28) that a policy is optimal iff it is greedy with respect to q^* . \square

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 10.15. *The optimal value functions v^* and q^* are a fixed-point of the so-called Bellman update,*

$$v^*(x) = \max_{a \in A} q^*(x, a), \quad (10.33)$$

$$= \max_{a \in A} r(x, a) + \gamma \mathbb{E}_{x'|x, a} [v^*(x')] \quad (10.34)$$

$$q^*(x, a) = r(x, a) + \gamma \mathbb{E}_{x'|x, a} \left[\max_{a' \in A} q^*(x', a') \right]. \quad (10.35)$$

using the definition of the q -function (10.9)

Proof. It follows from eq. (10.14) that

$$v^*(x) = \mathbb{E}_{a \sim \pi^*(x)} [q^*(x, a)]. \quad (10.36)$$

Thus, as π^* is greedy with respect to q^* , $v^*(x) = \max_{a \in A} q^*(x, a)$.

Equation (10.35) follows analogously from eq. (10.15). \square

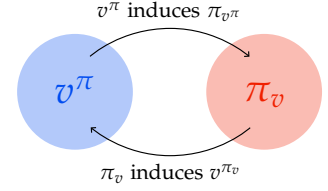


Figure 10.3: Cyclic dependency between **value function** and **greedy policy**.

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.

Exercise 10.16: Optimal policies

Recall the example of “becoming rich and famous” from fig. 10.2.

1. Show that the policy $\pi \equiv S$, which we considered in exercise 10.8, 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?

▷ *Solution*

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.

The two perspectives of Bellman's theorem naturally suggest two separate ways of finding the optimal policy. Policy iteration uses the perspective from eq. (10.32) of π^* as a fixed-point of the dependency between greedy policy and value function. In contrast, value iteration uses the perspective from eq. (10.33) 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.

10.3.3 Policy Iteration

Starting from an arbitrary initial policy, policy iteration as shown in alg. 10.17 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.

Let π_t be the policy after t iterations. We will now show that policy

Algorithm 10.17: 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

```

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 10.18 (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 (10.33) as the mapping $B^* : \mathbb{R}^n \rightarrow \mathbb{R}^n$,

$$(B^*v)(x) \doteq \max_{a \in A} q(x, a), \quad (10.37)$$

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. alg. 10.9) of $v^{\pi_{t+1}}$ initialized by v^{π_t} . We denote the iterates by \tilde{v}_t , in particular, we have that $\tilde{v}_0 = v^{\pi_t}$ and $\lim_{t \rightarrow \infty} \tilde{v}_t = v^{\pi_{t+1}}$.⁴ First, observe that for the first iteration of fixed-point iteration,

⁴ using the convergence of fixed-point iteration (10.25)

$$\begin{aligned}
 \tilde{v}_1(x) &= (B^*v^{\pi_t})(x) \\
 &= \max_{a \in A} q^{\pi_t}(x, a) \\
 &\geq q^{\pi_t}(x, \pi_t(x)) \\
 &= v^{\pi_t}(x) \\
 &= \tilde{v}_0(x).
 \end{aligned}$$

using that π_{t+1} is greedy wrt. v^{π_t}

using the definition of the Bellman update (10.37)

using (10.17)

Let us now consider a single iteration of fixed-point iteration. We have,

$$\tilde{v}_{t+1}(x) = r(x, \pi_{t+1}(x)) + \gamma \sum_{x' \in X} p(x' | x, \pi_{t+1}(x)) \cdot \tilde{v}_t(x').$$

using the definition of \tilde{v}_{t+1} (10.21)

Using an induction on t , we conclude,

$$\begin{aligned} &\geq r(x, \pi_{t+1}(x)) + \gamma \sum_{x' \in X} p(x' | x, \pi_{t+1}(x)) \cdot \tilde{v}_{t-1}(x') \\ &= \tilde{v}_t(x). \end{aligned}$$

using the induction hypothesis,
 $\tilde{v}_t(x') \geq \tilde{v}_{t-1}(x')$

This establishes the first claim,

$$v^{\pi_{t+1}} = \lim_{t \rightarrow \infty} \tilde{v}_t \geq \tilde{v}_0 = v^{\pi_t}. \quad (10.38)$$

For the second claim, recall from Bellman's theorem (10.33) 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 eq. (10.38) 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 eq. (10.39) that B^* is a contraction, implying that v^* is the unique fixed-point of B^* .

Theorem 10.19 (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 eq. (10.37),

$$(B^*v)(x) = \max_{a \in A} q(x, a),$$

where q was the state-action value function associated with the state value function v . The value iteration algorithm is shown in alg. 10.20.

We will now prove the convergence of value iteration using the fixed-point interpretation.

Theorem 10.21 (Convergence of value iteration). *Value iteration converges to an optimal policy.*

Algorithm 10.20: 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 (10.33) 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

$$\begin{aligned}
 \|B^* v - B^* v'\|_\infty &= \max_{x \in X} |(B^* v)(x) - (B^* v')(x)| \\
 &= \max_{x \in X} \left| \max_{a \in A} q(x, a) - \max_{a \in A} q'(x, a) \right| \\
 &\leq \max_{x \in X} \max_{a \in A} |q(x, a) - q'(x, a)| \\
 &\leq \gamma \max_{x \in X} \max_{a \in A} \sum_{x' \in X} p(x' | x, a) |v(x') - v'(x')| \\
 &\leq \gamma \|v - v'\|_\infty
 \end{aligned} \tag{10.39}$$

using the definition of the ℓ_∞ norm (10.23)

using the definition of the Bellman update (10.37)

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 $|v(x') - v'(x')| \leq \|v - v'\|_\infty$

where q and q' are the state-action value functions associated with v and v' , respectively. By eq. (10.22), B^* is a contraction and by Banach's fixed-point theorem v^* is its unique fixed-point. \square

Remark 10.22: 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.

Value iteration converges to an ϵ -optimal solution in a polynomial number of iterations. Unlike policy iteration, value iteration does not converge to an exact solution in general. 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 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.

⁶ Sparsity refers to the interconnectivity of the state space. When only few states are reachable from any state, we call an MDP sparse.

Exercise 10.23: 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

$$\|v^{\pi_t} - v^*\|_\infty \leq \gamma^t \|v^{\pi_0} - v^*\|_\infty \quad (10.40)$$

where v^{π_t} and v^* are vector representations of the functions v^{π_t} and v^* , respectively.

Hint: Recall from lemma 10.18 that $v^{\pi_{t+1}} \geq B^ v^{\pi_t} \geq v^{\pi_t}$.*

▷ *Solution*

Exercise 10.24: 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

$$f(x, x') \doteq \gamma \phi(x') - \phi(x), \quad \phi : X \rightarrow \mathbb{R}. \quad (10.41)$$

Show that the optimal policy remains unchanged under this definition of f .

▷ *Solution*

Readings

For a more extensive overview of finite Markov decision processes, refer to chapters 3 and 4 of “Reinforcement learning: An introduction” (Sutton and Barto, 2018) or chapter 17 of “Artificial intelligence: a modern approach” (Russell and Norvig, 2002).

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 Bayesian filtering and Kalman filters. 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 10.25 (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*

$$o(y | x) \doteq \mathbb{P}(Y_t = y | X_t = x). \quad (10.42)$$

Whereas MDPs are controlled Markov chains, POMDPs are controlled hidden Markov models.

Remark 10.26: 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 10.27 (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

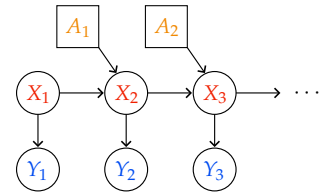


Figure 10.4: Directed graphical model of a partially observable Markov decision process with hidden states X_t , observables Y_t , and actions A_t .

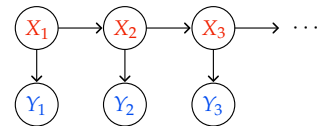


Figure 10.5: Directed graphical model of a hidden Markov model with hidden states X_t and observables Y_t .

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 fig. 10.5, its joint probability distribution factorizes into

$$\mathbb{P}(x_{1:t}, y_{1:t}) = \mathbb{P}(x_1) \cdot o(y_1 | x_1) \cdot \prod_{i=2}^t p(x_i | x_{i-1}) \cdot o(y_i | x_i). \quad (10.43)$$

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

$$\arg \max_{x_{1:t}} \mathbb{P}(x_{1:t} | y_{1:t}). \quad (10.44)$$

This task can be solved in linear time by a simple backtracking algorithm known as the *Viterbi algorithm*.

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.

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*,

$$b_t(x) \doteq \mathbb{P}(X_t = x | y_{1:t}, a_{1:t-1}), \quad (10.45)$$

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 *Bayesian 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,

$$\begin{aligned} b_{t+1}(x) &= \mathbb{P}(X_{t+1} = x | y_{1:t+1}, a_{1:t}) \\ &= \frac{1}{Z} \mathbb{P}(y_{t+1} | X_{t+1} = x) \mathbb{P}(X_{t+1} = x | y_{1:t}, a_{1:t}) \\ &= \frac{1}{Z} o(y_{t+1} | x) \mathbb{P}(X_{t+1} = x | y_{1:t}, a_{1:t}) \end{aligned}$$

by the definition of beliefs (10.45)

using Bayes' rule (1.59)

using the definition of observation probabilities (10.42)

$$\begin{aligned}
&= \frac{1}{Z} o(y_{t+1} | x) \sum_{x' \in X} p(x | x', a_t) \mathbb{P}(X_t = x' | y_{1:t}, a_{1:t-1}) && \text{by conditioning on the previous state } x', \\
&&& \text{noting } a_t \text{ does not influence } X_t \\
&= \frac{1}{Z} o(y_{t+1} | x) \sum_{x' \in X} p(x | x', a_t) b_t(x') && \text{using the definition of beliefs (10.45)}
\end{aligned} \tag{10.46}$$

where

$$Z \doteq \sum_{x \in X} o(y_{t+1} | x) \sum_{x' \in X} p(x | x', a_t) b_t(x'). \tag{10.47}$$

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 $(B_t)_{t \in \mathbb{N}_0}$,

$$B_t \doteq X_t | y_{1:t}, a_{1:t-1}, \tag{10.48}$$

where the (state-)space of all beliefs is the (infinite) space of all probability distributions over X ,⁸

⁸ This definition naturally extends to continuous state spaces \mathcal{X} .

$$\mathcal{B} \doteq \Delta^X \doteq \left\{ \mathbf{b} \in \mathbb{R}^{|X|} : \mathbf{b} \geq \mathbf{0}, \sum_{i=1}^{|X|} b(i) = 1 \right\}. \tag{10.49}$$

A Markov decision process, where every belief corresponds to a state is called a belief-state MDP.

Definition 10.28 (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*

$$\tau(b' | b, a) \doteq \mathbb{P}(B_{t+1} = b' | B_t = b, A_t = a), \tag{10.50}$$

- and *rewards*

$$\rho(b, a) \doteq \mathbb{E}_{x \sim b}[r(x, a)] = \sum_{x \in X} b(x) r(x, a). \tag{10.51}$$

It remains to derive the transition probabilities τ in terms of the original POMDP. We have,

$$\begin{aligned}
\tau(b_{t+1} | b_t, a_t) &= \mathbb{P}(b_{t+1} | b_t, a_t) \\
&= \sum_{y_{t+1} \in Y} \mathbb{P}(b_{t+1} | b_t, a_t, y_{t+1}) \mathbb{P}(y_{t+1} | b_t, a_t). && \text{by conditioning on } y_{t+1} \in Y
\end{aligned} \tag{10.52}$$

Using the Markovian structure of the belief updates, we naturally set,

$$\mathbb{P}(b_{t+1} \mid b_t, a_t, y_{t+1}) \doteq \begin{cases} 1 & \text{if } b_{t+1} \text{ matches the belief update of} \\ & \text{eq. (10.45) given } b_t, a_t, \text{ and } y_{t+1}, \\ 0 & \text{otherwise.} \end{cases} \quad (10.53)$$

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

$$\begin{aligned} \mathbb{P}(y_{t+1} \mid b_t, a_t) &= \mathbb{E}_{x \sim b_t} \left[\mathbb{E}_{x' \mid x, a_t} [\mathbb{P}(y_{t+1} \mid X_{t+1} = x')] \right] \\ &= \mathbb{E}_{x \sim b_t} \left[\mathbb{E}_{x' \mid x, a_t} [o(y_{t+1} \mid x')] \right] \\ &= \sum_{x \in X} b_t(x) \sum_{x' \in X} p(x' \mid x, a_t) \cdot o(y_{t+1} \mid x'). \end{aligned} \quad (10.54)$$

using the definition of observation probabilities (10.42)

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).

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. remark 1.10). 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.

Readings

For a more extensive overview of POMDPs and algorithms for solving POMDPs refer to section 17.4 of “Artificial intelligence: a modern approach” (Russell and Norvig, 2002).

Even though we focus on the fully observed setting throughout this course, 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.

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 and a continuous action space: 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.

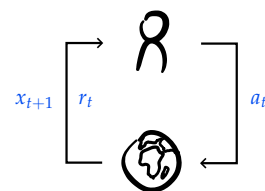


Figure 11.1: In reinforcement learning, an agent interacts with its environment in a sequence of rounds. After playing an action a_t , it observes rewards r_t and its new state x_{t+1} . The agent then uses this information to learn a model of the world.

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 (1.84) and Hoeffding's inequality (1.88) 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

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.

to some initial state. In contrast, in the *continuous setting* (or non-episodic / 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 $p(x' \mid x, a)$ 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 $x' \mid x, a$ 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 1.49, the MLE of the dynamics model coincides with the sample mean,

$$\hat{p}(x' \mid x, a) = \frac{N(x' \mid x, a)}{N(a \mid x)} \quad (11.4)$$

where $N(x' \mid x, a)$ counts the number of transitions from state x to state x' when playing action a and $N(a \mid x)$ 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),

$$\hat{r}(x, a) = \frac{1}{N(a \mid x)} \sum_{\substack{t=0 \\ x_t=x \\ a_t=a}}^{\infty} r_t. \quad (11.5)$$

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 (1.88) 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 eq. (10.32) 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  $u \leq \epsilon_t$  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,

$$\lim_{t \rightarrow \infty} \pi_t(a \mid x) = \mathbb{1}\{a = \arg \max_{a' \in A} Q_t^*(x, a')\} \quad (11.8)$$

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 $(\epsilon_t)_{t \in \mathbb{N}_0}$ satisfies the RM-conditions (1.109), e.g., if $\epsilon_t = 1/t$.

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.

³ see lemma 10.18 and theorem 10.19

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,

$$\pi_T(a \mid x) \propto \exp\left(\frac{1}{T} Q^*(x, a)\right), \quad (11.9)$$

which is the Gibbs distribution with temperature parameter $T > 0$. Observe that for $T \rightarrow 0$, softmax exploration corresponds to greedily maximizing the Q-function (i.e., greedy exploitation), whereas for $T \rightarrow \infty$, 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.

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** ∞ **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} .
-

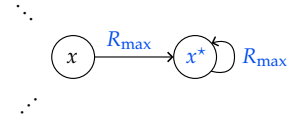


Figure 11.2: Illustration of the fairy-tale state of R_{\max} . If in doubt, the agent believes actions from the state x to lead to the fairy-tale state x^* with maximal rewards. This encourages the exploration of unknown states.

How many transitions are “enough”? We can use Hoeffding’s inequality (1.88) to get a rough idea! The key here, is our observation from

eq. (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

$$N(a \mid x) \geq \frac{R_{\max}^2}{2\epsilon^2} \log \frac{2}{\delta}. \quad (11.12) \quad \text{see (1.89)}$$

Lemma 11.7 (Exploration and exploitation of R_{\max}). *Every T time steps, with high probability, R_{\max} either*

- *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.

⁴Note that in the tabular setting, there are “only” polynomially many state-action pairs.

Theorem 11.8 (Convergence of R_{\max} , Brafman and Tennenholtz (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 $\hat{p}(x' \mid x, a)$ and $\hat{r}(x, a)$ in a table. This table already has $\mathcal{O}(n^2m)$ 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 $\mathcal{O}(nm)$ 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 therefore 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 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

$$\begin{aligned} v^\pi(x) &= r(x, \pi(x)) + \gamma \sum_{x' \in X} p(x' | x, \pi(x)) \cdot v^\pi(x') \\ &= \mathbb{E}_{R_0, X_1} [R_0 + \gamma v^\pi(X_1) | X_0 = x, A_0 = \pi(x)] \end{aligned} \quad (11.13)$$

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

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 !

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.

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 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 eq. (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 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 (1.109) 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 eq. (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 alg. 10.17). 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),

$$\begin{aligned}
 q^\pi(x, a) &= r(x, a) + \gamma \sum_{x' \in X} p(x' | x, a) \sum_{a' \in A} \pi(a' | x') q^\pi(x', a') \\
 &= \mathbb{E}_{R_0, X_1, A_1} [R_0 + \gamma q^\pi(X_1, A_1) | X_0 = x, A_0 = a] \quad (11.17)
 \end{aligned}$$

using Bellman's expectation equation (10.15)

interpreting the above expression as an expectation over R_0, X_1 and A_1

$$\approx r + \gamma q^\pi(x', a') \quad (11.18)$$

Monte Carlo approximation with a single sample

where the agent observed transitions (x, a, r, x') and (x', a', r', x'') .

The update rule from TD-learning is therefore adapted to⁶

$$Q^\pi(x, a) \leftarrow (1 - \alpha_t)Q^\pi(x, a) + \alpha_t(r + \gamma Q^\pi(x', a')). \quad (11.19)$$

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 .

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 (1.109) 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),

$$\begin{aligned} q^\pi(x, a) &= r(x, a) + \gamma \sum_{x' \in X} p(x' | x, a) \sum_{a' \in A} \pi(a' | x') q^\pi(x', a') \\ &= \mathbb{E}_{R_0, X_1} \left[R_0 + \gamma \sum_{a' \in A} \pi(a' | X_1) q^\pi(X_1, a') \mid X_0 = x, A_0 = a \right] \end{aligned} \quad (11.20)$$

using Bellman's expectation equation (10.15)

interpreting the above expression as an expectation over R_0 and X_1

$$\approx r + \gamma \sum_{a' \in A} \pi(a' | x') q^\pi(x', a') \quad (11.21)$$

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 (10.32) 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 (10.33) in place of Bellman's expectation equation (10.15), we obtain,

$$\begin{aligned} q^*(x, a) &= r(x, a) + \gamma \sum_{x' \in X} p(x' | x, a) \max_{a' \in A} q^*(x', a') \\ &= \mathbb{E}_{R_0, X_1} \left[R_0 + \gamma \max_{a' \in A} q^*(X_1, a') \mid X_0 = x, A_0 = a \right] \end{aligned} \quad (11.23)$$

$$\approx r + \gamma \max_{a' \in A} q^*(x', a') \quad (11.24)$$

using that the Q-function is a fixed-point of the Bellman update, see Bellman's theorem (10.33) interpreting the above expression as an expectation over R_0 and X_1 Monte Carlo approximation with a single sample

where the agent observed the transition (x, a, r, x') . Using a bootstrapping estimate Q^* for q^* , we obtain a structurally similar algorithm to 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 eq. (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 (1.109) 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).

Exercise 11.14: 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	x'	r
A	\downarrow	G_1	0	B	\leftarrow	A	0
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.

▷ *Solution*

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 alg. 11.15.

Algorithm 11.15: 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
    $(x, a, r, x')$ 
4    $Q^*(x, a) \leftarrow (1 - \alpha_t)Q^*(x, a) + \alpha_t(r + \gamma \max_{a' \in A} Q^*(x', a'))$ 
   // (11.27)
```

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 11.16. *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 $N(a | x) \leq T_{\text{init}}$,⁷*

⁷ $N(a | x)$ is the number of times action a is performed in state x .

$$Q_t^*(x, a) \geq V_{\max} \geq q^*(x, a). \quad (11.28)$$

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 eq. (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 11.17 (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 $\mathcal{O}(nm)$ 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 $\mathcal{O}(Tm)$. 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.

11.4.6 *Challenges of Model-free Approaches*

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.

Readings

For a more thorough treatment of tabular RL methods refer to

- chapter 22 of “Artificial intelligence: a modern approach” (Russell and Norvig, 2002) or
- chapter 6 of “Reinforcement learning: An introduction” (Sutton and Barto, 2018).

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 $\mathcal{O}(|\mathcal{X}|)$ space. If we learn the Q function, we even need $\mathcal{O}(|\mathcal{X}| \cdot |\mathcal{A}|)$ 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,

$$\bar{\ell}(\theta; x, r) \doteq \frac{1}{2}(v^\pi(x) - \theta(x))^2 \quad (12.2)$$

$$= \frac{1}{2}\left(r + \gamma \mathbb{E}_{x'|x, \pi(x)}[v^\pi(x')] - \theta(x)\right)^2 \quad (12.3)$$

using Bellman's expectation equation (10.11)

Note that this loss corresponds to a standard squared loss of the difference between the parameter $\theta(x)$ and the label $v^\pi(x)$ we want to learn.

We can now find the gradient of this loss. Elementary calculations yield,

$$\nabla_{\theta(x)} \bar{\ell}(\theta; x, r) = \theta(x) - \left(r + \gamma \mathbb{E}_{x'|x, \pi(x)}[v^\pi(x')]\right). \quad (12.4)$$

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 θ . To emphasize this, we write $V^\pi(x; \theta^{\text{old}}) \approx v^\pi(x)$ where $\theta^{\text{old}} = \theta$ but θ^{old} is treated as a constant with respect to θ .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.

¹ That is, the bootstrapping estimate $V^\pi(x; \theta^{\text{old}})$ is assumed to be constant with respect to $\theta(x)$ in the same way that $v^\pi(x)$ is constant with respect to $\theta(x)$.

If we were not using the bootstrapping estimate, the following derivation of the gradient of the loss would not be this simple.

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.

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.

Using the aforementioned shortcuts, let us define the loss ℓ after ob-

serving the single transition (x, a, r, x') ,

$$\ell(\theta; x, r, x') \doteq \frac{1}{2} \left(r + \gamma \theta^{\text{old}}(x') - \theta(x) \right)^2. \quad (12.5)$$

We define the gradient of this loss with respect to $\theta(x)$ as

$$\begin{aligned} \delta_{\text{TD}} &\doteq \nabla_{\theta(x)} \ell(\theta; x, r, x') \\ &= \theta(x) - \left(r + \gamma \theta^{\text{old}}(x') \right). \end{aligned} \quad (12.6)$$

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 (1.84) that Monte Carlo averages are unbiased.³ We therefore have,

$$\mathbb{E}_{x'|x, \pi(x)}[\delta_{\text{TD}}] = \nabla_{\theta(x)} \bar{\ell}(\theta; x, r). \quad (12.7)$$

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,

$$V^\pi(x; \theta) = \theta(x) \leftarrow \theta(x) - \alpha_t \delta_{\text{TD}} \quad (12.8)$$

$$\begin{aligned} &= (1 - \alpha_t) \theta(x) + \alpha_t \left(r + \gamma \theta^{\text{old}}(x') \right) \\ &= (1 - \alpha_t) V^\pi(x; \theta) + \alpha_t \left(r + \gamma V^\pi(x'; \theta^{\text{old}}) \right). \end{aligned} \quad (12.9)$$

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

³ Crucially, 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 alg. 1.60

using the definition of the temporal difference error (12.6)

substituting $V^\pi(x; \theta)$ for $\theta(x)$

⁴ An alternative interpretation is that TD-learning performs gradient descent with respect to the loss ℓ .

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.

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 (10.33) 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,

$$\theta \leftarrow \theta - \alpha_t \nabla_\theta \ell(\theta; x, a, r, x') \quad (12.13)$$

$$= \theta - \alpha_t \nabla_\theta \frac{1}{2} \left(r + \gamma \max_{a' \in \mathcal{A}} Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2$$

$$= \theta + \alpha_t \delta_B \nabla_\theta Q^*(x, a; \theta). \quad (12.14)$$

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.

⁵ That is, the value function takes similar values in “similar” states.

⁶ Note that often non-Bayesian neural networks (i.e., point estimates of the weights) are used. In the final chapter, chapter 13, we will explore the benefits of using Bayesian neural networks.

⁷ compare to eq. (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

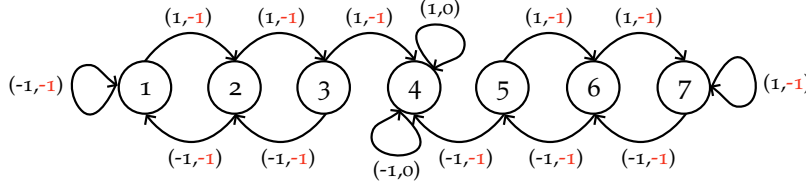


Figure 12.1: MDP studied in exercise 12.2. Each arrow marks a (deterministic) transition and is labeled with (action, reward).

Exercise 12.2: Q-learning and function approximation

Consider the MDP of fig. 12.1 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 $\alpha = 1/4$ to compute $\nabla_{\mathbf{w}} \ell(\mathbf{w}; \tau)$ and the updated weights $\mathbf{w}' = \mathbf{w} - \alpha \nabla_{\mathbf{w}} \ell(\mathbf{w}; \tau)$.

▷ *Solution*

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.

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} ,

$$\ell_{\text{DQN}}(\theta; \mathcal{D}) \doteq \frac{1}{2} \sum_{(x, a, r, x') \in \mathcal{D}} \left(r + \gamma \max_{a' \in \mathcal{A}} Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2. \quad \text{compare to the Q-learning loss (12.11)} \quad (12.16)$$

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 fig. 12.2. 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

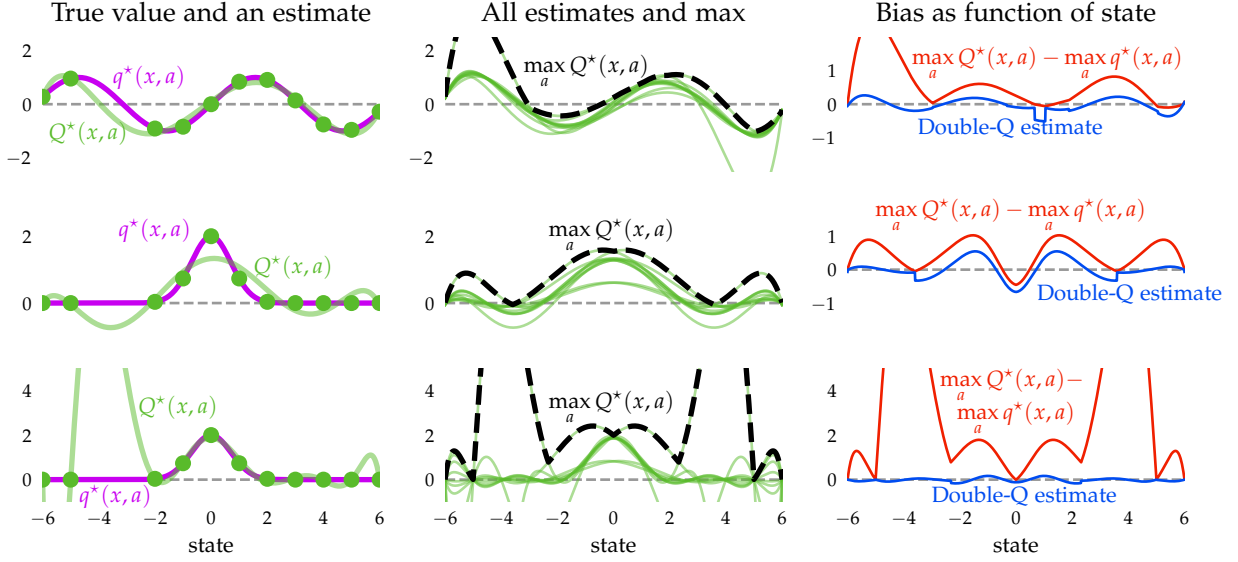


Figure 12.2: Illustration of overestimation during learning. In each state (x -axis), there are 10 actions. The left column shows the true values $v^*(x)$ (purple line). All true action values are defined by $q^*(x, a) \doteq v^*(x)$. The green line shows estimated values $Q^*(x, a)$ for one action as 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). Taken from “Deep reinforcement learning with double q-learning” (Van Hasselt et al., 2016).

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.

Readings

For more details on value function approximation, refer to chapter 9 of “Reinforcement learning: An introduction” (Sutton and Barto, 2018).

Optional Readings

- Mnih, Kavukcuoglu, Silver, Rusu, Veness, Bellemare, Graves, Riedmiller, Fidjeland, Ostrovski, et al. (2015).
Human-level control through deep reinforcement learning.
- Van Hasselt, Guez, and Silver (2016).
Deep reinforcement learning with double q-learning.

12.3 Policy Approximation

Q-learning defines a policy implicitly by

$$\pi^*(x) \doteq \arg \max_{a \in \mathcal{A}} Q^*(x, a). \quad (12.20)$$

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 that policy gradient methods fundamentally rely on randomized policies for exploration.

Remark 12.3: 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 $\pi_\varphi(a \mid x)$ the probability 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,

$$G_t = \sum_{m=0}^{\infty} \gamma^m R_{t+m}.$$

see eq. (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 define the policy value function.

Definition 12.4 (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 . We also define the bounded variant,

$$j_T(\pi) \doteq \mathbb{E}_\pi[G_{0:T}] = \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} \gamma^t R_t \right]. \quad (12.24)$$

For simplicity, we will abbreviate $j(\boldsymbol{\varphi}) \doteq j(\pi_{\boldsymbol{\varphi}})$.

Naturally, we want to maximize $j(\boldsymbol{\varphi})$. That is, we want to solve

$$\boldsymbol{\varphi}^* \doteq \arg \max_{\boldsymbol{\varphi}} j(\boldsymbol{\varphi}) \quad (12.25)$$

which is a non-convex optimization problem. Let us see how $j(\boldsymbol{\varphi})$ can be evaluated to understand the optimization problem better. We will again use a Monte Carlo estimate. Recall that a fixed $\boldsymbol{\varphi}$ 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,

$$\tau^{(1)}, \dots, \tau^{(m)} \sim \Pi_{\boldsymbol{\varphi}}, \quad (12.27)$$

where $\Pi_{\boldsymbol{\varphi}}$ is the distribution of all trajectories (i.e., rollouts) of the Markov chain induced by policy $\pi_{\boldsymbol{\varphi}}$. We denote the (bounded) discounted payoff of the i -th rollout by

$$g_{0:T}^{(i)} \doteq \sum_{t=0}^{T-1} \gamma^t r_t^{(i)} \quad (12.28)$$

where $r_t^{(i)}$ is the reward at time t of the i -th rollout. Using a Monte Carlo approximation, we can then estimate $j_T(\boldsymbol{\varphi})$. Moreover, due to the exponential discounting of future rewards, it is reasonable to approximate the policy value function using bounded trajectories,

$$j(\boldsymbol{\varphi}) \approx j_T(\boldsymbol{\varphi}) \approx \frac{1}{m} \sum_{i=1}^m g_{0:T}^{(i)}. \quad (12.29)$$

12.3.2 Policy Gradient

Let us now 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 estimates of its gradient,

$$\nabla_\varphi j(\varphi) \approx \nabla_\varphi j_T(\varphi) = \nabla_\varphi \mathbb{E}_{\tau \sim \Pi_\varphi}[G_0]. \quad (12.31)$$

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 done previously in remark 1.22. This should remind you of the reparameterization trick (see eq. (5.77)) that we used to solve a similar gradient in the context of variational inference. In this context, however, we cannot apply the reparameterization trick.⁹ 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 ,

$$\nabla_\varphi \log \Pi_\varphi(\tau). \quad (12.33)$$

Using the chain rule, the score function can be expressed as

$$\nabla_\varphi \log \Pi_\varphi(\tau) = \frac{\nabla_\varphi \Pi_\varphi(\tau)}{\Pi_\varphi(\tau)} \quad (12.34)$$

and by rearranging the terms, we obtain

$$\nabla_\varphi \Pi_\varphi(\tau) = \Pi_\varphi(\tau) \nabla_\varphi \log \Pi_\varphi(\tau). \quad (12.35)$$

This is called the *score function trick*.¹⁰

Now, assuming that state and action spaces are continuous, we obtain

$$\nabla_\varphi \mathbb{E}_{\tau \sim \Pi_\varphi}[G_0] = \nabla_\varphi \int \Pi_\varphi(\tau) \cdot G_0 d\tau$$

⁹ This 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 12.4.5 and 13.1.2.

¹⁰ We have already applied this “trick” in exercise 6.31.

using the definition of expectation (1.23b)

$$\begin{aligned}
&= \int \nabla_{\boldsymbol{\varphi}} \Pi_{\boldsymbol{\varphi}}(\tau) \cdot G_0 d\tau \\
&= \int G_0 \cdot \Pi_{\boldsymbol{\varphi}}(\tau) \nabla_{\boldsymbol{\varphi}} \log \Pi_{\boldsymbol{\varphi}}(\tau) d\tau \\
&= \mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}} [G_0 \nabla_{\boldsymbol{\varphi}} \log \Pi_{\boldsymbol{\varphi}}(\tau)]. \quad \square
\end{aligned}$$

using the regularity assumptions to swap gradient and integral
using the score function trick (12.35)
interpreting the integral as an expectation over $\Pi_{\boldsymbol{\varphi}}$

Intuitively, maximizing $j(\boldsymbol{\varphi})$ 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 $\nabla_{\boldsymbol{\varphi}} \log \Pi_{\boldsymbol{\varphi}}(\tau)$.

$$\begin{aligned}
&\nabla_{\boldsymbol{\varphi}} \log \Pi_{\boldsymbol{\varphi}}(\tau) \\
&= \nabla_{\boldsymbol{\varphi}} \left(\log p(x_0) + \sum_{t=0}^{T-1} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t | \mathbf{x}_t) + \sum_{t=0}^{T-1} \log p(x_{t+1} | \mathbf{x}_t, \mathbf{a}_t) \right) \\
&= \nabla_{\boldsymbol{\varphi}} \log p(x_0) + \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t | \mathbf{x}_t) + \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\varphi}} \log p(x_{t+1} | \mathbf{x}_t, \mathbf{a}_t) \\
&= \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t | \mathbf{x}_t). \tag{12.36}
\end{aligned}$$

using the definition of the distribution over trajectories $\Pi_{\boldsymbol{\varphi}}$
using that the first and third term are independent of $\boldsymbol{\varphi}$

When using a neural network for the parameterization of the policy π , we can use automatic differentiation to compute the gradients.

Exercise 12.6: Eligibility vector

The vector $\nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t | \mathbf{x}_t)$ 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

$$\pi_{\boldsymbol{\varphi}}(a | \mathbf{x}) \doteq \frac{\exp(h(\mathbf{x}, a, \boldsymbol{\varphi}))}{\sum_{b \in A} \exp(h(\mathbf{x}, b, \boldsymbol{\varphi}))} \tag{12.37}$$

with linear preferences $h(\mathbf{x}, a, \boldsymbol{\varphi}) \doteq \boldsymbol{\varphi}^\top \boldsymbol{\phi}(\mathbf{x}, a)$ where $\boldsymbol{\phi}(\mathbf{x}, a)$ is some feature vector, what is the form of the eligibility vector?

▷ *Solution*

The expectation of the score gradient estimator (12.32) can be approximated using Monte Carlo sampling,

$$\nabla_{\boldsymbol{\varphi}} j_T(\boldsymbol{\varphi}) \approx \frac{1}{m} \sum_{i=1}^m g_{0:T}^{(i)} \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t^{(i)} | \mathbf{x}_t^{(i)}). \tag{12.38}$$

However, typically the variance of these estimates is very large. Using so-called *baselines* can reduce the variance dramatically.

Lemma 12.7 (Score gradients with baselines). *We have,*

$$\mathbb{E}_{\tau \sim \Pi_\varphi} [G_0 \nabla_\varphi \log \Pi_\varphi(\tau)] = \mathbb{E}_{\tau \sim \Pi_\varphi} [(G_0 - b) \nabla_\varphi \log \Pi_\varphi(\tau)]. \quad (12.39)$$

Here, $b \in \mathbb{R}$ is called a baseline.

Proof. For the term to the right, we have due to linearity of expectation (1.24),

$$\begin{aligned} \mathbb{E}_{\tau \sim \Pi_\varphi} [(G_0 - b) \nabla_\varphi \log \Pi_\varphi(\tau)] &= \mathbb{E}_{\tau \sim \Pi_\varphi} [G_0 \nabla_\varphi \log \Pi_\varphi(\tau)] \\ &\quad - \mathbb{E}_{\tau \sim \Pi_\varphi} [b \cdot \nabla_\varphi \log \Pi_\varphi(\tau)]. \end{aligned}$$

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 (1.23b)

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

Exercise 12.8: 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.

1. 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 (1.24) 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

$$\text{Var}[f(\mathbf{X}) - g(\mathbf{X})] \leq \text{Var}[f(\mathbf{X})]. \quad (12.40)$$

2. Consider estimating $\nabla_\varphi j(\varphi)$. 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

$$\text{Var}_{\tau \sim \Pi_\varphi} [(G_0 - b) \nabla_\varphi \log \Pi_\varphi(\tau)] \leq \text{Var}_{\tau \sim \Pi_\varphi} [G_0 \nabla_\varphi \log \Pi_\varphi(\tau)]. \quad (12.41)$$

▷ *Solution*

Exercise 12.9: Score gradients with state-dependent baselines

One can even show, that we can subtract arbitrary baselines depending on *previous* states.

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

$$\begin{aligned} & \mathbb{E}_{\tau \sim \Pi_\varphi} [G_0 \nabla_\varphi \log \Pi_\varphi(\tau)] \\ &= \mathbb{E}_{\tau \sim \Pi_\varphi} \left[\sum_{t=0}^{T-1} (G_0 - b(\tau_{0:t-1})) \nabla_\varphi \log \pi_\varphi(\mathbf{a}_t \mid \mathbf{x}_t) \right] \end{aligned} \quad (12.42)$$

where we write $b(\tau_{0:-1}) = 0$.

▷ *Solution*

Example 12.10: Downstream returns

A commonly used baseline is

$$b(\tau_{0:t-1}) \doteq \sum_{m=0}^{t-1} \gamma^m r_m. \quad (12.43)$$

This baseline subtracts the returns of all actions before time t . Intuitively, using this baseline, the score gradient only considers downstream returns. Recall from eq. (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

$$G_0 - b(\tau_{0:t-1}) = \gamma^t G_{t:T}, \quad (12.44)$$

yielding the gradient estimator,

$$\begin{aligned} \nabla_\varphi j(\varphi) &\approx \nabla_\varphi j_T(\varphi) = \mathbb{E}_{\tau \sim \Pi_\varphi} [G_0 \nabla_\varphi \log \Pi_\varphi(\tau)] \\ &= \mathbb{E}_{\tau \sim \Pi_\varphi} \left[\sum_{t=0}^{T-1} \gamma^t G_{t:T} \nabla_\varphi \log \pi_\varphi(\mathbf{a}_t \mid \mathbf{x}_t) \right]. \end{aligned} \quad (12.45)$$

using the score gradient estimator
(12.32)

using a state-dependent baseline (12.42)

Performing stochastic gradient descent with the score gradient estimator and downstream returns is known as the *REINFORCE algorithm* (Williams, 1992) which is shown in alg. 12.11.

Algorithm 12.11: 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      $\varphi \leftarrow \varphi + \eta \gamma^t g_{t:T} \nabla_{\varphi} \log \pi_{\varphi}(a_t | x_t)$  // (12.46)
7 until converged

```

The variance of REINFORCE can be reduced further. A common technique is to subtract a term b_t from the downstream returns,

$$\nabla_{\varphi} j(\varphi) = \mathbb{E}_{\tau \sim \Pi_{\varphi}} \left[\sum_{t=0}^T \gamma^t (G_{t:T} - b_t) \nabla_{\varphi} \log \pi_{\varphi}(a_t | x_t) \right]. \quad (12.47)$$

For example, we can subtract the mean reward to go,

$$b_t \doteq \frac{1}{T} \sum_{t=0}^{T-1} G_{t:T}. \quad (12.48)$$

Exercise 12.12: Policy gradients with downstream returns

Suppose we are training an agent to solve a computer game. There are only two possible actions, specifically:

1. *do nothing*; and
2. *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 $\tau \sim \Pi_{\theta}$.

▷ *Solution*

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.¹¹

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.

¹¹ This is because the score gradient estimator is used to obtain gradients of the policy value function with respect to the *current* policy.

Readings

For more details on policy gradient methods, refer to chapter 13 (up to section 13.4) of “Reinforcement learning: An introduction” (Sutton and Barto, 2018).

12.4 Actor-Critic Methods

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 12.13 (Advantage function). Given a policy π , the *advantage function*,

$$a^\pi(x, a) \doteq q^\pi(x, a) - v^\pi(x) \quad (12.49)$$

$$= q^\pi(x, a) - \mathbb{E}_{a' \sim \pi(x)} [q^\pi(x, a')], \quad (12.50) \quad \text{using eq. (10.14)}$$

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 eq. (12.50) 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,

$$\max_{a \in \mathcal{A}} a^\pi(x, a) \geq 0. \quad (12.51)$$

Moreover, it follows directly from Bellman’s theorem (10.32) that

$$\pi \text{ is optimal} \iff \forall x \in \mathcal{X}, a \in \mathcal{A} : a^\pi(x, a) \leq 0. \quad (12.52)$$

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

$$\pi_q(\mathbf{x}) \doteq \arg \max_{a \in \mathcal{A}} q(\mathbf{x}, a) \quad (12.53)$$

since

$$\arg \max_{a \in \mathcal{A}} q(\mathbf{x}, a) = \arg \max_{a \in \mathcal{A}} q(\mathbf{x}, a) - v(\mathbf{x}) = \arg \max_{a \in \mathcal{A}} q(\mathbf{x}, a),$$

as $v(\mathbf{x})$ is independent of a . This coincides with our initial definition of greedy policies in eq. (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 (12.45) that we had introduced in the previous section,

$$\nabla_{\boldsymbol{\varphi}} j_T(\boldsymbol{\varphi}) = \mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}} \left[\sum_{t=0}^{T-1} \gamma^t G_{t:T} \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t \mid \mathbf{x}_t) \right].$$

Previously, we have approximated the policy value function $j(\boldsymbol{\varphi})$ by the bounded policy value function $j_T(\boldsymbol{\varphi})$. 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 $j(\boldsymbol{\varphi})$.

$$\begin{aligned} \nabla_{\boldsymbol{\varphi}} j(\boldsymbol{\varphi}) &= \lim_{T \rightarrow \infty} \nabla_{\boldsymbol{\varphi}} j_T(\boldsymbol{\varphi}) \\ &= \sum_{t=0}^{\infty} \mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}} [\gamma^t G_t \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t \mid \mathbf{x}_t)]. \end{aligned} \quad (12.54)$$

substituting the score gradient estimator with downstream returns (12.45) and using linearity of expectation (1.24)

Observe that because the expectations only consider downstream returns, we can disregard all data from the trajectory prior to time t . Let us define

$$\tau_{t:\infty} \doteq ((\mathbf{x}_t, \mathbf{a}_t, r_t, \mathbf{x}_{t+1}), (\mathbf{x}_{t+1}, \mathbf{a}_{t+1}, r_{t+1}, \mathbf{x}_{t+2}), \dots), \quad (12.55)$$

as the trajectory from time step t . Then,

$$= \sum_{t=0}^{\infty} \mathbb{E}_{\tau_{t:\infty} \sim \Pi_{\boldsymbol{\varphi}}} [\gamma^t G_t \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t \mid \mathbf{x}_t)].$$

We now condition on x_t and a_t ,

$$= \sum_{t=0}^{\infty} \mathbb{E}_{x_t, a_t} [\gamma^t \mathbb{E}_{r_t, \tau_{t+1}:\infty} [G_t \mid x_t, a_t] \nabla_{\varphi} \log \pi_{\varphi}(a_t \mid x_t)].$$

using that $\pi_{\varphi}(a_t \mid x_t)$ is a constant given x_t and 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 eq. (10.6),

$$\begin{aligned} &= \sum_{t=0}^{\infty} \mathbb{E}_{x_t, a_t} [\gamma^t \mathbb{E}_{\pi_{\varphi}} [G_t \mid x_t, a_t] \nabla_{\varphi} \log \pi_{\varphi}(a_t \mid x_t)] \\ &= \sum_{t=0}^{\infty} \mathbb{E}_{x_t, a_t} [\gamma^t q^{\pi_{\varphi}}(x_t, a_t) \nabla_{\varphi} \log \pi_{\varphi}(a_t \mid x_t)]. \end{aligned} \quad (12.56)$$

using the definition of the Q-function (10.8)

It turns out that $\mathbb{E}_{x_t, a_t} [q^{\pi_{\varphi}}(x_t, a_t)]$ exhibits much less variance than our previous estimator $\mathbb{E}_{x_t, a_t} \mathbb{E}_{\pi_{\varphi}} [G_t \mid x_t, a_t]$. Equation (12.56) 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*,¹²

$$\rho_{\varphi}^{\infty}(x) \doteq (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}_{\pi_{\varphi}}(\mathbf{X}_t = x). \quad (12.57)$$

The factor $(1 - \gamma)$ ensures that ρ_{φ}^{∞} is a probability density as

$$\int \rho_{\varphi}^{\infty}(x) dx = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \int \mathbb{P}_{\pi_{\varphi}}(\mathbf{X}_t = x) dx = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t = 1.$$

Intuitively, $\rho_{\varphi}^{\infty}(x)$ measures how often we visit state x when following policy π_{φ} . It can be thought of as a “discounted frequency”.

Theorem 12.14 (Policy gradient theorem in terms of ρ_{φ}^{∞}). *Policy gradients can be represented in terms of the Q-function,*

$$\nabla_{\varphi} j(\varphi) \propto \mathbb{E}_{x \sim \rho_{\varphi}^{\infty}} \mathbb{E}_{a \sim \pi_{\varphi}(\cdot \mid x)} [q^{\pi_{\varphi}}(x, a) \nabla_{\varphi} \log \pi_{\varphi}(a \mid x)]. \quad (12.58)$$

Proof. The right hand side of eq. (12.56) can be expressed as

$$\begin{aligned} &\sum_{t=0}^{\infty} \int \mathbb{P}_{\pi_{\varphi}}(\mathbf{X}_t = x) \cdot \mathbb{E}_{a_t \sim \pi_{\varphi}(\cdot \mid x)} [\gamma^t q^{\pi_{\varphi}}(x, a_t) \nabla_{\varphi} \log \pi_{\varphi}(a_t \mid x)] dx \\ &= \frac{1}{1 - \gamma} \int \rho_{\varphi}^{\infty}(x) \cdot \mathbb{E}_{a \sim \pi_{\varphi}(\cdot \mid x)} [q^{\pi_{\varphi}}(x, a) \nabla_{\varphi} \log \pi_{\varphi}(a \mid x)] dx \end{aligned}$$

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.

¹² 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.

Exercise 12.15: Policy gradient with an exponential family

1. Suppose, we can choose between two actions $a \in \{0, 1\}$ in each state. A natural stochastic policy is induced by a Bernoulli distribution,

$$a \sim \text{Bern}(\sigma(f_{\varphi}(x))), \quad (12.59)$$

where σ is the logistic function from eq. (5.10). First, write down the expression for $\pi_{\varphi}(a | x)$. Then, derive the expression for $\nabla_{\varphi} j(\varphi)$ in terms of $q^{\pi_{\varphi}}$, $\sigma(f_{\varphi}(x))$, and $\nabla_{\varphi} f_{\varphi}(x)$ using the policy gradient theorem.

2. 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

$$\pi_{\varphi}(a | x) = h(a) \exp(af_{\varphi}(x) - A(f_{\varphi}(x))) \quad (12.60)$$

where h , f , and A are known functions.¹³ Derive the expression of the policy gradient $\nabla_{\varphi} j(\varphi)$ for such a distribution.

3. Can you relate the results of the previous two exercises (1) and (2)? What are h and A in case of the Bernoulli distribution?
4. The Gaussian distribution with unit variance $\mathcal{N}(f_{\varphi}(x), 1)$ is of the same canonical form with

$$A(f_{\varphi}(x)) = \frac{f_{\varphi}(x)^2}{2}. \quad (12.61)$$

Determine the policy gradient $\nabla_{\varphi} j(\varphi)$.

5. For a Gaussian policy, can we instead apply the reparameterization trick (5.80) that we have seen in the context of variational inference? If yes, how? If not, why?

▷ *Solution*

Observe that we cannot use the policy gradient to calculate the gradient exactly, as we do not know $q^{\pi_{\varphi}}$. Instead, we will use bootstrapping estimates $Q^{\pi_{\varphi}}$ of $q^{\pi_{\varphi}}$.

12.4.3 On-policy Actor-Critics

Actor-Critic methods consist of two components:

- a parameterized policy, $\pi(a | x; \varphi) \doteq \pi_{\varphi}$, which is called *actor*; and (12.62)
- 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 $Q^{\pi_{\varphi}}$ by Q . (12.63)

¹³ Observe that this form is equivalent to the form introduced in eq. (5.57) where $f_{\varphi}(x)$ is the natural parameter, and we let $A(f) = \log Z(f)$.

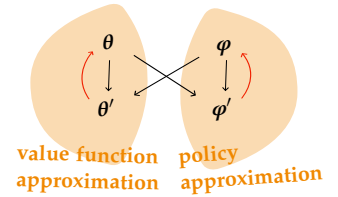


Figure 12.3: Illustration of one iteration of actor-critic methods. The dependencies between the actors and critics are shown as arrows. Methods differ in the exact order in which actor and critic are updated.

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

$$\nabla_{\boldsymbol{\varphi}} J(\boldsymbol{\varphi}) \approx \sum_{t=0}^{\infty} \mathbb{E}_{(x_t, a_t) \sim \pi_{\boldsymbol{\varphi}}} [\gamma^t Q(x_t, a_t; \boldsymbol{\theta}) \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(a_t | x_t)] \quad (12.64) \quad \text{see eq. (12.56)}$$

where Q is a bootstrapping estimate of $q^{\pi_{\boldsymbol{\varphi}}}$. This algorithm is known as *online actor-critic* or *Q actor-critic* and shown in alg. 12.16.

Algorithm 12.16: Online actor-critic

```

1 initialize parameters  $\boldsymbol{\varphi}$  and  $\boldsymbol{\theta}$ 
2 repeat
3   use  $\pi_{\boldsymbol{\varphi}}$  to obtain transition  $(x, a, r, x')$ 
4    $\delta = r + \gamma Q(x', \pi_{\boldsymbol{\varphi}}(x'); \boldsymbol{\theta}) - Q(x, a; \boldsymbol{\theta})$ 
      // actor update
5    $\boldsymbol{\varphi} \leftarrow \boldsymbol{\varphi} + \eta \gamma^t Q(x, a; \boldsymbol{\theta}) \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(a | x)$  // (12.65)
      // critic update
6    $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \delta \nabla_{\boldsymbol{\theta}} Q(x, a; \boldsymbol{\theta})$  // (12.66)
7 until converged

```

Comparing to the derivation for TD-learning from eq. (12.8), we observe that eq. (12.66) corresponds to the SARSA update rule.¹⁴ Due to the use of SARSA for learning the critic, this algorithm is fundamentally on-policy.

¹⁴ The gradient with respect to θ_Q appears analogously to our derivation of approximate Q-learning (12.15).

Crucially, by neglecting the dependence of the bootstrapping estimate Q on the policy parameters $\boldsymbol{\varphi}$, we introduce bias in the gradient estimates. In other words, using the bootstrapping estimate Q means that 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.

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,

$$\boldsymbol{\varphi} \leftarrow \boldsymbol{\varphi} + \eta_t \gamma^t (Q(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta}) - V(\mathbf{x}; \boldsymbol{\theta})) \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a} \mid \mathbf{x}) \quad (12.67)$$

$$= \boldsymbol{\varphi} + \eta_t \gamma^t A(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a} \mid \mathbf{x}) \quad (12.68)$$

using the definition of the advantage function (12.49)

where $A(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta})$ is a bootstrapped estimate of the advantage function $a^{\pi_{\boldsymbol{\varphi}}}$. 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 for effectively trading bias and variance. This leads to algorithms such as *generalized advantage estimation* (GAE) (Schulman et al., 2015b).

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 12.5). 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., 2015a). TRPO

uses multiple iterations, where in each iteration a fixed critic is used to optimize the policy.¹⁵ During iteration k , we have,

$$\boldsymbol{\varphi}_{k+1} \leftarrow \arg \max_{\boldsymbol{\varphi}} J(\boldsymbol{\varphi}) \quad \text{subject to } \mathbb{E}_{x \sim \rho_{\boldsymbol{\varphi}_k}^{\infty}} \text{KL}(\pi_{\boldsymbol{\varphi}_k}(\cdot | x) \| \pi_{\boldsymbol{\varphi}}(\cdot | x)) \leq \delta \quad (12.69)$$

where

$$J(\boldsymbol{\varphi}) \doteq \mathbb{E}_{x \sim \rho_{\boldsymbol{\varphi}_k}^{\infty}, a \sim \pi_{\boldsymbol{\varphi}_k}(\cdot | x)} \left[\frac{\pi_{\boldsymbol{\varphi}}(a | x)}{\pi_{\boldsymbol{\varphi}_k}(a | x)} A^{\pi_{\boldsymbol{\varphi}_k}}(x, a) \right]. \quad (12.70)$$

Notably, J is an expectation with respect to the *previous* policy $\pi_{\boldsymbol{\varphi}_k}$ and the previous critic $A^{\pi_{\boldsymbol{\varphi}_k}}$. TRPO is an example of *importance weighting* where the importance weights,

$$\frac{\pi_{\boldsymbol{\varphi}}(a | x)}{\pi_{\boldsymbol{\varphi}_k}(a | x)},$$

are used to correct for using the previous policy. 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

$$\mathbb{E}_{x \sim \rho_{\boldsymbol{\varphi}_k}^{\infty}} \text{KL}(\pi_{\boldsymbol{\varphi}_k}(\cdot | x) \| \pi_{\boldsymbol{\varphi}}(\cdot | x)) \leq \delta$$

to optimize only in the “neighborhood” of the current policy. This constraint is also necessary for the importance weights not to blow up.

Taking the expectation with respect to the previous policy $\pi_{\boldsymbol{\varphi}_k}$ 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 heuristic variant of TRPO that often works well in practice (Schulman et al., 2017; Wang et al., 2020). It is based on controlling the importance weights directly rather than imposing a constraint on the KL-divergence. PPO is used, for example, to train large-scale language models such as GPT (Stiennon et al., 2020; OpenAI, 2023).

12.4.4 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,¹⁶ instead of score gradient estimators.

¹⁵ Intuitively, each iteration performs a collection of gradient ascent steps.

¹⁶ see section 5.5.1

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.¹⁷

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 12.4.5. Recall that our initial motivation to consider policy gradient methods and then actor-critic methods was the intractability of the DQN loss

$$\ell_{\text{DQN}}(\theta; \mathcal{D}) = \frac{1}{2} \sum_{(x, a, r, x') \in \mathcal{D}} \left(r + \gamma \max_{a' \in \mathcal{A}} Q^*(x', a'; \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2 \quad \text{see eq. (12.16)}$$

when the action space \mathcal{A} is large. What if we simply replace the exact maximum over actions by a parameterized policy?

$$\ell_{\text{DQN}}(\theta; \mathcal{D}) \approx \frac{1}{2} \sum_{(x, a, r, x') \in \mathcal{D}} \left(r + \gamma Q^*(x', \pi_\varphi(x'); \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2. \quad (12.71)$$

We want to train our parameterized policy to learn the maximization over action, that is, to approximate the greedy policy¹⁸

$$\pi_\varphi(x) \approx \pi^*(x) = \arg \max_{a \in \mathcal{A}} Q^*(x, a; \theta). \quad (12.72)$$

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

$$\varphi^* = \arg \max_{\varphi} \mathbb{E}_{x \sim \mu} [Q^*(x, \pi_\varphi(x); \theta)] \quad (12.73)$$

where $\mu(x) > 0$ is an *exploration distribution* over states with full support.¹⁹ We refer to this expectation by

$$J_\mu(\varphi; \theta) \doteq \mathbb{E}_{x \sim \mu} [Q^*(x, \pi_\varphi(x); \theta)]. \quad (12.74)$$

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 J_μ with respect to φ :

$$\nabla_\varphi J_\mu(\varphi; \theta) = \mathbb{E}_{x \sim \mu} [\nabla_\varphi Q^*(x, \pi_\varphi(x); \theta)]. \quad (12.75)$$

¹⁷ Policy evaluation is at the core of policy iteration. See alg. 10.17 for the definition of policy iteration and section 10.2 for a summary of policy evaluation in the context of Markov decision processes.

¹⁸ Here, we already apply the improvement of DDQN to use the most-recent estimate of the Q-function for action selection (see eq. (12.17)).

¹⁹ We require full support to ensure that all states are explored.

as we have seen in remark 1.22

Analogously to on-policy actor-critics (see eq. (12.64)), 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,

$$\nabla_\varphi Q^*(x, \pi_\varphi(x); \theta) = D_\varphi \pi_\varphi(x) \cdot \nabla_a Q^*(x, a; \theta)|_{a=\pi_\varphi(x)}. \quad (12.76)$$

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 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”*.²⁰ This corresponds to an algorithm called *deep deterministic policy gradients* (Lillicrap et al., 2015) shown in alg. 12.17. This algorithm is essentially equivalent to Q-learning with function approximation (e.g., DQN),²¹ with the only exception that we replace the maximization over actions with the learned policy π_φ .

²⁰ Intuitively, this adds “additional randomness” to the policy π_φ .

²¹ see eq. (12.16)

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.

12.4.5 Off-Policy Actor Critics with 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,

$$\ell_{\text{DQN}}(\theta; \mathcal{D}) \approx \frac{1}{2} \left(r + \gamma Q^*(x', \pi(x'; \varphi); \theta^{\text{old}}) - Q^*(x, a; \theta) \right)^2,$$

refer to the squared loss of Q-learning (12.11)

Algorithm 12.17: 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  $a = \pi_{\varphi}(x) + \epsilon$  for  $\epsilon \sim \mathcal{N}(\mathbf{0}, \lambda I)$ 
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              $\varphi^{\text{old}} \leftarrow (1 - \rho) \varphi^{\text{old}} + \rho \varphi$ 

```

which only considers the fixed action $\pi(\mathbf{x}'; \boldsymbol{\varphi}^{\text{old}})$ with an expected squared loss,

$$\ell_{\text{DQN}}(\boldsymbol{\theta}; \mathcal{D}) \approx \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{x}', \boldsymbol{\varphi})} \left[\frac{1}{2} \left(r + \gamma Q^*(\mathbf{x}', \mathbf{a}'; \boldsymbol{\theta}^{\text{old}}) - Q^*(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta}) \right)^2 \right], \quad (12.77)$$

which considers a distribution over actions.

It turns out that we can still compute gradients of this expectation.

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{x}', \boldsymbol{\varphi})} \left[\frac{1}{2} \left(r + \gamma Q^*(\mathbf{x}', \mathbf{a}'; \boldsymbol{\theta}^{\text{old}}) - Q^*(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta}) \right)^2 \right] \\ = \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{x}', \boldsymbol{\varphi})} \left[\nabla_{\boldsymbol{\theta}} \frac{1}{2} \left(r + \gamma Q^*(\mathbf{x}', \mathbf{a}'; \boldsymbol{\theta}^{\text{old}}) - Q^*(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta}) \right)^2 \right]. \end{aligned} \quad \text{as we have seen in remark 1.22}$$

Similarly to our definition of the Bellman error (12.12), we define by

$$\delta_{\text{B}}(\mathbf{a}') \doteq r + \gamma Q^*(\mathbf{x}', \mathbf{a}'; \boldsymbol{\theta}^{\text{old}}) - Q^*(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta}), \quad (12.78)$$

the *Bellman error* for a fixed action \mathbf{a}' . Using the chain rule, we obtain

$$= \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{x}', \boldsymbol{\varphi})} [\delta_{\text{B}}(\mathbf{a}') \nabla_{\boldsymbol{\theta}} Q^*(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta})]. \quad (12.79)$$

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 $\nabla_{\boldsymbol{\theta}} Q^*(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta})$. 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

$$J_{\mu}(\boldsymbol{\varphi}; \boldsymbol{\theta}) \doteq \mathbb{E}_{\mathbf{x} \sim \mu} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{x}; \boldsymbol{\varphi})} [Q^*(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta})].$$

of which we can obtain gradients via

$$\nabla_{\boldsymbol{\varphi}} J_{\mu}(\boldsymbol{\varphi}; \boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim \mu} \nabla_{\boldsymbol{\varphi}} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{x}; \boldsymbol{\varphi})} [Q^*(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta})]. \quad (12.80) \quad \text{see remark 1.22}$$

Note that the inner expectation is with respect to a measure that depends on the parameters $\boldsymbol{\varphi}$, 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.²² Earlier, in chapter 5 on variational inference we have already seen reparameterization gradients.²³ Here, if our policy is reparameterizable, we can use the *reparameterization trick* from theorem 5.31!

²² see eq. (12.32)

²³ see (5.78)

Example 12.18: 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

$$a = g(\epsilon; x, \varphi) \doteq \Sigma^{1/2}(x; \varphi)\epsilon + \mu(x; \varphi), \quad \epsilon \sim \mathcal{N}(0, I) \quad (12.81)$$

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 5.32.

As we have seen, not only Gaussians are reparameterizable. In general, we called a distribution (in this context, a policy) reparameterizable iff $a \sim \pi(x; \varphi)$ is such that $a = g(\epsilon; x, \varphi)$, where $\epsilon \sim \phi$ is an independent random variable.

Then, we have,

$$\begin{aligned} \nabla_{\varphi} \mathbb{E}_{a \sim \pi(x; \varphi)} [Q^*(x, a; \theta)] \\ = \mathbb{E}_{\epsilon \sim \phi} [\nabla_{\varphi} Q^*(x, g(\epsilon; x, \varphi); \theta)] \end{aligned} \quad (12.82)$$

using the reparameterization trick (5.78)

$$= \mathbb{E}_{\epsilon \sim \phi} \left[\nabla_a Q^*(x, a; \theta) \Big|_{a=g(\epsilon; x, \varphi)} \cdot D_{\varphi} g(\epsilon; x, \varphi) \right]. \quad (12.83)$$

using the chain rule analogously to eq. (12.76)

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 eq. (12.79) to obtain gradients for the critic and reparameterization gradients for the actor is called *stochastic value gradients* (SVG) (Heess et al., 2015).

Readings

For more details on actor-critic methods, refer to chapter 13 (from section 13.5) of “Reinforcement learning: An introduction” (Sutton and Barto, 2018).

12.5 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.²⁴ 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

²⁴ see section 5.4

$$j_\lambda(\boldsymbol{\varphi}) \doteq j(\boldsymbol{\varphi}) + \lambda H[\Pi_\varphi] \quad (12.84)$$

$$= \sum_{t=1}^T \mathbb{E}_{(x_t, a_t) \sim \Pi_\varphi} [r(x_t, a_t) + \lambda H[\pi_\varphi(\cdot | x_t)]], \quad (12.85)$$

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):

$$j(\boldsymbol{\varphi}) = \sum_{t=1}^T \mathbb{E}_{(x_t, a_t) \sim \Pi_\varphi} [r(x_t, a_t)]. \quad (12.86)$$

Here, for notational convenience, we begin the sum with $t = 1$ rather than $t = 0$.

12.5.1 Entropy Regularization as Probabilistic Inference

The entropy-regularized objective from eq. (12.85) 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_\star the distribution over trajectories τ under the optimal policy π^\star . 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 optimal we can derive Π_\star 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 .

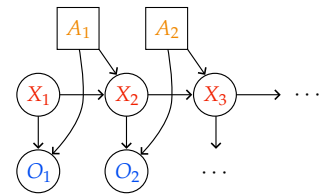


Figure 12.4: Directed graphical model of the underlying hidden Markov model with hidden states X_t , optimality variables O_t , and actions A_t .

We consider the HMM defined by the Gibbs distribution

$$p(\mathcal{O}_t \mid \mathbf{x}_t, \mathbf{a}_t) \propto \exp\left(\frac{1}{\lambda} r(\mathbf{x}_t, \mathbf{a}_t)\right), \quad \text{with } \lambda > 0 \quad (12.87)$$

which is a natural choice as we have seen in exercise 6.25 that the Gibbs distribution maximizes entropy subject to $\mathbb{E}[\mathcal{O}_t \cdot r(\mathbf{x}_t, \mathbf{a}_t) \mid \mathbf{x}_t, \mathbf{a}_t] < \infty$.

The distribution over trajectories conditioned on optimality of actions (i.e., conditioned on $\mathcal{O}_{1:T}$) is given by

$$\Pi_\star(\tau) \doteq p(\tau \mid \mathcal{O}_{1:T}) = p(\mathbf{x}_1) \prod_{t=1}^T p(\mathbf{a}_t \mid \mathbf{x}_t, \mathcal{O}_t) p(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{a}_t). \quad (12.88)$$

It remains to determine $p(\mathbf{a}_t \mid \mathbf{x}_t, \mathcal{O}_t)$ which corresponds to the optimal policy $\pi^\star(\mathbf{a}_t \mid \mathbf{x}_t)$. It is generally useful to think of the situation where the prior policy $p(\mathbf{a}_t \mid \mathbf{x}_t)$ is uniform on \mathcal{A} ,²⁵ in which case by Bayes' rule (1.59), $p(\mathbf{a}_t \mid \mathbf{x}_t, \mathcal{O}_t) \propto p(\mathcal{O}_t \mid \mathbf{x}_t, \mathbf{a}_t)$, so

$$\Pi_\star(\tau) \propto \left[p(\mathbf{x}_1) \prod_{t=1}^T p(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{a}_t) \right] \exp\left(\frac{1}{\lambda} \sum_{t=1}^T r(\mathbf{x}_t, \mathbf{a}_t)\right). \quad (12.89)$$

Recall that our fundamental goal is to approximate Π_\star with a distribution over trajectories Π_φ under the parameterized policy π_φ . It is therefore a natural idea to minimize $\text{KL}(\Pi_\varphi \parallel \Pi_\star)$:²⁶

$$\begin{aligned} & \arg \min_{\varphi} \text{KL}(\Pi_\varphi \parallel \Pi_\star) \\ &= \arg \min_{\varphi} H[\Pi_\varphi \parallel \Pi_\star] - H[\Pi_\varphi] \\ &= \arg \max_{\varphi} \mathbb{E}_{\tau \sim \Pi_\varphi} [\log \Pi_\star(\tau) - \log \Pi_\varphi(\tau)] \\ &= \arg \max_{\varphi} \mathbb{E}_{\tau \sim \Pi_\varphi} \left[\sum_{t=1}^T r(\mathbf{x}_t, \mathbf{a}_t) - \lambda \log \pi_\varphi(\mathbf{a}_t \mid \mathbf{x}_t) \right] \\ &= \arg \max_{\varphi} \sum_{t=1}^T \mathbb{E}_{(\mathbf{x}_t, \mathbf{a}_t) \sim \Pi_\varphi} [r(\mathbf{x}_t, \mathbf{a}_t) + \lambda H[\pi_\varphi(\cdot \mid \mathbf{x}_t)]]. \end{aligned} \quad (12.90)$$

That is, entropy regularization is equivalent to minimizing the KL-divergence from Π_\star to Π_φ . This highlights a very natural tradeoff between exploration and exploitation, wherein $H[\Pi_\varphi \parallel \Pi_\star]$ encourages exploitation and $H[\Pi_\varphi]$ encourages exploration.

It can be shown that a “softmax” version of the Bellman optimality equation (10.35) can be obtained for eq. (12.90):

$$q^\star(\mathbf{x}, \mathbf{a}) = \frac{1}{\lambda} r(\mathbf{x}, \mathbf{a}) + \mathbb{E}_{\mathbf{x}' \sim \mathbf{x}, \mathbf{a}} \left[\log \int_{\mathcal{A}} \exp(q^\star(\mathbf{x}', \mathbf{a}')) d\mathbf{a}' \right] \quad (12.91)$$

Using eq. (12.30). We assume here that the dynamics and initial state distribution are “fixed”, that is, we assume $p(\mathbf{x}_1 \mid \mathcal{O}_{1:T}) = p(\mathbf{x}_1)$ and $p(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{a}_t, \mathcal{O}_{1:T}) = p(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{a}_t)$.

²⁵ This is not a restriction as any non-uniform prior can be pushed into eq. (12.87).

²⁶ Observe that we cannot easily minimize forward-KL as we cannot sample from Π_\star . 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 (5.44)

using the definition of cross-entropy (5.40) and entropy (5.32)

using eqs. (12.30) and (12.89) and simplifying

using the definition of entropy (5.32) and linearity of expectation (1.24)

and the optimal policy has the form $\pi^*(a | x) \propto \exp(q^*(x, a))$. Here, q^* is called a *soft value function*. The second term of eq. (12.91) quantifies downstream rewards. In comparison to the “standard” Bellman optimality equation (10.35), the soft value function is less greedy which tends to encourage robustness.

Exercise 12.19: Soft value function

In this exercise, we derive the soft value function from eq. (12.91).

1. We let $\beta(a_t | x_t) \doteq p(a_t | x_t, \mathcal{O}_t)$, $Z(x) \doteq \int_{\mathcal{A}} \beta(a | x) da$, and denote by $\hat{\pi}(\cdot | x)$ the policy $\beta(\cdot | x)/Z(x)$. Show that

$$\begin{aligned} \text{KL}(\Pi_{\varphi} \| \Pi_{\star}) \\ = \sum_{t=1}^T \mathbb{E}_{x_t \sim \Pi_{\varphi}} [\text{KL}(\pi_{\varphi}(\cdot | x_t) \| \hat{\pi}(\cdot | x_t)) - \log Z(x_t)]. \end{aligned} \quad (12.92)$$

2. Show that if the space of policies parameterized by φ is sufficiently expressive, $\pi^*(a | x) \propto \exp(q^*(x, a))$ solves eq. (12.90).

▷ *Solution*

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 and 12.4 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.

We can also express this optimization in terms of an evidence lower bound. The evidence lower bound for the observations $\mathcal{O}_{1:T}$ is

$$L(\Pi_{\varphi}, \Pi_{\star}; \mathcal{O}_{1:T}) = \mathbb{E}_{\tau \sim \Pi_{\varphi}} [\log p(\mathcal{O}_{1:T} | \tau) + \log \Pi(\tau) - \log \Pi_{\varphi}(\tau)] \quad \text{using the definition of the ELBO (5.63c)} \quad (12.93)$$

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

$$-L(\Pi_{\varphi}, \Pi_{\star}; \mathcal{O}_{1:T}) = \mathbb{E}_{\tau \sim \Pi_{\varphi}} [S[p(\mathcal{O}_{1:T} | \tau)]] + \text{KL}(\Pi_{\varphi} \| \Pi) \quad (12.94)$$

$$= \underbrace{S[p(\mathcal{O}_{1:T})]}_{\text{“extrinsic” value}} + \underbrace{\text{KL}(\Pi_{\varphi} \| \Pi_{\star})}_{\text{“epistemic” value}}. \quad (12.95) \quad \text{using that } \Pi_{\star}(\tau) = p(\tau | \mathcal{O}_{1:T})$$

which we already encountered in eq. (5.68) in the context of variational

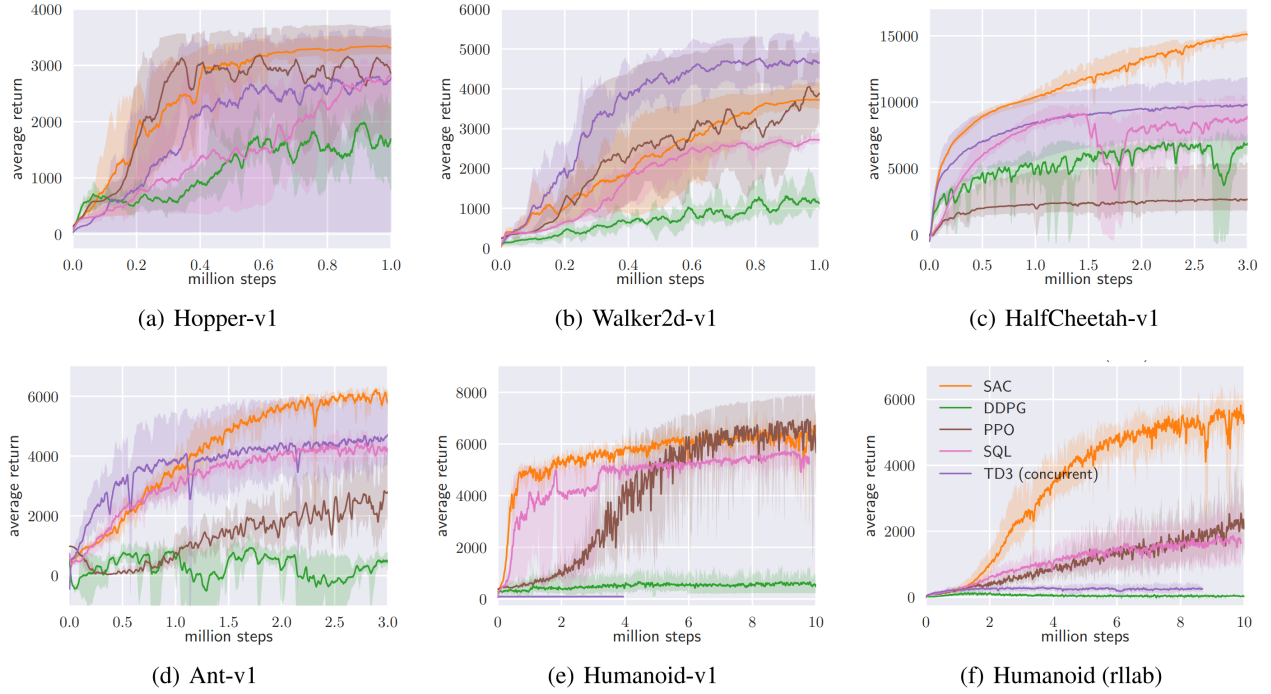


Figure 12.5: Comparison of training curves of a selection of on-policy and off-policy policy gradient methods. Taken from “Soft actor-critic algorithms and applications” (Haarnoja et al., 2018b).

inference. To summarize, we have seen that

$$\arg \max_{\boldsymbol{\varphi}} j_{\lambda}(\boldsymbol{\varphi}) = \arg \min_{\boldsymbol{\varphi}} \text{KL}(\Pi_{\boldsymbol{\varphi}} \| \Pi_{\star}) = \arg \max_{\boldsymbol{\varphi}} L(\Pi_{\boldsymbol{\varphi}}, \Pi_{\star}; \mathcal{O}_{1:T}).$$

Recall that the free energy $-L(\Pi_{\boldsymbol{\varphi}}, \Pi_{\star}; \mathcal{O}_{1:T})$ is a variational upper bound to the surprise about observations $S[\Pi_{\star}(\mathcal{O}_{1:T})]$ when following an optimal policy.²⁷ 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.²⁸

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*.

²⁷ see remark 5.27

²⁸ This is because, a path leading to a low reward state will include sub-optimal actions.

Optional Readings

For more details, read “Reinforcement learning and control as probabilistic inference: Tutorial and review” (Levine, 2018).

Remark 12.20: Towards active inference

Maximum entropy reinforcement learning makes one fundamental assumption, namely, that the “biased” distribution about observations specifying the underlying HMM is

$$p(\mathcal{O}_t \mid \mathbf{x}_t, \mathbf{a}_t) \propto \exp\left(\frac{1}{\lambda} r(\mathbf{x}_t, \mathbf{a}_t)\right).$$

This assumption is how the task of optimal control (in an unknown MDP) with a certain reward function is framed 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:

$$p(\mathcal{O}_t \mid \mathbf{x}_t, \mathbf{a}_t) \propto \exp\left(\frac{1}{\lambda} Q(\mathbf{x}_t, \mathbf{a}_t; \theta)\right).$$

Crucially, the choice of $p(\mathcal{O}_t \mid \mathbf{x}_t, \mathbf{a}_t)$ 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 HMMs) is known as *active inference* (Friston et al., 2015; Millidge et al., 2020, 2021).

12.5.2 Summary

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. 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.

We then 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,

and SAC that are widely used and work quite well in practice. Importantly, they are significantly more sample efficient than on-policy policy gradient methods, which often results in much faster learning of a near-optimal policy.

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 (2015b).
High-dimensional continuous control using generalized advantage estimation.
- **TRPO:** Schulman, Levine, Abbeel, Jordan, and Moritz (2015a).
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 (2015).
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.

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 Bayesian learning, 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 .

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  $\pi$  and no (or some) initial data  $\mathcal{D}$ 
2 for several episodes do
3   roll out policy  $\pi$  to collect data
4   learn a model of the dynamics  $f$  and rewards  $r$  from data
5   plan a new policy  $\pi$  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 non-linear 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 non-linear 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,

$$x_{t+1} = f(x_t, a_t). \quad (13.1)$$

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,

$$\max_{a_{0:\infty}} \sum_{t=0}^{\infty} \gamma^t r(x_t, a_t) \quad \text{such that } x_{t+1} = f(x_t, a_t). \quad (13.2)$$

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  $x_t$ 
3   plan over a finite horizon  $H$ ,
      
$$\max_{a_{t:t+H-1}} \sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} r(x_\tau, a_\tau) \quad \text{such that } x_{\tau+1} = f(x_\tau, a_\tau) \quad (13.3)$$

4   carry out action  $a_t$ 

```

Observe that the state x_τ can be interpreted as a deterministic function $x_\tau(a_{t:\tau-1})$, which depends on all actions from time t to time $\tau - 1$ and the state x_t . To solve the optimization problem of a single iteration, we therefore need to maximize,

$$J_H(a_{t:t+H-1}) \doteq \sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} r(x_\tau(a_{t:\tau-1}), a_\tau). \quad (13.4)$$

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 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.

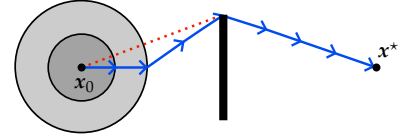


Figure 13.1: Illustration of model predictive control in a deterministic transition model. The agent starts in position x_0 and wants to reach x^* despite the black obstacle. We use the reward function $r(x) = -\|x - x^*\|$. The gray concentric circles represent the length of a single step. We plan with 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 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 naïve 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*, 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

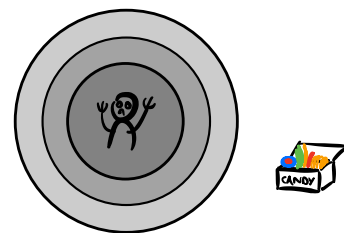


Figure 13.2: Illustration of finite-horizon planning with sparse rewards. When the finite time horizon does not suffice to “reach” a reward, the agent has no signal to follow.

future rewards.

$$J_H(\mathbf{a}_{t:t+H-1}) \doteq \underbrace{\sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} r(\mathbf{x}_\tau(\mathbf{a}_{t:\tau-1}), \mathbf{a}_\tau)}_{\text{short-term}} + \underbrace{\gamma^H V(\mathbf{x}_{t+H})}_{\text{long-term}}. \quad (13.6)$$

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,

$$\mathbf{a}_t = \arg \max_{\mathbf{a} \in \mathcal{A}} J_1(\mathbf{a}) = \pi_V(\mathbf{x}_t). \quad (13.7)$$

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  $\mathbf{x}_t$ 
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  $\mathbf{a}_t$ 

```

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 non-linear 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 12.4.5). 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 $\mathbf{x}_{t+1} \sim f(\mathbf{x}_t, \mathbf{a}_t)$ is such that $\mathbf{x}_{t+1} = g(\epsilon; \mathbf{x}_t, \mathbf{a}_t)$, where $\epsilon \sim \phi$ is a random variable that is independent of \mathbf{x}_t and \mathbf{a}_t . We have already seen in example 12.18 (in the context of stochastic policies) how a Gaussian transition model can be reparameterized.

In this case, \mathbf{x}_{τ} is determined recursively by $\mathbf{a}_{t:\tau-1}$ and $\epsilon_{t:\tau-1}$,

$$\begin{aligned} \mathbf{x}_{\tau} &= \mathbf{x}_{\tau}(\epsilon_{t:\tau-1}; \mathbf{a}_{t:\tau-1}) \\ &\doteq g(\epsilon_{\tau-1}; g(\dots; (\epsilon_{t+1}; g(\epsilon_t; \mathbf{x}_t, \mathbf{a}_t)), \mathbf{a}_{t+1}), \dots), \mathbf{a}_{\tau-1}). \end{aligned} \quad (13.9)$$



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 5.31 and eq. (12.82)

This allows us to obtain unbiased estimates of J_H using Monte Carlo approximation,

$$J_H(\mathbf{a}_{t:t+H-1}) \approx \frac{1}{m} \sum_{i=1}^m \sum_{\tau=t}^{t+H-1} \left(\gamma^{\tau-t} r(\mathbf{x}_\tau(\boldsymbol{\epsilon}_{t:\tau-1}^{(i)}; \mathbf{a}_{t:\tau-1}), \mathbf{a}_\tau) + \gamma^H V(\mathbf{x}_{t+H}) \right) \quad (13.10)$$

where $\boldsymbol{\epsilon}_{t:t+H-1}^{(i)} \stackrel{\text{iid}}{\sim} \phi$ 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,

$$\pi^*(\mathbf{x}) = \arg \max_{\mathbf{a} \in \mathcal{A}} Q^*(\mathbf{x}, \mathbf{a}; \boldsymbol{\theta}), \quad \text{see eq. (12.20)}$$

and therefore had to solve an optimization problem over all actions. We accelerated this by learning an approximate policy that “mimicked” this optimization,

$$\boldsymbol{\varphi}^* = \arg \max_{\boldsymbol{\varphi}} \underbrace{\mathbb{E}_{\mathbf{x} \sim \mu} [Q^*(\mathbf{x}, \pi_{\boldsymbol{\varphi}}(\mathbf{x}); \boldsymbol{\theta})]}_{J_{\mu}(\boldsymbol{\varphi}; \boldsymbol{\theta})} \quad \text{see eq. (12.73)}$$

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,

$$J_{\mu,H}(\boldsymbol{\varphi}; \boldsymbol{\theta}) \doteq \mathbb{E}_{x_0 \sim \mu, x_{1:H} | \pi_{\boldsymbol{\varphi},f}} \left[\sum_{\tau=0}^{H-1} \gamma^{\tau} r_{\tau} + \gamma^H Q^*(x_H, \pi_{\boldsymbol{\varphi}}(x_H); \boldsymbol{\theta}) \right] \quad (13.12)$$

This approach naturally extends to randomized policies using reparameterization gradients, which we have discussed in section 12.4.5. 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.

Remark 13.6: Planning as inference

We have been using *Monte Carlo rollouts* to estimate the expectation of eq. (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 course. 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 eq. (13.12). A prominent example of this approach is *probabilistic inference for learning control* (PILCO).
- In section 12.5, 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 eq. (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 Bayesian Learning

In the following, we will discuss how we can use the techniques from Bayesian learning, 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,

$$x_{t+1} \sim f(x_t, a_t; \psi). \quad (13.14)$$

Example 13.7: Conditional Gaussian dynamics

We could, for example, use a conditional Gaussian for the transition model,

$$x_{t+1} \sim \mathcal{N}(\mu(x_t, a_t; \psi), \Sigma(x_t, a_t; \psi)). \quad (13.15)$$

As we have seen in eq. (1.134), we can rewrite the covariance matrix Σ as a product of a lower-triangular matrix \mathcal{L} and its transpose using the Cholesky decomposition $\Sigma = \mathcal{L}\mathcal{L}^\top$ of Σ . 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 2.3 that epistemic uncertainty corresponds to our uncertainty about the model, $p(f \mid \mathcal{D})$, while aleatoric uncertainty corresponds to the uncertainty of the transitions in the underlying Markov decision process (which can be thought of as “irreducible” noise), $p(x_{t+1} \mid f, x_t, a_t)$.

Intuitively, Bayesian learning 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 (as a Bayesian prior) about the transition and rewards model.

Example 13.9: Inference with conditional Gaussian dynamics

Let us revisit inference with our conditional Gaussian dynamics model from eq. (13.15),

$$x_{t+1} \sim \mathcal{N}(\mu(x_t, a_t; \psi), \Sigma(x_t, a_t; \psi)).$$

Recall that in the setting of Bayesian neural networks, most ap-

⁴ Quantifying the uncertainty of an estimate is a problem that we have spent the first few chapters exploring. Notably, refer to

- section 2.3 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.

proximate inference techniques represented the approximate posterior using some form of a mixture of Gaussians,⁵

$$p(\mathbf{x}_{t+1} \mid \mathcal{D}, \mathbf{x}_t, \mathbf{a}_t) \approx \frac{1}{m} \sum_{i=1}^m \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_t, \mathbf{a}_t; \boldsymbol{\psi}^{(i)}), \boldsymbol{\Sigma}(\mathbf{x}_t, \mathbf{a}_t; \boldsymbol{\psi}^{(i)})). \quad (13.16)$$

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 J_H ,

$$J_H(\mathbf{a}_{t:t+H-1}) \approx \frac{1}{m} \sum_{i=1}^m J_H(\mathbf{a}_{t:t+H-1}; \mathbf{f}^{(i)}) \quad \text{where} \quad (13.17)$$

$$J_H(\mathbf{a}_{t:t+H-1}; f) \doteq \sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} r(\mathbf{x}_\tau(\boldsymbol{\epsilon}_{t:\tau-1}^{(i)}; \mathbf{a}_{t:\tau-1}, f), \mathbf{a}_\tau) + \gamma^H V(\mathbf{x}_{t+H}). \quad (13.18)$$

Here, $\mathbf{f}^{(i)} \stackrel{\text{iid}}{\sim} p(\mathbf{f} \mid \mathcal{D})$ are independent samples of the transition model, and $\boldsymbol{\epsilon}_{t:t+H-1}^{(i)} \stackrel{\text{iid}}{\sim} \boldsymbol{\phi}$ parameterizes the dynamics analogously to eq. (13.9):

$$\begin{aligned} \mathbf{x}_\tau(\boldsymbol{\epsilon}_{t:\tau-1}; \mathbf{a}_{t:\tau-1}, f) \\ \doteq f(\boldsymbol{\epsilon}_{\tau-1}; f(\dots; (\boldsymbol{\epsilon}_{t+1}; f(\boldsymbol{\epsilon}_t; \mathbf{x}_t, \mathbf{a}_t), \mathbf{a}_{t+1}), \dots), \mathbf{a}_{\tau-1}). \end{aligned} \quad (13.19)$$

Observe that the epistemic and aleatoric uncertainty are treated differently. Within a particular MDP f , we ensure that randomness (i.e., aleatoric 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 eq. (13.17) we are now using

⁵ see

- eq. (7.18) for variational inference;
- eq. (7.21a) for Marko chain Monte Carlo;
- eq. (7.28) for dropout regularization; and
- eq. (7.29) for probabilistic ensembles.

⁶ see section 7.3.1

⁷ see eq. (13.10)

Monte Carlo trajectory sampling as a subroutine and average over an “ensemble” of Markov decision processes.



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 $J_H(\pi)$ instead of $J_H(a_{t:t+H-1})$.

This leads us to a first greedily exploitative algorithm for model-based reinforcement learning, which is shown in alg. 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 $p(f) = p(f \mid \mathcal{D})$
 - 2 **for** several episodes **do**
 - 3 plan a new policy π to (approximately) maximize,

$$\max_{\pi} \mathbb{E}_{f \sim p(\cdot \mid \mathcal{D})} [J_H(\pi; f)] \quad (13.20)$$
 - 4 roll out policy π to collect more data
 - 5 update posterior $p(f \mid \mathcal{D})$
-

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.

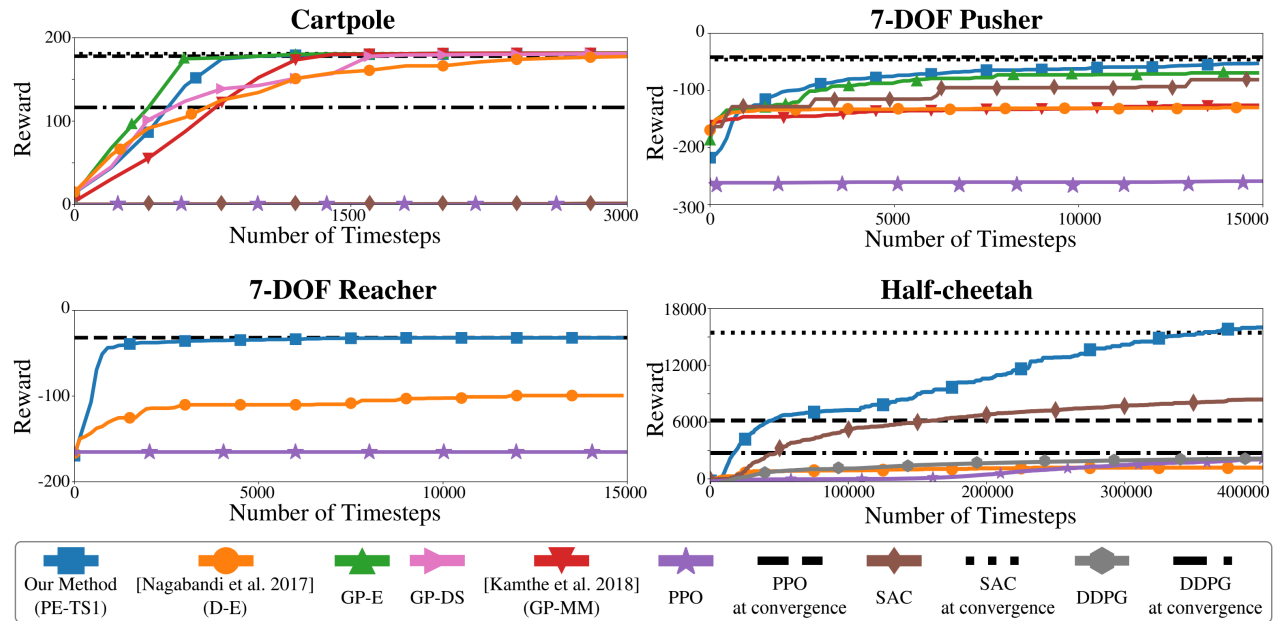


Figure 13.5: Sample efficiency of model-free and model-based reinforcement learning. DDPG and SAC are shown as constant (black) lines, because they take an order of magnitude more time steps before learning a good model. They also compare the PETS algorithm (blue) to deterministic ensembles (orange), where the transition models are assumed to be deterministic (or to have covariance 0). Taken from “Deep reinforcement learning in a handful of trials using probabilistic dynamics models” (Chua et al., 2018).

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.

Optional Readings

- Deisenroth and Rasmussen (2011).
PILCO: A model-based and data-efficient approach to policy search.
- Chua, Calandra, McAllister, and Levine (2018).
Deep reinforcement learning in a handful of trials using probabilistic dynamics models.

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 alg. 13.11 which is an immediate adaptation of greedy exploitation and straightforward to implement.

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 explo-

Algorithm 13.11: Thompson sampling

```

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   sample a model  $f \sim p(\cdot \mid \mathcal{D})$ 
4   plan a new policy  $\pi$  to (approximately) maximize,
      
$$\max_{\pi} J_H(\pi; f) \quad (13.21)$$

5   roll out policy  $\pi$  to collect more data
6   update posterior  $p(f \mid \mathcal{D})$ 

```

ration. 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,

$$f_i(\mathbf{x}, \mathbf{a}) \in \mathcal{C}_i \doteq [\mu_i(\mathbf{x}, \mathbf{a} \mid \mathcal{D}) - \beta_i \sigma_i(\mathbf{x}, \mathbf{a} \mid \mathcal{D}), \mu_i(\mathbf{x}, \mathbf{a} \mid \mathcal{D}) + \beta_i \sigma_i(\mathbf{x}, \mathbf{a} \mid \mathcal{D})], \quad (13.22)$$

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 alg. 13.13 takes the optimal step with respect to the most optimistic model among all transition models that are consistent with the data.

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

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  $\pi$  to (approximately) maximize,
      
$$\max_{\pi} \max_{f \in \mathcal{M}(\mathcal{D})} J_H(\pi; f) \quad (13.24)$$

4   roll out policy  $\pi$  to collect more data
5   update posterior  $p(f \mid \mathcal{D})$ 

```

choose which of the plausible dynamics it follows. The “luck” of the agent can be considered as additional decision variables. Consider the optimization problem,

$$\pi_{t+1} \doteq \arg \max_{\pi} \max_{\eta(\cdot) \in [-1,1]^d} J_H(\pi; \hat{f}_t) \quad (13.25)$$

with “optimistic” dynamics

$$\hat{f}_{t,i}(\mathbf{x}, \mathbf{a}) \doteq \mu_{t,i}(\mathbf{x}, \mathbf{a}) + \beta_{t,i} \eta_i(\mathbf{x}, \mathbf{a}) \sigma_{t,i}(\mathbf{x}, \mathbf{a}). \quad (13.26)$$

Here the decision variables η_i control the variance of an action. That is, within the confidence bounds of the transition model, the agent can freely choose the state that is reached by playing an action \mathbf{a} from state \mathbf{x} . Essentially, this corresponds to maximizing expected reward in an augmented (optimistic) MDP with *known* dynamics \hat{f} and with a larger action space that also includes the decision variables η . This is a known MDP for which we can use our toolbox for planning which we developed in section 13.1.

The algorithm that maximizes expected reward in this augmented MDP is called *hallucinated upper confidence reinforcement learning* (H-UCRL) (Curi et al., 2020; Treven et al., 2023). H-UCRL can be seen as the natural extension of the UCB acquisition function from Bayesian optimization to reinforcement learning. An illustration of the algorithm is given in fig. 13.6.

Intuitively, the agent has the tendency to believe that it can achieve much more than it actually can. As more data is collected, the confidence bounds shrink and the optimistic policy rewards converge to the actual rewards. Yet, crucially, we only collect data in regions of the state-action space that are more promising than the regions that we have already explored. That is, we only collect data in the most promising regions.

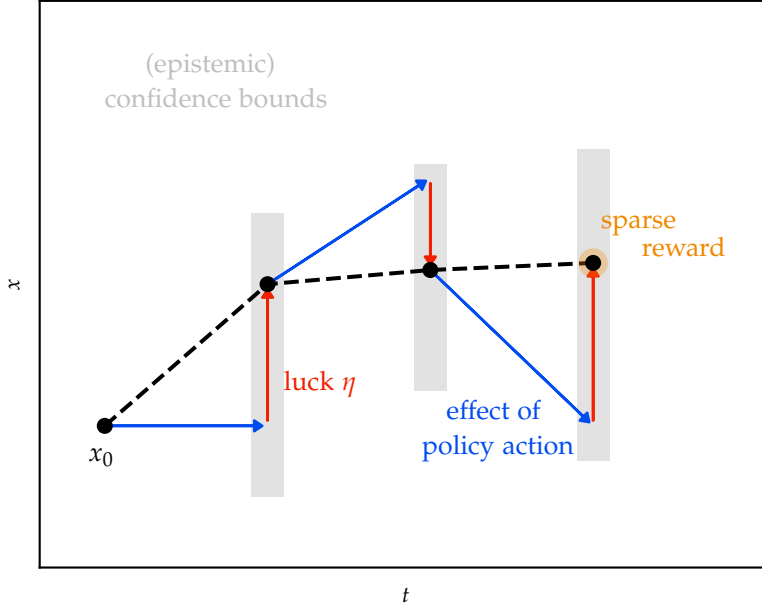


Figure 13.6: Illustration of H-UCRL in a one-dimensional state space. The agent “hallucinates” that it takes the black trajectory when, in reality, the outcomes of its actions are as shown in blue. The agent can hallucinate to land anywhere within the gray confidence regions (i.e., the epistemic uncertainty in the model) using the luck decision variables η . This allows agents to discover long sequences of actions leading to sparse rewards.

Exercise 13.14: Bounding the regret of H-UCRL

Analogously to our analysis of the UCB acquisition function for Bayesian optimization, we can use optimism to bound the regret of H-UCRL. We assume for simplicity that $d = 1$ and drop the index i in the following. Let us denote by $x_{t,k}$ the k -th state visited during episode t when following policy π_t . We denote by $\hat{x}_{t,k}$ the corresponding state but under the optimistic dynamics \hat{f} rather than the true dynamics f . The instantaneous regret during episode t is given by

$$r_t = J_H(\pi^*; f) - J_H(\pi_t; f) \quad (13.27)$$

where we take the objective to be $J_H(\pi; f) = \sum_{k=0}^{H-1} r(x_{t,k}, \pi(x_{t,k}))$.

You may assume that $f(x, \pi(x))$, $r(x, \pi(x))$, and $\sigma(x, \pi(x))$ are Lipschitz continuous in x .

1. Show by induction that for any $k \geq 0$, with high probability,

$$\|\hat{x}_{t,k} - x_{t,k}\| \leq 2\beta_t \sum_{l=0}^{k-1} \alpha_t^{k-1-l} \sigma_{t-1}(x_{t,l}, \pi_t(x_{t,l})) \quad (13.28)$$

where α_t depends on the Lipschitz constants and β_t .

2. Assuming w.l.o.g. that $\alpha_t \geq 1$, show that with high probability,

$$r_t \leq \mathcal{O}\left(\beta_t H \alpha_t^{H-1} \sum_{k=0}^{H-1} \sigma_{t-1}(x_{t,k}, \pi_t(x_{t,k}))\right) \quad (13.29)$$

3. Let $\Gamma_T \doteq \max_{\pi_1, \dots, \pi_T} \sum_{t=1}^T \sum_{k=0}^{H-1} \sigma_{t-1}^2(x_{t,k}, \pi_t(x_{t,k}))$. Analogously to exercise 9.7, it can be derived that $\Gamma_T \leq \mathcal{O}(H\gamma_T)$ if the dynamics are modeled by a Gaussian process.⁹ Bound the cumulative regret $R_T = \sum_{t=1}^T r_t \leq \mathcal{O}\left(\beta_T H^{\frac{3}{2}} \alpha_T^{H-1} \sqrt{T\Gamma_T}\right)$.

Thus, if the dynamics are modeled by a Gaussian process with kernel such that γ_T is sublinear, the regret of H-UCRL is sublinear.

▷ *Solution*

Optimistic exploration yields the strongest effects for hard exploration tasks, for example, settings with large penalties associated with performing certain actions and settings with sparse rewards.¹⁰ In those settings, most other strategies (i.e., those that are not sufficiently explorative), learn not to act at all. However, even in settings of “ordinary rewards”, optimistic exploration often learns good policies faster.

13.3.3 Constrained Exploration

Besides making exploration more efficient, another use of uncertainty is to make exploration more safe. Today, reinforcement learning is still far away from being deployed directly to the real world. In practice, reinforcement learning is almost always used together with a simulator, in which the agent can “safely” train and explore. Yet, in many domains, it is not possible to simulate the training process, either because we are lacking a perfect model of the environment, or because simulating such a model is too computationally inefficient. This is where we can make use of uncertainty estimates of our model to avoid unsafe states.

Let us denote by $\mathcal{X}_{\text{unsafe}}$ the subset of unsafe states of our state space \mathcal{X} . A natural idea is to perform planning using confidence bounds of our epistemic uncertainty. This allows us to pessimistically forecast the plausible consequences of our actions, given what we have already learned about the transition model. As we collect more data, the confidence bounds will shrink, requiring us to be less conservative over time. This idea is also at the heart of fields like *robust control*.

A general formalism for planning under constraints is the notion of *constrained Markov decision processes* (Altman, 1999). Given dynamics f , planning in constrained MDPs amounts to the following optimization problem:

$$\max_{\pi} J_{\mu}(\pi; f) \doteq \mathbb{E}_{x_0 \sim \mu, x_{1:\infty} | \pi, f} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \quad (13.30a)$$

$$\text{subject to } J_{\mu}^c(\pi; f) \leq \delta \quad (13.30b)$$

⁹ For details, see appendix A of Treven et al. (2023).

¹⁰ Action penalties are often used to discourage the agent from exhibiting unwanted behavior. However, increasing the action penalty increases the difficulty of exploration. Therefore, optimistic exploration is especially useful in settings where we want to practically disallow many actions by attributing large penalties to them.

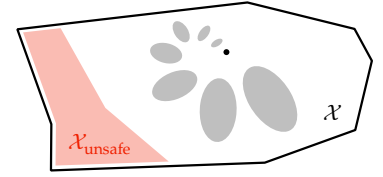


Figure 13.7: Illustration of planning with confidence bounds. The unsafe set of states is shown as the red region. The agent starts at the position denoted by the black dot and plans a sequence of actions. The confidence bounds on the subsets of states that are reached by this action sequence are shown in gray.

where μ is a distribution over the initial state and

$$J_\mu^c(\pi; f) \doteq \mathbb{E}_{x_0 \sim \mu, x_{1:\infty} | \pi, f} \left[\sum_{t=0}^{\infty} \gamma^t c(x_t) \right] \quad (13.31)$$

are expected discounted costs with respect to a cost function $c : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$.¹¹ Observe that for the cost function $c(x) \doteq \mathbb{1}\{x \in \mathcal{X}_{\text{unsafe}}\}$, the value $J_\mu^c(\pi; f)$ can be interpreted as an upper bound on the (discounted) probability of visiting unsafe states under dynamics f ,¹² and hence, the constraint $J_\mu^c(\pi; f) \leq \delta$ bounds the probability of visiting an unsafe state when following policy π .

The *augmented Lagrangian method* can be used to solve the optimization problem of eq. (13.30).¹³ Thereby, one solves

$$\max_{\pi} \min_{\lambda \geq 0} J_\mu(\pi; f) - \lambda(J_\mu^c(\pi; f) - \delta) \quad (13.32)$$

$$= \max_{\pi} \begin{cases} J_\mu(\pi; f) & \text{if } \pi \text{ is feasible} \\ -\infty & \text{otherwise.} \end{cases} \quad (13.33)$$

Observe that if π is feasible, then $J_\mu^c(\pi; f) \leq \delta$ and so the minimum over λ is satisfied if $\lambda = 0$. Conversely, if π is infeasible, λ can be made arbitrarily large to solve the optimization problem. In practice, an additional penalty term is added to smooth the objective when transitioning between feasible and infeasible policies.

Note that solving constrained optimization problems such as eq. (13.30) yields an optimal safe policy. However, it is not ensured that constraints are not violated during the search for the optimal policy. Generally, exterior penalty methods such as the augmented Lagrangian method allow for generating infeasible points during the search, and are therefore not suitable when constraints have to be strictly enforced at all times. Thus, this method is more applicable in the episodic setting (e.g. when an agent is first “trained” in a simulated environment and then “deployed” to the actual task) rather than in the continuous setting where the agent has to operate in the “real world” from the beginning and cannot easily be reset.

Remark 13.15: Barrier functions

The augmented Lagrangian method is merely one example of optimizing a joint objective of maximizing rewards and minimizing costs to find a policy with bounded costs. Another example of this approach are so-called *Barrier functions* which augment the reward objective with a smoothed approximation of the boundary

¹¹ It is straightforward to extend this framework to allow for multiple constraints.

¹² This follows from a simple union bound (1.1).

¹³ For an introduction, read chapter 17 of “Numerical optimization” (Wright et al., 1999).

of a set $\mathcal{X}_{\text{unsafe}}$. That is, one solves

$$\max_{\pi} J_{\mu}(\pi; f) - \lambda \cdot B_{\mu}^c(\pi; f) \quad (13.34)$$

for some $\lambda > 0$ where the barrier function $B_{\mu}^c(\pi; f)$ goes to infinity as a state on the boundary of $\mathcal{X}_{\text{unsafe}}$ is approached. Examples for barrier functions are $-\log c(x)$ and $\frac{1}{c(x)}$.

Barrier functions are an interior penalty method which ensure that points are feasible during the search for the optimal solution.

So far, we have assumed knowledge of the true dynamics to solve the optimization problem of eq. (13.30). If we do not know the true dynamics but instead have access to a set of plausible models $\mathcal{M}(\mathcal{D})$ given data \mathcal{D} (cf. section 13.3.2) which encompasses the true dynamics, then a natural strategy is to be **optimistic** with respect to future rewards and **pessimistic** with respect to future constraint violations. More specifically, we solve the optimization problem,

$$\max_{\pi} \max_{f \in \mathcal{M}(\mathcal{D})} J_{\mu}(\pi; f) \quad (13.35a) \quad \text{optimistic}$$

$$\text{subject to } \max_{f \in \mathcal{M}(\mathcal{D})} J_{\mu}^c(\pi; f) \leq \delta. \quad (13.35b) \quad \text{pessimistic}$$

Intuitively, jointly maximizing $J_{\mu}(\pi; f)$ with respect to π and (plausible) f can lead the agent to try behaviors with potentially high reward due to *optimism* (i.e., the agent “dreams” about potential outcomes). On the other hand, being *pessimistic* with respect to constraint violations enforces the safety constraints (i.e., the agent has “nightmares” about potential outcomes).

If the policy values $J_{\mu}(\pi; f)$ and $J_{\mu}^c(\pi; f)$ are modeled as a distribution (e.g., using Bayesian neural networks), then the inner maximizations over plausible dynamics can be approximated using samples from the posterior distributions. Thus, the augmented Lagrangian method can also be used to solve the general optimization problem of eq. (13.35). The resulting algorithm is known as *Lagrangian model-based agent* (LAMBDA) (As et al., 2022).

13.3.4 Safe Exploration

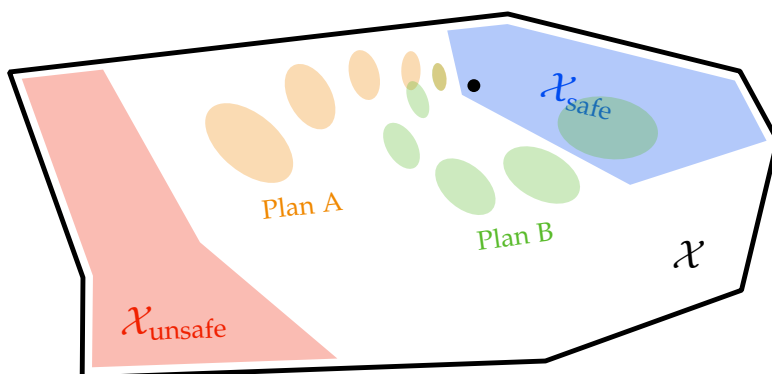
As noted in the previous section, in many settings we do not only want to eventually find a safe policy, but we also want to ensure that we act safely while searching for an optimal policy. To this end, recall the general approach outlined in the beginning of the previous section wherein we plan using (pessimistic) confidence bounds of the plausible consequences of our actions.

One key challenge of this approach is that we need to forecast plausible trajectories. The confidence bounds of such trajectories cannot be nicely represented anymore. One approach is to over-approximate the confidence bounds over reachable states from trajectories.

Theorem 13.16 (Koller et al. (2018)). *For conditional Gaussian dynamics, the reachable states of trajectories can be over-approximated with high probability.*

Another key challenge is that actions might have consequences that exceed the time horizon used for planning. In other words, by performing an action now, our agent might put itself into a state that is not yet unsafe, but out of which it cannot escape and which will eventually lead to an unsafe state. You may think of a car driving towards a wall. When a crash with the wall is designated as an unsafe state, there are quite a few states in advance at which it is already impossible to avoid the crash. Thus, looking ahead a finite number of steps is not sufficient to prevent entering unsafe states.

It turns out that it is still possible to use the epistemic uncertainty about the model and short-term plausible behaviors to make guarantees about certain long-term consequences. One idea is to also consider a set of safe states $\mathcal{X}_{\text{safe}}$ alongside the set of unsafe states $\mathcal{X}_{\text{unsafe}}$, for which our agent knows how to stay inside (i.e., remain safe). In other words, for states $x \in \mathcal{X}_{\text{safe}}$, we know that our agent can always behave in such a way that it will not reach an unsafe state.¹⁴ An illustration of this approach is shown in fig. 13.9.



The problem is that this set of safe states might be very conservative. That is to say, it is likely that rewards are mostly attained *outside* of the set of safe states. The key idea is to plan two sequences of actions, instead of only one. “Plan A” (the *performance plan*) is planned with the objective to solve the task, that is, attain maximum reward. “Plan

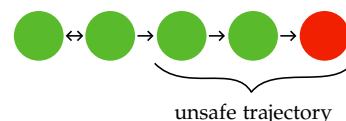


Figure 13.8: Illustration of long-term consequence when planning a finite number of steps. Green dots are to denote safe states and the red dot is to denote an unsafe state. After performing the first action, the agent is still able to return to the previous state. Yet, after reaching the third state, the agent is already guaranteed to end in an unsafe state. When using only a finite horizon of $H = 2$ for planning, the agent might make this transition regardless.

¹⁴ In the example of driving a car, the set of safe states corresponds to those states where we know that we can still safely break before hitting the wall.

Figure 13.9: Illustration of long-term consequence when planning a finite number of steps. The unsafe set of states is shown in red and the safe set of states is shown in blue. The confidence intervals corresponding to actions of the performance plan and safety plan are shown in orange and green, respectively.

B'' (the *safety plan*) is planned with the objective to return to the set of safe states $\mathcal{X}_{\text{safe}}$. In addition, we enforce that both plans must agree on the first action to be played.

Under the assumption that the confidence bounds are conservative estimates, this guarantees that after playing this action, our agent will still be in a state of which it can return to the set of safe states. In this way, we can gradually increase the set of states that are safe to explore. It can be shown that under suitable conditions, returning to the safe set can be guaranteed (Koller et al., 2018).

13.3.5 Safety Filters

An alternative approach is to (slightly) adjust a potentially unsafe policy π to obtain a policy $\hat{\pi}$ which avoids entering unsafe states with high probability.

Following our interpretation of constrained policy optimization in terms of optimism and pessimism from section 13.3.3, to pessimistically evaluate the safety of a policy with respect to a cost function c given a set of plausible models $\mathcal{M}(\mathcal{D})$, we can use

$$C^\pi(x) \doteq \max_{f \in \mathcal{M}(\mathcal{D})} J_{\delta_x}^c(\pi; f) \quad (13.36) \quad \delta_x \text{ is the point density at } x$$

where the initial-state distribution δ_x is to denote that the initial state is x . Observe that eq. (13.36) permits a reparameterization in terms of additional decision variables η which is analogous to our discussion in section 13.3.2. Concretely, we have

$$C^\pi(x) = \max_{\eta} J_{\delta_x}^c(\pi; \hat{f}) \quad (13.37)$$

where \hat{f} are the adjusted dynamics (13.26) which are based on the “luck” variables η . In the context of estimating costs which we aim to minimize (as opposed to rewards which we aim to maximize), η can be interpreted as the “bad luck” of the agent.

When our only objective is to act safely, that is, we only aim to minimize cost and are indifferent to rewards, then this reparameterization allows us to find a “maximally safe” policy,

$$\pi_{\text{safe}} \doteq \arg \min_{\pi} \mathbb{E}_{x \sim \mu} [C^\pi(x)] = \arg \min_{\pi} \max_{\eta} J_{\mu}^c(\pi; \hat{f}). \quad (13.38)$$

Under some conditions π_{safe} can be shown to satisfy $\pi_{\text{safe}}(x) = \pi_x(x)$ for any x where $\pi_x \doteq \arg \min_{\pi} C^\pi(x)$ (Curi et al., 2022, proposition 4.2). On its own, the policy π_{safe} is rather useless for exploring the state space. In particular, when in a state that we already deem safe, following policy π_{safe} , the agent aims simply to “stay” within the set of

safe states which means that it has no incentive to explore / maximize reward.

Instead, one can interpret π_{safe} as a “backup” policy in case we realize upon exploring that a certain trajectory is too dangerous, akin to our notion of the “safety plan” in section 13.3.4. That is, given any (potentially explorative and dangerous) policy π , we can find the adjusted policy,

$$\hat{\pi}(\mathbf{x}) \doteq \arg \min_{a \in \mathcal{A}} d(\pi(\mathbf{x}), a) \quad (13.39a)$$

$$\text{subject to } \max_{\eta} C^{\pi_{\text{safe}}}(\hat{f}(\mathbf{x}, a)) \leq \delta \quad (13.39b)$$

for some metric $d(\cdot, \cdot)$ on \mathcal{A} . The constraint ensures that the pessimistic next state $\hat{f}(\mathbf{x}, a)$ is recoverable by following policy π_{safe} . In this way,

$$\tilde{\pi}(\mathbf{x}) \doteq \begin{cases} \hat{\pi}(\mathbf{x}) & \text{if } C^{\pi_{\text{safe}}}(\mathbf{x}) \leq \delta \\ \pi_{\text{safe}}(\mathbf{x}) & \text{if } C^{\pi_{\text{safe}}}(\mathbf{x}) > \delta \end{cases} \quad (13.40)$$

is the policy “closest” to π which is δ -safe with respect to the pessimistic cost estimates $C^{\pi_{\text{safe}}}$ (Curi et al., 2022).¹⁵ This is also called a *safety filter*. Similar approaches using backup policies also allow safe exploration in the model-free setting (Sukhija et al., 2023; Widmer et al., 2023).

¹⁵ The policy $\tilde{\pi}$ is required as, a priori, it is not guaranteed that the state \mathbf{x}_t will satisfy $C^{\pi_{\text{safe}}}(\mathbf{x}_t) \leq \delta$ for all t unless $\tilde{\pi}$ is replanned after every step.

Optional Readings

- Curi, Berkenkamp, and Krause (2020).
Efficient model-based reinforcement learning through optimistic policy search and planning.
- Berkenkamp, Turchetta, Schoellig, and Krause (2017).
Safe model-based reinforcement learning with stability guarantees.
- Koller, Berkenkamp, Turchetta, and Krause (2018).
Learning-based model predictive control for safe exploration.
- As, Usmanova, Curi, and Krause (2022).
Constrained Policy Optimization via Bayesian World Models.
- Curi, Lederer, Hirche, and Krause (2022).
Safe Reinforcement Learning via Confidence-Based Filters.
- Turchetta, Berkenkamp, and Krause (2019).
Safe exploration for interactive machine learning.

A

Mathematical Background

A.1 Useful Matrix Identities and Inequalities

- *Woodbury's matrix identity* states that for any matrices $A \in \mathbb{R}^{n \times n}$, $C \in \mathbb{R}^{m \times m}$, and $U, V^\top \in \mathbb{R}^{n \times m}$ where A and C are invertible,

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}. \quad (\text{A.1})$$

The following identity is a direct consequence:

$$(A + xx^\top)^{-1} = A^{-1} - A^{-1}x(1 + x^\top A^{-1}x)^{-1}x^\top A^{-1} \quad (\text{A.2a})$$

$$= A^{-1} - \frac{(A^{-1}x)(A^{-1}x)^\top}{1 + x^\top A^{-1}x} \quad (\text{A.2b})$$

for any symmetric and invertible matrix $A \in \mathbb{R}^{n \times n}$ and $x \in \mathbb{R}^n$.

- The *matrix inversion lemma* states that for matrices $A, B \in \mathbb{R}^{n \times n}$,

$$(A + B)^{-1} = A^{-1} - A^{-1}(A^{-1} + B^{-1})^{-1}A^{-1}. \quad (\text{A.3})$$

- *Hadamard's inequality* states that

$$\det(M) \leq \prod_{i=0}^d M(i, i) \quad (\text{A.4})$$

for any positive definite matrix $M \in \mathbb{R}^{d \times d}$.

- The *Weinstein-Aronszajn identity* for positive definite matrices $A \in \mathbb{R}^{d \times n}$ and $B \in \mathbb{R}^{n \times d}$ states that

$$\det(I_{d \times d} + AB) = \det(I_{n \times n} + BA). \quad (\text{A.5})$$

B

Solutions

Fundamentals

Solution to exercise 1.5 (Properties of probability).

1. By the third axiom and $B = A \cup (B \setminus A)$,

$$\mathbb{P}(B) = \mathbb{P}(A) + \mathbb{P}(B \setminus A).$$

Noting from the first axiom that $\mathbb{P}(B \setminus A) \geq 0$ completes the proof.

2. By the second axiom,

$$\mathbb{P}(A \cup \overline{A}) = \mathbb{P}(\Omega) = 1$$

and by the third axiom,

$$\mathbb{P}(A \cup \overline{A}) = \mathbb{P}(A) + \mathbb{P}(\overline{A}).$$

Reorganizing the equations completes the proof.

3. Define the countable sequence of events

$$A'_i \doteq A_i \setminus \left(\bigcup_{j=1}^{i-1} A_j \right).$$

Note that the sequence of events $\{A'_i\}_i$ is disjoint. Thus, we have by the third axiom and then using (1) that

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

Solution to exercise 1.19 (Random walks on graphs). We show that any vertex v is visited eventually with probability 1.

We denote by $w \rightarrow v$ the event that the random walk starting at vertex w visits the vertex v eventually, we denote by $\Gamma(w)$ the neighborhood

of w , and we write $\deg(w) \doteq |\Gamma(w)|$. We have,

$$\begin{aligned}\mathbb{P}(w \rightarrow v) &= \sum_{w' \in \Gamma(w)} \mathbb{P}(\text{the random walk first visits } w') \cdot \mathbb{P}(w' \rightarrow v) && \text{using the law of total probability (1.15)} \\ &= \frac{1}{\deg(w)} \sum_{w' \in \Gamma(w)} \mathbb{P}(w' \rightarrow v). && \text{using that the random walk moves to a} \\ &&& \text{neighbor uniformly at random}\end{aligned}$$

Take u to be the vertex such that $\mathbb{P}(u \rightarrow v)$ is minimized. Then,

$$\mathbb{P}(u \rightarrow v) = \frac{1}{\deg(u)} \sum_{u' \in \Gamma(u)} \underbrace{\mathbb{P}(u' \rightarrow v)}_{\geq \mathbb{P}(u \rightarrow v)} \geq \mathbb{P}(u \rightarrow v).$$

That is, for all neighbors u' of u , $\mathbb{P}(u \rightarrow v) = \mathbb{P}(u' \rightarrow v)$. Using that the graph is connected and finite, we conclude $\mathbb{P}(u \rightarrow v) = \mathbb{P}(w \rightarrow v)$ for any vertex w . Finally, note that $\mathbb{P}(v \rightarrow v) = 1$, and hence, the random walk starting at any vertex u visits the vertex v eventually with probability 1.

Solution to exercise 1.26 (Law of total expectation). Let $\mathbb{1}\{A_i\}$ be the indicator random variable for the event A_i . Then,

$$\begin{aligned}\mathbb{E}[\mathbf{X} \cdot \mathbb{1}\{A_i\}] &= \mathbb{E}[\mathbb{E}[\mathbf{X} \cdot \mathbb{1}\{A_i\} \mid \mathbb{1}\{A_i\}]] && \text{by the tower rule (1.33)} \\ &= \mathbb{E}[\mathbf{X} \cdot \mathbb{1}\{A_i\} \mid \mathbb{1}\{A_i\} = 1] \cdot \mathbb{P}(\mathbb{1}\{A_i\} = 1) && \text{expanding the outer expectation} \\ &\quad + 0 \cdot \mathbb{P}(\mathbb{1}\{A_i\} = 0) \\ &= \mathbb{E}[\mathbf{X} \mid A_i] \cdot \mathbb{P}(A_i). && \text{the event } \mathbb{1}\{A_i\} = 1 \text{ is equivalent to } A_i\end{aligned}$$

Summing up for all i , the left-hand side becomes

$$\mathbb{E}[\mathbf{X}] = \sum_{i=1}^k \mathbb{E}[\mathbf{X} \cdot \mathbb{1}\{A_i\}].$$

Solution to exercise 1.28 (Covariance matrices are positive semi-definite). Let $\Sigma \doteq \text{Var}[\mathbf{X}]$ be a covariance matrix of the random vector \mathbf{X} and fix any $\mathbf{z} \in \mathbb{R}^n$. Then,

$$\begin{aligned}\mathbf{z}^\top \Sigma \mathbf{z} &= \mathbf{z}^\top \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top] \mathbf{z} && \text{using the definition of variance (1.44)} \\ &= \mathbb{E}[\mathbf{z}^\top (\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top \mathbf{z}]. && \text{using linearity of expectation (1.24)}\end{aligned}$$

Define the random variable $Z \doteq \mathbf{z}^\top (\mathbf{X} - \mathbb{E}[\mathbf{X}])$. Then,

$$= \mathbb{E}[Z^2] \geq 0.$$

Solution to exercise 1.33 (Bayes' rule). Let us start by defining some events that we will reason about later. For ease of writing, let us call the person in question X .

$D = X$ has the disease,

P = The test shows a positive response.

Now we can translate the information in the question to formal statements in terms of D and P ,

$$\begin{aligned}\mathbb{P}(D) &= 10^{-4} \\ \mathbb{P}(P \mid D) &= \mathbb{P}(\bar{P} \mid \bar{D}) = 0.99.\end{aligned}$$

the disease is rare

the test is accurate

We want to determine $\mathbb{P}(D \mid P)$. One can find this probability by using Bayes' rule (1.59),

$$\mathbb{P}(D \mid P) = \frac{\mathbb{P}(P \mid D) \cdot \mathbb{P}(D)}{\mathbb{P}(P)}.$$

From the quantities above, we have everything except for $\mathbb{P}(P)$. This, however, we can compute using the law of total probability,

$$\begin{aligned}\mathbb{P}(P) &= \mathbb{P}(P \mid D) \cdot \mathbb{P}(D) + \mathbb{P}(P \mid \bar{D}) \cdot \mathbb{P}(\bar{D}) \\ &= 0.99 \cdot 10^{-4} + 0.01 \cdot (1 - 10^{-4}) \\ &= 0.010098.\end{aligned}$$

Hence, $\mathbb{P}(D \mid P) = 0.99 \cdot 10^{-4} / 0.010098 \approx 0.0098 = 0.98\%$.

Solution to exercise 1.38 (Sample variance). Using that X is zero-mean, we have that

$$\begin{aligned}\mathbb{E}[\bar{X}_n^2] &= \text{Var}[\bar{X}_n] = \frac{1}{n} \text{Var}[X] \quad \text{and} \\ \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X^{(i)2}\right] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X^{(i)2}] = \text{Var}[X].\end{aligned}$$

Thus,

$$\mathbb{E}[S_n^2] = \frac{n}{n-1} \left(\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X^{(i)2}\right] - \mathbb{E}[\bar{X}_n^2] \right) = \text{Var}[X].$$

Solution to exercise 1.43 (Simple concentration inequalities).

1. W.l.o.g. we assume that X is continuous. We have

$$\begin{aligned}\mathbb{E}X &= \int_0^\infty x \cdot p(x) dx \\ &= \underbrace{\int_0^\epsilon x \cdot p(x) dx}_{\geq 0} + \int_\epsilon^\infty x \cdot p(x) dx \geq \epsilon \underbrace{\int_\epsilon^\infty p(x) dx}_{\mathbb{P}(X \geq \epsilon)}.\end{aligned}$$

using the definition of expectation (1.23b)

2. Consider the non-negative random variable $Y \doteq (X - \mathbb{E}X)^2$. We have

$$\begin{aligned}\mathbb{P}(|X - \mathbb{E}X| \geq \epsilon) &= \mathbb{P}((X - \mathbb{E}X)^2 \geq \epsilon^2) \\ &\leq \frac{\mathbb{E}[(X - \mathbb{E}X)^2]}{\epsilon^2} \\ &= \frac{\text{Var}X}{\epsilon^2}.\end{aligned}$$

using Markov's inequality (1.82)

using the definition of variance (1.44)

Solution to exercise 1.45 (Weak law of large numbers). Fix any $\epsilon > 0$. Applying Chebyshev's inequality and noting that $\mathbb{E}\bar{X}_n = \mathbb{E}X$, we obtain

$$\mathbb{P}(|\bar{X}_n - \mathbb{E}X| \geq \epsilon) \leq \frac{\text{Var}\bar{X}_n}{\epsilon^2}.$$

We further have for the variance of the sample mean that

$$\text{Var}\bar{X}_n = \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] = \frac{\text{Var}X}{n}.$$

Thus,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mathbb{E}X| \geq \epsilon) \leq \lim_{n \rightarrow \infty} \frac{\text{Var}X}{\epsilon^2 n} = 0$$

which is precisely the definition of $\bar{X}_n \xrightarrow{\mathbb{P}} \mathbb{E}X$.

Solution to exercise 1.55 (Directional derivatives). Fix a $\lambda > 0$. Using a first-order expansion, we have

$$f(x + \lambda d) = f(x) + \lambda \nabla f(x)^\top d + o(\lambda \|d\|_2).$$

Dividing by λ yields,

$$\frac{f(x + \lambda d) - f(x)}{\lambda} = \nabla f(x)^\top d + \underbrace{\frac{o(\lambda \|d\|_2)}{\lambda}}_{\rightarrow 0}.$$

Taking the limit $\lambda \rightarrow 0$ gives the desired result.

Solution to exercise 1.65 (Product of Gaussian PDFs). We need to find $\mu \in \mathbb{R}^n$ and $\Sigma \in \mathbb{R}^{n \times n}$ such that

$$(x - \mu)^\top \Sigma^{-1} (x - \mu) \propto (x - \mu_1)^\top \Sigma_1^{-1} (x - \mu_1) + (x - \mu_2)^\top \Sigma_2^{-1} (x - \mu_2). \quad (\text{B.1})$$

for any $x \in \mathbb{R}^n$.

The left-hand side of eq. (B.1) is equal to

$$x^\top \Sigma^{-1} x - 2x^\top \Sigma^{-1} \mu + \mu^\top \Sigma^{-1} \mu.$$

The right-hand side of eq. (B.1) is equal to

$$\begin{aligned} & \left(x^\top \Sigma_1^{-1} x + x^\top \Sigma_2^{-1} x \right) - 2 \left(x^\top \Sigma_1^{-1} \mu_1 + x^\top \Sigma_2^{-1} \mu_2 \right) \\ & \quad + \left(\mu_1^\top \Sigma_1^{-1} \mu_1 + \mu_2^\top \Sigma_2^{-1} \mu_2 \right) \\ & = x^\top \left(\Sigma_1^{-1} + \Sigma_2^{-1} \right) x - 2x^\top \left(\Sigma_1^{-1} \mu_1 + \Sigma_2^{-1} \mu_2 \right) \\ & \quad + \left(\mu_1^\top \Sigma_1^{-1} \mu_1 + \mu_2^\top \Sigma_2^{-1} \mu_2 \right). \end{aligned}$$

We observe that both sides are equal if

$$\Sigma^{-1} = \Sigma_1^{-1} + \Sigma_2^{-1} \quad \text{and} \quad \Sigma^{-1} \mu = \Sigma_1^{-1} \mu_1 + \Sigma_2^{-1} \mu_2.$$

Solution to exercise 1.67 (Marginal / conditional distribution of a Gaussian). Let $\mathbf{x} \sim \mathbf{X}$. The joint distribution can be expressed as

$$p(\mathbf{x}) = p(\mathbf{x}_A, \mathbf{x}_B) \\ = \frac{1}{Z} \exp \left(-\frac{1}{2} \begin{bmatrix} \mathbf{x}_A - \boldsymbol{\mu}_A \\ \mathbf{x}_B - \boldsymbol{\mu}_B \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\Lambda}_{AA} & \boldsymbol{\Lambda}_{AB} \\ \boldsymbol{\Lambda}_{BA} & \boldsymbol{\Lambda}_{BB} \end{bmatrix} \begin{bmatrix} \mathbf{x}_A - \boldsymbol{\mu}_A \\ \mathbf{x}_B - \boldsymbol{\mu}_B \end{bmatrix} \right)$$

where Z denotes the normalizing constant. To simplify the notation, we write

$$\begin{bmatrix} \boldsymbol{\Delta}_A \\ \boldsymbol{\Delta}_B \end{bmatrix} \doteq \begin{bmatrix} \mathbf{x}_A - \boldsymbol{\mu}_A \\ \mathbf{x}_B - \boldsymbol{\mu}_B \end{bmatrix}.$$

1. We obtain

$$\begin{aligned} p(\mathbf{x}_A) &= \frac{1}{Z} \int \exp \left(-\frac{1}{2} \begin{bmatrix} \boldsymbol{\Delta}_A \\ \boldsymbol{\Delta}_B \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\Lambda}_{AA} & \boldsymbol{\Lambda}_{AB} \\ \boldsymbol{\Lambda}_{BA} & \boldsymbol{\Lambda}_{BB} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Delta}_A \\ \boldsymbol{\Delta}_B \end{bmatrix} \right) d\mathbf{x}_B && \text{using the sum rule (1.10)} \\ &= \frac{1}{Z} \exp \left(-\frac{1}{2} \left[\boldsymbol{\Delta}_A^\top (\boldsymbol{\Lambda}_{AA} - \boldsymbol{\Lambda}_{AB} \boldsymbol{\Lambda}_{BB}^{-1} \boldsymbol{\Lambda}_{BA}) \boldsymbol{\Delta}_A \right] \right) && \text{using the first hint} \\ &\quad \cdot \int \exp \left(-\frac{1}{2} \left[(\boldsymbol{\Delta}_B + \boldsymbol{\Lambda}_{BB}^{-1} \boldsymbol{\Lambda}_{BA} \boldsymbol{\Delta}_A)^\top \boldsymbol{\Lambda}_{BB} (\boldsymbol{\Delta}_B + \boldsymbol{\Lambda}_{BB}^{-1} \boldsymbol{\Lambda}_{BA} \boldsymbol{\Delta}_A) \right] \right) d\mathbf{x}_B. \end{aligned}$$

Observe that the integrand is an unnormalized Gaussian PDF, and hence, the integral evaluates to $\sqrt{\det(2\pi\boldsymbol{\Lambda}_{BB}^{-1})}$ (so is constant with respect to \mathbf{x}_A). We can therefore simplify to

$$\begin{aligned} &= \frac{1}{Z'} \exp \left(-\frac{1}{2} \left[\boldsymbol{\Delta}_A^\top (\boldsymbol{\Lambda}_{AA} - \boldsymbol{\Lambda}_{AB} \boldsymbol{\Lambda}_{BB}^{-1} \boldsymbol{\Lambda}_{BA}) \boldsymbol{\Delta}_A \right] \right) \\ &= \frac{1}{Z'} \exp \left(-\frac{1}{2} \left[\boldsymbol{\Delta}_A^\top \boldsymbol{\Sigma}_{AA}^{-1} \boldsymbol{\Delta}_A \right] \right). \end{aligned} \quad \text{using the second hint}$$

2. We obtain

$$\begin{aligned} p(\mathbf{x}_B | \mathbf{x}_A) &= \frac{p(\mathbf{x}_A, \mathbf{x}_B)}{p(\mathbf{x}_A)} && \text{using the definition of conditional distributions (1.13)} \\ &= \frac{1}{Z'} \exp \left(-\frac{1}{2} \begin{bmatrix} \boldsymbol{\Delta}_A \\ \boldsymbol{\Delta}_B \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\Lambda}_{AA} & \boldsymbol{\Lambda}_{AB} \\ \boldsymbol{\Lambda}_{BA} & \boldsymbol{\Lambda}_{BB} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Delta}_A \\ \boldsymbol{\Delta}_B \end{bmatrix} \right) && \text{noting that } p(\mathbf{x}_A) \text{ is constant with respect to } \mathbf{x}_B \\ &= \frac{1}{Z'} \exp \left(-\frac{1}{2} \left[\boldsymbol{\Delta}_A^\top (\boldsymbol{\Lambda}_{AA} - \boldsymbol{\Lambda}_{AB} \boldsymbol{\Lambda}_{BB}^{-1} \boldsymbol{\Lambda}_{BA}) \boldsymbol{\Delta}_A \right] \right) && \text{using the first hint} \\ &\quad \cdot \exp \left(-\frac{1}{2} \left[(\boldsymbol{\Delta}_B + \boldsymbol{\Lambda}_{BB}^{-1} \boldsymbol{\Lambda}_{BA} \boldsymbol{\Delta}_A)^\top \boldsymbol{\Lambda}_{BB} (\boldsymbol{\Delta}_B + \boldsymbol{\Lambda}_{BB}^{-1} \boldsymbol{\Lambda}_{BA} \boldsymbol{\Delta}_A) \right] \right) \\ &= \frac{1}{Z''} \exp \left(-\frac{1}{2} \left[(\boldsymbol{\Delta}_B + \boldsymbol{\Lambda}_{BB}^{-1} \boldsymbol{\Lambda}_{BA} \boldsymbol{\Delta}_A)^\top \boldsymbol{\Lambda}_{BB} (\boldsymbol{\Delta}_B + \boldsymbol{\Lambda}_{BB}^{-1} \boldsymbol{\Lambda}_{BA} \boldsymbol{\Delta}_A) \right] \right). \end{aligned} \quad \text{observing that the first exponential is constant with respect to } \mathbf{x}_B$$

Finally, observe that

$$\Delta_B + \Lambda_{BB}^{-1} \Lambda_{BA} \Delta_A = \mu_B - \Lambda_{BB}^{-1} \Lambda_{BA} (x_A - \mu_A) = \mu_{B|A} \quad \text{and}$$

using the third hint

$$\Lambda_{BB}^{-1} = \Sigma_{BB} - \Sigma_{BA} \Sigma_{AA}^{-1} \Sigma_{AB} = \Sigma_{B|A}.$$

using the second hint

Thus, $x_B \mid x_A \sim \mathcal{N}(\mu_{B|A}, \Sigma_{B|A})$.

Solution to exercise 1.69 (Closedness properties of Gaussians).

1. Let $\mathbf{Y} \doteq \mathbf{A}\mathbf{X} + \mathbf{b}$ and define $\mathbf{s} \doteq \mathbf{A}^\top \mathbf{t}$. We have for the MGF of \mathbf{Y} that

$$\begin{aligned} \varphi_{\mathbf{Y}}(\mathbf{t}) &= \mathbb{E} \exp(\mathbf{t}^\top \mathbf{Y}) \\ &= \mathbb{E} \exp(\mathbf{t}^\top \mathbf{A}\mathbf{X}) \cdot \exp(\mathbf{t}^\top \mathbf{b}) \\ &= \mathbb{E} \exp(\mathbf{s}^\top \mathbf{X}) \cdot \exp(\mathbf{t}^\top \mathbf{b}) \\ &= \varphi_{\mathbf{X}}(\mathbf{s}) \cdot \exp(\mathbf{t}^\top \mathbf{b}) \\ &= \exp\left(\mathbf{s}^\top \boldsymbol{\mu} + \frac{1}{2} \mathbf{s}^\top \boldsymbol{\Sigma} \mathbf{s}\right) \cdot \exp(\mathbf{t}^\top \mathbf{b}) \\ &= \exp\left(\mathbf{t}^\top (\mathbf{A}\boldsymbol{\mu} + \mathbf{b}) + \frac{1}{2} \mathbf{t}^\top \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top \mathbf{t}\right), \end{aligned}$$

which implies that $\mathbf{Y} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$.

2. We have for the MGF of $\mathbf{Y} \doteq \mathbf{X} + \mathbf{X}'$ that

$$\begin{aligned} \varphi_{\mathbf{Y}}(\mathbf{t}) &= \mathbb{E} \exp(\mathbf{t}^\top \mathbf{Y}) \\ &= \mathbb{E} \exp(\mathbf{t}^\top \mathbf{X}) \cdot \mathbb{E} \exp(\mathbf{t}^\top \mathbf{X}') \\ &= \varphi_{\mathbf{X}}(\mathbf{t}) \cdot \varphi_{\mathbf{X}'}(\mathbf{t}) \\ &= \exp\left(\mathbf{t}^\top \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t}\right) \cdot \exp\left(\mathbf{t}^\top \boldsymbol{\mu}' + \frac{1}{2} \mathbf{t}^\top \boldsymbol{\Sigma}' \mathbf{t}\right) \\ &= \exp\left(\mathbf{t}^\top (\boldsymbol{\mu} + \boldsymbol{\mu}') + \frac{1}{2} \mathbf{t}^\top (\boldsymbol{\Sigma} + \boldsymbol{\Sigma}') \mathbf{t}\right), \end{aligned}$$

using the independence of \mathbf{X} and \mathbf{X}'

which implies that $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu} + \boldsymbol{\mu}', \boldsymbol{\Sigma} + \boldsymbol{\Sigma}')$.

Solution to exercise 1.72 (Expectation and Variance of Gaussians).

Note that if $Y \sim \mathcal{N}(0, 1)$ is a (univariate) standard normal random variable, then

$$\begin{aligned} \mathbb{E}[Y] &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} y \cdot \exp\left(-\frac{1}{2}y^2\right) dy \\ &= \frac{1}{\sqrt{2\pi}} \left(\int_{-\infty}^0 y \cdot \exp\left(-\frac{1}{2}y^2\right) dy + \int_0^{\infty} y \cdot \exp\left(-\frac{1}{2}y^2\right) dy \right) \\ &= \frac{1}{\sqrt{2\pi}} \left(\left[-\exp\left(-\frac{1}{2}y^2\right) \right]_{-\infty}^0 + \left[-\exp\left(-\frac{1}{2}y^2\right) \right]_0^{\infty} \right) \end{aligned}$$

using the PDF of the standard normal distribution (1.6)

$$\begin{aligned}
&= \frac{1}{\sqrt{2\pi}}([-1+0] + [0+1]) = 0, \\
\text{Var}[Y] &= \mathbb{E}[Y^2] - \underbrace{\mathbb{E}[Y]^2}_0 \\
&= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} y^2 \cdot \exp\left(-\frac{1}{2}y^2\right) dy && \text{using the definition of variance (1.45)} \\
&= \frac{1}{\sqrt{2\pi}} \left(\int_{-\infty}^0 y^2 \cdot \exp\left(-\frac{1}{2}y^2\right) dy + \int_0^{\infty} y^2 \cdot \exp\left(-\frac{1}{2}y^2\right) dy \right) && \text{using the PDF of the standard normal distribution (1.6)} \\
&= \frac{1}{\sqrt{2\pi}} \left(\left[-y \cdot \exp\left(-\frac{1}{2}y^2\right) \right]_{-\infty}^0 + \int_{-\infty}^0 \exp\left(-\frac{1}{2}y^2\right) dy \right. \\
&\quad \left. + \left[-y \cdot \exp\left(-\frac{1}{2}y^2\right) \right]_0^{\infty} + \int_0^{\infty} \exp\left(-\frac{1}{2}y^2\right) dy \right) && \text{integrating by parts} \\
&= \frac{1}{\sqrt{2\pi}} \left(\int_{-\infty}^0 \exp\left(-\frac{1}{2}y^2\right) dy + \int_0^{\infty} \exp\left(-\frac{1}{2}y^2\right) dy \right) \\
&= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2}y^2\right) dy = 1. && \text{a PDF integrates to 1}
\end{aligned}$$

Recall from eq. (1.125) that we can express $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ as

$$\mathbf{X} = \boldsymbol{\Sigma}^{1/2} \mathbf{Y} + \boldsymbol{\mu}$$

where $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Using that \mathbf{Y} is a vector of independent univariate standard normal random variables, we conclude that $\mathbb{E}[\mathbf{Y}] = \mathbf{0}$ and $\text{Var}[\mathbf{Y}] = \mathbf{I}$. We obtain

$$\begin{aligned}
\mathbb{E}[\mathbf{X}] &= \mathbb{E}[\boldsymbol{\Sigma}^{1/2} \mathbf{Y} + \boldsymbol{\mu}] = \boldsymbol{\Sigma}^{1/2} \underbrace{\mathbb{E}[\mathbf{Y}]}_{\mathbf{0}} + \boldsymbol{\mu} = \boldsymbol{\mu}, \quad \text{and} && \text{using linearity of expectation (1.24)} \\
\text{Var}[\mathbf{X}] &= \text{Var}[\boldsymbol{\Sigma}^{1/2} \mathbf{Y} + \boldsymbol{\mu}] = \boldsymbol{\Sigma}^{1/2} \underbrace{\text{Var}[\mathbf{Y}]}_{\mathbf{I}} \boldsymbol{\Sigma}^{1/2 \top} = \boldsymbol{\Sigma}. && \text{using eq. (1.47)}
\end{aligned}$$

Bayesian Linear Regression

Solution to exercise 2.2 (Closed-form linear regression).

1. We begin by deriving the gradient and Hessian of the least squares and ridge losses. For least squares,

$$\begin{aligned}
\nabla_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 &= \nabla_{\mathbf{w}} (\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \|\mathbf{y}\|_2^2) \\
&= 2(\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y}), \\
\mathbf{H}_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 &= 2\mathbf{X}^\top \mathbf{X}.
\end{aligned}$$

Similarly, for ridge regression,

$$\begin{aligned}
\nabla_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 &= 2(\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y}) + \lambda \mathbf{w}, \\
\mathbf{H}_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 &= 2\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}.
\end{aligned}$$

From the assumption that the Hessians are positive definite, we know that any minimizer is a unique globally optimal solution (due to strict convexity), and that $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}$ are invertible.

Using the first-order optimality condition for convex functions, we attain the solutions to least squares and ridge regression by setting the respective gradient to $\mathbf{0}$.

2. We choose $\hat{\mathbf{w}}_{\text{ls}}$ such that $\mathbf{X}\hat{\mathbf{w}}_{\text{ls}} = \Pi_{\mathbf{X}}\mathbf{y}$. This implies that $\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}_{\text{ls}} \perp \mathbf{X}\mathbf{w}$ for all $\mathbf{w} \in \mathbb{R}^d$. In other words, $(\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}_{\text{ls}})^\top \mathbf{X}\mathbf{w} = 0$ for all \mathbf{w} , which implies that $(\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}_{\text{ls}})^\top \mathbf{X} = \mathbf{0}$. By simple algebraic manipulation it can be seen that this condition is equivalent to the gradient condition of the previous exercise.

Solution to exercise 2.4 (Variance of least squares around training data). In the two-dimensional setting, i.e., data is of the form $\mathbf{x}_i = [1 \ x_i]^\top$ ($x_i \in \mathbb{R}$), we have

$$\mathbf{X}^\top \mathbf{X} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}.$$

Thus,

$$\text{Var}[\hat{\mathbf{w}}_{\text{ls}}] = \sigma_n^2 (\mathbf{X}^\top \mathbf{X})^{-1} = \frac{\sigma_n^2}{Z} \begin{bmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{bmatrix} \quad \text{using eq. (2.9)}$$

where $Z \doteq n(\sum x_i^2) - (\sum x_i)^2$.

Therefore, the predictive variance at a point $[1 \ x^*]^\top$ is

$$\begin{bmatrix} 1 & x^* \end{bmatrix} \text{Var}[\hat{\mathbf{w}}_{\text{ls}}] \begin{bmatrix} 1 \\ x^* \end{bmatrix} = \frac{\sigma_n^2}{Z} \sum x_i^2 - 2x_i x^* + (x^*)^2 = \frac{\sigma_n^2}{Z} \sum (x_i - x^*)^2.$$

Thus, the predictive variance is minimized for $x^* = \frac{1}{n} \sum x_i$.

Solution to exercise 2.8 (Aleatoric and epistemic uncertainty of BLR).

The law of total variance (2.22) yields the following decomposition of the predictive variance,

$$\begin{aligned} \text{Var}_{y^*}[y^* \mid \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}] &= \mathbb{E}_{\mathbf{w}}[\text{Var}_{y^*}[y^* \mid \mathbf{x}^*, \mathbf{w}] \mid \mathbf{x}_{1:n}, y_{1:n}] \\ &\quad + \text{Var}_{\mathbf{w}}[\mathbb{E}_{y^*}[y^* \mid \mathbf{x}^*, \mathbf{w}] \mid \mathbf{x}_{1:n}, y_{1:n}] \end{aligned}$$

wherein the first term corresponds to the **aleatoric uncertainty** and the second term corresponds to the **epistemic uncertainty**.

The aleatoric uncertainty is given by

$$\mathbb{E}_{\mathbf{w}}[\text{Var}_{y^*}[y^* \mid \mathbf{x}^*, \mathbf{w}] \mid \mathbf{x}_{1:n}, y_{1:n}] = \mathbb{E}_{\mathbf{w}}[\sigma_n^2 \mid \mathbf{x}_{1:n}, y_{1:n}] = \sigma_n^2 \quad \text{using the definition of } \sigma_n^2 \text{ (2.7)}$$

For the epistemic uncertainty,

$$\begin{aligned}\text{Var}_w[\mathbb{E}_{y^*}[y^* \mid \mathbf{x}^*, \mathbf{w}] \mid \mathbf{x}_{1:n}, y_{1:n}] &= \text{Var}_w[\mathbf{w}^\top \mathbf{x}^* \mid \mathbf{x}_{1:n}, y_{1:n}] && \text{using that } y^* = \mathbf{w}^\top \mathbf{x}^* + \epsilon \text{ where } \epsilon \text{ is} \\ &= \mathbf{x}^{*\top} \text{Var}_w[\mathbf{w} \mid \mathbf{x}_{1:n}, y_{1:n}] \mathbf{x}^* && \text{zero-mean noise} \\ &= \mathbf{x}^{*\top} \boldsymbol{\Sigma} \mathbf{x}^* && \text{using eq. (1.47)} \\ & && \text{using eq. (2.13)}\end{aligned}$$

where $\boldsymbol{\Sigma}$ is the posterior covariance matrix.

Solution to exercise 2.7 (Bayesian linear regression).

1. Recall from section 2.1.1 that the MLE and least squares estimate coincide if the noise is zero-mean Gaussian. We therefore have,

$$\hat{\mathbf{w}}_{\text{MLE}} = \hat{\mathbf{w}}_{\text{ls}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \begin{bmatrix} 0.63 \\ 1.83 \end{bmatrix}.$$

2. Recall from eq. (2.13) that $\mathbf{w} \mid \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with

$$\begin{aligned}\boldsymbol{\Sigma} &= (\sigma_n^{-2} \mathbf{X}^\top \mathbf{X} + \sigma_p^{-2} \mathbf{I})^{-1} \quad \text{and} \\ \boldsymbol{\mu} &= \sigma_n^{-2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y}.\end{aligned}$$

Inserting $\sigma_n^2 = 0.1$ and $\sigma_p^2 = 0.05$ yields,

$$\boldsymbol{\Sigma} = \begin{bmatrix} 0.019 & -0.014 \\ -0.014 & 0.019 \end{bmatrix}.$$

Then,

$$\hat{\mathbf{w}}_{\text{MAP}} = \boldsymbol{\mu} = \begin{bmatrix} 0.91 \\ 1.31 \end{bmatrix}.$$

3. Recall from eq. (2.21) that $y^* \mid \mathbf{X}, \mathbf{y}, \mathbf{x}^* \sim \mathcal{N}(\mu_{y^*}^*, \sigma_{y^*}^2)$ with

$$\mu_{y^*}^* = \boldsymbol{\mu}^\top \mathbf{x}^* \quad \text{and} \quad \sigma_{y^*}^2 = \mathbf{x}^{*\top} \boldsymbol{\Sigma} \mathbf{x}^* + \sigma_n^2.$$

Inserting $\mathbf{x}^* = [3 \ 3]^\top$, and $\sigma_n^2, \boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ from above yields,

$$\mu_{y^*}^* = 6.66 \quad \text{and} \quad \sigma_{y^*}^2 = 0.186.$$

4. One would have to let $\sigma_p^2 \rightarrow \infty$.

Solution to exercise 2.9 (Online Bayesian linear regression). We denote by \mathbf{X}_t the design matrix and by \mathbf{y}_t the vector of observations including the first t data points.

1. Note that

$$\mathbf{X}_t^\top \mathbf{X}_t = \sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^\top \quad \text{and} \quad \mathbf{X}_t^\top \mathbf{y}_t = \sum_{i=1}^t y_i \mathbf{x}_i.$$

This means that after observing the $(t+1)$ -st data point, we have that

$$\begin{aligned}\mathbf{X}_{t+1}^\top \mathbf{X}_{t+1} &= \mathbf{X}_t^\top \mathbf{X}_t + \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \quad \text{and} \\ \mathbf{X}_{t+1}^\top \mathbf{y}_{t+1} &= \mathbf{X}_t^\top \mathbf{y}_t + y_{t+1} \mathbf{x}_{t+1}.\end{aligned}$$

Hence, by just keeping $\mathbf{X}_t^\top \mathbf{X}_t$ (which is a $d \times d$ matrix) and $\mathbf{X}_t^\top \mathbf{y}_t$ (which is a vector in \mathbb{R}^d) in memory, and updating them as above, we do not need to keep the whole data in memory.

2. One has to compute $(\sigma_n^{-2} \mathbf{X}_t^\top \mathbf{X}_t + \sigma_p^{-2} \mathbf{I})^{-1}$ for finding $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ in every round. We can write

$$\begin{aligned}(\sigma_n^{-2} \mathbf{X}_{t+1}^\top \mathbf{X}_{t+1} + \sigma_p^{-2} \mathbf{I})^{-1} &= \sigma_n^2 (\mathbf{X}_{t+1}^\top \mathbf{X}_{t+1} + \sigma_n^2 \sigma_p^{-2} \mathbf{I})^{-1} \\ &= \sigma_n^2 (\underbrace{\mathbf{X}_t^\top \mathbf{X}_t + \sigma_n^2 \sigma_p^{-2} \mathbf{I}}_{\mathbf{A}_t} + \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top)^{-1}\end{aligned}$$

where $\mathbf{A}_t \in \mathbb{R}^{d \times d}$. Using the Woodbury matrix identity (A.2) and that we know the inverse of \mathbf{A}_t (from the previous iteration), the computation of the inverse of $(\mathbf{A}_t + \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top)$ is of $\mathcal{O}(d^2)$, which is much better than computing the inverse of $(\mathbf{A}_t + \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top)$ from scratch.

Kalman Filters

Solution to exercise 3.6 (Kalman update). Recall from eq. (3.10) that

$$p(x_{t+1} | y_{1:t+1}) = \frac{1}{Z} p(x_{t+1} | y_{1:t}) p(y_{t+1} | x_{t+1}). \quad (\text{B.2})$$

Using the sensor model (3.14b),

$$p(y_{t+1} | x_{t+1}) = \frac{1}{Z'} \exp\left(-\frac{1}{2} \frac{(y_{t+1} - x_{t+1})^2}{\sigma_y^2}\right).$$

It remains to compute the predictive distribution,

$$\begin{aligned}p(x_{t+1} | y_{1:t}) &= \int p(x_{t+1} | x_t) p(x_t | y_{1:t}) dx_t \\ &= \frac{1}{Z''} \int \exp\left(-\frac{1}{2} \left[\frac{(x_{t+1} - x_t)^2}{\sigma_x^2} + \frac{(x_t - \mu_t)^2}{\sigma_t^2} \right]\right) dx_t \\ &= \frac{1}{Z''} \int \exp\left(-\frac{1}{2} \left[\frac{\sigma_t^2 (x_{t+1} - x_t)^2 + \sigma_x^2 (x_t - \mu_t)^2}{\sigma_t^2 \sigma_x^2} \right]\right) dx_t.\end{aligned}$$

using eq. (3.11)

using the motion model (3.14a) and previous update

The exponent is the sum of two expressions that are quadratic in x_t . *Completing the square* allows rewriting any quadratic $ax^2 + bx + c$ as the sum of a squared term $a(x + \frac{b}{2a})^2$ and a residual term $c - \frac{b^2}{4a}$ that is independent of x . In this case, we have $a = (\sigma_t^2 + \sigma_x^2)/(\sigma_t^2 \sigma_x^2)$, $b = -2(\sigma_t^2 x_{t+1} + \sigma_x^2 \mu_t)/(\sigma_t^2 \sigma_x^2)$, and $c = (\sigma_t^2 x_{t+1}^2 + \sigma_x^2 \mu_t^2)/(\sigma_t^2 \sigma_x^2)$. The residual term can be taken outside the integral, giving

$$= \frac{1}{Z''} \exp\left(-\frac{1}{2}\left[c - \frac{b^2}{4a}\right]\right) \int \exp\left(-\frac{a}{2}\left[x_t + \frac{b}{2a}\right]^2\right) dx_t.$$

The integral is simply the integral of a Gaussian over its entire support, and thus evaluates to 1. We are therefore left with only the residual term from the quadratic. Plugging back in the expressions for a, b , and c and simplifying, we obtain

$$= \frac{1}{Z''} \exp\left(-\frac{1}{2} \frac{(x_{t+1} - \mu_t)^2}{\sigma_t^2 + \sigma_x^2}\right).$$

That is, $X_{t+1} | y_{1:t} \sim \mathcal{N}(\mu_t, \sigma_t^2 + \sigma_x^2)$.

Plugging our results back into eq. (B.2), we obtain

$$p(x_{t+1} | y_{1:t+1}) = \frac{1}{Z'''} \exp\left(-\frac{1}{2} \left[\frac{(x_{t+1} - \mu_t)^2}{\sigma_t^2 + \sigma_x^2} + \frac{(y_{t+1} - x_{t+1})^2}{\sigma_y^2} \right]\right).$$

Completing the square analogously to our derivation of the predictive distribution yields,

$$= \frac{1}{Z'''} \exp\left(-\frac{1}{2} \frac{\left(x_{t+1} - \frac{(\sigma_t^2 + \sigma_x^2)y_{t+1} + \sigma_y^2 \mu_t}{\sigma_t^2 + \sigma_x^2 + \sigma_y^2}\right)^2}{\frac{(\sigma_t^2 + \sigma_x^2)\sigma_y^2}{\sigma_t^2 + \sigma_x^2 + \sigma_y^2}}\right).$$

Hence, $X_{t+1} | y_{1:t+1} \sim \mathcal{N}(\mu_{t+1}, \sigma_{t+1}^2)$ as defined in eq. (3.15).

Solution to exercise 3.7 (Parameter estimation using Kalman filters).

1. The given parameter estimation problem can be formulated as a Kalman filter in the following way:

$$\begin{aligned} x_t &= \pi & \forall t, \\ y_t &= x_t + \eta_t & \eta_t \sim \mathcal{N}(0, \sigma_y^2). \end{aligned}$$

Thus, in terms of Kalman filters, this yields $f = h = 1, \epsilon_t = \sigma_x^2 = 0$. Using eq. (3.16), the *Kalman gain* is given by

$$k_{t+1} = \frac{\sigma_t^2}{\sigma_t^2 + \sigma_y^2},$$

whereas the *variance of the estimation error* σ_t^2 satisfies

$$\sigma_{t+1}^2 = \sigma_y^2 k_{t+1} = \frac{\sigma_t^2 \sigma_y^2}{\sigma_t^2 + \sigma_y^2}. \quad \text{using eq. (3.19)}$$

To get the closed form, observe that

$$\frac{1}{\sigma_{t+1}^2} = \frac{1}{\sigma_t^2} + \frac{1}{\sigma_y^2} = \frac{1}{\sigma_{t-1}^2} + \frac{2}{\sigma_y^2} = \cdots = \frac{1}{\sigma_0^2} + \frac{t+1}{\sigma_y^2},$$

yielding,

$$\sigma_{t+1}^2 = \frac{\sigma_0^2 \sigma_y^2}{(t+1)\sigma_0^2 + \sigma_y^2} \quad \text{and} \quad k_{t+1} = \frac{\sigma_0^2}{(t+1)\sigma_0^2 + \sigma_y^2}.$$

2. When $t \rightarrow \infty$, we get $k_{t+1} \rightarrow 0$ and $\sigma_{t+1}^2 \rightarrow 0$, giving

$$\mu_{t+1} = \mu_t + k_{t+1}(y_{t+1} - \mu_t) = \mu_t, \quad \text{using eq. (3.18)}$$

thus resulting in a stationary sequence.

3. Observe that $\sigma_0^2 \rightarrow \infty$ implies $k_{t+1} = \frac{1}{t+1}$. Therefore,

$$\begin{aligned} \mu_{t+1} &= \mu_t + \frac{1}{t+1}(y_{t+1} - \mu_t) && \text{using eq. (3.18)} \\ &= \frac{t}{t+1}\mu_t + \frac{y_{t+1}}{t+1} \\ &= \frac{t}{t+1}\left(\frac{t-1}{t}\mu_{t-1} + \frac{y_t}{t}\right) + \frac{y_{t+1}}{t+1} && \text{using eq. (3.18)} \\ &= \frac{t-1}{t+1}\mu_{t-1} + \frac{y_t + y_{t+1}}{t+1} \\ &\vdots \\ &= \frac{y_1 + \cdots + y_{t+1}}{t+1}, \end{aligned}$$

which is simply the *sample mean*.

Gaussian Processes

Solution to exercise 4.4 (Feature space of Gaussian kernel).

1. We have for any $j \in \mathbb{N}_0$ and $x, y \in \mathbb{R}$,

$$\frac{1}{\sqrt{j!}} e^{-\frac{x^2}{2}} x^j \cdot \frac{1}{\sqrt{j!}} e^{-\frac{y^2}{2}} y^j = e^{-\frac{x^2+y^2}{2}} \frac{(xy)^j}{j!}.$$

Summing over all j , we obtain

$$\boldsymbol{\phi}(x)^\top \boldsymbol{\phi}(y) = \sum_{j=0}^{\infty} \phi_j(x) \phi_j(y)$$

$$\begin{aligned}
&= e^{-\frac{x^2+y^2}{2}} \sum_{j=0}^{\infty} \frac{(xy)^j}{j!} \\
&= e^{-\frac{x^2+y^2}{2}} e^{xy} \\
&= e^{-\frac{(x-y)^2}{2}} \\
&= k(x, y).
\end{aligned}$$

using the Taylor series expansion for the exponential function

2. As we have seen in section 2.6, the effective dimension is n . The crucial difference of kernelized regression (e.g., Gaussian processes) to linear regression is that the effective dimension *grows* with the sample size, whereas it is fixed for linear regression. Models where the effective dimension may depend on the sample size are called *non-parametric* models, and models where the effective dimension is fixed are called *parametric* models.

Solution to exercise 4.5 (A Kalman filter as a Gaussian process).

First we look at the mean,

$$\mu(t) = \mathbb{E}[X_t] = \mathbb{E}[X_{t-1} + \epsilon_{t-1}] = \mathbb{E}[X_{t-1}] = \mu(t-1).$$

Knowing that $\mu(0) = 0$, we can derive that $\mu(t) = 0$ ($\forall t$).

Now we look at the variance of X_t ,

$$\text{Var}[X_t] = \text{Var}\left[X_0 + \sum_{\tau=0}^{t-1} \epsilon_{\tau}\right] = \sigma_0^2 + t\sigma_x^2.$$

using that the noise is independent

Finally, we look at the distribution of $[f(t) \ f(t')]^{\top}$ for arbitrary $t \leq t'$.

$$\begin{bmatrix} X_t \\ X_{t'} \end{bmatrix} = \begin{bmatrix} X_t \\ X_t \end{bmatrix} + \sum_{\tau=t}^{t'-1} \begin{bmatrix} 0 \\ \epsilon_{\tau} \end{bmatrix}.$$

Therefore, we get that

$$\begin{aligned}
\begin{bmatrix} X_t \\ X_{t'} \end{bmatrix} &\sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \text{Var}[X_t] \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + (t' - t) \begin{bmatrix} 0 & 0 \\ 0 & \sigma_x^2 \end{bmatrix}\right) \\
&= \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \text{Var}[X_t] & \text{Var}[X_t] \\ \text{Var}[X_t] & \text{Var}[X_t] + (t' - t)\sigma_x^2 \end{bmatrix}\right).
\end{aligned}$$

We take the kernel $k_{\text{KF}}(t, t')$ to be the covariance between $f(t)$ and $f(t')$, which is $\text{Var}[X_t] = \sigma_0^2 + \sigma_x^2 t$. Notice, however, that we assumed $t \leq t'$. Thus, overall, the kernel is described by

$$k_{\text{KF}}(t, t') = \sigma_0^2 + \sigma_x^2 \min\{t, t'\}.$$

Solution to exercise 4.9 (Reproducing property and RKHS norm).

1. As $f \in \mathcal{H}_k(\mathcal{X})$ we can express $f(x)$ for some $\beta_i \in \mathbb{R}$ and $x_i \in \mathcal{X}$ as

$$\begin{aligned} f(x) &= \sum_{i=1}^n \beta_i k(x, x_i) \\ &= \sum_{i=1}^n \beta_i \langle k(x_i, \cdot), k(x, \cdot) \rangle_k \\ &= \left\langle \sum_{i=1}^n \beta_i k(x_i, \cdot), k(x, \cdot) \right\rangle_k \\ &= \langle f(\cdot), k(x, \cdot) \rangle_k. \end{aligned}$$

2. By applying Cauchy-Schwarz,

$$\begin{aligned} |f(x) - f(y)| &= |\langle f, k(x, \cdot) - k(y, \cdot) \rangle_k| \\ &\leq \|f\|_k \|k(x, \cdot) - k(y, \cdot)\|_k \end{aligned}$$

Solution to exercise 4.11 (MAP estimate of Gaussian processes).

1. By the representer theorem, $\hat{f}(x) = \hat{\mathbf{a}}^\top \mathbf{k}_{x,A}$. In particular, we have $\mathbf{f} = \mathbf{K}\hat{\mathbf{a}}$ and therefore

$$\begin{aligned} -\log p(y_{1:n} \mid x_{1:n}, f) &= \frac{1}{2\sigma_n^2} \|\mathbf{y} - \mathbf{f}\|_2^2 + \text{const} \\ &= \frac{1}{2\sigma_n^2} \|\mathbf{y} - \mathbf{K}\hat{\mathbf{a}}\|_2^2 + \text{const}. \end{aligned}$$

The regularization term simplifies to

$$\|f\|_k^2 = \langle f, f \rangle_k = \hat{\mathbf{a}}^\top \mathbf{K} \hat{\mathbf{a}} = \|\hat{\mathbf{a}}\|_K^2.$$

Combining, we have that

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a} \in \mathbb{R}^n} \frac{1}{2\sigma_n^2} \|\mathbf{y} - \mathbf{K}\mathbf{a}\|_2^2 + \frac{1}{2} \|\mathbf{a}\|_K^2$$

as desired. It follows by multiplying through with $2\sigma_n^2$ that $\lambda = \sigma_n^2$.

2. Expanding the objective determined in (1), we are looking for the minimizer of

$$\mathbf{a}^\top (\sigma_n^2 \mathbf{K} + \mathbf{K}^2) \mathbf{a} - 2\mathbf{y}^\top \mathbf{K} \mathbf{a} + \mathbf{y}^\top \mathbf{y}.$$

Differentiating with respect to the coefficients \mathbf{a} , we obtain the minimizer $\hat{\mathbf{a}} = (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$. Thus, the prediction at a point x^* is $\mathbf{k}_{x^*,A}^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$ which coincides with the MAP estimate.

Solution to exercise 4.13 (Gradient of the marginal likelihood).

1. Applying the two hints eqs. (4.39) and (4.40) to eq. (4.37) yields,

$$\frac{\partial}{\partial \theta_j} \log p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) = \frac{1}{2} \mathbf{y}^\top \mathbf{K}_{\mathbf{y},\boldsymbol{\theta}}^{-1} \frac{\partial \mathbf{K}_{\mathbf{y},\boldsymbol{\theta}}}{\partial \theta_j} \mathbf{K}_{\mathbf{y},\boldsymbol{\theta}}^{-1} \mathbf{y} - \frac{1}{2} \text{tr} \left(\mathbf{K}_{\mathbf{y},\boldsymbol{\theta}}^{-1} \frac{\partial \mathbf{K}_{\mathbf{y},\boldsymbol{\theta}}}{\partial \theta_j} \right).$$

We can simplify to

$$\begin{aligned}
&= \frac{1}{2} \text{tr} \left(\mathbf{y}^\top \mathbf{K}_{y,\theta}^{-1} \frac{\partial \mathbf{K}_{y,\theta}}{\partial \theta_j} \mathbf{K}_{y,\theta}^{-1} \mathbf{y} \right) - \frac{1}{2} \text{tr} \left(\mathbf{K}_{y,\theta}^{-1} \frac{\partial \mathbf{K}_{y,\theta}}{\partial \theta_j} \right) && \text{using that } \mathbf{y}^\top \mathbf{K}_{y,\theta}^{-1} \frac{\partial \mathbf{K}_{y,\theta}}{\partial \theta_j} \mathbf{K}_{y,\theta}^{-1} \mathbf{y} \text{ is a scalar} \\
&= \frac{1}{2} \text{tr} \left(\mathbf{K}_{y,\theta}^{-1} \mathbf{y} \mathbf{y}^\top \mathbf{K}_{y,\theta}^{-1} \frac{\partial \mathbf{K}_{y,\theta}}{\partial \theta_j} - \mathbf{K}_{y,\theta}^{-1} \frac{\partial \mathbf{K}_{y,\theta}}{\partial \theta_j} \right) && \text{using the cyclic property and linearity of the trace} \\
&= \frac{1}{2} \text{tr} \left(\mathbf{K}_{y,\theta}^{-1} \mathbf{y} (\mathbf{K}_{y,\theta}^{-1} \mathbf{y})^\top \frac{\partial \mathbf{K}_{y,\theta}}{\partial \theta_j} - \mathbf{K}_{y,\theta}^{-1} \frac{\partial \mathbf{K}_{y,\theta}}{\partial \theta_j} \right) && \text{using that } \mathbf{K}_{y,\theta}^{-1} \text{ is symmetric} \\
&= \frac{1}{2} \text{tr} \left((\alpha \alpha^\top - \mathbf{K}_{y,\theta}^{-1}) \frac{\partial \mathbf{K}_{y,\theta}}{\partial \theta_j} \right).
\end{aligned}$$

2. We denote by $\tilde{\mathbf{K}}$ the covariance matrix of \mathbf{y} for the covariance function \tilde{k} , so $\mathbf{K}_{y,\theta} = \theta_0 \tilde{\mathbf{K}}$. Then,

$$\frac{\partial}{\partial \theta_0} \log p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) = \frac{1}{2} \text{tr} \left((\theta_0^{-2} \tilde{\mathbf{K}}^{-1} \mathbf{y} (\tilde{\mathbf{K}}^{-1} \mathbf{y})^\top - \theta_0^{-1} \tilde{\mathbf{K}}^{-1}) \tilde{\mathbf{K}} \right). \quad \text{using eq. (4.38)}$$

Simplifying the terms and using linearity of the trace, we obtain that

$$\frac{\partial}{\partial \theta_0} \log p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) = 0 \quad \Longleftrightarrow \quad \theta_0 = \frac{1}{n} \text{tr}(\mathbf{y} \mathbf{y}^\top \tilde{\mathbf{K}}^{-1}).$$

If we define $\tilde{\mathbf{A}} \doteq \tilde{\mathbf{K}}^{-1}$ as the precision matrix associated to \mathbf{y} for the covariance function \tilde{k} , we can express θ_0^* in closed form as

$$\theta_0^* = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \tilde{\mathbf{A}}(i, j) y_i y_j. \quad (\text{B.3})$$

3. We immediately see from eq. (B.3) that θ_0^* scales by s^2 if \mathbf{y} is scaled by s .

Solution to exercise 4.19 (Uniform convergence of Fourier features).

1. We have that $s(\cdot) \in [-\sqrt{2}, \sqrt{2}]$, and hence, $s(\Delta_i)$ is $\sqrt{2}$ -sub-Gaussian.¹ It then follows from Hoeffding's inequality (1.88) that

$$\mathbb{P}(|f(\Delta_i)| \geq \epsilon) \leq 2 \exp \left(-\frac{m\epsilon^2}{4} \right).$$

2. We can apply Markov's inequality (1.82) to obtain

$$\begin{aligned}
\mathbb{P} \left(\|\nabla f(\Delta^*)\|_2 \geq \frac{\epsilon}{2r} \right) &= \mathbb{P} \left(\|\nabla f(\Delta^*)\|_2^2 \geq \left(\frac{\epsilon}{2r} \right)^2 \right) \\
&\leq \left(\frac{2r \mathbb{E} \|\nabla f(\Delta^*)\|_2}{\epsilon} \right)^2.
\end{aligned}$$

It remains to bound the expectation. We have

$$\mathbb{E} \|\nabla f(\Delta^*)\|_2^2 = \mathbb{E} \|\nabla s(\Delta^*) - \nabla k(\Delta^*)\|_2^2$$

¹ see example 1.47

$$= \mathbb{E} \|\nabla s(\Delta^*)\|_2^2 - 2 \nabla k(\Delta^*)^\top \mathbb{E} \nabla s(\Delta^*) + \mathbb{E} \|\nabla k(\Delta^*)\|_2^2. \quad \text{using linearity of expectation (1.24)}$$

Note that $\mathbb{E} \nabla s(\Delta) = \nabla k(\Delta)$ using eq. (1.28) and using that $s(\Delta)$ is an unbiased estimator of $k(\Delta)$. Therefore,

$$\begin{aligned} &= \mathbb{E} \|\nabla s(\Delta^*)\|_2^2 - \mathbb{E} \|\nabla k(\Delta^*)\|_2^2 \\ &\leq \mathbb{E} \|\nabla s(\Delta^*)\|_2^2 \\ &\leq \mathbb{E}_{\omega \sim p} \|\omega\|_2^2 = \sigma_p^2. \end{aligned}$$

using that s is the cos of a linear function in ω

3. Using a union bound (1.1) and then the result of (1),

$$\mathbb{P} \left(\bigcup_{i=1}^T |f(\Delta_i)| \geq \frac{\epsilon}{2} \right) \leq \sum_{i=1}^T \mathbb{P} \left(|f(\Delta_i)| \geq \frac{\epsilon}{2} \right) \leq 2T \exp \left(-\frac{m\epsilon^2}{16} \right).$$

4. First, note that by contraposition,

$$\sup_{\Delta \in \mathcal{M}_\Delta} |f(\Delta)| \geq \epsilon \implies \exists i. |f(\Delta_i)| \geq \frac{\epsilon}{2} \text{ or } \|\nabla f(\Delta^*)\|_2 \geq \frac{\epsilon}{2r},$$

and therefore,

$$\begin{aligned} \mathbb{P} \left(\sup_{\Delta \in \mathcal{M}_\Delta} |f(\Delta)| \geq \epsilon \right) &\leq \mathbb{P} \left(\bigcup_{i=1}^T |f(\Delta_i)| \geq \frac{\epsilon}{2} \cup \|\nabla f(\Delta^*)\|_2 \geq \frac{\epsilon}{2r} \right) \\ &\leq \mathbb{P} \left(\bigcup_{i=1}^T |f(\Delta_i)| \geq \frac{\epsilon}{2} \right) + \mathbb{P} \left(\|\nabla f(\Delta^*)\|_2 \geq \frac{\epsilon}{2r} \right) \quad \text{using a union bound (1.1)} \\ &\leq 2T \exp \left(-\frac{m\epsilon^2}{2^4} \right) + \left(\frac{2r\sigma_p}{\epsilon} \right)^2 \quad \text{using the results from (2) and (3)} \\ &\leq \alpha r^{-d} + \beta r^2 \quad \text{using } T \leq (4 \text{ diam}(\mathcal{M})/r)^d \end{aligned}$$

with $\alpha = 2(4 \text{ diam}(\mathcal{M}))^d \exp \left(-\frac{m\epsilon^2}{2^4} \right)$ and $\beta = \left(\frac{2\sigma_p}{\epsilon} \right)^2$. Using the hint, we obtain

$$\begin{aligned} &= 2\beta^{\frac{d}{d+2}} \alpha^{\frac{2}{d+2}} \\ &= 2 \left(2 \left(\frac{2^3 \sigma_p \text{diam}(\mathcal{M})}{\epsilon} \right)^{2d} \right)^{\frac{1}{d+2}} \exp \left(-\frac{m\epsilon^2}{2^3(d+2)} \right) \\ &\leq 2^8 \left(\frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon} \right)^2 \exp \left(-\frac{m\epsilon^2}{2^3(d+2)} \right) \quad \text{using } \frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon} \geq 1 \end{aligned}$$

5. We have

$$\begin{aligned} \sigma_p^2 &= \mathbb{E}_{\omega \sim p} [\omega^\top \omega] \\ &= \int \omega^\top \omega \cdot p(\omega) d\omega \\ &= \int \omega^\top \omega \cdot e^{i\omega^\top \mathbf{0}} p(\omega) d\omega. \end{aligned}$$

Now observe that $\frac{\partial^2}{\partial \Delta_j^2} \int e^{i\omega^\top \Delta} p(\omega) d\omega = - \int \omega_j^2 e^{i\omega^\top \Delta} p(\omega) d\omega$. Thus,

$$\begin{aligned} &= -\text{tr} \left(H_\Delta \int p(\omega) e^{i\omega^\top \Delta} d\omega \Big|_{\Delta=0} \right) \\ &= -\text{tr}(H_\Delta k(0)). \end{aligned}$$

using that p is the Fourier transform of k (4.48)

Finally, we have for the Gaussian kernel that

$$\frac{\partial^2}{\partial \Delta_j^2} \exp \left(-\frac{\Delta^\top \Delta}{2\ell^2} \right) \Big|_{\Delta_j=0} = -\frac{1}{\ell^2}.$$

Solution to exercise 4.22 (Subset of regressors).

1. We write $\tilde{f} \doteq [f \ f^*]^\top$. From the definition of SoR (4.59), we gather

$$q_{\text{SoR}}(\tilde{f} \mid \mathbf{u}) = \mathcal{N} \left(\tilde{f}; \begin{bmatrix} K_{AU} K_{UU}^{-1} \\ K_{*U} K_{UU}^{-1} \end{bmatrix} \mathbf{u}, \mathbf{0} \right). \quad (\text{B.4})$$

We know that \tilde{f} and \mathbf{u} are jointly Gaussian, and hence, the marginal distribution of \tilde{f} is also Gaussian. We have for the mean and variance that

$$\begin{aligned} \mathbb{E}[\tilde{f}] &= \mathbb{E}[\mathbb{E}[\tilde{f} \mid \mathbf{u}]] && \text{using the tower rule (1.33)} \\ &= \mathbb{E}_{\mathbf{u}} \left[\begin{bmatrix} K_{AU} K_{UU}^{-1} \\ K_{*U} K_{UU}^{-1} \end{bmatrix} \mathbf{u} \right] && \text{using eq. (B.4)} \\ &= \begin{bmatrix} K_{AU} K_{UU}^{-1} \\ K_{*U} K_{UU}^{-1} \end{bmatrix} \mathbb{E}[\mathbf{u}] = \mathbf{0} && \text{using linearity of expectation (1.24)} \\ \text{Var}[\tilde{f}] &= \underbrace{\mathbb{E} \text{Var}[\tilde{f} \mid \mathbf{u}]}_{\mathbf{0}} + \text{Var} \mathbb{E}[\tilde{f} \mid \mathbf{u}] && \text{using the law of total variance (1.50)} \\ &= \text{Var}_{\mathbf{u}} \left[\begin{bmatrix} K_{AU} K_{UU}^{-1} \\ K_{*U} K_{UU}^{-1} \end{bmatrix} \mathbf{u} \right] && \text{using eq. (B.4)} \\ &= \begin{bmatrix} K_{AU} K_{UU}^{-1} \\ K_{*U} K_{UU}^{-1} \end{bmatrix} \underbrace{\text{Var}[\mathbf{u}]}_{K_{UU}} \begin{bmatrix} K_{AU} K_{UU}^{-1} \\ K_{*U} K_{UU}^{-1} \end{bmatrix}^\top && \text{using eq. (1.47)} \\ &= \begin{bmatrix} Q_{AA} & Q_{A*} \\ Q_{*A} & Q_{**} \end{bmatrix}. \end{aligned}$$

Having determined $q_{\text{SoR}}(f, f^*)$, $q_{\text{SoR}}(f^* \mid \mathbf{y})$ follows directly using the formulas for finding the Gaussian process predictive posterior (4.6).

2. The given covariance function follows directly from inspecting the derived covariance matrix, $\text{Var}[f] = Q_{AA}$.

Variational Inference

Solution to exercise 5.3 (Logistic loss).

1. We have

$$\begin{aligned}
 \nabla_w \ell_{\log}(w^\top x; y) &= \nabla_w \log(1 + \exp(-yw^\top x)) && \text{using the definition of the logistic loss (5.14)} \\
 &= \nabla_w \log \frac{1 + \exp(yw^\top x)}{\exp(yw^\top x)} \\
 &= \nabla_w \log(1 + \exp(yw^\top x)) - yw^\top x \\
 &= \frac{1}{1 + \exp(-yw^\top x)} \cdot \exp(-yw^\top x) \cdot yx - yx && \text{using the chain rule} \\
 &= -yx \cdot \left(1 - \frac{\exp(yw^\top x)}{1 + \exp(yw^\top x)}\right) \\
 &= -yx \cdot \sigma(-yw^\top x). && \text{using the definition of the logistic function (5.10)}
 \end{aligned}$$

2. As suggested in the hint, we compute the first derivative of σ ,

$$\begin{aligned}
 \sigma'(z) &\doteq \frac{d}{dz} \sigma(z) = \frac{d}{dz} \frac{1}{1 + \exp(-z)} && \text{using the definition of the logistic function (5.10)} \\
 &= \frac{\exp(-z)}{(1 + \exp(-z))^2} && \text{using the quotient rule} \\
 &= \sigma(z) \cdot (1 - \sigma(z)). && \text{using the definition of the logistic function (5.10)}
 \end{aligned}$$

We get for the Hessian of ℓ_{\log} ,

$$\begin{aligned}
 H_w \ell_{\log}(w^\top x; y) &= D_w \nabla_w \ell_{\log}(w^\top x; y) && \text{using the definition of a Hessian (1.106) and symmetry} \\
 &= -yx D_w \sigma(-yw^\top x). && \text{using the gradient of the logistic loss from (1)}
 \end{aligned}$$

We have $D_w - yw^\top x = -yx^\top$ and $D_z \sigma(z) = \sigma'(z)$, and therefore, using the chain rule of multivariate calculus (5.20),

$$\begin{aligned}
 &= -yx \cdot \sigma'(-yw^\top x) \cdot (-yx^\top) \\
 &= xx^\top \cdot \sigma'(-yw^\top x). && \text{using } y^2 = 1
 \end{aligned}$$

Finally, recall that

$$\begin{aligned}
 \sigma(w^\top x) &= \mathbb{P}(y = +1 \mid x, w) \quad \text{and} \\
 \sigma(-w^\top x) &= \mathbb{P}(y = -1 \mid x, w).
 \end{aligned}$$

Thus,

$$\begin{aligned}
 \sigma'(-yw^\top x) &= \mathbb{P}(Y \neq y \mid x, w) \cdot (1 - \mathbb{P}(Y \neq y \mid x, w)) \\
 &= \mathbb{P}(Y \neq y \mid x, w) \cdot \mathbb{P}(Y = y \mid x, w) \\
 &= (1 - \mathbb{P}(Y = y \mid x, w)) \cdot \mathbb{P}(Y = y \mid x, w) \\
 &= \sigma'(w^\top x).
 \end{aligned}$$

3. By the second-order characterization of convexity (cf. remark 1.57), a twice differentiable function f is convex if and only if its Hessian is positive semi-definite. Writing $c \doteq \sigma'(\mathbf{w}^\top \mathbf{x})$, we have for any $\delta \in \mathbb{R}^n$ that

$$\delta^\top \mathbf{H}_w \ell_{\log}(\mathbf{w}^\top \mathbf{x}; y) \delta = \delta^\top (c \cdot \mathbf{x} \mathbf{x}^\top) \delta = c(\delta^\top \mathbf{x})^2 \geq 0, \quad \text{using } c \geq 0 \text{ and } (\cdot)^2 \geq 0$$

and hence, the logistic loss is convex in \mathbf{w} .

Solution to exercise 5.5 (Gaussian process classification).

1. Using the law of total probability (1.15),

$$\begin{aligned} p(y^* = +1 \mid \mathbf{x}_{1:n}, y_{1:n}, \mathbf{x}^*) \\ &= \int p(y^* = +1 \mid f^*) p(f^* \mid \mathbf{x}_{1:n}, y_{1:n}, \mathbf{x}^*) df^* \\ &= \int \sigma(f^*) p(f^* \mid \mathbf{x}_{1:n}, y_{1:n}, \mathbf{x}^*) df^*. \end{aligned} \quad \text{using } y^* \sim \text{Bern}(\sigma(f^*))$$

Due to the non-Gaussian likelihood, the integral is analytically intractable. However, as the integral is one-dimensional, numerical approximations such as the Gauss-Legendre quadrature can be used.

2. (a) According to Bayes' rule (1.59), we know that

$$\begin{aligned} \psi(f) &= \log p(f \mid \mathbf{x}_{1:n}, y_{1:n}) && \text{using eq. (5.2)} \\ &= \log p(y_{1:n} \mid f) + \log p(f \mid \mathbf{x}_{1:n}) - \log p(y_{1:n} \mid \mathbf{x}_{1:n}) \\ &= \log p(y_{1:n} \mid f) + \log p(f \mid \mathbf{x}_{1:n}) + \text{const.} && \text{using Bayes' rule (1.59)} \end{aligned}$$

Note that $p(f \mid \mathbf{x}_{1:n}) = \mathcal{N}(f; \mathbf{0}, \mathbf{K}_{AA})$. Plugging in this closed-form Gaussian distribution of the GP prior gives

$$= \log p(y_{1:n} \mid f) - \frac{1}{2} \mathbf{f}^\top \mathbf{K}_{AA}^{-1} \mathbf{f} + \text{const}$$

Differentiating with respect to f yields

$$\begin{aligned} \nabla \psi(f) &= \nabla \log p(y_{1:n} \mid f) - \mathbf{K}_{AA}^{-1} \mathbf{f} \\ \mathbf{H} \psi(f) &= \mathbf{H}_f \log p(y_{1:n} \mid f) - \mathbf{K}_{AA}^{-1}. \end{aligned}$$

Hence, $\mathbf{\Lambda} = \mathbf{K}_{AA}^{-1} + \mathbf{W}$ where $\mathbf{W} \doteq -\mathbf{H}_f \log p(y_{1:n} \mid f)|_{f=\hat{f}}$. It remains to derive \mathbf{W} for the probit likelihood. Using independence of the training examples,

$$\log p(y_{1:n} \mid f) = \sum_{i=1}^n \log p(y_i \mid f_i),$$

and hence, the Hessian of this expression is diagonal. Using the symmetry of $\Phi(z; 0, \sigma_n^2)$ around zero, we can write

$$\log p(y_i \mid f_i) = \log \Phi(y_i f_i; 0, \sigma_n^2).$$

In the following, we write $\mathcal{N}(z) \doteq \mathcal{N}(z; 0, \sigma_n^2)$ and $\Phi(z) \doteq \Phi(z; 0, \sigma_n^2)$ to simplify the notation. Differentiating with respect to f_i , we obtain

$$\begin{aligned}\frac{\partial}{\partial f_i} \log \Phi(y_i f_i) &= \frac{y_i \mathcal{N}(f_i)}{\Phi(y_i f_i)} \\ \frac{\partial^2}{\partial f_i^2} \log \Phi(y_i f_i) &= -\frac{\mathcal{N}(f_i)^2}{\Phi(y_i f_i)^2} - \frac{y_i f_i \mathcal{N}(f_i)}{\sigma_n^2 \Phi(y_i f_i)},\end{aligned}$$

using $\mathcal{N}(y_i f_i) = \mathcal{N}(f_i)$ since \mathcal{N} is an even function and $y_i \in \{\pm 1\}$

$$\text{and } \mathbf{W} = -\text{diag}\left\{\frac{\partial^2}{\partial f_i^2} \log \Phi(y_i f_i)\right\}_{i=1}^n \Big|_{f=\hat{f}}.$$

- (b) Note that Λ' is a precision matrix over weights w and $f = Xw$, so by eq. (1.47) the corresponding variance over latent values f is $X\Lambda'^{-1}X^\top$. The two precision matrices are therefore equivalent if $\Lambda^{-1} = X\Lambda'^{-1}X^\top$.

Analogously to example 5.2 and exercise 5.3, we have that

$$\begin{aligned}\mathbf{W} &= -\mathbf{H}_f \log p(y_{1:n} | f) \Big|_{f=\hat{f}} \\ &= \text{diag}_{i \in [n]} \{\sigma(\hat{f}_i)(1 - \sigma(\hat{f}_i))\} \\ &= \text{diag}_{i \in [n]} \{\pi_i(1 - \pi_i)\}.\end{aligned}$$

By Woodbury's matrix identity (A.1),

$$\begin{aligned}\Lambda'^{-1} &= (\mathbf{I} + \mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \\ &= \mathbf{I} - \mathbf{X}^\top (\mathbf{W}^{-1} + \mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X}.\end{aligned}$$

Thus,

$$\begin{aligned}X\Lambda'^{-1}X^\top &= \mathbf{X} \mathbf{X}^\top - \mathbf{X} \mathbf{X}^\top (\mathbf{W}^{-1} + \mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{X}^\top \\ &= \mathbf{K}_{AA} - \mathbf{K}_{AA} (\mathbf{W}^{-1} + \mathbf{K}_{AA})^{-1} \mathbf{K}_{AA} \\ &= (\mathbf{K}_{AA}^{-1} + \mathbf{W})^{-1} \\ &= \Lambda^{-1}.\end{aligned}$$

using that $\mathbf{K}_{AA} = \mathbf{X} \mathbf{X}^\top$
by the matrix inversion lemma (A.3)

- (c) Using the formulas for the conditional distribution of a Gaussian (1.118), evaluating a conditional GP at a test point x^* yields,

$$f^* | x^*, f \sim \mathcal{N}(\mu^*, k^*), \quad \text{where} \quad (\text{B.5a})$$

$$\mu^* \doteq k_{x^*, A}^\top \mathbf{K}_{AA}^{-1} f, \quad (\text{B.5b})$$

$$k^* \doteq k(x^*, x^*) - k_{x^*, A}^\top \mathbf{K}_{AA}^{-1} k_{x^*, A}. \quad (\text{B.5c})$$

Then, using the tower rule (1.33),

$$\begin{aligned}\mathbb{E}_q[f^* | x^*, x_{1:n}, y_{1:n}] &= \mathbb{E}_{f \sim q} \left[\mathbb{E}_{f^* \sim p}[f^* | x^*, f] \mid x_{1:n}, y_{1:n} \right] \\ &= k_{x^*, A}^\top \mathbf{K}_{AA}^{-1} \mathbb{E}_q[f | x_{1:n}, y_{1:n}]\end{aligned}$$

using eq. (B.5b) and linearity of expectation (1.24)

$$= \mathbf{k}_{\mathbf{x}^*, A}^\top \mathbf{K}_{AA}^{-1} \hat{f}. \quad (\text{B.6})$$

Observe that any maximum \hat{f} of $\psi(f)$ needs to satisfy $\nabla \psi(f) = \mathbf{0}$. Hence, $\hat{f} = \mathbf{K}_{AA}(\nabla_f \log p(y_{1:n} | f))$, and the expectation simplifies to

$$= \mathbf{k}_{\mathbf{x}^*, A}^\top (\nabla_f \log p(y_{1:n} | f)).$$

Using the law of total variance (1.50),

$$\begin{aligned} & \text{Var}_q[f^* | \mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}] \\ &= \mathbb{E}_{f \sim q} \left[\text{Var}_{f^* \sim p} [f^* | \mathbf{x}^*, \mathbf{x}_{1:n}, f] \mid \mathbf{x}_{1:n}, y_{1:n} \right] \\ & \quad + \text{Var}_{f \sim q} \left[\mathbb{E}_{f^* \sim p} [f^* | \mathbf{x}^*, \mathbf{x}_{1:n}, f] \mid \mathbf{x}_{1:n}, y_{1:n} \right] \\ &= \mathbb{E}_q[k^* | \mathbf{x}_{1:n}, y_{1:n}] + \text{Var}_q \left[\mathbf{k}_{\mathbf{x}^*, A}^\top \mathbf{K}_{AA}^{-1} f \mid \mathbf{x}_{1:n}, y_{1:n} \right] \\ &= k^* + \mathbf{k}_{\mathbf{x}^*, A}^\top \mathbf{K}_{AA}^{-1} \text{Var}_q[f | \mathbf{x}_{1:n}, y_{1:n}] \mathbf{K}_{AA}^{-1} \mathbf{k}_{\mathbf{x}^*, A}. \end{aligned}$$

using eqs. (B.5b) and (B.5c)

using that k^* is independent of f ,
eq. (1.47), and symmetry of \mathbf{K}_{AA}^{-1}

Recall from (a) that $\text{Var}_q[f | \mathbf{x}_{1:n}, y_{1:n}] = (\mathbf{K}_{AA} + \mathbf{W}^{-1})^{-1}$, so

$$\begin{aligned} &= k(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{k}_{\mathbf{x}^*, A}^\top \mathbf{K}_{AA}^{-1} \mathbf{k}_{\mathbf{x}^*, A} \\ & \quad + \mathbf{k}_{\mathbf{x}^*, A}^\top \mathbf{K}_{AA}^{-1} (\mathbf{K}_{AA} + \mathbf{W}^{-1})^{-1} \mathbf{K}_{AA}^{-1} \mathbf{k}_{\mathbf{x}^*, A} \\ &= k(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{k}_{\mathbf{x}^*, A}^\top (\mathbf{K}_{AA} + \mathbf{W}^{-1})^{-1} \mathbf{k}_{\mathbf{x}^*, A}. \end{aligned} \quad (\text{B.7})$$

plugging in the expression for k^* (B.5c)

using the matrix inversion lemma (A.3)

Notice the similarity of eqs. (B.6) and (B.7) to eqs. (B.5b) and (B.5c).

The latter is the conditional mean and variance if f is known whereas the former is the conditional mean and variance given the noisy observation $y_{1:n}$ of f . The matrix \mathbf{W} quantifies the noise in the observations.

(d) Recall that

$$\begin{aligned} p(y^* = +1 | \mathbf{x}_{1:n}, y_{1:n}, \mathbf{x}^*) &\approx \int \sigma(f^*) q(f^* | \mathbf{x}_{1:n}, y_{1:n}, \mathbf{x}^*) df^* \\ &= \mathbb{E}_q[\sigma(f^*)]. \end{aligned}$$

using the Laplace-approximated latent predictive posterior

using LOTUS (1.30)

This quantity can be interpreted as the *averaged prediction* over all latent predictions f^* . In contrast, $\sigma(\mathbb{E}_q[f^*])$ can be understood as the “MAP” prediction, which is obtained using the MAP estimate of f^* .² As σ is non-linear, the two quantities are not identical, and generally the averaged prediction is preferred.

² As q is a Gaussian, its mode (i.e., the MAP estimate) and its mean coincide.

3. We have

$$\begin{aligned} \mathbb{E}_\epsilon \mathbb{1}\{f(\mathbf{x}) + \epsilon \geq 0\} &= \mathbb{P}_\epsilon(f(\mathbf{x}) + \epsilon \geq 0) \\ &= \mathbb{P}_\epsilon(-\epsilon \leq f(\mathbf{x})) \\ &= \mathbb{P}_\epsilon(\epsilon \leq f(\mathbf{x})) \\ &= \Phi(f(\mathbf{x}); 0, \sigma_n^2). \end{aligned}$$

using $\mathbb{E}[X] = p$ if $X \sim \text{Bern}(p)$

using that the distribution of ϵ is
symmetric around 0

Solution to exercise 5.17 (Binary cross-entropy loss). If $y = 1$ then

$$\ell_{\text{bce}}(\hat{y}; y) = -\log \hat{y} = \log(1 + e^{-\hat{f}}) = \ell_{\log}(\hat{f}; y).$$

If $y = -1$ then

$$\ell_{\text{bce}}(\hat{y}; y) = -\log(1 - \hat{y}) = \log(1 + e^{\hat{f}}) = \ell_{\log}(\hat{f}; y).$$

Here the second equality follows from the simple algebraic fact

$$1 - \frac{1}{1 + e^{-z}} = 1 - \frac{e^z}{e^z + 1} = \frac{1}{1 + e^z}. \quad \text{multiplying by } \frac{e^z}{e^z}$$

Solution to exercise 5.13 (Jensen's inequality).

1. Recall that as f is convex,

$$\forall x_1, x_2, \forall \lambda \in [0, 1] : \quad f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2).$$

We prove the statement by induction on k . The base case, $k = 2$, follows trivially from the convexity of f . For the induction step, suppose that the statement holds for some fixed $k \geq 2$ and assume w.l.o.g. that $\theta_{k+1} \in (0, 1)$. We then have,

$$\begin{aligned} \sum_{i=1}^{k+1} \theta_i f(x_i) &= (1 - \theta_{k+1}) \left(\sum_{i=1}^k \frac{\theta_i}{1 - \theta_{k+1}} f(x_i) \right) + \theta_{k+1} f(x_{k+1}) \\ &\geq (1 - \theta_{k+1}) \cdot f \left(\sum_{i=1}^k \frac{\theta_i}{1 - \theta_{k+1}} x_i \right) + \theta_{k+1} f(x_{k+1}) && \text{using the induction hypothesis} \\ &\geq f \left(\sum_{i=1}^{k+1} \theta_i x_i \right). && \text{using the convexity of } f \end{aligned}$$

2. Noting that \log_2 is concave, we have by Jensen's inequality,

$$\begin{aligned} H[p] &= \mathbb{E}_{x \sim p} \left[\log_2 \left(\frac{1}{p(x)} \right) \right] && \text{by definition of entropy (5.33a)} \\ &\leq \log_2 \mathbb{E}_{x \sim p} \left[\frac{1}{p(x)} \right] && \text{by Jensen's inequality (5.38)} \\ &= \log_2 n. \end{aligned}$$

Solution to exercise 5.19 (Gibbs' inequality).

1. Let p and q be two (continuous) distributions. The KL-divergence between p and q is

$$\begin{aligned} \text{KL}(p \| q) &= \mathbb{E}_{x \sim p} \left[\log \frac{p(x)}{q(x)} \right] && \text{using the definition of KL-divergence (5.44)} \\ &= \mathbb{E}_{x \sim p} \left[\text{S} \left[\frac{q(x)}{p(x)} \right] \right]. \end{aligned}$$

Note that the surprise $S[u] = -\log u$ is a convex function in u , and hence,

$$\begin{aligned} &\geq S\left[\mathbb{E}_{x \sim p}\left[\frac{q(x)}{p(x)}\right]\right] \\ &= S\left[\int q(x) dx\right] \\ &= S[1] = 0. \end{aligned}$$

using Jensen's inequality (5.37)

a probability density integrates to 1

2. We observe from the derivation of (1) that $\text{KL}(p\|q) = 0$ iff equality holds for Jensen's inequality. Now, if p and q are discrete with finite and identical support, we can follow from the hint that Jensen's inequality degenerates to an equality iff p and q are point wise identical.

Solution to exercise 5.20 (Maximum entropy principle).

1. We have

$$\begin{aligned} H[f\|g] &= -\int_{\mathbb{R}} f(x) \log g(x) dx \\ &= -\int_{\mathbb{R}} f(x) \cdot \left(\log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \frac{(x-\mu)^2}{2\sigma^2}\right) dx \\ &= -\log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) \underbrace{\int_{\mathbb{R}} f(x) dx}_1 + \frac{1}{2\sigma^2} \int_{\mathbb{R}} f(x)(x-\mu)^2 dx \\ &= \log(\sigma\sqrt{2\pi}) + \frac{1}{2\sigma^2} \underbrace{\mathbb{E}_{x \sim f}[(x-\mu)^2]}_{\sigma^2} \\ &= \log(\sigma\sqrt{2\pi}) + \frac{1}{2} \\ &= H[g]. \end{aligned}$$

using the definition of cross-entropy (5.40)

using that $g(x) = \mathcal{N}(x; \mu, \sigma^2)$

using the entropy of Gaussians (5.35)

2. We have shown that

$$H[g] - H[f] = \text{KL}(f\|g) \geq 0,$$

and hence, $H[g] \geq H[f]$. That is, for a fixed mean μ and variance σ^2 , the distribution that has maximum entropy among all distributions on \mathbb{R} is the normal distribution.

Solution to exercise 5.22 (KL-divergence of Gaussians). We can rewrite the KL-divergence as

$$\begin{aligned} \text{KL}(p\|q) &= \mathbb{E}_{x \sim p}[\log p(x) - \log q(x)] \\ &= \mathbb{E}_{x \sim p} \left[\frac{1}{2} \log \frac{\det(\Sigma_q)}{\det(\Sigma_p)} - \frac{1}{2} (x - \mu_p)^\top \Sigma_p^{-1} (x - \mu_p) \right. \\ &\quad \left. + \frac{1}{2} (x - \mu_q)^\top \Sigma_q^{-1} (x - \mu_q) \right] \end{aligned}$$

using the definition of KL-divergence (5.44)

using the Gaussian PDF (1.115)

$$= \frac{1}{2} \log \frac{\det(\Sigma_q)}{\det(\Sigma_p)} - \frac{1}{2} \mathbb{E}_{x \sim p} \left[(x - \mu_p)^\top \Sigma_p^{-1} (x - \mu_p) \right] + \frac{1}{2} \mathbb{E}_{x \sim p} \left[(x - \mu_q)^\top \Sigma_q^{-1} (x - \mu_q) \right]$$

using linearity of expectation (1.24)

As $(x - \mu_p)^\top \Sigma_p^{-1} (x - \mu_p) \in \mathbb{R}$, we can rewrite the second term as

$$\begin{aligned} & \frac{1}{2} \mathbb{E}_{x \sim p} \left[\text{tr} \left((x - \mu_p)^\top \Sigma_p^{-1} (x - \mu_p) \right) \right] \\ &= \frac{1}{2} \mathbb{E}_{x \sim p} \left[\text{tr} \left((x - \mu_p)(x - \mu_p)^\top \Sigma_p^{-1} \right) \right] \\ &= \frac{1}{2} \text{tr} \left(\mathbb{E}_{x \sim p} \left[(x - \mu_p)(x - \mu_p)^\top \right] \Sigma_p^{-1} \right) \\ &= \frac{1}{2} \text{tr} \left(\Sigma_p \Sigma_p^{-1} \right) \\ &= \frac{1}{2} \text{tr}(I) = \frac{d}{2}. \end{aligned}$$

using the cyclic property of the trace

using linearity of the trace and linearity of expectation (1.24)

using the definition of the covariance matrix (1.44)

For the third term, we use the hint (5.51) to obtain

$$\begin{aligned} & \frac{1}{2} \mathbb{E}_{x \sim p} \left[(x - \mu_q)^\top \Sigma_q^{-1} (x - \mu_q) \right] \\ &= \frac{1}{2} \left((\mu_p - \mu_q)^\top \Sigma_q^{-1} (\mu_p - \mu_q) + \text{tr} \left(\Sigma_q^{-1} \Sigma_p \right) \right). \end{aligned}$$

Putting all terms together we get

$$\begin{aligned} \text{KL}(p \| q) &= \frac{1}{2} \left(\log \frac{\det(\Sigma_q)}{\det(\Sigma_p)} - d + (\mu_p - \mu_q)^\top \Sigma_q^{-1} (\mu_p - \mu_q) \right. \\ &\quad \left. + \text{tr} \left(\Sigma_q^{-1} \Sigma_p \right) \right). \end{aligned}$$

Solution to exercise 5.25 (Forward vs reverse KL).

1. Let p and q be discrete distributions. The derivation is analogous if p and q are taken to be continuous. First, we write the KL-divergence between p and q as

$$\begin{aligned} \text{KL}(p \| q) &= \sum_x \sum_y p(x, y) \log_2 \frac{p(x, y)}{q(x)q(y)} \\ &= \sum_x \sum_y p(x, y) \log_2 p(x, y) - \sum_x \sum_y p(x, y) \log_2 q(x) \\ &\quad - \sum_x \sum_y p(x, y) \log_2 q(y) \\ &= \sum_x \sum_y p(x, y) \log_2 p(x, y) - \sum_x p(x) \log_2 q(x) - \sum_y p(y) \log_2 q(y) \\ &= -H[p(x, y)] + H[p(x) \| q(x)] + H[p(y) \| q(y)] \\ &= -H[p(x, y)] + H[p(x)] + H[p(y)] + \text{KL}(p(x) \| q(x)) \\ &\quad + \text{KL}(p(y) \| q(y)) \end{aligned}$$

using the definition of KL-divergence (5.44)

using the sum rule (1.10)

using the definitions of entropy (5.32) and cross-entropy (5.40)

using eq. (5.41)

$$= \text{KL}(p(x)\|q(x)) + \text{KL}(p(y)\|q(y)) + \text{const.}$$

Hence, to minimize $\text{KL}(p\|q)$ with respect to the variational distributions $q(x)$ and $q(y)$ we should set $\text{KL}(p(x)\|q(x)) = 0$ and $\text{KL}(p(y)\|q(y)) = 0$, respectively. This is obtained when

$$q(x) = p(x) \quad \text{and} \quad q(y) = p(y).$$

2. The reverse KL-divergence $\text{KL}(q\|p)$ on the finite domain $x, y \in \{1, 2, 3, 4\}$ is defined as

$$\text{KL}(q\|p) = \sum_x \sum_y q(x)q(y) \log_2 \frac{q(x)q(y)}{p(x,y)}.$$

We can easily observe from the above formula that the support of q must be a subset of the support of p . In other words, if $q(x, y)$ is positive outside the support of p (i.e., when $p(x, y) = 0$) then $\text{KL}(q\|p) = \infty$. Hence, the reverse KL-divergence has an infinite value except when the support of q is either $\{1, 2\} \times \{1, 2\}$ or $\{(3, 3)\}$ or $\{(4, 4)\}$. Thus, it has three local minima.

For the first case, the minimum is achieved when $q(x) = q(y) = (\frac{1}{2}, \frac{1}{2}, 0, 0)$. The corresponding KL-divergence is $\text{KL}(q\|p) = \log_2 2 = 1$. For the second case and the third case, $q(x) = q(y) = (0, 0, 1, 0)$ and $q(x) = q(y) = (0, 0, 0, 1)$, respectively. The KL-divergence in both cases is $\text{KL}(q\|p) = \log_2 4 = 2$.

3. Let us compute $p(x = 4)$ and $p(y = 1)$:

$$p(x = 4) = \sum_y p(x = 4, y) = \frac{1}{4},$$

$$p(y = 1) = \sum_x p(x, y = 1) = \frac{1}{4}.$$

Hence, $q(x = 4, y = 1) = p(x = 4)p(y = 1) = \frac{1}{16}$, however, $p(x = 4, y = 1) = 0$. We therefore have for the reverse KL-divergence that $\text{KL}(q\|p) = \infty$.

Solution to exercise 5.29 (Gaussian VI vs Laplace approximation).

- Recall from section 5.4.6 that q^* matches the first and second *moments* of p . In contrast, the Laplace approximation matches the *mode* of p and the second derivative of $-\log p$. In general, the mean is different from the mode, and so is the second moment from the second derivative.
- We have

$$\begin{aligned} & \arg \min_{q \in \mathcal{Q}} \text{KL}(q\|p(\cdot \mid \mathcal{D})) \\ &= \arg \max_{q \in \mathcal{Q}} L(q, p; \mathcal{D}) \end{aligned}$$

using eq. (5.62)

$$= \arg \max_{q \in \mathcal{Q}} \mathbb{E}_{\theta \sim q} [\log p(y_{1:n}, \theta \mid \mathbf{x}_{1:n})] + H[q] \quad \text{using eq. (5.63a)}$$

$$= \arg \max_{q \in \mathcal{Q}} \mathbb{E}_{\theta \sim q} [\log p(y_{1:n}, \theta \mid \mathbf{x}_{1:n})] + \frac{n}{2} \log(2\pi e) + \frac{1}{2} \log \det(\Sigma). \quad \text{using eq. (5.36)}$$

Differentiating with respect to μ and Σ , we have that \tilde{q} must satisfy

$$\begin{aligned} \mathbf{0} &= \nabla_{\mu} \mathbb{E}_{\theta \sim \tilde{q}} [\log p(y_{1:n}, \theta \mid \mathbf{x}_{1:n})] \\ \Sigma^{-1} &= -2 \nabla_{\Sigma} \mathbb{E}_{\theta \sim \tilde{q}} [\log p(y_{1:n}, \theta \mid \mathbf{x}_{1:n})]. \end{aligned} \quad \text{using the first hint}$$

The result follows by eq. (5.72).

Solution to exercise 5.30 (Gradient of reverse-KL). To simplify the notation, we write $\Sigma \doteq \text{diag}\{\sigma_1^2, \dots, \sigma_d^2\}$. The reverse KL-divergence can be expressed as

$$\begin{aligned} \text{KL}(q_{\lambda} \| p(\cdot)) &= \frac{1}{2} \left(\text{tr}(\sigma_p^{-2} \Sigma) + \sigma_p^{-2} \mu^{\top} \mu - d + \log \frac{(\sigma_p^2)^d}{\det(\Sigma)} \right) \\ &= \frac{1}{2} \left(\sigma_p^{-2} \sum_{i=1}^d \sigma_i^2 + \sigma_p^{-2} \mu^{\top} \mu - d + d \log \sigma_p^2 - \sum_{i=1}^d \log \sigma_i^2 \right). \end{aligned} \quad \text{using the expression for the KL-divergence of Gaussians (5.50)}$$

It follows immediately that $\nabla_{\mu} \text{KL}(q_{\lambda} \| p(\cdot)) = \sigma_p^{-2} \mu$. Moreover,

$$\frac{\partial}{\partial \sigma_i} \text{KL}(q_{\lambda} \| p(\cdot)) = \frac{1}{2} \left(\underbrace{\sigma_p^{-2} \frac{\partial}{\partial \sigma_i} \sigma_i^2}_{2\sigma_i} - \underbrace{\frac{\partial}{\partial \sigma_i} \log \sigma_i^2}_{2/\sigma_i} \right) = \frac{\sigma_i}{\sigma_p^2} - \frac{1}{\sigma_i}.$$

Solution to exercise 5.33 (Reparameterizable distributions).

1. Let $Y \sim \text{Unif}([0, 1])$. Then, using the two hints,

$$(b - a)Y + a \sim \text{Unif}([a, b]).$$

2. Let $Z_1 \sim \mathcal{N}(\mu, \sigma^2)$ and $Z_2 \sim \mathcal{N}(0, 1)$. We have, $X = e^{Z_1}$. Recall from eq. (1.125) that Z_1 can equivalently be expressed in terms of Z_2 as $Z_1 = \sigma Z_2 + \mu$. This yields,

$$X = e^{Z_1} = e^{\sigma Z_2 + \mu}.$$

3. Denote by F the CDF of $\text{Cauchy}(0, 1)$. Observe that F is invertible with inverse $F^{-1}(y) = \tan(\pi(y - \frac{1}{2}))$. Let $Y \sim \text{Unif}([0, 1])$ and write $X = F^{-1}(Y)$. Then,

$$\begin{aligned} P_X(x) &= \mathbb{P}(X \leq x) \\ &= \mathbb{P}(F^{-1}(Y) \leq x) \\ &= \mathbb{P}(Y \leq F(x)) \\ &= F(x). \end{aligned}$$

This reparameterization works for any distribution with invertible CDF (not just Cauchy) and is known as the *universality of the uniform* (cf. section 1.1.5). The universality of the uniform is commonly used in pseudo-random number generators as it allows “lifting” samples from a uniform distribution to countless other well-known distributions.

4. The derivative of $\text{ReLU}(z)$ is $\mathbb{1}\{z > 0\}$. Applying the reparameterization trick gives

$$\begin{aligned}
 \frac{d}{d\mu} \mathbb{E}_{x \sim \mathcal{N}(\mu, 1)} \text{ReLU}(wx) &= \mathbb{E}_{\epsilon \sim \mathcal{N}(0, 1)} \frac{d}{d\mu} \text{ReLU}(w(\mu + \epsilon)) \\
 &= w \mathbb{E}_{\epsilon} \mathbb{1}\{w(\mu + \epsilon) > 0\} && \text{using the chain rule} \\
 &= w \mathbb{E}_{\epsilon} \mathbb{1}\{\mu + \epsilon > 0\} \\
 &= w \mathbb{P}_{\epsilon}(\mu + \epsilon > 0) \\
 &= w \mathbb{P}_{\epsilon}(\epsilon < \mu) \\
 &= w \mathbb{P}_{\epsilon}(\epsilon \leq \mu) \\
 &= w \Phi(\mu).
 \end{aligned}$$

Markov Chain Monte Carlo Methods

Solution to exercise 6.3 (Markov chain update). We have

$$q_{t+1}(x') = \mathbb{P}(X_{t+1} = x') = \sum_x \underbrace{\mathbb{P}(X_t = x)}_{q_t(x)} \underbrace{\mathbb{P}(X_{t+1} = x' \mid X_t = x)}_{p(x'|x)}.$$

using the sum rule (1.10)

Noting that $p(x' \mid x) = P(x, x')$, we conclude $q_{t+1} = q_t P$.

Solution to exercise 6.4 (k -step transitions). It follows directly from the definition of matrix multiplication that

$$\begin{aligned}
 P^k(x, x') &= \sum_{x_1, \dots, x_{k-1}} P(x, x_1) \cdot P(x_1, x_2) \cdots P(x_{k-1}, x') \\
 &= \sum_{x_1, \dots, x_{k-1}} \mathbb{P}(X_1 = x_1 \mid X_0 = x) \cdots \mathbb{P}(X_k = x' \mid X_{k-1} = x_{k-1}) && \text{using the definition of the transition matrix (6.8)} \\
 &= \sum_{x_1, \dots, x_{k-1}} \mathbb{P}(X_1 = x_1, \dots, X_{k-1} = x_{k-1}, X_k = x' \mid X_0 = x) && \text{using the product rule (1.14)} \\
 &= \mathbb{P}(X_k = x' \mid X_0 = x). && \text{using the sum rule (1.10)}
 \end{aligned}$$

Solution to exercise 6.14 (Finding stationary distributions). We consider the transition matrix

$$P = \begin{bmatrix} 0.60 & 0.30 & 0.10 \\ 0.50 & 0.25 & 0.25 \\ 0.20 & 0.40 & 0.40 \end{bmatrix}.$$

We note that the entries of P are all different from 0, thus the Markov chain corresponding to this transition matrix is ergodic.³ Thus, there exists a unique stationary distribution π to which the Markov chain converges irrespectively of the distribution over initial states q_0 .

³ All elements of the transition matrix being strictly greater than 0 is a sufficient, but not necessary, condition for ergodicity.

We know that $P^\top \pi = \pi$ (where we write π as a column vector), therefore, to find the stationary distribution π , we need to find the normalized eigenvector associated with eigenvalue 1 of the matrix P^\top . That is, we want to solve $(P^\top - I)\pi = \mathbf{0}$ for π . We obtain the linear system of equations,

$$\begin{aligned} -0.40\pi_1 + 0.50\pi_2 + 0.20\pi_3 &= 0 \\ 0.30\pi_1 - 0.75\pi_2 + 0.40\pi_3 &= 0 \\ 0.10\pi_1 + 0.25\pi_2 - 0.60\pi_3 &= 0. \end{aligned}$$

Note that the left hand side of equation i corresponds to the probability of entering state i at stationarity minus π_i . Quite intuitively, this difference should be 0, that is, after one iteration the random walk is at state i with the same probability as before the iteration.

Solving this system of equations, for example, using the Gaussian elimination algorithm, we obtain the normalized eigenvector

$$\pi = \frac{1}{72} \begin{bmatrix} 35 \\ 22 \\ 15 \end{bmatrix}.$$

Thus, we conclude that in the long run, the percentage of news days that will be classified as “good” is $35/72$.

Solution to exercise 6.23 (Practical examples of Gibbs sampling).

1. We have to compute the conditional distributions. Notice that for $x \in \{0, \dots, n\}$ and $y \in [0, 1]$,

$$p(x, y) = \binom{n}{x} y^x (1-y)^{n-x} \cdot y^{\alpha-1} (1-y)^{\beta-1} = \text{Bin}(x; n, y) \cdot C_y,$$

where $\text{Bin}(n, y)$ is the PMF of the binomial distribution (1.66) with n trials and success probability y , and C_y is a constant depending on y . It is clear that

$$\begin{aligned} p(x | y) &= \frac{p(x, y)}{p(y)} \\ &= \frac{\text{Bin}(x; n, y) \cdot C_y}{p(y)} \\ &= \text{Bin}(x; n, y). \end{aligned}$$

using the definition of conditional probability (1.11)

using that $p(x | y)$ is a probability distribution over x and $\text{Bin}(x; n, y)$ already sums to 1, so $C_y = p(y)$

So in short, sampling from $p(x | y)$ is equivalent to sampling from a binomial distribution, which amounts to n times throwing a coin with bias y , and outputting the number of heads.

For the other conditional distribution, recall the PMF of the beta distribution with parameters α, β ,

$$\text{Beta}(y; \alpha, \beta) = C \cdot y^{\alpha-1} (1-y)^{\beta-1}$$

where C is some constant depending on α and β only. We then have

$$p(x, y) = \text{Beta}(y; x + \alpha, n - x + \beta) \cdot C_x,$$

where C_x is some constant depending on x, α, β . This shows (analogously to above), that

$$p(y | x) = \text{Beta}(y; x + \alpha, n - x + \beta).$$

So for sampling y given x , one can sample from the beta distribution. There are several methods for sampling from a beta distribution, and we refer the reader to the corresponding Wikipedia page.

2. We first derive the posterior distribution of μ . We have

$$\begin{aligned} \log p(\mu | \lambda, x_{1:n}) &= \log p(\mu) + \log p(x_{1:n} | \mu, \lambda) + \text{const} \\ &= -\frac{\lambda_0}{2}(\mu - \mu_0)^2 - \frac{\lambda}{2} \sum_{i=1}^n (x_i - \mu)^2 + \text{const} \\ &= -\frac{1}{2} \underbrace{(\lambda_0 + n\lambda)}_{l_\lambda} \mu^2 + \underbrace{\left(\lambda_0 \mu_0 + \lambda \sum_{i=1}^n x_i \right)}_{m_\lambda l_\lambda} \mu + \text{const}. \end{aligned} \quad \text{by expanding the squares}$$

It follows that $l_\lambda = \lambda_0 + n\lambda$, $m_\lambda = (\lambda_0 \mu_0 + \lambda \sum_{i=1}^n x_i) / l_\lambda$, and $\mu | \lambda, x_{1:n} \sim \mathcal{N}(\mu_\lambda, l_\lambda^{-1})$. That is, the posterior precision is the sum of prior precision and the precisions of each observation, and the posterior mean is a weighted average of prior mean and observations (where the weights are the precisions).

For the posterior distribution of λ , we have

$$\begin{aligned} p(\lambda | \mu, x_{1:n}) &\propto p(\lambda) \cdot p(x_{1:n} | \mu, \lambda) \\ &= \lambda^{\alpha-1} e^{-\beta\lambda} \cdot \lambda^{\frac{n}{2}} e^{-\frac{\lambda}{2} \sum_{i=1}^n (x_i - \mu)^2} \\ &= \lambda^{\alpha + \frac{n}{2} - 1} e^{-\lambda(\beta + \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2)} \\ &= \text{Gamma}(\lambda; a_\mu, b_\mu) \end{aligned}$$

where $a_\mu = \alpha + \frac{n}{2}$ and $b_\mu = \beta + \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2$.

3. We have

$$p(\alpha, c | x_{1:n}) \propto p(\alpha, c) \cdot p(x_{1:n} | \alpha, c)$$

$$\begin{aligned}
& \propto \mathbb{1}\{\alpha, c > 0\} \prod_{i=1}^n \frac{\alpha c^\alpha}{x_i^{\alpha+1}} \mathbb{1}\{x_i \geq c\} \\
& = \frac{\alpha^n c^{n\alpha}}{(\prod_{i=1}^n x_i)^{\alpha+1}} \mathbb{1}\{c < x^*\} \mathbb{1}\{\alpha, c > 0\}
\end{aligned}$$

where $x^* \doteq \min\{x_1, \dots, x_n\}$. It is not clear how one could sample from this distribution directly. Instead, we use Gibbs sampling.

For the posterior distribution of α , we have

$$\begin{aligned}
p(\alpha \mid c, x_{1:n}) & \propto p(\alpha, c \mid x_{1:n}) \\
& \propto \frac{\alpha^n c^{n\alpha}}{(\prod_{i=1}^n x_i)^{\alpha+1}} \mathbb{1}\{\alpha > 0\} \\
& = \alpha^n \exp\left(-\alpha \left(\sum_{i=1}^n \log x_i - n \log c\right)\right) \mathbb{1}\{\alpha > 0\} \\
& \propto \text{Gamma}(\alpha; a, b)
\end{aligned}$$

where $a \doteq n + 1$ and $b \doteq \sum_{i=1}^n \log x_i - n \log c$. On the other hand,

$$p(c \mid \alpha, x_{1:n}) \propto p(\alpha, c \mid x_{1:n}) \propto c^{n\alpha} \mathbb{1}\{0 < c < x^*\}.$$

It remains to show that it is easy to sample from a random variable X with $p_X(x) \propto x^a \mathbb{1}\{0 < x < b\}$ for $a, b > 0$. We have that the normalizing constant is given by $\int_0^b x^a dx = \frac{b^{a+1}}{a+1}$. Therefore, the CDF of X is

$$\begin{aligned}
P_X(x) & = \int_0^x p(y) dy \\
& = \frac{a+1}{b^{a+1}} \int_0^x y^a dy \\
& = \frac{a+1}{b^{a+1}} \frac{x^{a+1}}{a+1} \\
& = \left(\frac{x}{b}\right)^{a+1}.
\end{aligned}$$

The inverse CDF is given by $P_X^{-1}(y) = by^{\frac{1}{a+1}}$. Therefore, we can sample from X using inverse transform sampling (cf. section 1.1.5) by sampling $Y \sim \text{Unif}([0, 1])$ and setting $X = bY^{\frac{1}{a+1}}$.

Solution to exercise 6.25 (Maximum entropy property of Gibbs distribution).

1. Our goal is to solve the optimization problem

$$\begin{aligned}
& \max_{p \in \Delta^{\mathcal{T}}} - \sum_{x \in \mathcal{T}} p(x) \log_2 p(x) \\
& \text{subject to } \sum_{x \in \mathcal{T}} p(x) f(x) = \mu
\end{aligned} \tag{B.8}$$

for some $\mu \in \mathbb{R}$. The Lagrangian with dual variables λ_0 and λ_1 is given by

$$\begin{aligned} L(p, \lambda_0, \lambda_1) &= - \sum_{x \in \mathcal{T}} p(x) \log_2 p(x) + \lambda_0 \left(1 - \sum_{x \in \mathcal{T}} p(x) \right) + \\ &\quad \lambda_1 \left(\mu - \sum_{x \in \mathcal{T}} p(x) f(x) \right) \\ &= - \sum_{x \in \mathcal{T}} p(x) (\log_2 p(x) + \lambda_0 + \lambda_1 f(x)) + \text{const.} \end{aligned}$$

Note that $L(p, \lambda_0, \lambda_1) = H[p]$ if the constraints are satisfied. Thus, we simply need to solve the (dual) optimization problem

$$\max_{p \geq 0} \min_{\lambda_0, \lambda_1 \in \mathbb{R}} L(p, \lambda_0, \lambda_1).$$

We have

$$\frac{\partial}{\partial p(x)} L(p, \lambda_0, \lambda_1) = -\log_2 p(x) - \lambda_0 - \lambda_1 f(x) - 1.$$

Setting the partial derivatives to zero, we obtain

$$\begin{aligned} p(x) &= 2 \exp(-\lambda_0 - \lambda_1 f(x) - 1) \\ &\propto \exp(-\lambda_1 f(x)). \end{aligned} \quad \text{using } \log_2(\cdot) = \frac{\log(\cdot)}{\log(2)}$$

Clearly, p is a valid probability mass function when normalized (i.e., for an appropriate choice of λ_0). We complete the proof by setting $T \doteq \frac{1}{\lambda_1}$.

2. As $T \rightarrow \infty$ (and $\lambda_1 \rightarrow 0$), the optimization problem reduces to picking the maximum entropy distribution without the first-moment constraint. This distribution is the uniform distribution over \mathcal{T} . Conversely, as $T \rightarrow 0$ (and $\lambda_1 \rightarrow \infty$), the Gibbs distribution reduces to a point density around its mode.

Solution to exercise 6.26 (Energy function of Bayesian logistic regression). First, note that the sum of convex functions is convex, hence, we consider each term individually.

The Hessian of the regularization term is λI , and thus, by the second-order characterization of convexity, this term is convex in w .

Finally, note that the second term is a sum of logistic losses ℓ_{\log} (5.14), and we have seen in exercise 5.3 that ℓ_{\log} is convex in w .

Solution to exercise 6.28 (Energy reduction of Gibbs sampling). Recall the following two facts:

1. Gibbs sampling is an instance of Metropolis-Hastings with proposal distribution

$$r(\mathbf{x}' | \mathbf{x}) = \begin{cases} p(x'_i | \mathbf{x}'_{-i}) & \text{if } \mathbf{x}' \text{ differs from } \mathbf{x} \text{ only in entry } i \\ 0 & \text{otherwise} \end{cases}$$

and acceptance distribution $\alpha(\mathbf{x}' | \mathbf{x}) = 1$.

2. The acceptance distribution of Metropolis-Hastings where the stationary distribution p is a Gibbs distribution with energy function f is

$$\alpha(\mathbf{x}' | \mathbf{x}) = \min \left\{ 1, \frac{r(\mathbf{x} | \mathbf{x}')}{r(\mathbf{x}' | \mathbf{x})} \exp(f(\mathbf{x}) - f(\mathbf{x}')) \right\}.$$

We therefore know that

$$\frac{r(\mathbf{x} | \mathbf{x}')}{r(\mathbf{x}' | \mathbf{x})} \exp(f(\mathbf{x}) - f(\mathbf{x}')) \geq 1.$$

We remark that this inequality even holds with equality using our derivation of theorem 6.21. Taking the logarithm and reorganizing the terms, we obtain

$$f(\mathbf{x}') \leq f(\mathbf{x}) + \log r(\mathbf{x} | \mathbf{x}') - \log r(\mathbf{x}' | \mathbf{x}). \quad (\text{B.9})$$

By the definition of the proposal distribution of Gibbs sampling,

$$r(\mathbf{x}' | \mathbf{x}) = p(x'_i | \mathbf{x}_{-i}) \quad \text{and} \quad r(\mathbf{x} | \mathbf{x}') = p(x_i | \mathbf{x}_{-i}). \quad \text{using } \mathbf{x}'_{-i} = \mathbf{x}_{-i}$$

Taking the expectation of eq. (B.9),

$$\begin{aligned} \mathbb{E}_{x'_i \sim p(\cdot | \mathbf{x}_{-i})} [f(\mathbf{x}')] &\leq f(\mathbf{x}) + \log p(x_i | \mathbf{x}_{-i}) - \mathbb{E}_{x'_i \sim p(\cdot | \mathbf{x}_{-i})} [\log p(x'_i | \mathbf{x}_{-i})] \\ &= f(\mathbf{x}) - \mathbb{S}[p(x_i | \mathbf{x}_{-i})] + \mathbb{H}[p(\cdot | \mathbf{x}_{-i})]. \end{aligned}$$

That is, the energy is expected to decrease if the expected surprise of the new sample $x'_i | \mathbf{x}_{-i}$ is smaller than the surprise of the current sample $x_i | \mathbf{x}_{-i}$.

Solution to exercise 6.31 (Mixing time of Langevin dynamics).

1. We take the minimum of eq. (6.52) with respect to \mathbf{y} on both sides. The minimum of the left-hand side is $f(\mathbf{0}) = 0$. To find the minimum of the right-hand side, we differentiate with respect to \mathbf{y} :

$$\frac{\partial}{\partial \mathbf{y}} \left[f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\alpha}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \right] = \nabla f(\mathbf{x}) + \alpha(\mathbf{y} - \mathbf{x}).$$

Setting the partial derivative to zero, we obtain

$$\mathbf{y} = \mathbf{x} - \frac{1}{\alpha} \nabla f(\mathbf{x}).$$

Plugging this \mathbf{y} into eq. (6.52), we have

$$\begin{aligned} 0 &\geq f(\mathbf{x}) - \frac{1}{\alpha} \|\nabla f(\mathbf{x})\|_2^2 + \frac{1}{2\alpha} \|\nabla f(\mathbf{x})\|_2^2 \\ &= f(\mathbf{x}) - \frac{1}{2\alpha} \|\nabla f(\mathbf{x})\|_2^2. \end{aligned}$$

2. Using the chain rule,

$$\frac{d}{dt}f(\mathbf{x}_t) = \nabla f(\mathbf{x}_t)^\top \frac{d}{dt}\mathbf{x}_t.$$

Note that $\frac{d}{dt}\mathbf{x}_t = -\nabla f(\mathbf{x}_t)$ by eq. (6.51), so

$$\begin{aligned} &= -\|\nabla f(\mathbf{x}_t)\|_2^2 \\ &\leq -2\alpha f(\mathbf{x}_t) \end{aligned}$$

where the last step follows from the PL-inequality (6.53).

3. Follows directly from (2) and Grönwall's inequality (6.54) by letting $g(t) = f(\mathbf{x}_t)$ and noting that $-\int_0^t 2\alpha ds = -2\alpha t$.
4. It suffices to show $\nabla q_t = q_t \nabla \log q_t$. By the chain rule,

$$\nabla \log q_t = \frac{\nabla q_t}{q_t}.$$

We obtain the desired result by rearranging the terms.

5. We have

$$\frac{d}{dt}\text{KL}(q_t\|p) = \frac{d}{dt} \int q_t \log \frac{q_t}{p} d\boldsymbol{\theta}.$$

By the chain rule,

$$= \int \frac{\partial q_t}{\partial t} \log \frac{q_t}{p} d\boldsymbol{\theta} + \int q_t \frac{\partial}{\partial t} \log \frac{q_t}{p} d\boldsymbol{\theta}.$$

For the second term, we have

$$\int q_t \frac{\partial}{\partial t} \log \frac{q_t}{p} d\boldsymbol{\theta} = \int \frac{\partial q_t}{\partial t} d\boldsymbol{\theta} = \frac{d}{dt} \underbrace{\int q_t d\boldsymbol{\theta}}_1 = 0.$$

Plugging in the Fokker-Planck equation (6.56) into the first term, we obtain

$$= \int \nabla \cdot \left(q_t \nabla \log \frac{q_t}{p} \right) \log \frac{q_t}{p} d\boldsymbol{\theta}.$$

Letting $\varphi \doteq \log \frac{q_t}{p}$ and $\mathbf{F} \doteq \nabla \varphi$, and then applying the hint, we have

$$\begin{aligned} &= \int (\nabla \cdot q_t \mathbf{F}) \varphi d\boldsymbol{\theta} \\ &= - \int q_t \|\nabla \varphi\|_2^2 d\boldsymbol{\theta} \end{aligned}$$

$$\begin{aligned}
&= -\mathbb{E}_{\theta \sim q_t} \left[\left\| \nabla \log \frac{q_t(\theta)}{p(\theta)} \right\|_2^2 \right] \\
&= -J(q_t \| p).
\end{aligned}$$

6. Noting that p satisfies the LSI with constant α and combining with (5), we have that $\frac{d}{dt} \text{KL}(q_t \| p) \leq -2\alpha \text{KL}(q_t \| p)$. Observe that this result is analogous to the result derived in (2) and the LSI can intuitively be seen as the PL-inequality, but in the space of distributions. Analogously to (3), we obtain the desired convergence result by applying Grönwall's inequality (6.54).
7. By Pinsker's inequality (6.21), $\|q_t - p\|_{\text{TV}} \leq e^{-\alpha t} \sqrt{2\text{KL}(q_0 \| p)}$. It follows from elementary algebra that $\|q_t - p\|_{\text{TV}} \leq \epsilon$ if $t \geq \frac{1}{\alpha} \log \frac{C}{\epsilon}$ where $C \doteq \sqrt{2\text{KL}(q_0 \| p)}$. Thus, $\tau_{\text{TV}}(\epsilon) \in \mathcal{O}(\log(1/\epsilon))$.

Solution to exercise 6.33 (Hamiltonian Monte Carlo).

1. We show that under the dynamics (6.62), the Hamiltonian $H(\mathbf{x}, \mathbf{y})$ is a constant. In particular, $H(\mathbf{x}', \mathbf{y}') = H(\mathbf{x}, \mathbf{y})$. This directly implies that

$$\alpha((\mathbf{x}', \mathbf{y}') \mid (\mathbf{x}, \mathbf{y})) = \min\{1, \exp(H(\mathbf{x}', \mathbf{y}') - H(\mathbf{x}, \mathbf{y}))\} = 1.$$

To see why $H(\mathbf{x}, \mathbf{y})$ is constant, we compute

$$\begin{aligned}
\frac{d}{dt} H(\mathbf{x}, \mathbf{y}) &= \nabla_{\mathbf{x}} H \cdot \frac{d\mathbf{x}}{dt} + \nabla_{\mathbf{y}} H \cdot \frac{d\mathbf{y}}{dt} && \text{using the chain rule} \\
&= \nabla_{\mathbf{x}} H \cdot \nabla_{\mathbf{y}} H - \nabla_{\mathbf{x}} H \cdot \nabla_{\mathbf{y}} H && \text{using the Hamiltonian dynamics (6.62)} \\
&= 0.
\end{aligned}$$

2. By applying one Leapfrog step, we have for \mathbf{x} ,

$$\begin{aligned}
\mathbf{x}(t + \tau) &= \mathbf{x}(t) + \tau \mathbf{y}(t + \tau/2) && \text{using eq. (6.63b)} \\
&= \mathbf{x}(t) + \tau \left(\mathbf{y}(t) - \frac{\tau}{2} \nabla_{\mathbf{x}} f(\mathbf{x}(t)) \right) && \text{using eq. (6.63a)} \\
&= \mathbf{x}(t) - \frac{\tau^2}{2} \nabla_{\mathbf{x}} f(\mathbf{x}(t)) + \tau \mathbf{y}(t).
\end{aligned}$$

Now observe that $\mathbf{y}(t)$ is a Gaussian random variable, independent of \mathbf{x} (because we sample \mathbf{y} freshly at the beginning of the L Leapfrog steps, and we are doing just one Leapfrog step). By renaming $\tau' \doteq \tau^2/2$ and $\epsilon \doteq \mathbf{y}(t)$, we get

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \tau' \nabla_{\mathbf{x}} f(\mathbf{x}_t) + \sqrt{2\tau'} \epsilon$$

which coincides with the proposal distribution of Langevin Monte Carlo (6.47).

Bayesian Deep Learning

Solution to exercise 7.2 (Softmax is a generalization of the logistic function). We have

$$\begin{aligned}\sigma_1(f) &= \frac{\exp(f_1)}{\exp(f_0) + \exp(f_1)} \\ &= \frac{1}{\exp(f_0 - f_1) + 1} \\ &= \sigma(-(f_0 - f_1)).\end{aligned}$$

using the definition of the softmax function (7.5)

multiplying by $\frac{\exp(-f_1)}{\exp(-f_1)}$

using the definition of the logistic function (5.10)

$\sigma_0(f) = 1 - \sigma(f)$ follows from the fact that $\sigma_0(f) + \sigma_1(f) = 1$.

Active Learning

Solution to exercise 8.4 (Mutual information and KL divergence).

We have

$$\begin{aligned}I(\mathbf{X}; \mathbf{Y}) &= H[\mathbf{X}] - H[\mathbf{X} | \mathbf{Y}] \\ &= \mathbb{E}_x[-\log p(x)] - \mathbb{E}_{(x,y)}[-\log p(x | y)] \\ &= \mathbb{E}_{(x,y)}[-\log p(x)] - \mathbb{E}_{(x,y)}[-\log p(x | y)] \\ &= \mathbb{E}_{(x,y)} \left[\log \frac{p(x | y)}{p(x)} \right].\end{aligned}$$

using the definition of mutual information (8.9)

using the definitions of entropy (5.32) and conditional entropy (8.2)

using the law of total expectation (1.15)

From this we get directly that

$$\begin{aligned}I(\mathbf{X}; \mathbf{Y}) &= \mathbb{E}_y \mathbb{E}_{x|y} \left[\log \frac{p(x | y)}{p(x)} \right] \\ &= \mathbb{E}_y [\text{KL}(p(x | y) \| p(x))],\end{aligned}$$

using the definition of KL-divergence (5.44)

and we also conclude

$$\begin{aligned}I(\mathbf{X}; \mathbf{Y}) &= \mathbb{E}_{(x,y)} \left[\log \frac{p(x, y)}{p(x)p(y)} \right] \\ &= \text{KL}(p(x, y) \| p(x)p(y)).\end{aligned}$$

using the definition of conditional probability (1.11)

using the definition of KL-divergence (5.44)

Solution to exercise 8.7 (Non-monotonicity of cond. mutual information). Symmetry of conditional mutual information (8.19) and eq. (8.18) give the following relationship,

$$I(\mathbf{X}; \mathbf{Y}, \mathbf{Z}) = I(\mathbf{X}; \mathbf{Y}) + I(\mathbf{X}; \mathbf{Z} | \mathbf{Y}) = I(\mathbf{X}; \mathbf{Z}) + I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}). \quad (\text{B.10})$$

1. $\mathbf{X} \perp \mathbf{Z}$ implies $I(\mathbf{X}; \mathbf{Z}) = 0$. Thus, eq. (B.10) simplifies to

$$I(\mathbf{X}; \mathbf{Y}) + I(\mathbf{X}; \mathbf{Z} | \mathbf{Y}) = I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}).$$

Using that $I(\mathbf{X}; \mathbf{Z} | \mathbf{Y}) \geq 0$, we conclude $I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) \geq I(\mathbf{X}; \mathbf{Y})$.

2. $\mathbf{X} \perp \mathbf{Z} \mid \mathbf{Y}$ implies $I(\mathbf{X}; \mathbf{Z} \mid \mathbf{Y}) = 0$. Equation (B.10) simplifies to

$$I(\mathbf{X}; \mathbf{Y}) = I(\mathbf{X}; \mathbf{Z}) + I(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z}).$$

Using that $I(\mathbf{X}; \mathbf{Z}) \geq 0$, we conclude $I(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z}) \leq I(\mathbf{X}; \mathbf{Y})$.

3. Again, eq. (B.10) simplifies to

$$I(\mathbf{X}; \mathbf{Y}) = I(\mathbf{X}; \mathbf{Z}) + I(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z}).$$

Using that $I(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z}) \geq 0$, we conclude $I(\mathbf{X}; \mathbf{Z}) \leq I(\mathbf{X}; \mathbf{Y})$.

Solution to exercise 8.8 (Interaction information).

1. Expanding the definition of interaction information, one obtains

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y}; \mathbf{Z}) &= (H[\mathbf{X}] + H[\mathbf{Y}] + H[\mathbf{Z}]) \\ &\quad - (H[\mathbf{X}, \mathbf{Y}] + H[\mathbf{X}, \mathbf{Z}] + H[\mathbf{Y}, \mathbf{Z}]) \\ &\quad + H[\mathbf{X}, \mathbf{Y}, \mathbf{Z}], \end{aligned}$$

and hence, interaction information is symmetric.

2. Conditional on either one of X_1 or X_2 , the distribution of Y remains unchanged, and hence $I(Y; X_1 \mid X_2) = I(Y; X_2 \mid X_1) = 0$. Conversely, conditional on both X_1 and X_2 , Y is fully determined, and hence $I(Y; X_1, X_2) = 1$ noting that Y encodes one bit worth of information. Thus, $I(Y; X_1; X_2) = -1$ meaning that there is synergy between X_1 and X_2 with respect to Y .

Solution to exercise 8.10 (Marginal gain of maximizing mutual information). As suggested, we derive the result in two steps.

1. First, we have

$$\begin{aligned} \Delta_I(\mathbf{x} \mid A) &= I(A \cup \{\mathbf{x}\}) - I(A) \\ &= I(f_{A \cup \{\mathbf{x}\}}; \mathbf{y}_A, \mathbf{y}_x) - I(f_A; \mathbf{y}_A) \\ &= I(f_{A \cup \{\mathbf{x}\}}; \mathbf{y}_A, \mathbf{y}_x) - I(f_{A \cup \{\mathbf{x}\}}; \mathbf{y}_A) \\ &= H[f_{A \cup \{\mathbf{x}\}} \mid \mathbf{y}_A] - H[f_{A \cup \{\mathbf{x}\}} \mid \mathbf{y}_A, \mathbf{y}_x] \\ &= I(f_{A \cup \{\mathbf{x}\}}; \mathbf{y}_x \mid \mathbf{y}_A) \\ &= I(f_A; \mathbf{y}_x \mid f_x, \mathbf{y}_A) + I(f_x; \mathbf{y}_x \mid \mathbf{y}_A) \\ &= I(f_x; \mathbf{y}_x \mid \mathbf{y}_A). \end{aligned}$$

using the definition of marginal gain (8.24)

using the definition of I (8.23)

using $\mathbf{y}_A \perp f_x \mid f_A$

using the definition of MI (8.9)

using the definition of cond. MI (8.16)

using eq. (8.18)

using $I(f_A; \mathbf{y}_x \mid f_x, \mathbf{y}_A) = 0$ as $\mathbf{y}_x \perp f_A \mid f_x$

2. For the second part, we get

$$\begin{aligned} I(f_x; \mathbf{y}_x \mid \mathbf{y}_A) &= I(\mathbf{y}_x; f_x \mid \mathbf{y}_A) \\ &= H[\mathbf{y}_x \mid \mathbf{y}_A] - H[\mathbf{y}_x \mid f_x, \mathbf{y}_A] \\ &= H[\mathbf{y}_x \mid \mathbf{y}_A] - H[\mathbf{y}_x \mid f_x] \\ &= H[\mathbf{y}_x \mid \mathbf{y}_A] - H[\epsilon_x]. \end{aligned}$$

using symmetry of conditional MI (8.19)

using the definition of cond. MI (8.16)

using that $\mathbf{y}_x \perp \epsilon_A$ so $\mathbf{y}_x \perp \mathbf{y}_A \mid f_x$

given f_x , the only randomness in \mathbf{y}_x originates from ϵ_x

Solution to exercise 8.14 (Submodularity means no synergy). We have

$$\begin{aligned} I(f_x; y_x; y_{B \setminus A} \mid y_A) &= I(f_x; y_x \mid y_A) - I(f_x; y_x \mid y_B) \\ &= \Delta_I(x \mid A) - \Delta_I(x \mid B) \\ &\geq 0 \end{aligned}$$

using the definition of interaction information (8.21)
using eq. (8.25)

where the final inequality follows from submodularity of I .

Bayesian Optimization

Solution to exercise 9.2 (Convergence to the static optimum). First, observe that

$$\begin{aligned} \lim_{T \rightarrow \infty} \frac{R_T}{T} &= \max_x f^*(x) - \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T f^*(x_t) \\ &= \max_x f^*(x) - \lim_{t \rightarrow \infty} f^*(x_t). \end{aligned} \quad (\text{B.11})$$

using the definition of regret (9.4)

using the Cesàro mean (9.7)

Now, suppose that the algorithm converges to the static optimum,

$$\lim_{t \rightarrow \infty} f^*(x_t) = \max_x f^*(x).$$

Together with eq. (B.11) we conclude that the algorithm achieves sub-linear regret.

For the other direction, we prove the contrapositive. That is, we assume that the algorithm does not converge to the static optimum and show that it has (super-)linear regret. We distinguish between two cases.

1. First, suppose that $\lim_{t \rightarrow \infty} f^*(x_t)$ does not exist. Then, a sub-optimal x_t is sampled infinitely often, and hence, R_T grows at least linearly with T .
2. On the other hand, if $\lim_{t \rightarrow \infty} f^*(x_t)$ exists, our assumption is formalized by

$$\lim_{t \rightarrow \infty} f^*(x_t) < \max_x f^*(x).$$

Together with eq. (B.11) we conclude $\lim_{T \rightarrow \infty} R_T/T > 0$.

Solution to exercise 9.5 (Bayesian confidence intervals).

1. As suggested by the hint, we let $Z \sim \mathcal{N}(0, 1)$ and $c > 0$, and bound

$$\begin{aligned} \mathbb{P}(Z > c) &= \mathbb{P}(Z - c > 0) \\ &= \int_0^\infty \frac{1}{\sqrt{2\pi}} e^{-(z+c)^2/2} dz. \end{aligned}$$

using the PDF of the univariate normal distribution (1.6)

Note that for $z \geq 0$,

$$\frac{(z+c)^2}{2} = \frac{z^2}{2} + zc + \frac{c^2}{2} \geq \frac{z^2}{2} + \frac{c^2}{2}.$$

Thus,

$$\begin{aligned} &\leq e^{-c^2/2} \int_0^\infty \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz \\ &= e^{-c^2/2} \mathbb{P}(Z > 0) \\ &= \frac{1}{2} e^{-c^2/2}. \end{aligned} \tag{B.12}$$

using symmetry of the standard normal distribution around 0

Since we made the Bayesian assumption $f^*(\mathbf{x}) \sim \mathcal{N}(\mu_0(\mathbf{x}), \sigma_0^2(\mathbf{x}))$ and assumed Gaussian observation noise $y_t \sim \mathcal{N}(f^*(\mathbf{x}_t), \sigma_n^2)$, the posterior is also Gaussian:

$$f^*(\mathbf{x}) \mid \mathbf{x}_{1:t}, y_{1:t} \sim \mathcal{N}(\mu_t(\mathbf{x}), \sigma_t^2(\mathbf{x})).$$

Hence, writing $\mathbb{P}_t(\cdot) \doteq \mathbb{P}(\cdot \mid \mathbf{x}_{1:t}, y_{1:t})$,

$$\begin{aligned} \mathbb{P}_{t-1}(f^*(\mathbf{x}) \notin \mathcal{C}_t(\mathbf{x})) &= \mathbb{P}_{t-1}(|f^*(\mathbf{x}) - \mu_{t-1}(\mathbf{x})| > \beta_t \sigma_{t-1}(\mathbf{x})) \\ &= 2\mathbb{P}_{t-1}\left(\frac{f^*(\mathbf{x}) - \mu_{t-1}(\mathbf{x})}{\sigma_{t-1}(\mathbf{x})} > \beta_t\right) \\ &\leq e^{-\beta_t^2/2}. \end{aligned}$$

using symmetry of the Gaussian distribution

using eq. (B.12) with $c = \beta_t$

2. We have

$$\begin{aligned} \mathbb{P}\left(\bigcup_{\mathbf{x} \in \mathcal{X}} \bigcup_{t \geq 1} f^*(\mathbf{x}) \notin \mathcal{C}_t(\mathbf{x})\right) &\leq \sum_{\mathbf{x} \in \mathcal{X}} \sum_{t \geq 1} \mathbb{P}(f^*(\mathbf{x}) \notin \mathcal{C}_t(\mathbf{x})) \\ &\leq |\mathcal{X}| \sum_{t \geq 1} e^{-\beta_t^2/2}. \end{aligned}$$

using a union bound (1.1)

using (1)

Letting $\beta_t^2 \doteq 2 \log(|\mathcal{X}| (\pi t)^2 / (6\delta))$, we get

$$\begin{aligned} &= \frac{6\delta}{\pi^2} \sum_{t \geq 1} \frac{1}{t^2} \\ &= \delta. \end{aligned}$$

using $\sum_{t \geq 1} \frac{1}{t^2} = \frac{\pi^2}{6}$

Solution to exercise 9.7 (Regret of GP-UCB).

1. We denote the static optimum by \mathbf{x}^* . By the definition of \mathbf{x}_t ,

$$\begin{aligned} \mu_{t-1}(\mathbf{x}_t) + \beta_t \sigma_{t-1}(\mathbf{x}_t) &\geq \mu_{t-1}(\mathbf{x}^*) + \beta_t \sigma_{t-1}(\mathbf{x}^*) \\ &\geq f^*(\mathbf{x}^*). \end{aligned}$$

using eq. (9.9)

Thus,

$$\begin{aligned} r_t &= f^*(\mathbf{x}^*) - f^*(\mathbf{x}_t) \\ &\leq \beta_t \sigma_{t-1}(\mathbf{x}_t) + \mu_{t-1}(\mathbf{x}_t) - f^*(\mathbf{x}_t) \\ &\leq 2\beta_t \sigma_{t-1}(\mathbf{x}_t). \end{aligned}$$

again using eq. (9.9)

2. We have for any fixed T ,

$$\begin{aligned}
 I(f_{T+1}; y_{T+1}) &= H[y_{T+1}] - H[\epsilon_{T+1}] \\
 &= H[y_T] - H[\epsilon_T] + H[y_{x_{T+1}} | y_T] - H[\epsilon_{x_{T+1}}] \\
 &= I(f_T; y_T) + I(f_{x_{T+1}}; y_{x_{T+1}} | y_T) \\
 &= I(f_T; y_T) + \frac{1}{2} \log \left(1 + \frac{\sigma_T^2(x_{T+1})}{\sigma_n^2} \right).
 \end{aligned}$$

analogously to eq. (8.26)

using the chain rule for entropy (8.5)
and the mutual independence of ϵ_{T+1}
using the definition of MI (8.9)

using eq. (8.15)

Note that $I(f_0; y_0) = 0$. The result then follows by induction.

3. By Cauchy-Schwarz, $R_T^2 \leq T \sum_{t=1}^T r_t^2$, and hence, it suffices to show $\sum_{t=1}^T r_t^2 \leq \mathcal{O}(\beta_T^2 \gamma_T)$. We have

$$\begin{aligned}
 \sum_{t=1}^T r_t^2 &\leq 4\beta_T^2 \sum_{t=1}^T \sigma_{t-1}^2(x_t) \\
 &= 4\sigma_n^2 \beta_T^2 \sum_{t=1}^T \frac{\sigma_{t-1}^2(x_t)}{\sigma_n^2}.
 \end{aligned}$$

using part (1)

Observe that $\sigma_{t-1}^2(x_t)/\sigma_n^2$ is bounded by $M \doteq \max_{x \in \mathcal{X}} \sigma_0^2(x)/\sigma_n^2$ as variance is monotonically decreasing (cf. section 1.6). Applying the hint, we obtain

$$\begin{aligned}
 &\leq 4C\sigma_n^2 \beta_T^2 \sum_{t=1}^T \log \left(1 + \frac{\sigma_{t-1}^2(x_t)}{\sigma_n^2} \right) \\
 &= 8C\sigma_n^2 \beta_T^2 I(f_T; y_T) \\
 &\leq 8C\sigma_n^2 \beta_T^2 \gamma_T.
 \end{aligned}$$

using part (2)

using the definition of γ_T (9.13)

Solution to exercise 9.9 (Sublinear regret of GP-UCB for a linear kernel).

1. Let $S \subseteq \mathcal{X}$ be such that $|S| \leq T$. Recall from eq. (8.15) that $I(f_S; y_S) = \frac{1}{2} \log \det(I + \sigma_n^{-2} K_{SS})$. Using that the kernel is linear we can rewrite $K_{SS} = X_S^\top X_S$. Using Weinstein-Aronszajn's identity (A.5) we have

$$I(f_S; y_S) = \frac{1}{2} \log \det(I + \sigma_n^{-2} X_S^\top X_S) = \frac{1}{2} \log \det(I + \sigma_n^{-2} X_S X_S^\top).$$

If we define $M \doteq I + \sigma_n^{-2} X_S X_S^\top$ as a sum of symmetric positive definite matrices, M itself is symmetric positive definite. Thus, we have from Hadamard's inequality (A.4),

$$\begin{aligned}
 \det(M) &\leq \det(\text{diag}\{I + \sigma_n^{-2} X_S X_S^\top\}) \\
 &= \det(I + \sigma_n^{-2} \text{diag}\{X_S X_S^\top\}).
 \end{aligned}$$

$\text{diag}\{A\}$ refers to the diagonal matrix whose elements are those of A

Note that

$$\text{diag}\{X_S X_S^\top\}(i, i) = \sum_{t=1}^{|S|} x_t(i)^2$$

$$\leq \sum_{i=1}^d \sum_{t=1}^{|S|} \mathbf{x}_t(i)^2 = \sum_{t=1}^{|S|} \underbrace{\|\mathbf{x}_t\|_2^2}_{\leq 1} \leq |S| \leq T.$$

If we denote by $\lambda \leq T$ the largest term of $\text{diag}\{\mathbf{X}_S \mathbf{X}_S^\top\}$ then we have

$$\det(\mathbf{M}) \leq (1 + \sigma_n^{-2} \lambda)^d \leq (1 + \sigma_n^{-2} T)^d,$$

yielding,

$$\mathbf{I}(f_S; \mathbf{y}_S) \leq \frac{d}{2} \log(1 + \sigma_n^{-2} T)$$

implying that $\gamma_T = \mathcal{O}(d \log T)$.

2. Using the regret bound (cf. theorem 9.6) and the Bayesian confidence intervals (cf. theorem 9.4), and then $\gamma_T = \mathcal{O}(d \log T)$, we have

$$R_T = \mathcal{O}(\beta_T \sqrt{\gamma_T T}) = \tilde{\mathcal{O}}(\sqrt{dT}),$$

and hence, $\lim_{T \rightarrow \infty} \frac{R_T}{T} = 0$.

Solution to exercise 9.12 (Closed-form expected improvement).

1. Note that f is a Gaussian process, and hence, our posterior distribution after round t is entirely defined by the mean function μ_t and the covariance function k_t . Reparameterizing the posterior distribution using a standard Gaussian (1.125), we obtain

$$f(\mathbf{x}) \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:t} = \mu_t(\mathbf{x}) + \sigma_t(\mathbf{x})\epsilon$$

for $\epsilon \sim \mathcal{N}(0, 1)$. We get

$$\begin{aligned} \text{EI}_t(\mathbf{x}) &= \mathbb{E}_{f(\mathbf{x}) \sim \mathcal{N}(\mu_t(\mathbf{x}), \sigma_t^2(\mathbf{x}))} [I_t(\mathbf{x})] \\ &= \mathbb{E}_{f(\mathbf{x}) \sim \mathcal{N}(\mu_t(\mathbf{x}), \sigma_t^2(\mathbf{x}))} [(f(\mathbf{x}) - \hat{f}_t)_+] \\ &= \mathbb{E}_{\epsilon \sim \mathcal{N}(0, 1)} [(\mu_t(\mathbf{x}) + \sigma_t(\mathbf{x})\epsilon - \hat{f}_t)_+] \\ &= \int_{-\infty}^{+\infty} (\mu_t(\mathbf{x}) + \sigma_t(\mathbf{x})\epsilon - \hat{f}_t)_+ \cdot \phi(\epsilon) d\epsilon. \end{aligned}$$

using the definition of expected improvement (9.23)

using the definition of improvement (9.19)

using the reparameterization

For $\epsilon < \frac{\hat{f}_t - \mu_t(\mathbf{x})}{\sigma_t(\mathbf{x})} \doteq z_t(\mathbf{x})$ we have $(\mu_t(\mathbf{x}) + \sigma_t(\mathbf{x})\epsilon - \hat{f}_t)_+ = 0$. Thus, we obtain

$$\text{EI}_t(\mathbf{x}) = \int_{z_t(\mathbf{x})}^{+\infty} (\mu_t(\mathbf{x}) + \sigma_t(\mathbf{x})\epsilon - \hat{f}_t) \cdot \phi(\epsilon) d\epsilon. \quad (\text{B.13})$$

2. By splitting the integral from eq. (B.13) into two distinct terms, we obtain

$$\begin{aligned} \text{EI}_t(\mathbf{x}) &= (\mu_t(\mathbf{x}) - \hat{f}_t) \int_{z_t(\mathbf{x})}^{+\infty} \phi(\epsilon) d\epsilon \\ &\quad - \sigma_t(\mathbf{x}) \int_{z_t(\mathbf{x})}^{+\infty} (-\epsilon) \cdot \phi(\epsilon) d\epsilon. \end{aligned}$$

For the first term, we use the symmetry of $\mathcal{N}(0, 1)$ around 0 to write the integral in terms of the CDF. For the second term, we notice that $(-\epsilon) \cdot \phi(\epsilon) = \frac{1}{\sqrt{2\pi}} \frac{d}{d\epsilon} e^{-\epsilon^2/2}$. Thus, we can derive this integral directly,

$$= (\mu_t(\mathbf{x}) - \hat{f}_t) \Phi(-z_t(\mathbf{x})) - \sigma_t(\mathbf{x}) \left(\lim_{\epsilon \rightarrow \infty} \phi(\epsilon) - \phi(z_t(\mathbf{x})) \right).$$

Using the symmetry of ϕ around 0, we obtain

$$\text{EI}_t(\mathbf{x}) = (\mu_t(\mathbf{x}) - \hat{f}_t) \Phi\left(\frac{\mu_t(\mathbf{x}) - \hat{f}_t}{\sigma_t(\mathbf{x})}\right) + \sigma_t(\mathbf{x}) \phi\left(\frac{\mu_t(\mathbf{x}) - \hat{f}_t}{\sigma_t(\mathbf{x})}\right).$$

Solution to exercise 9.16 (Regret of IDS).

1. As $\mathbf{x}_t^{\text{UCB}}$ is the UCB action,

$$\hat{\Delta}_t(\mathbf{x}_t^{\text{UCB}}) = u_t(\mathbf{x}_t^{\text{UCB}}) - l_t(\mathbf{x}_t^{\text{UCB}}) = 2\beta_{t+1}\sigma_t(\mathbf{x}_t^{\text{UCB}}).$$

2. We first bound $I_t(\mathbf{x}_t^{\text{UCB}})$:

$$\begin{aligned} I_t(\mathbf{x}_t^{\text{UCB}}) &= \mathbb{I}\left(f_{\mathbf{x}_t^{\text{UCB}}}; y_{\mathbf{x}_t^{\text{UCB}}} \mid \mathbf{x}_{1:t}, y_{1:t}\right) \\ &= \frac{1}{2} \log \left(1 + \frac{\sigma_t^2(\mathbf{x}_t^{\text{UCB}})}{\sigma_n^2} \right). \end{aligned} \quad \text{using eq. (8.15)}$$

Note that $\sigma_t^2(\mathbf{x})/\sigma_n^2 \leq C$ for some constant C since variance is decreasing monotonically. So, applying the hint,

$$\geq \frac{\sigma_t^2(\mathbf{x}_t^{\text{UCB}})}{2C\sigma_n^2}.$$

Combining this with (1), we obtain

$$\hat{\Psi}_t(\mathbf{x}_t^{\text{UCB}}) = \frac{\hat{\Delta}_t(\mathbf{x}_t^{\text{UCB}})^2}{I_t(\mathbf{x}_t^{\text{UCB}})} \leq 8C\sigma_n^2\beta_{t+1}^2.$$

3. With high probability,

$$\Psi_t(\mathbf{x}_{t+1}) \leq \hat{\Psi}_t(\mathbf{x}_{t+1}) \leq \hat{\Psi}_t(\mathbf{x}_t^{\text{UCB}}) \leq 8C\sigma_n^2\beta_{t+1}^2,$$

where the first inequality is due to $\Delta(\mathbf{x}) \leq \hat{\Delta}_t(\mathbf{x})$ with high probability, the second inequality due to the definition of the IDS algorithm (9.35), and the third inequality is from (2). Invoking theorem 9.13 with $\bar{\Psi}_T \doteq 8C\sigma_n^2\beta_T^2$, we obtain that

$$R_T \leq \sqrt{\gamma_T \bar{\Psi}_T T} = \beta_T \sqrt{8C\sigma_n^2 \gamma_T T}$$

with high probability.

Markov Decision Processes

Solution to exercise 10.8 (Value functions). We can use eq. (10.15) to write the state-action values as a linear system of equations (i.e., as a “table”). This linear system can be solved, for example, using Gaussian elimination to yield the desired result.

Solution to exercise 10.13 (Greedy policies). It follows directly from the definition of the state-action value function (10.9) that

$$\arg \max_{a \in A} q(x, a) = \arg \max_{a \in A} r(x, a) + \gamma \sum_{x' \in X} p(x' | x, a) \cdot v(x').$$

Solution to exercise 10.16 (Optimal policies).

1. Recall from Bellman’s theorem (10.32) that a policy is optimal iff it is greedy with respect to its state-action value function. Now, observe that in the “poor, unknown” state, the policy π is not greedy.
2. Analogously to exercise 10.8, we write the state-action values as a linear system of equations and solve the system using, e.g., Gaussian elimination.
3. Observe from the result of (2) that π' is greedy with respect to its state-action value function, and hence, it follows from Bellman’s theorem that π' is optimal.

Solution to exercise 10.23 (Linear convergence of policy iteration).

Using the hint and $v^* \geq v^\pi$ for any policy π ,

$$\begin{aligned} \|v^{\pi_t} - v^*\|_\infty &\leq \|B^* v^{\pi_{t-1}} - v^*\|_\infty \\ &= \|B^* v^{\pi_{t-1}} - B^* v^*\|_\infty \\ &\leq \gamma \|v^{\pi_{t-1}} - v^*\|_\infty \\ &\leq \gamma^t \|v^{\pi_0} - v^*\|_\infty. \end{aligned}$$

using that v^* is a fixed-point of B^* , that is, $B^* v^* = v^*$
using that B^* is a contraction, see theorem 10.21
by induction

Solution to exercise 10.24 (Reward modification).

1. Recall that the value function v_M^π for an MDP M is defined as $v_M^\pi(x) = \mathbb{E}_\pi[\sum_{t=0}^\infty \gamma^t R_t | X_0 = x]$. Given an optimal policy π^* for M and any policy π , we know that for any $x \in X$,

$$\begin{aligned} v_M^{\pi^*}(x) &\geq v_M^\pi(x) \\ \iff \mathbb{E}_{\pi^*}[\sum_{t=0}^\infty \gamma^t R_t | X_0 = x] &\geq \mathbb{E}_\pi[\sum_{t=0}^\infty \gamma^t R_t | X_0 = x] \\ \iff \mathbb{E}_{\pi^*}[\sum_{t=0}^\infty \gamma^t \alpha R_t | X_0 = x] &\geq \mathbb{E}_\pi[\sum_{t=0}^\infty \gamma^t \alpha R_t | X_0 = x] \\ \iff \mathbb{E}_{\pi^*}[\sum_{t=0}^\infty \gamma^t R'_t | X_0 = x] &\geq \mathbb{E}_\pi[\sum_{t=0}^\infty \gamma^t R'_t | X_0 = x] \\ \iff v_{M'}^{\pi^*}(x) &\geq v_{M'}^\pi(x). \end{aligned}$$

multiplying both sides by α

Thus, π^* is an optimal policy for M' .

2. We give an example where the optimal policies differ when rewards are shifted.

Consider an MDP with three states $\{1, 2, 3\}$ where 1 is the initial state and 3 is a terminal state. If one plays action A in states 1 or 2 one transitions directly to the terminal state. Additionally, in state 1 one can play action B which leads to state 2. Let every transition give a deterministic reward of $r \doteq -1$. Then it is optimal to traverse the shortest path to the terminal state, in particular, to choose action A when in state 1.

If we consider the reward $r' \doteq r + 2 = 1$, then it is optimal to traverse the longest path to the terminal state, in particular, to choose action B when in state 1.

3. For an MDP M , we know that its optimal state-action value function satisfies Bellman's optimality equation (10.35),

$$q_M^*(x, a) = \mathbb{E}_{x'|x, a} \left[r(x, x') + \gamma \max_{a' \in A} q_M^*(x', a') \right].$$

For the MDP M' , we have

$$\begin{aligned} q_{M'}^*(x, a) &= \mathbb{E}_{x'|x, a} \left[r'(x, x') + \gamma \max_{a' \in A} q_{M'}^*(x', a') \right] \\ &= \mathbb{E}_{x'|x, a} \left[r(x, x') + f(x, x') + \gamma \max_{a' \in A} q_{M'}^*(x', a') \right] \\ &= \mathbb{E}_{x'|x, a} \left[r(x, x') + \gamma \phi(x') - \phi(x) + \gamma \max_{a' \in A} q_{M'}^*(x', a') \right]. \end{aligned}$$

Reorganizing the terms, we obtain

$$q_{M'}^*(x, a) + \phi(x) = \mathbb{E}_{x'|x, a} \left[r(x, x') + \gamma \max_{a' \in A} (q_{M'}^*(x', a') + \phi(x')) \right].$$

If we now define $q(x, a) \doteq q_{M'}^*(x, a) + \phi(x)$, we have

$$q(x, a) = \mathbb{E}_{x'|x, a} \left[r(x, x') + \gamma \max_{a' \in A} q(x', a') \right].$$

This is exactly Bellman's optimality equation for the MDP M with reward function r , and hence, $q \equiv q_M^*$.

If we take π^* to be an optimal policy for M , then it satisfies

$$\begin{aligned} \pi^*(x) &\in \arg \max_{a \in A} q_M^*(x, a) \\ &= \arg \max_{a \in A} q_{M'}^*(x, a) + \phi(x) \\ &= \arg \max_{a \in A} q_{M'}^*(x, a). \end{aligned}$$

using the above characterization of q_M^*

using that $\phi(x)$ is independent of a

Tabular Reinforcement Learning

Solution to exercise 11.14 (Q-learning).

1. $Q^*(A, \downarrow) = 1.355$, $Q^*(G_1, \text{exit}) = 5.345$, $Q^*(G_2, \text{exit}) = 0.5$

2. Repeating the given episodes infinitely often will not lead to convergence to the optimal Q-function because not all state-action pairs are visited infinitely often.

Let us assume we observe the following episode instead of the first episode.

Episode 3			
x	a	x'	r
A	\rightarrow	B	0
B	\downarrow	G_2	0
G_2	exit		1

If we repeat episodes 2 and 3 infinitely often, Q-learning will converge to the optimal Q-function as all state-action pairs will be visited infinitely often.

3. First, recall that Q-learning is an off-policy algorithm, and hence, even if episodes are obtained off-policy, Q-learning will still converge to the optimal Q-function (if the other convergence conditions are met). Note that it only matters which state-action pairs are observed and not which policies were followed to obtain these observations.

The “closer” the initial Q-values are to the optimal Q-function, the faster the convergence of Q-learning. However, if the convergence conditions are met, Q-learning will converge to the optimal Q-function regardless of the initial Q-values.

4. $v^*(A) = 10$, $v^*(B) = 10$, $v^*(G_1) = 10$, $v^*(G_2) = 1$

Model-free Reinforcement Learning

Solution to exercise 12.2 (Q-learning and function approximation).

1. We have to show that

$$v^*(x) = \max_{a \in A} r(x, a) + \gamma \mathbb{E}_{x'|x, a} [v^*(x')]$$

for every $x \in \{1, 2, \dots, 7\}$. We give a derivation here for $x = 1$ and $x = 2$.

- For $x = 1$,

$$\begin{aligned} v^*(1) &= -3 \\ \max_{a \in A} r(1, a) + \gamma \mathbb{E}_{x'|1, a} [v^*(x')] &= -3 \end{aligned}$$

since

$$r(1, a) + \gamma \mathbb{E}_{x'|1, a} [v^*(x')] = \begin{cases} -1 + -2 = -3 & \text{if } a = 1 \\ -1 + -3 = -4 & \text{if } a = -1. \end{cases}$$

- Likewise, for $x = 2$,

$$\begin{aligned} v^*(2) &= -2 \\ \max_{a \in A} r(2, a) + \gamma \mathbb{E}_{x'|2,a} [v^*(x')] &= -2 \end{aligned}$$

since

$$r(2, a) + \gamma \mathbb{E}_{x'|2,a} [v^*(x')] = \begin{cases} -1 + -1 = -2 & \text{if } a = 1 \\ -1 + -3 = -4 & \text{if } a = -1. \end{cases}$$

2. We have

$$\begin{aligned} Q(3, -1) &= 0 + \frac{1}{2} \left(-1 + \max_{a' \in A} Q(2, a') \right) = \frac{1}{2} (-1 + 0) = -\frac{1}{2} \\ Q(2, 1) &= 0 + \frac{1}{2} \left(-1 + \max_{a' \in A} Q(3, a') \right) = \frac{1}{2} (-1 + 0) = -\frac{1}{2} \\ Q(3, 1) &= 0 + \frac{1}{2} \left(-1 + \max_{a' \in A} Q(4, a') \right) = \frac{1}{2} (-1 + 0) = -\frac{1}{2} \\ Q(4, 1) &= 0 + \frac{1}{2} \left(0 + \max_{a' \in A} Q(4, a') \right) = \frac{1}{2} (0 + 0) = 0. \end{aligned}$$

3. We compute

$$\begin{aligned} \nabla_w \ell(w; \tau) &= - \left(r + \gamma \max_{a' \in A} Q(x', a'; w^{\text{old}}) - Q(x, a; w) \right) \begin{bmatrix} x \\ a \\ 1 \end{bmatrix} && \text{using the derivation of eq. (12.15)} \\ &= - \left(-1 + \max_{a' \in A} \{1 - a' - 2\} - (-2 - 1 + 1) \right) \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} -4 \\ 2 \\ -2 \end{bmatrix}. \end{aligned}$$

This gives

$$\begin{aligned} w' &= w - \alpha \nabla_w \ell(w; \tau) \\ &= \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} - \frac{1}{4} \begin{bmatrix} -4 \\ 2 \\ -2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1/2 \\ 3/2 \end{bmatrix}. \end{aligned}$$

Solution to exercise 12.6 (Eligibility vector). We have

$$\begin{aligned} \nabla_{\phi} \log \pi_{\phi}(a | x) &= \frac{\nabla_{\phi} \pi_{\phi}(a | x)}{\pi_{\phi}(a | x)} && \text{using the chain rule} \\ &= \phi(x, a) - \frac{\sum_{b \in A} \phi(x, b) \exp(\phi^{\top} \phi(x, b))}{\sum_{b \in A} \exp(\phi^{\top} \phi(x, b))} && \text{using elementary calculus} \\ &= \phi(x, a) - \sum_{b \in A} \pi_{\phi}(b | x) \cdot \phi(x, b). \end{aligned}$$

Solution to exercise 12.8 (Variance of score gradients with baselines).

1. The result follows directly using that

$$\text{Var}[f(\mathbf{X}) - g(\mathbf{X})] = \text{Var}[f(\mathbf{X})] + \text{Var}[g(\mathbf{X})] - 2\text{Cov}[f(\mathbf{X}), g(\mathbf{X})]. \quad \text{using eq. (1.48)}$$

2. Denote by $r(\tau)$ the discounted rewards attained by trajectory τ . Let $f(\tau) \doteq r(\tau) \nabla_{\boldsymbol{\varphi}} \Pi_{\boldsymbol{\varphi}}(\tau)$ and $g(\tau) \doteq b \nabla_{\boldsymbol{\varphi}} \Pi_{\boldsymbol{\varphi}}(\tau)$. Recall that $\mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}}[g(\tau)] = 0$, implying that

$$\begin{aligned} \text{Var}[g(\tau)] &= \text{Var}[b \nabla_{\boldsymbol{\varphi}} \Pi_{\boldsymbol{\varphi}}(\tau)] \\ &= \mathbb{E}[(b \nabla_{\boldsymbol{\varphi}} \Pi_{\boldsymbol{\varphi}}(\tau))^2]. \end{aligned} \quad \text{using the definition of variance (1.44)}$$

On the other hand,

$$\begin{aligned} \text{Cov}[f(\tau), g(\tau)] &= \mathbb{E}[(f(\tau) - \mathbb{E}[f(\tau)])g(\tau)] \\ &= \mathbb{E}[f(\tau)g(\tau)] - \mathbb{E}[f(\tau)] \underbrace{\mathbb{E}[g(\tau)]}_0 \\ &= \mathbb{E}[b \cdot r(\tau) \cdot (\nabla_{\boldsymbol{\varphi}} \Pi_{\boldsymbol{\varphi}}(\tau))^2]. \end{aligned} \quad \begin{array}{l} \text{using the definition of covariance (1.35)} \\ \text{using linearity of expectation (1.24)} \end{array}$$

Therefore, 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 the result follows from eq. (12.40).

Solution to exercise 12.9 (Score gradients with state-dependent baselines). First, observe that

$$\mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}}[G_0 \nabla_{\boldsymbol{\varphi}} \log \Pi_{\boldsymbol{\varphi}}(\tau)] = \mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}}\left[\sum_{t=0}^{T-1} G_0 \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t \mid \mathbf{x}_t)\right], \quad \text{using eq. (12.30)}$$

and hence, it suffices to show

$$\begin{aligned} &\mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}}\left[\sum_{t=0}^{T-1} G_0 \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t \mid \mathbf{x}_t)\right] \\ &= \mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}}\left[\sum_{t=0}^{T-1} (G_0 - b(\tau_{0:t-1})) \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t \mid \mathbf{x}_t)\right]. \end{aligned} \quad (\text{B.14})$$

We prove eq. (B.14) with an induction on T . The base case ($T = 0$) is satisfied trivially. Fixing any T and assuming eq. (B.14) holds for T , we have,

$$\begin{aligned} &\mathbb{E}_{\tau \sim \Pi_{\boldsymbol{\varphi}}}\left[\sum_{t=0}^T (G_0 - b(\tau_{0:t-1})) \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t \mid \mathbf{x}_t)\right] \\ &= \mathbb{E}_{\tau_{0:T-1}}\left[\mathbb{E}_{\tau_T}\left[\sum_{t=0}^T (G_0 - b(\tau_{0:t-1})) \nabla_{\boldsymbol{\varphi}} \log \pi_{\boldsymbol{\varphi}}(\mathbf{a}_t \mid \mathbf{x}_t) \mid \tau_{0:T-1}\right]\right] \end{aligned} \quad \text{using the tower rule (1.33)}$$

$$\begin{aligned}
&= \mathbb{E}_{\tau_{0:T-1}} \left[\sum_{t=0}^{T-1} G_0 \nabla_{\varphi} \log \pi_{\varphi}(a_t | x_t) \right. \\
&\quad \left. + \mathbb{E}_{\tau_T} [(G_0 - b(\tau_{0:T-1})) \nabla_{\varphi} \log \pi_{\varphi}(a_T | x_T) | \tau_{0:T-1}] \right].
\end{aligned}$$

using the induction hypothesis

Using the score function trick for the score function $\nabla_{\varphi} \log \pi_{\varphi}(a_t | x_t)$ analogously to the proof of lemma 12.7, we have,

$$\begin{aligned}
&\mathbb{E}_{\tau_T} [b(\tau_{0:T-1}) \nabla_{\varphi} \log \pi_{\varphi}(a_T | x_T) | \tau_{0:T-1}] \\
&= \mathbb{E}_{a_T} [b(\tau_{0:T-1}) \nabla_{\varphi} \log \pi_{\varphi}(a_T | x_T) | \tau_{0:T-1}] \\
&= b(\tau_{0:T-1}) \int \pi_{\varphi}(a_T | x_T) \nabla_{\varphi} \log \pi_{\varphi}(a_T | x_T) da_T \\
&= b(\tau_{0:T-1}) \int \nabla_{\varphi} \pi_{\varphi}(a_T | x_T) da_T \\
&= 0.
\end{aligned}$$

using the score function trick,

$$\nabla_{\varphi} \log \pi_{\varphi}(a_t | x_t) = \nabla_{\varphi} \pi_{\varphi}(a_t | x_t) / \pi_{\varphi}(a_t | x_t)$$

Thus,

$$\begin{aligned}
&= \mathbb{E}_{\tau_{0:T-1}} \left[\mathbb{E}_{\tau_T} \left[\sum_{t=0}^T G_0 \nabla_{\varphi} \log \pi_{\varphi}(a_t | x_t) \mid \tau_{0:T-1} \right] \right] \\
&= \mathbb{E}_{\tau \sim \Pi_{\varphi}} \left[\sum_{t=0}^T G_0 \nabla_{\varphi} \log \pi_{\varphi}(a_t | x_t) \right].
\end{aligned}$$

using the tower rule (1.33) again

Solution to exercise 12.12 (Policy gradients with downstream returns). Each trajectory τ is described by four transitions,

$$\tau = (x_0, a_0, r_0, x_1, a_1, r_1, x_2, a_2, r_2, x_3, a_3, r_3, x_4).$$

Moreover, we have given

$$\begin{aligned}
\pi_{\theta}(2 | x) &= \theta, & \frac{\partial \pi_{\theta}(2 | x)}{\partial \theta} &= +1 \\
\pi_{\theta}(1 | x) &= 1 - \theta, & \frac{\partial \pi_{\theta}(1 | x)}{\partial \theta} &= -1.
\end{aligned}$$

We first compute the downstream returns for the given episode,

$$\begin{aligned}
G_{0:4} &= r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 = 1 + \frac{1}{2} \cdot 0 + \frac{1}{4} \cdot 1 + \frac{1}{8} \cdot 1 = \frac{11}{8} \\
G_{1:4} &= r_1 + \gamma r_2 + \gamma^2 r_3 = 0 + \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 1 = \frac{3}{4} \\
G_{2:4} &= r_2 + \gamma r_3 = 1 + \frac{1}{2} \cdot 1 = \frac{3}{2} \\
G_{3:4} &= r_3 = 1.
\end{aligned}$$

Lastly, we can combine them to compute the policy gradient,

$$\nabla_{\theta} j(\theta) \approx \sum_{t=0}^3 \gamma^t G_{t:4} \nabla_{\theta} \log \pi_{\theta}(a_t | x_t)$$

using Monte Carlo approximation of eq. (12.45) with a single sample

$$= 1 \cdot \frac{11}{8} - 1 \cdot \frac{1}{2} \cdot \frac{3}{4} + 1 \cdot \frac{1}{4} \cdot \frac{3}{2} - 1 \cdot \frac{1}{8} \cdot 1 = \frac{5}{4}.$$

Solution to exercise 12.15 (Policy gradient with an exponential family).

1. The distribution on actions is given by

$$\pi_{\varphi}(a | x) = \sigma(f_{\varphi}(x))^a \cdot (1 - \sigma(f_{\varphi}(x)))^{1-a}.$$

To simplify the notation, we write \mathbb{E} for $\mathbb{E}_{(x,a) \sim \pi_{\varphi}}$ and q for $q^{\pi_{\varphi}}$. We get

$$\begin{aligned} \nabla_{\varphi} j(\varphi) &= \mathbb{E} [q(x, a) \nabla_{\varphi} (a \log \sigma(f_{\varphi}(x)) + (1 - a) \log(1 - \sigma(f_{\varphi}(x))))] && \text{using the policy gradient theorem (12.58)} \\ &= \mathbb{E} \left[q(x, a) \nabla_f (a \log \sigma(f) + (1 - a) \log(1 - \sigma(f))) \nabla_{\varphi} f_{\varphi}(x) \right] && \text{using the chain rule} \\ &= \mathbb{E} \left[q(x, a) \nabla_f \left(-a \log(1 + e^{-f}) + (1 - a) \log \left(\frac{e^{-f}}{1 + e^{-f}} \right) \right) \nabla_{\varphi} f_{\varphi}(x) \right] && \text{using the definition of the logistic function (5.10)} \\ &= \mathbb{E} \left[q(x, a) \nabla_f (-f + af - \log(1 + e^{-f})) \nabla_{\varphi} f_{\varphi}(x) \right] \\ &= \mathbb{E} \left[q(x, a) \left(a - 1 + \frac{e^{-f}}{1 + e^{-f}} \right) \nabla_{\varphi} f_{\varphi}(x) \right] \\ &= \mathbb{E} [q(x, a)(a - \sigma(f)) \nabla_{\varphi} f_{\varphi}(x)]. && \text{using the definition of the logistic function (5.10)} \end{aligned}$$

The term $a - \sigma(f)$ can be understood as a residual as it corresponds to the difference between the target action a and the expected action $\sigma(f)$.

2. We have

$$\begin{aligned} \nabla_{\varphi} j(\varphi) &= \mathbb{E} [q(x, a) \nabla_f \log \pi_f(a | x) \nabla_{\varphi} f(x)] && \text{using eq. (12.58) and the chain rule} \\ &= \mathbb{E} [q(x, a) \nabla_f (\log h(a) + af - A(f)) \nabla_{\varphi} f(x)] \\ &= \mathbb{E} [q(x, a)(a - \nabla_f A(f)) \nabla_{\varphi} f(x)]. \end{aligned}$$

3. We have $\nabla_f A(f) = \sigma(f)$. We are therefore looking for a function $A(f)$ whose derivative is $\sigma(f) = \frac{1}{1+e^{-f}} = \frac{e^f}{1+e^f}$. With this equality of the sigmoid we can find the integral, and we have $A(f) = \log(1 + e^f) + c$. Let us confirm that this gives us the Bernoulli distribution with $c = 0$:

$$\begin{aligned} \pi_f(a | x) &= h(a) \exp(af - \log(1 + e^f)) \\ &= h(a) \frac{e^{af}}{1 + e^f} \\ &= h(a) \begin{cases} \sigma(f) & \text{if } a = 1 \\ 1 - \sigma(f) & \text{if } a = 0. \end{cases} \end{aligned}$$

This is the Bernoulli distribution with parameter $\sigma(f)$ where we have $h(a) = 1$.

4. Using that $\nabla_f A(f) = f$, we immediately get

$$\nabla_{\varphi} j(\varphi) = \mathbb{E}[q(x, a)(a - f) \nabla_{\varphi} f(x)].$$

5. No, we cannot use the reparameterization trick since we do not know how the states x depend on action a . These dependencies are determined by the unknown dynamics of the environment. Nonetheless, we can apply it after sampling an episode according to a policy and then updating policy parameters in hindsight. This is for example done by the soft actor-critic (SAC) algorithm (cf. section 12.5).

Solution to exercise 12.19 (Soft value function).

1. Equation (12.90) can be written as

$$\begin{aligned} & \text{KL}(\Pi_{\varphi} \| \Pi_{\star}) \\ &= \sum_{t=1}^T \mathbb{E}_{(x_t, a_t) \sim \Pi_{\varphi}} [\text{S}[\beta(a_t | x_t)] - \text{H}[\pi_{\varphi}(\cdot | x_t)]] . \end{aligned}$$

Adding and subtracting $\log Z(x_t)$ gives

$$\begin{aligned} &= \sum_{t=1}^T \mathbb{E}_{(x_t, a_t) \sim \Pi_{\varphi}} [\text{S}[\hat{\pi}(a_t | x_t)] - \log Z(x_t) - \text{H}[\pi_{\varphi}(\cdot | x_t)]] \\ &= \sum_{t=1}^T \mathbb{E}_{x_t \sim \Pi_{\varphi}} [\text{H}[\pi_{\varphi}(\cdot | x_t) \| \hat{\pi}(\cdot | x_t)] - \log Z(x_t) - \text{H}[\pi_{\varphi}(\cdot | x_t)]] && \text{using the definition of cross-entropy (5.40)} \\ &= \sum_{t=1}^T \mathbb{E}_{x_t \sim \Pi_{\varphi}} [\text{KL}(\pi_{\varphi}(\cdot | x_t) \| \hat{\pi}(\cdot | x_t)) - \log Z(x_t)] . && \text{using the definition of KL-divergence (5.44)} \end{aligned}$$

2. We prove the statement by (reverse) induction on t . For the base case, note that the term

$$\mathbb{E}_{x_T \sim \Pi_{\varphi}} [\text{KL}(\pi_{\varphi}(\cdot | x_T) \| \hat{\pi}(\cdot | x_T)) - \log Z(x_T)]$$

is minimized for $\pi_{\varphi} \equiv \hat{\pi}$. The KL-divergence then evaluates to zero, and we are left only with the $\log Z(x_T)$ term.

For the inductive step, fix any $1 \leq t < T$ and $\pi^*(a_t | x_t)$ must minimize the two terms

$$\begin{aligned} & \mathbb{E}_{x_t \sim \Pi_{\varphi}} [\text{KL}(\pi_{\varphi}(\cdot | x_t) \| \hat{\pi}(\cdot | x_t)) - \log Z(x_t)] \\ &+ \mathbb{E}_{(x_t, a_t) \sim \Pi_{\varphi}} [\mathbb{E}_{x_{t+1} \sim p(\cdot | x_t, a_t)} [-\log Z(x_{t+1})]] \end{aligned}$$

where the first term stems directly from the objective (12.92) and the second term represents the contribution of $\pi_\varphi(a_t | x_t)$ to all subsequent terms. Letting $\beta^*(a | x) \doteq \exp(q^*(x, a))$, $Z^*(x) \doteq \int_{\mathcal{A}} \beta^*(a | x) da$, and recalling that we denote by $\pi^*(\cdot | x)$ the policy $\beta^*(\cdot | x)/Z^*(x)$, this objective can be reexpressed as

$$= \mathbb{E}_{x_t \sim \Pi_\varphi} [\text{KL}(\pi_\varphi(\cdot | x_t) \| \pi^*(\cdot | x_t)) - \log Z^*(x_t)]$$

which is minimized for $\pi_\varphi \equiv \pi^*$, leaving only the $\log Z^*(x_t)$ term. It remains only to observe that $\beta^*(a_T | x_T) = \beta(a_T | x_T)$, so in the final state x_T , $\log Z^*(x_T) = \log Z(x_T)$ and the policies π^* and $\hat{\pi}$ coincide.

Model-based Reinforcement Learning

Solution to exercise 13.14 (Bounding the regret of H-UCRL).

1. To simplify the notation, we use $z_{t,k}$ as shorthand for $(x_{t,k}, \pi_t(x_{t,k}))$ (and similarly $\hat{z}_{t,k}$ for $(\hat{x}_{t,k}, \pi_t(\hat{x}_{t,k}))$). The base case is implied trivially. For the induction step, assume that eq. (13.28) holds at iteration k . We have

$$\begin{aligned} \|\hat{x}_{t,k+1} - x_{t,k+1}\| &= \|\hat{f}_{t-1}(\hat{z}_{t,k}) - f(z_{t,k})\| \\ &\leq \|\hat{f}_{t-1}(\hat{z}_{t,k}) - f(\hat{z}_{t,k})\| + \|f(\hat{z}_{t,k}) - f(z_{t,k})\| \\ &\leq 2\beta_t \sigma_{t-1}(\hat{z}_{t,k}) + L_1 \|\hat{x}_{t,k} - x_{t,k}\| \end{aligned}$$

adding and subtracting $f(\hat{z}_{t,k})$ and using Cauchy-Schwarz

where the final inequality follows with high probability from the assumption that the confidence intervals are well-calibrated (cf. eq. (13.22)) and the assumed Lipschitzness.

$$\begin{aligned} &= 2\beta_t [\sigma_{t-1}(z_{t,k}) + \sigma_{t-1}(\hat{z}_{t,k}) - \sigma_{t-1}(z_{t,k})] \\ &\quad + L_1 \|\hat{x}_{t,k} - x_{t,k}\| \end{aligned}$$

Once more using Lipschitz continuity, we obtain

$$\begin{aligned} &\leq 2\beta_t [\sigma_{t-1}(z_{t,k}) + L_2 \|\hat{x}_{t,k} - x_{t,k}\|] \\ &\quad + L_1 \|\hat{x}_{t,k} - x_{t,k}\| \\ &= 2\beta_t \sigma_{t-1}(z_{t,k}) + \alpha_t \|\hat{x}_{t,k} - x_{t,k}\| \end{aligned}$$

where $\alpha_t \doteq (L_1 + 2\beta_t L_2)$. By the induction hypothesis,

$$\leq 2\beta_t \sum_{l=0}^k \alpha_t^{k-l} \sigma_{t-1}(z_{t,l}).$$

This is identical to the analysis of UCB from exercise 9.7, only that here errors compound along the trajectory.

2. The assumption that $\alpha_t \geq 1$ implies that

$$\|\hat{x}_{t,k} - x_{t,k}\| \leq 2\beta_t \alpha_t^{H-1} \sum_{l=0}^{k-1} \sigma_{t-1}(z_{t,l}). \quad (\text{B.15})$$

Moreover, by definition of π_t , we have with high probability that $J_H(\pi_t; \hat{f}) \geq J_H(\pi^*; f)$. This is because π_t maximizes reward under *optimistic* dynamics. Thus,

$$\begin{aligned} r_t &= \sum_{k=0}^{H-1} r(x_{t,k}, \pi^*(x_{t,k})) - r(x_{t,k}, \pi_t(x_{t,k})) \\ &\leq \sum_{k=0}^{H-1} r(\hat{x}_{t,k}, \pi_t(\hat{x}_{t,k})) - r(x_{t,k}, \pi_t(x_{t,k})) \\ &\leq L_3 \sum_{k=0}^{H-1} \|\hat{x}_{t,k} - x_{t,k}\| && \text{using Lipschitz-continuity of } r \\ &\leq 2\beta_t \alpha_t^{H-1} L_3 \sum_{k=0}^{H-1} \sum_{l=0}^{k-1} \sigma_{t-1}(z_{t,l}) && \text{using eq. (B.15)} \\ &\leq 2\beta_t H \alpha_t^{H-1} L_3 \sum_{k=0}^{H-1} \sigma_{t-1}(z_{t,k}). \end{aligned}$$

3. Let us first bound R_T^2 . By the Cauchy-Schwarz inequality,

$$\begin{aligned} R_T^2 &\leq T \sum_{t=1}^T r_t^2 \\ &\leq \mathcal{O} \left(T \sum_{t=1}^T \beta_t^2 H^2 \alpha_t^{2(H-1)} \left(\sum_{k=0}^{H-1} \sigma_{t-1}(z_{t,k}) \right)^2 \right) && \text{using (2)} \\ &\leq \mathcal{O} \left(T \beta_T^2 H^3 \alpha_T^{2(H-1)} \sum_{t=1}^T \sum_{k=0}^{H-1} \sigma_{t-1}^2(z_{t,k}) \right). && \text{using Cauchy-Schwarz and assuming} \\ &&& \text{w.l.o.g. that } \beta_t \text{ is monotonically increasing} \end{aligned}$$

Taking the square root, we obtain

$$\begin{aligned} R_T &\leq \mathcal{O} \left(\beta_T H^{\frac{3}{2}} \alpha_T^{H-1} \sqrt{T \sum_{t=1}^T \sum_{k=0}^{H-1} \sigma_{t-1}^2(z_{t,k})} \right) \\ &\leq \mathcal{O} \left(\beta_T H^{\frac{3}{2}} \alpha_T^{H-1} \sqrt{T \Gamma_T} \right). \end{aligned}$$

Bibliography

- Eitan Altman. *Constrained Markov decision processes: stochastic modeling*. Routledge, 1999.
- Yarden As, Ilnura Usmanova, Sebastian Curi, and Andreas Krause. Constrained policy optimization via bayesian world models. *International Conference on Learning Representations*, 2022.
- Dominique Bakry and Michel Émery. Diffusions hypercontractives. In *Séminaire de Probabilités XIX 1983/84: Proceedings*, pages 177–206. Springer, 2006.
- Felix Berkenkamp, Matteo Turchetta, Angela Schoellig, and Andreas Krause. Safe model-based reinforcement learning with stability guarantees. *Advances in neural information processing systems*, 30, 2017.
- Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International conference on machine learning*, pages 1613–1622. PMLR, 2015.
- Ilija Bogunovic, Jonathan Scarlett, Andreas Krause, and Volkan Cevher. Truncated variance reduction: A unified approach to bayesian optimization and level-set estimation. *Advances in neural information processing systems*, 29, 2016.
- Léon Bottou. Online learning and stochastic approximations. *On-line learning in neural networks*, 17(9):142, 1998.
- Ronen I Brafman and Moshe Tennenholtz. R-max-a general polynomial time algorithm for near-optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct):213–231, 2002.
- Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.

- Kathryn Chaloner and Isabella Verdinelli. Bayesian experimental design: A review. *Statistical Science*, pages 273–304, 1995.
- Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient hamiltonian monte carlo. In *International conference on machine learning*, pages 1683–1691. PMLR, 2014.
- Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International Conference on Machine Learning*, pages 844–853. PMLR, 2017.
- Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. *Advances in neural information processing systems*, 31, 2018.
- Thomas Cover and Joy Thomas. *Elements of information theory*. Wiley-Interscience, 2006.
- Sebastian Curi, Felix Berkenkamp, and Andreas Krause. Efficient model-based reinforcement learning through optimistic policy search and planning. *Advances in Neural Information Processing Systems*, 33:14156–14170, 2020.
- Sebastian Curi, Armin Lederer, Sandra Hirche, and Andreas Krause. Safe reinforcement learning via confidence-based filters. In *2022 IEEE 61st Conference on Decision and Control (CDC)*, pages 3409–3415. IEEE, 2022.
- Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*, pages 465–472. Citeseer, 2011.
- Marc Peter Deisenroth, A Aldo Faisal, and Cheng Soon Ong. *Mathematics for machine learning*. Cambridge University Press, 2020.
- Simon Duane, Anthony D Kennedy, Brian J Pendleton, and Duncan Roweth. Hybrid monte carlo. *Physics letters B*, 195(2):216–222, 1987.
- David Duvenaud. *Automatic model construction with Gaussian processes*. PhD thesis, University of Cambridge, 2014.
- David Duvenaud and Ryan P Adams. Black-box stochastic variational inference in five lines of python. In *NIPS Workshop on Black-box Learning and Inference*, 2015.
- Eyal Even-Dar and Yishay Mansour. Convergence of optimistic and incremental q-learning. *Advances in neural information processing systems*, 14, 2001.

- Eyal Even-Dar, Yishay Mansour, and Peter Bartlett. Learning rates for q-learning. *Journal of machine learning Research*, 5(1), 2003.
- Matthew Fellows, Anuj Mahajan, Tim GJ Rudner, and Shimon Whiteson. Virel: A variational inference framework for reinforcement learning. *Advances in neural information processing systems*, 32, 2019.
- Karl Friston. The free-energy principle: a unified brain theory? *Nature reviews neuroscience*, 11(2):127–138, 2010.
- Karl Friston, Francesco Rigoli, Dimitri Ognibene, Christoph Mathys, Thomas Fitzgerald, and Giovanni Pezzulo. Active inference and epistemic value. *Cognitive neuroscience*, 6(4):187–214, 2015.
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pages 1587–1596. PMLR, 2018.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Insights and applications. In *Deep Learning Workshop, ICML*, volume 1, page 2, 2015.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *International Conference on Machine Learning*, pages 1183–1192. PMLR, 2017.
- Roman Garnett. *Bayesian optimization*. Cambridge University Press, 2023.
- Daniel Golovin, Benjamin Solnik, Subhodeep Moitra, Greg Kochanski, John Karro, and David Sculley. Google vizier: A service for black-box optimization. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1487–1495, 2017.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. PMLR, 2018a.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta,

- Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018b.
- Jouni Hartikainen and Simo Särkkä. Kalman filtering and smoothing solutions to temporal gaussian process regression models. In *2010 IEEE international workshop on machine learning for signal processing*, pages 379–384. IEEE, 2010.
- Nicolas Heess, Gregory Wayne, David Silver, Timothy Lillicrap, Tom Erez, and Yuval Tassa. Learning continuous control policies by stochastic value gradients. *Advances in neural information processing systems*, 28, 2015.
- Philipp Hennig and Christian J Schuler. Entropy search for information-efficient global optimization. *Journal of Machine Learning Research*, 13(6), 2012.
- James Hensman, Alexander Matthews, and Zoubin Ghahramani. Scalable variational gaussian process classification. In *Artificial Intelligence and Statistics*, pages 351–360. PMLR, 2015.
- José Miguel Hernández-Lobato, Matthew W Hoffman, and Zoubin Ghahramani. Predictive entropy search for efficient global optimization of black-box functions. *Advances in neural information processing systems*, 27, 2014.
- Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- Matthew W Hoffman and Zoubin Ghahramani. Output-space predictive entropy search for flexible global optimization. In *NIPS workshop on Bayesian Optimization*, pages 1–5, 2015.
- Thomas Hofmann. *Computational Intelligence Lab*. ETH Zurich, 2022.
- Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. *Conference on Uncertainty in Artificial Intelligence*, 2018.
- Tommi Jaakkola, Michael Jordan, and Satinder Singh. Convergence of stochastic iterative dynamic programming algorithms. *Advances in neural information processing systems*, 6, 1993.
- Edwin T Jaynes. *Probability theory: The logic of science*. Cambridge university press, 2003.

- Kirthevasan Kandasamy, Akshay Krishnamurthy, Jeff Schneider, and Barnabás Póczos. Parallelised bayesian optimisation via thompson sampling. In *International Conference on Artificial Intelligence and Statistics*, pages 133–142. PMLR, 2018.
- Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in neural information processing systems*, 30, 2017.
- Johannes Kirschner and Andreas Krause. Information directed sampling and bandits with heteroscedastic noise. In *Conference On Learning Theory*, pages 358–384. PMLR, 2018.
- Torsten Koller, Felix Berkenkamp, Matteo Turchetta, and Andreas Krause. Learning-based model predictive control for safe exploration. In *2018 IEEE conference on decision and control (CDC)*, pages 6059–6066. IEEE, 2018.
- Andreas Krause and Daniel Golovin. Submodular function maximization. *Tractability*, 3:71–104, 2014.
- Andreas Krause and Fanny Yang. *Introduction to Machine Learning*. ETH Zurich, 2022.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- David A Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and review. *arXiv preprint arXiv:1805.00909*, 2018.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Yi-An Ma, Yuansi Chen, Chi Jin, Nicolas Flammarion, and Michael I Jordan. Sampling can be faster than optimization. *Proceedings of the National Academy of Sciences*, 116(42):20881–20885, 2019.
- David JC MacKay. Information-based objective functions for active data selection. *Neural computation*, 4(4):590–604, 1992.
- Beren Millidge, Alexander Tschantz, Anil K Seth, and Christopher L

- Buckley. On the relationship between active inference and control as inference. In *Active Inference: First International Workshop, IWAI 2020, Co-located with ECML/PKDD 2020, Ghent, Belgium, September 14, 2020, Proceedings 1*, pages 3–11. Springer, 2020.
- Beren Millidge, Alexander Tschantz, and Christopher L Buckley. Whence the expected free energy? *Neural Computation*, 33(2):447–482, 2021.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
- Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte carlo gradient estimation in machine learning. *J. Mach. Learn. Res.*, 21(132):1–62, 2020.
- Kevin P Murphy. Conjugate bayesian analysis of the gaussian distribution. *def*, 1(202):16, 2007.
- Jayakrishnan Nair, Adam Wierman, and Bert Zwart. *The Fundamentals of Heavy Tails*, volume 53. Cambridge University Press, 2022.
- George L Nemhauser, Laurence A Wolsey, and Marshall L Fisher. An analysis of approximations for maximizing submodular set functions. *Mathematical programming*, 14:265–294, 1978.
- Vu Nguyen, Sunil Gupta, Santu Rana, Cheng Li, and Svetha Venkatesh. Regret for expected improvement over the best-observed value and stopping condition. In *Asian Conference on Machine Learning*, pages 279–294. PMLR, 2017.
- OpenAI. Gpt-4 technical report, 2023.
- Manfred Opper and Cédric Archambeau. The variational gaussian approximation revisited. *Neural computation*, 21(3):786–792, 2009.
- Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7(15):510, 2008.
- Joaquin Quinonero-Candela and Carl Edward Rasmussen. A unifying

- view of sparse approximate gaussian process regression. *The Journal of Machine Learning Research*, 6, 2005.
- Maxim Raginsky, Alexander Rakhlin, and Matus Telgarsky. Non-convex learning via stochastic gradient langevin dynamics: a nonasymptotic analysis. In *Conference on Learning Theory*, pages 1674–1703. PMLR, 2017.
- Ali Rahimi, Benjamin Recht, et al. Random features for large-scale kernel machines. In *NIPS*, 2007.
- Rajesh Ranganath, Sean Gerrish, and David Blei. Black box variational inference. In *Artificial intelligence and statistics*, pages 814–822. PMLR, 2014.
- Gareth O Roberts and Jeffrey S Rosenthal. General state space markov chains and mcmc algorithms. 2004.
- Philip A Romero, Andreas Krause, and Frances H Arnold. Navigating the protein fitness landscape with gaussian processes. *Proceedings of the National Academy of Sciences*, 110(3):E193–E201, 2013.
- Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- Stuart Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Prentice Hall, 2002.
- Daniel Russo and Benjamin Van Roy. Learning to optimize via information-directed sampling. *Advances in Neural Information Processing Systems*, 27, 2014.
- Daniel Russo and Benjamin Van Roy. An information-theoretic analysis of thompson sampling. *The Journal of Machine Learning Research*, 17(1):2442–2471, 2016.
- Grant Sanderson. But what is the fourier transform? a visual introduction, 2018. URL <https://www.youtube.com/watch?v=spUNpyF58BY>.
- Bernhard Schölkopf, Ralf Herbrich, and Alex J Smola. A generalized representer theorem. In *International conference on computational learning theory*, pages 416–426. Springer, 2001.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015a.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015b.

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Guy Shani, Joelle Pineau, and Robert Kaplow. A survey of point-based pomdp solvers. *Autonomous Agents and Multi-Agent Systems*, 27(1): 1–51, 2013.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharmashan Kumaran, Thore Graepel, et al. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*, 2017.
- Satinder Singh, Tommi Jaakkola, Michael L Littman, and Csaba Szepesvári. Convergence results for single-step on-policy reinforcement-learning algorithms. *Machine learning*, 38(3):287–308, 2000.
- Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *International Conference on Machine Learning*, pages 1015–1022. ICML, 2010.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15 (1):1929–1958, 2014.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.
- Bhavya Sukhija, Matteo Turchetta, David Lindner, Andreas Krause, Sebastian Trimpe, and Dominik Baumann. Gosafeopt: Scalable safe exploration for global optimization of dynamical systems. *Artificial Intelligence*, 320:103922, 2023.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

- Yee Whye Teh, Alexandre H Thiery, and Sebastian J Vollmer. Consistency and fluctuations for stochastic gradient langevin dynamics. *Journal of Machine Learning Research*, 17, 2016.
- Michalis Titsias and Miguel Lázaro-Gredilla. Doubly stochastic variational bayes for non-conjugate inference. In *International conference on machine learning*, pages 1971–1979. PMLR, 2014.
- Lenart Treven, Jonas Hübottter, Bhavya Sukhija, Florian Dörfler, and Andreas Krause. Efficient exploration in continuous-time model-based reinforcement learning. *Advances in neural information processing systems*, 2023.
- Matteo Turchetta, Felix Berkenkamp, and Andreas Krause. Safe exploration for interactive machine learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Sattar Vakili, Kia Khezeli, and Victor Picheny. On information gain and regret bounds in gaussian process bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 82–90. PMLR, 2021.
- Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
- Santosh Vempala and Andre Wibisono. Rapid convergence of the unadjusted langevin algorithm: Isoperimetry suffices. *Advances in neural information processing systems*, 32, 2019.
- Li Wan, Matthew Zeiler, Sixin Zhang, Yann Le Cun, and Rob Fergus. Regularization of neural networks using dropconnect. In *International conference on machine learning*, pages 1058–1066. PMLR, 2013.
- Yuhui Wang, Hao He, and Xiaoyang Tan. Truly proximal policy optimization. In *Uncertainty in Artificial Intelligence*, pages 113–122. PMLR, 2020.
- Zi Wang and Stefanie Jegelka. Max-value entropy search for efficient bayesian optimization. In *International Conference on Machine Learning*, pages 3627–3635. PMLR, 2017.
- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688. Citeseer, 2011.
- Daniel Widmer, Dongho Kang, Bhavya Sukhija, Jonas Hübottter, Andreas Krause, and Stelian Coros. Tuning legged locomotion controllers via safe bayesian optimization. *arXiv preprint arXiv:2306.07092*, 2023.

Christopher K Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*, volume 2. MIT press, 2006.

Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.

Stephen Wright, Jorge Nocedal, et al. Numerical optimization. *Springer Science*, 35(67-68):7, 1999.

Pan Xu, Jinghui Chen, Difan Zou, and Quanquan Gu. Global convergence of langevin dynamics based algorithms for nonconvex optimization. *Advances in Neural Information Processing Systems*, 31, 2018.

Yinyu Ye. The simplex and policy-iteration methods are strongly polynomial for the markov decision problem with a fixed discount rate. *Mathematics of Operations Research*, 36(4):593–603, 2011.

Acronyms

<i>a.s.</i>	almost surely, with high probability, with probability 1	<i>LOTP</i>	law of total probability
<i>A2C</i>	advantage actor-critic	<i>LOTUS</i>	law of the unconscious statistician
<i>BALD</i>	Bayesian active learning by disagreement	<i>LOTV</i>	law of total variance
<i>BLR</i>	Bayesian linear regression	<i>LSI</i>	log-Sobolev inequality
<i>CDF</i>	cumulative distribution function	<i>MAB</i>	multi-armed bandits
<i>DDPG</i>	deep deterministic policy gradients	<i>MALA</i>	Metropolis adjusted Langevin algorithm
<i>DDQN</i>	double deep Q-networks	<i>MAP</i>	maximum a posteriori
<i>DQN</i>	deep Q-networks	<i>MC</i>	Monte Carlo
<i>ECE</i>	expected calibration error	<i>MCE</i>	maximum calibration error
<i>EI</i>	expected improvement	<i>MCMC</i>	Markov chain Monte Carlo
<i>ELBO</i>	evidence lower bound	<i>MDP</i>	Markov decision process
<i>ES</i>	entropy search	<i>MERL</i>	maximum entropy reinforcement learning
<i>FITC</i>	fully independent training conditional	<i>MES</i>	max-value entropy search
<i>GAE</i>	generalized advantage estimation	<i>MGF</i>	moment-generating function
<i>GD</i>	gradient descent	<i>MI</i>	mutual information
<i>GLIE</i>	greedy in the limit with infinite exploration	<i>MLE</i>	maximum likelihood estimate
<i>GP</i>	Gaussian process	<i>MLL</i>	marginal log likelihood
<i>GPC</i>	Gaussian process classification	<i>MPC</i>	model predictive control
<i>GRV</i>	Gaussian random vector	<i>MSE</i>	mean squared error
<i>H-UCRL</i>	hallucinated upper confidence reinforcement learning	<i>NLL</i>	negative log likelihood
<i>HMC</i>	Hamiltonian Monte Carlo	<i>ODE</i>	ordinary differential equation
<i>HMM</i>	hidden Markov model	<i>OPES</i>	output-space predictive entropy search
<i>i.i.d.</i>	independent and identically distributed	<i>PBPI</i>	point-based policy iteration
<i>IDS</i>	information-directed sampling	<i>PBVI</i>	point-based value iteration
<i>iff</i>	if and only if	<i>PDF</i>	probability density function
<i>KL</i>	Kullback-Leibler	<i>PES</i>	predictive entropy search
<i>LAMBDA</i>	Lagrangian model-based agent	<i>PETS</i>	probabilistic ensembles with trajectory sampling
<i>LASSO</i>	least absolute shrinkage and selection operator	<i>PI</i>	probability of improvement
<i>LD</i>	Langevin dynamics	<i>PI</i>	policy iteration
<i>LMC</i>	Langevin Monte Carlo	<i>PILCO</i>	probabilistic inference for learning control
<i>LOTE</i>	law of total expectation	<i>PL</i>	Polyak-Łojasiewicz
		<i>PMF</i>	probability mass function
		<i>POMDP</i>	partially observable Markov decision process

<i>PPO</i>	proximal policy optimization	<i>SVG</i>	stochastic value gradients
<i>RBF</i>	radial basis function	<i>SVGD</i>	Stein variational gradient descent
<i>ReLU</i>	rectified linear unit	<i>SWA</i>	stochastic weight averaging
<i>RHC</i>	receding horizon control	<i>SWAG</i>	stochastic weight averaging-Gaussian
<i>RKHS</i>	reproducing kernel Hilbert space	<i>Tanh</i>	hyperbolic tangent
<i>RL</i>	reinforcement learning	<i>TD</i>	temporal difference
<i>RM</i>	Robbins-Monro	<i>TD3</i>	twin delayed DDPG
<i>SAA</i>	stochastic average approximation	<i>TRPO</i>	trust-region policy optimization
<i>SAC</i>	soft actor-critic	<i>UCB</i>	upper confidence bound
<i>SARSA</i>	state-action-reward-state-action	<i>ULA</i>	unadjusted Langevin algorithm
<i>SDE</i>	stochastic differential equation	<i>VI</i>	variational inference
<i>SG-HMC</i>	stochastic gradient Hamiltonian Monte Carlo	<i>VI</i>	value iteration
<i>SGD</i>	stochastic gradient descent	<i>w.l.o.g.</i>	without loss of generality
<i>SGLD</i>	stochastic gradient Langevin dynamics	<i>w.r.t.</i>	with respect to
<i>SLLN</i>	strong law of large numbers	<i>WLLN</i>	weak law of large numbers
<i>SoR</i>	subset of regressors		

Index

- R_{\max} algorithm, 221
- ℓ_1 -regularization, 36
- ℓ_2 -regularization, 36
- ℓ_∞ norm, 201
- ϵ -greedy, **219**, 277
- σ -algebra, 2

- a-optimal design, 175
- acceptance distribution, 130
- acquisition function, **182**, 277
- activation, 148
- activation function, 148
- active inference, 261
- actor, 248
- actor-critic method, 248
- Adagrad, 37
- Adam, 37
- adaptive learning rate, 37
- additive white Gaussian noise, 49
- advantage actor-critic, 250
- advantage function, 245
- affine transformation, 38
- aleatoric uncertainty, 54
- almost sure convergence, 26
- alpha-beta pruning, 266
- aperiodic Markov chain, 125
- augmented Lagrangian method, 282
- automatic differentiation (auto-diff), 150

- backpropagation, 150
- bagging, 157
- Banach fixed-point theorem, 201
- Barrier function, 282

- baseline, 241
- batch, 35
- batch size, 36
- Bayes by Backprop, 153
- Bayes' rule, 20
- Bayes' rule for entropy, 164
- Bayesian active learning by disagreement, 173
- Bayesian experimental design, 174
- Bayesian filtering, 61, **63**, 210
- Bayesian linear regression, 53
- Bayesian logistic regression, 93
- Bayesian networks, 11
- Bayesian neural network, 151
- Bayesian optimization, 183
- Bayesian reasoning, 1
- Bayesian smoothing, 64
- belief, 62, **210**
- belief-state Markov decision process, 211
- Bellman error, **234**, 255
- Bellman expectation equation, 198
- Bellman optimality equation, 204
- Bellman update, 203
- Bellman's optimality principle, 204
- Bellman's theorem, 203
- Bernoulli distribution, 5
- best-arm identification, 175
- bias, 23
- bias-variance tradeoff, **24**, 250
- binary cross-entropy loss, 104
- binomial distribution, 5
- black box stochastic variational inference, 118
- Bochner's theorem, 84
- Boltzmann distribution, 134

- Boltzmann exploration, 220
- bootstrapping, 157, **223**
- Borel σ -algebra, 3
- bounded discounted payoff, 239
- Brownian motion, 138
- burn-in time, 129
- calibration, 157
- categorical distribution, 5
- Cauchy distribution, 118
- central limit theorem, 7
- Cesàro mean, 182
- chain rule for entropy, 164
- chain rule of mutual information, 191
- chain rule of probability, 9
- change of variables formula, 17
- Chebyshev's inequality, 26
- Cholesky decomposition, 42
- classification, 18
- closed-loop control, 269
- competitive ratio, 182
- completing the square, 299
- computation graph, 147
- concave function, 32
- conditional distribution, 9
- conditional entropy, 164
- conditional expectation, 14
- conditional independence, 10
- conditional likelihood, 20
- conditional linear Gaussian, 42
- conditional mutual information, 166
- conditional probability, 8
- conditional variance, 16
- conditioning, 9
- conjugate prior, 22
- consistent estimator, 24
- constrained Markov decision processes, 281
- continuous setting, 217
- contraction, 201
- control, 182
- convergence in distribution, 26
- convergence in probability, 26
- convergence with probability 1, 26
- convex body chasing, 182
- convex function, 32
- convex function chasing, 182
- correlation, 15
- covariance, 14
- covariance function, 69
- covariance matrix, 16
- critic, 248
- cross-entropy, 103
- cross-entropy loss, 150
- cross-entropy method, 266
- curse of dimensionality, 186
- d-optimal design, 174
- d-separation, 12
- data processing inequality, 167
- deep deterministic policy gradients, **254**, 269
- deep Q-networks, 236
- deep reinforcement learning, 234
- density, 6
- design, 174
- design criterion, 174
- design matrix, 47
- detailed balance equation, 128
- differentiation under the integral sign, 13
- diffusion process, 139
- Dirac delta function, 9
- directed graphical model, 11
- directional derivative, 32
- discounted payoff, **196**, 239, 265
- discounted state occupancy measure, 247
- double DQN, 236
- downstream return, 243
- dropconnect regularization, 155
- dropout regularization, 155
- dynamics model, 195
- e-optimal design, 175
- eligibility vector, 241
- empirical Bayes, 80
- empirical risk, 19
- energy function, 134
- entropy, 50, **101**
- entropy regularization, 257
- entropy search, 176

- episode, 216
- episodic setting, 216
- epistemic uncertainty, 54
- epoch, 36
- equilibrium, 139
- ergodic theorem, 129
- ergodicity, 125
- estimator, 23
- Euler's formula, 84
- event, 2
- event space, 2
- evidence, 112
- exclusive KL-divergence, 107
- expectation, 12
- expected calibration error, 159
- expected improvement, 188
- experience replay, 236
- experimental design, 174
- exploration distribution, 252
- exploration noise, 277
- exploration-exploitation dilemma, 114, **179**
- exponential family, 110, 248
- exponential kernel, 73
- feature space, 56
- feed-forward neural network, 150
- finite measure, **122**
- finite-horizon, 197
- first-order characterization of convexity, 32
- first-order expansion, 31
- first-order optimality condition, 31
- fixed-point iteration, 200
- Fokker-Planck equation, 140
- forward KL-divergence, 107
- forward sampling, 71
- forward-backward algorithm, 65
- Fourier transform, 84
- free energy, 113
- free energy principle, **113**, 260
- fully independent training conditional, 90
- fully observable environment, 196
- function class, 18
- function-space view, 56
- fundamental theorem of ergodic Markov chains, 126
- gamma distribution, 133
- Gaussian, 6, 37
- Gaussian kernel, 72
- Gaussian Markov Process, 88
- Gaussian noise "dithering", 253, **277**
- Gaussian process, 69
- Gaussian process classification, 93, **96**
- Gaussian random vector, 37
- generalized advantage estimation, 250
- generative model, **20**
- Gibbs distribution, 134
- Gibbs sampling, 132
- Gibbs' inequality, 105
- global optimum, 31
- gradient, 13
- gradient flow, 139
- greedy in the limit with infinite exploration, 220
- greedy policy, 202
- Grönwall's inequality, 140
- Gumbel-max trick, 256
- Hadamard's inequality, 287
- hallucinated upper confidence reinforcement learning, 279
- Hamiltonian, 143
- Hamiltonian Monte Carlo, 142
- heavy-tailed distribution, 25
- Hessian, 34
- heteroscedastic, **18**
- heteroscedastic noise, 152
- hidden layer, 147
- hidden Markov model, 209
- histogram binning, 159
- Hoeffding's inequality, 28
- Hoeffding's lemma, 28
- homoscedastic, **18**
- homoscedastic noise, 152
- hyperbolic tangent, 148
- hyperprior, **82**, 192
- importance weighting, 251

- improper prior, 21
- improvement, 187
- inclusive KL-divergence, 107
- independence, 10
- inducing points method, 88
- inference, 19
- information gain, 165
- information never hurts, **164**, 165
- information-directed sampling, 191
- input layer, 147
- inputs, 18
- instantaneous, 197
- instantaneous regret, 181
- interaction information, 167
- inverse Fourier transform, 84
- inverse transform sampling, 7
- irreducible Markov chain, 124
- isotonic regression, 159
- isotropic Gaussian, 38
- isotropic kernel, 75

- Jacobian, 18
- Jensen's inequality, 102
- joint distribution, 8
- joint entropy, 164
- joint likelihood, 20

- Kalman filter, 61
- Kalman gain, 65, **66**
- Kalman smoothing, 64
- Kalman update, 66
- kernel function, 57, 69, **71**
- kernel matrix, 57
- kernel ridge regression, 78
- kernel trick, 57
- kernelization, 57
- Kolmogorov axioms, 2
- Kullback-Leibler divergence, 104

- labels, 18
- Lagrangian model-based agent, 283
- Langevin dynamics, 139
- Langevin Monte Carlo, 138
- Laplace approximation, 92
- Laplace distribution, 25
- Laplace kernel, 73
- Laplace's method, 91
- lasso, 52
- law of large numbers, 27
- law of the unconscious statistician, 13
- law of total expectation, 14
- law of total probability, 9
- law of total variance, 17
- layer, 147
- Leapfrog method, 144
- learning, 19
- least squares estimator, 48
- length scale, 72
- light-tailed distribution, 25
- likelihood, 20
- linear kernel, 72
- linear regression, 47
- linearity of expectation, 12
- local optimum, 31
- log-concave distribution, 134
- log-likelihood, 29
- log-normal distribution, 25
- log-prior, 30
- log-Sobolev inequality, 141
- logistic function, 93
- logistic loss, 94
- logits, 147
- loss function, 30
- Lyapunov function, 139

- Mahalanobis norm, 42
- marginal gain, 169
- marginal likelihood, 20
- marginalization, 8
- Markov chain, 122
- Markov decision process, 195
- Markov property, 62, **122**
- Markov's inequality, 26
- matrix inversion lemma, 287
- Matérn kernel, 73
- max-value entropy search, 177
- maximization bias, 236
- maximizing the marginal likelihood, 79

- maximum a posteriori estimate, 30
- maximum calibration error, 159
- maximum entropy principle, 50
- maximum entropy reinforcement learning, 257
- maximum likelihood estimate, 29
- mean function, 69
- mean payoff, 197
- mean square continuity, 74
- mean square convergence, 74
- mean square differentiability, 74
- mean squared error, **24**, 149
- mean-field distribution, 99
- Mercer's theorem, 72
- method of moments, 110
- metrical task systems, 182
- Metropolis adjusted Langevin algorithm, 138
- Metropolis-Hastings algorithm, 130
- Metropolis-Hastings theorem, 130
- mixing time, 127
- mode, 6
- model predictive control, 265
- model-based reinforcement learning, **217**, 264
- model-free reinforcement learning, 217
- moment matching, **110**, 271
- moment-generating function, 40
- momentum, 37
- monotone submodularity, 169
- monotonicity of mutual information, 191
- Monte Carlo approximation, **26**, 271
- Monte Carlo control, 219
- Monte Carlo rollouts, 270
- Monte Carlo trajectory sampling, 268
- Monte Carlo tree search, 266
- most likely explanation, 210
- multi-armed bandits, 175, **180**
- multicollinearity, 48
- multilayer perceptron, 147
- multinomial distribution, 5
- mutual information, 163, **165**

- natural parameters, 110
- negative log-likelihood, 29, **150**
- neural fitted Q-iteration, 236
- neural network, 147

- normal distribution, 6, **37**
- normalizing constant, 21

- objective function, 30
- Occam's razor, 50
- off-policy, 217
- on-policy, 217
- online, 62
- online actor-critic, 249
- online algorithms, 182
- online Bayesian linear regression, 55
- online learning, 180
- open-loop control, 269
- optimal control, 265
- optimal design, 174
- optimism in the face of uncertainty, **181**, 183, 189, 221, 228, 277
- optimistic exploration, 279
- optimistic Q-learning, 228
- output layer, 147
- output scale, 72, **75**
- output-space predictive entropy search, 177
- overfitting, 19

- Pareto distribution, 133
- partially observable Markov decision process, 209
- partition function, 110
- performance plan, 284
- Pinsker's inequality, 126
- planning, 114, **195**, 215
- plate notation, 12
- Platt scaling, 159
- plausible model, 278
- point density, 9
- point-based policy iteration, 212
- point-based value iteration, 212
- pointwise convergence, 87
- policy, 196
- policy gradient method, 238
- policy gradient theorem, 247
- policy iteration, **205**, 252
- policy search method, 238
- policy value function, 239
- Polyak averaging, 236

- Polyak-Łojasiewicz inequality, 140
- population risk, 19
- positive definite, 42
- positive definite kernel, 72
- positive semi-definite, 42
- posterior, 20, **21**
- precision matrix, 37
- predictive entropy search, 176
- predictive posterior, 21
- prior, 20
- probabilistic ensembles with trajectory sampling, 276
- probabilistic inference for learning control, 271, **275**
- probability, 3
- probability density function, 6
- probability matching, 189
- probability measure, 2
- probability of improvement, 187
- probability simplex, 5, 212
- probability space, 3
- probit likelihood, 97
- product rule, 9
- proposal distribution, 130
- proximal policy optimization, 251
- pseudo-random number generators, 7
- pushforward measure, 18

- Q actor-critic, 249
- Q-function, 197
- Q-learning, 227
- quadratic form, 42
- quantile function, 7

- radial basis function kernel, 72
- random Fourier features, 84
- random shooting methods, 266
- random variable, 4
- random vector, 8
- rapidly mixing Markov chain, 127
- realizations of a random variable, 4
- receding horizon control, 265
- rectified linear unit, 148
- recursive Bayesian estimation, 61
- redundancy, 167
- regression, 18
- regret, 181
- regret-information ratio, 190
- Reichenbach's common cause principle, 11
- REINFORCE algorithm, 244
- reinforcement learning, 215
- relative entropy, 104
- relative Fisher information, 141
- reliability diagram, 158
- reparameterizable distribution, 117
- reparameterization trick, **116**, 255, 268
- replay buffer, 236
- representer theorem, 77
- reproducing kernel Hilbert space, 76
- reproducing property, 77
- reverse KL-divergence, 107
- reversible Markov chain, 128
- reward shaping, 208
- reward to go, 243
- ridge regression, 48
- Robbins-Monro algorithm, 35
- Robbins-Monro conditions, 35
- robust control, 281
- rollout, 239

- saddle point, 31
- safety filter, 286
- safety plan, 285
- sample covariance matrix, 23
- sample mean, 23
- sample space, 2
- sample variance, 23
- SARSA, 225
- score function, **116**, 240
- score function trick, 240
- score gradient estimator, 116, **240**
- second-order characterization of convexity, 34
- second-order expansion, 34
- semi-conjugate prior, 133
- sharply concentrated, 24
- shift-invariant kernel, 75
- small tails, 25
- snapshot, 154

- soft actor critic, **259**, 271
- soft Q-learning, 259
- soft value function, 259
- softmax exploration, 220
- softmax function, 148
- spectral density, 85
- square root, 43
- squared exponential kernel, 72
- standard deviation, 15
- standard normal distribution, 6, **37**
- state of a random variable, 4
- state space model, 61
- state value function, 197
- state-action value function, 197
- state-dependent baseline, 243
- stationary distribution, 124
- stationary kernel, 75
- stationary point, 31
- stochastic average approximation, 267
- stochastic environment, 195
- stochastic gradient descent, 35
- stochastic gradient Langevin dynamics, 142
- stochastic matrix, 123
- stochastic process, 122
- stochastic semi-gradient descent, 233
- stochastic value gradients, **256**, 270
- stochastic weight averaging, 155
- stochastic weight averaging-Gaussian, 155
- strict convexity, 32
- sub-Gaussian random variable, 27
- sublinear regret, 181
- submodularity, 169
- subsampling, 154
- subset of regressors, 89
- sufficient statistic, 110, 154, **167**
- sum rule, 8
- supervised learning, **18**, 163
- support, 5
- supremum norm, 201
- surprise, 100
- synergy, 167
- tabular setting, 215
- tail distribution function, 25
- tail index, 134
- target space, 4
- temperature scaling, 159
- temporal models, 67
- temporal-difference error, 233
- temporal-difference learning, 224
- testing conditional, 89
- Thompson sampling, **189**, 277
- time-homogeneous process, 122
- total variation distance, 126
- tower rule, 14
- training conditional, 89
- trajectory, 216
- trajectory sampling, 267
- transition, 216
- transition function, 122
- transition graph, 123
- transition matrix, 123
- trust-region policy optimization, 250
- twin delayed DDPG, 253
- two-filter smoothing, 65
- unadjusted Langevin algorithm, 138
- unbiased estimator, 23
- uncertainty sampling, 172
- uncorrelated, 15
- uniform convergence, 87
- uniform distribution, 5, 7
- union bound, 3
- universal approximation theorem, 148
- universality of the uniform, 7, 315
- upper confidence bound, 183
- value iteration, 207
- variance, 16
- variational family, 99
- variational inference, **98**, 271
- variational parameters, 91
- variational posterior, 91
- Viterbi algorithm, 210
- weak law of large numbers, 27
- weight decay, 36
- weight-space view, 50

Weinstein-Aronszajn identity, 287

well-calibrated confidence interval, 184

Wiener process, 74, 122, **138**

Woodbury matrix identity, 287

Errata

October 6, 2023	fixed typo in eq. (1.61)
October 19, 2023	revised structure of chapter 5 and section 6.3
October 23, 2023	fixed typo in example 3.5
October 29, 2023	ensure existence of limit in eq. (5.56)
October 31, 2023	clarify index i in theorem 6.21
November 4, 2023	revised structure of chapters 2 and 7
November 11, 2023	fig. 8.4 was wrongly stating that uncertainty sampling picks the point maximizing “total” uncertainty
November 14, 2023	section 7.3.3 was incorrectly referring to “dropout” rather than “dropconnect”
November 14, 2023	fixed typo in eq. (9.14)