

# Decision Making Under Uncertainty

April 28, 2014



# Contents

<b>1</b>	<b>Introduction</b>	<b>7</b>
1.1	Exploration-exploitation . . . . .	8
1.1.1	Introduction . . . . .	8
<b>2</b>	<b>Probability concepts</b>	<b>13</b>
2.1	Probability . . . . .	14
2.1.1	Sets, experiments and sample spaces . . . . .	15
2.1.2	Events, measure and probability . . . . .	16
2.1.3	Conditioning and independence . . . . .	19
2.2	Random variables . . . . .	21
2.2.1	Discrete and continuous random variables . . . . .	23
2.2.2	Random vectors . . . . .	23
2.2.3	Moments . . . . .	25
2.3	Conclusion . . . . .	26
<b>3</b>	<b>Subjective probability and utility</b>	<b>29</b>
3.1	Subjective probability . . . . .	30
3.1.1	Relative likelihood . . . . .	30
3.1.2	Subjective probability assumptions . . . . .	31
3.1.3	Assigning unique probabilities . . . . .	32
3.1.4	Conditional likelihoods . . . . .	33
3.1.5	Probability elicitation . . . . .	34
3.2	Utility theory . . . . .	35
3.2.1	Rewards and preferences . . . . .	35
3.2.2	Preferences among distributions . . . . .	37
3.2.3	Utility . . . . .	37
3.2.4	Measuring utility . . . . .	39
3.2.5	Convex and concave utility functions . . . . .	40
3.3	Summary . . . . .	40
<b>4</b>	<b>Decision problems</b>	<b>43</b>
4.1	Introduction . . . . .	44
4.2	Rewards that depend on the outcome of an experiment . . . . .	44
4.2.1	Formalisation of the problem setting . . . . .	45
4.2.2	Decision diagrams . . . . .	47
4.2.3	Statistical estimation* . . . . .	48
4.3	Bayes decisions . . . . .	49
4.3.1	Convexity of the Bayes-optimal utility* . . . . .	50

4.4	Decision problems with observations . . . . .	52
4.4.1	Calculating posteriors . . . . .	56
<b>5</b>	<b>Estimation</b>	<b>59</b>
5.1	Calculation of posterior distributions . . . . .	60
5.2	Sufficient statistics . . . . .	60
5.2.1	Formalisation of sufficient statistics . . . . .	61
5.2.2	Exponential families . . . . .	62
5.3	Conjugate priors . . . . .	63
5.3.1	Bernoulli-Beta conjugate pair . . . . .	63
5.4	Credible intervals . . . . .	66
5.5	Concentration inequalities . . . . .	69
5.6	Approximate Bayesian approaches . . . . .	73
5.6.1	Approximate Bayesian Computation . . . . .	74
5.7	Other conjugate families . . . . .	74
5.7.1	Conjugates for the normal distribution . . . . .	74
5.7.2	Conjugates for multivariate distributions . . . . .	77
5.8	Summary . . . . .	80
<b>6</b>	<b>Hypothesis testing</b>	<b>83</b>
6.1	Decision problems . . . . .	84
6.2	Unconditional and conditional errors . . . . .	87
6.3	Two-point test with Bernoulli trials . . . . .	89
6.4	Null hypothesis tests . . . . .	92
6.5	The fallacy of $P$ -values . . . . .	95
<b>7</b>	<b>Sequential sampling</b>	<b>99</b>
7.1	Gains from sequential sampling . . . . .	100
7.1.1	An example sequential problem . . . . .	101
7.2	Sequential decision procedures . . . . .	104
7.3	Calculating the expected utility of a sequential decision procedure	105
7.4	Backwards induction . . . . .	107
7.5	Unbounded sequential decision procedures . . . . .	107
7.6	The sequential probability ratio test . . . . .	108
7.6.1	Wald's theorem . . . . .	112
7.7	Martingales . . . . .	112
7.7.1	Doob martingales . . . . .	113
7.7.2	The Azuma-Hoeffding inequality . . . . .	113
7.8	Markov processes . . . . .	114
<b>8</b>	<b>Experiment design and Markov decision processes</b>	<b>117</b>
8.1	Introduction . . . . .	118
8.1.1	Experiment design: examples . . . . .	118
8.2	Bandit problems . . . . .	118
8.2.1	Bernoulli bandits . . . . .	119
8.2.2	Decision-theoretic bandit process . . . . .	120
8.3	Markov decision processes and reinforcement learning . . . . .	121
8.3.1	Value functions . . . . .	124
8.4	Finite horizon, undiscounted problems . . . . .	125
8.4.1	Policy evaluation . . . . .	125

8.4.2	Monte-Carlo policy evaluation . . . . .	126
8.4.3	Finite horizon backwards induction . . . . .	128
8.5	Infinite-horizon . . . . .	129
8.5.1	Examples . . . . .	129
8.5.2	Markov chain theory for discounted problems . . . . .	132
8.5.3	Optimality equations . . . . .	134
8.5.4	MDP Algorithms . . . . .	136
8.6	Further reading . . . . .	142
<b>9</b>	<b>Reinforcement learning and stochastic approximation</b>	<b>145</b>
9.1	Introduction . . . . .	146
9.1.1	Bandit problems . . . . .	146
9.1.2	Estimation and Robbins-Monro approximation . . . . .	147
9.1.3	The theory of the approximation . . . . .	149
9.2	Dynamic problems . . . . .	150
9.2.1	Monte-Carlo policy evaluation and iteration . . . . .	152
9.2.2	Temporal difference methods . . . . .	154
9.2.3	Stochastic value iteration methods . . . . .	155
<b>10</b>	<b>Approximate dynamic programming</b>	<b>163</b>
10.1	Introduction . . . . .	164
10.1.1	Error bounds . . . . .	164
10.1.2	Features . . . . .	165
10.2	Approximate policy iteration . . . . .	166
10.2.1	Estimation building blocks . . . . .	167
10.2.2	The value estimation step . . . . .	168
10.2.3	Policy estimation . . . . .	169
10.2.4	Rollout-based policy iteration methods . . . . .	170
10.2.5	Least Squares Methods . . . . .	171
10.3	Approximate Value Iteration . . . . .	172
10.3.1	Approximate backwards induction . . . . .	173
10.3.2	State aggregation . . . . .	173
10.3.3	Representative states . . . . .	174
<b>11</b>	<b>Bayesian reinforcement learning</b>	<b>177</b>
11.1	Introduction . . . . .	178
11.2	Bayesian reinforcement learning . . . . .	178
11.2.1	Updating the belief . . . . .	180
11.3	Finding Bayes-optimal policies . . . . .	181
11.3.1	The expected MDP heuristic . . . . .	181
11.3.2	The maximum MDP heuristic . . . . .	183
11.3.3	Bounds on the expected utility . . . . .	184
11.3.4	Tighter lower bounds . . . . .	185
11.3.5	Stochastic branch and bound . . . . .	189
11.4	Partially observable Markov decision processes . . . . .	189

<b>12 Distribution-free reinforcement learning</b>	<b>193</b>
12.1 Introduction . . . . .	194
12.2 Bandit problems . . . . .	194
12.2.1 UCB . . . . .	196
12.3 Structured bandit problems . . . . .	197
12.4 Reinforcement learning problems . . . . .	198
12.4.1 Optimality Criteria . . . . .	198
12.4.2 UCRL . . . . .	200
12.4.3 Bibliographical remarks . . . . .	202
.1 Symbols . . . . .	203
.2 Index . . . . .	204
.3 Glossary . . . . .	206

# Chapter 1

## Introduction

The purpose of this book is to collect all fundamental results for decision making under uncertainty in one place, much as the book by Puterman on Markov decision processes did for Markov decision process theory. In particular, the aim is to give a unified account of algorithms and theory for sequential problems and reinforcement learning. Starting from elementary statistical decision theory, we progress to the reinforcement learning problem, its formalisation and various solution methods. The end of the book focuses on the current state of the art in approximations.

The book may be read as follows . . .

### *The problem of decision making under uncertainty*

- *Modelling* our **uncertainty** about the world  $\Rightarrow$  learning
- *Optimising* our **decisions** given our knowledge  $\Rightarrow$  planning

### *Applications and related problems*

- Optimisation: robust decisions, efficient search, planning.
- AI: modelling, learning from interaction and/or demonstration.
- Economics: Mechanism design, behavioural modelling.
- Security: Cryptography, Biometrics, Intrusion detection and response
- Biology and Medicine: Automatic experiment design, clinical trials, cognitive science.

Planning and learning are connected through the exploration-exploitation trade-off.

## 1.1 Exploration-exploitation

### 1.1.1 Introduction

#### The exploration-exploitation trade-off

**Example 1.1.1** (Selecting a restaurant). *Consider the problem of selecting the restaurant to go to.*

- *You usually go to Les Epinards. The food there is usually to your taste and satisfactory.*
- *You heard that King's Arm is really good! It is tempting to try it out. But there is a risk involved.*



- *It's Friday. Do you:*
  - *Go to Les Epinards?*
  - *Call King's Arm to reserve?*
  - *Check the menu of King's Arm and then decide?*

### ***The exploration-exploitation trade-off***

---

- **Exploit** knowledge about the world to gain a *known* reward.
- **Explore** the world to **learn**, *potentially* getting less or more reward.
- Arises when data collection is **interactive**.

### **Formalising decision problems**

- How do our decisions depend on what we want?
- How do we weigh alternatives?
- Is there a good concept of rationality?

### **Beliefs, learning and planning**

- How can we express belief and how does belief change?
- How might we make decisions according to our beliefs?
- What if our decisions can affect our beliefs?

### **Why decision theory?**

- Formalising trade-offs makes problems well-posed.
- Better overall solutions could be found.
- We may ignore non-essential aspects.

Reinforcement learning is the problem of learning to act in an unknown environment, solely through interaction and some limited feedback. The learner does not necessarily have access to detailed instructions on how to perform a

task. Instead, it performs *actions*, which affect the environment and obtains some *feedback*. Sometimes the actions result in *rewards*, which represent the agent's desires. The learning problem is then formulated as the problem of learning how to act to maximise total reward.

This is a fundamental problem in artificial intelligence, since frequently we can tell robots, computers and cars only *what* we want them to do, but have no idea how we can program them. We would like to simply give them a description of our goals and then let them *explore* the environment on their own to find a good solution. Since the world is partially unknown, the learner always has to decide between two general types of decisions. It can either decide to do something which will give it some known reward, or it could take a risk and try something new. This is called the exploration-exploitation tradeoff.

Of course, animals and humans do learn in this way. Through imitation, exploration and reward signals, humans shape their behaviour to achieve their goals. In fact, it has been known since the 90s that there is some connection between some reinforcement learning algorithms and mechanisms in the basal ganglia.

Finally, there are many connections between reinforcement learning and other fields. Some algorithms used for reinforcement learning are also equivalent to fundamental algorithms in other fields. In particular, the general Bayes-optimal reinforcement learning algorithm is formally identical to Bayes-optimal algorithms for the automatic design of experiments and clinical trials. In addition, algorithms developed for more specialised problems have found application in optimisation of stochastic functions and game-theoretic problems.

## Outline

- \* Probability refresher. Measure theory; Axiomatic definition of probability; Conditional probability; Bayes' theorem; Random variables; Expectation
- 1. Subjective probability and utility. The notion of subjective probability; eliciting priors; the concept of utility; expected utility
- 2. Decision problems. maximising expected utility; maximin utility; regret;
- 3. Estimation. Estimation as conditioning; families of distributions that are closed under conditioning; conjugate priors; concentration inequalities; PAC and high-probability bounds; (\*) Markov Chain Monte Carlo; (\*) ABC estimation
- \* Hypothesis testing. Bayesian hypothesis testing; point hypotheses; null hypothesis testing; unconditional and conditional hypothesis test procedures; frequentist testing; the fallacy of  $p$ -values
- 4. Sequential sampling and optimal stopping. Sequential sampling problems; the cost of sampling; optimal stopping; Martingales
- 5. Automatic experiment design and bandit problems. Belief and information state; bandit problems; Markov decision processes; backwards induction

6. Reinforcement learning I: Markov decision processes and fundamental algorithms. simple learning strategies for bandit problems; Markov decision processes; value iteration; policy iteration; policy search; model-free methods
7. Reinforcement learning II: Stochastic and approximation algorithms  $Q$ -learning; approximate value iteration; approximate policy iteration
8. Reinforcement learning III: Generalised problems. continuous case; partially observable case; (\*) multi-agent case
9. Project meeting.
10. Reinforcement learning IV: Bayesian algorithms Bounds on the utility; Thompson sampling; Stochastic Branch and Bound; Sparse sampling; Rollout sampling;
11. Reinforcement learning V: Bandit algorithms and regret Tree and metric bandits; UCRL; (\*) Bounds for Thompson sampling
12. Project meeting.
13. Learning with expert advice Regret; Prediction with expert advice; Prediction with incomplete information; Prediction with side-information; Connections with game theory
14. Learning by demonstration; Preference Elicitation

**Assessment****Exercises and feedback: 40%**

- Exercises after every unit.
- Exercise sets include feedback form.
- Necessary for a good project!

**Participation: 10%**

- Active participation in the course.
- Corrections on course notes.

**Project: 50%**

Competition, presentation and report.

- Team competition using rl-glue socket API.
- Each team codes:
  - An environment (test-bed).
  - An agent.
- Agents are evaluated on all environments.

**Themes**

- Models for representing belief and preferences.
- Algorithms for decision making.
- Fast optimisation.
- Applications in finance.
- Decision making in animals.
- Inferring preferences and beliefs.
- Automatic design of experiments.

## Chapter 2

# Probability concepts

## 2.1 Probability

### Two notions of probability

While probability is a simple mathematical construction, philosophically it has had at least two different meanings. In the classical sense, a probability distribution is a description for a truly random event. In the subjectivist sense, probability is merely a description for uncertainty which may or may not be due to randomness.

#### Classical Probability

- A *random experiment* is performed, with a given set  $\Omega$  of possible outcomes. A simple example is the 2-slit experiment in physics, where a particle is generated and which can go through either one of two slits. According to our current understanding of quantum theory, it is impossible to predict which slit the particle will go through. There, the set of possible events correspond to the particle passing through one or the other slit.
- We care about the probability that the particle will go through one of the two slots in the experiment. Does it depend on where the other particles have passed through? In the 2-slit experiment, the probabilities of either event can be actually accurately calculated. However, which slit the particle will go through is fundamentally unpredictable.

Such quantum experiments are the only ones that are currently thought of as truly random (though some people disagree about that too). Any other procedure, such as tossing a coin or casting a die, is inherently deterministic and only *appears* random due to the difficulty in predicting the outcome. That is, modelling a coin toss as a random process is usually the best approximation we can make in practice, given our uncertainty about the complex dynamics involved.

#### Subjective Probability

- We assume that  $\Omega$  is a set of possible *worlds* or realities. This set can be quite large and include anything imaginable. For example, it may include worlds where dragons are real. However, in practice one only cares about certain aspects of the world.
- We can interpret the probability of a world in  $\Omega$  as a belief that it is the true world.

In such a setting there is an actual true world  $\omega^* \in \Omega$ , which is simply unknown. This could have been set by Nature to an arbitrary value deterministically. The probability only reflects our lack of knowledge.

### 2.1.1 Sets, experiments and sample spaces

#### Set theory definitions

A very useful way to describe a set  $A$  is as follows

$$A \triangleq \{x \mid x \text{ have property } Y\}$$

for example

$$B(c, r) \triangleq \{x \in \mathbb{R}^n \mid \|x - c\| \leq r\}$$

describes the set of points enclosed in an  $n$ -dimensional sphere of radius  $r$  with center  $c \in \mathbb{R}^n$ .

- If an element  $x$  *belongs* to a set  $A$ , we write  $x \in A$ .
- Let the *sample space*  $\Omega$  be a set such that  $\omega \in \Omega$  always.
- We say that  $A$  is a *subset* of  $B$  or that  $B$  *contains*  $A$ , and write  $A \subset B$ , iff,  $x \in B$  for any  $x \in A$ .
- Let  $B \setminus A \triangleq \{x \mid x \in B \wedge x \notin A\}$  be the set difference.
- Let  $A \triangle B \triangleq (B \setminus A) \cup (A \setminus B)$  be the symmetric set difference.
- The *complement* of any  $A \subset \Omega$  is  $A^c \triangleq \Omega \setminus A$ .
- The *empty set* is  $\emptyset = \Omega^c$ .
- The *union* of  $n$  sets:  $A_1, \dots, A_n$  is  $\bigcup_{i=1}^n A_i = A_1 \cup \dots \cup A_n$ .
- The *intersection* of  $n$  sets  $A_1, \dots, A_n$  is  $\bigcap_{i=1}^n A_i = A_1 \cap \dots \cap A_n$ .
- $A$  and  $B$  are *disjoint* if  $A \cap B = \emptyset$ .

#### Experiments and sample spaces

##### Experiments

The set of possible experimental outcomes of an experiment is called the *sample space*  $\Omega$ .

- $\Omega$  must contain all possible outcomes.
- Each statistician  $i$  may consider a different  $\Omega_i$  for the same experiment.

**Example 2.1.1.** *Experiment: give medication to a patient.*

- $\Omega_1 = \{\text{Recovery within a day, No recovery after a day}\}$ .
- $\Omega_2 = \{\text{The medication has side-effects, No side-effect}\}$ .
- $\Omega_3 = \text{all combinations of the above.}$

**Product spaces**

- We perform  $n$  experiments.
- Assume that the  $i$ -th experiment has sample space  $\Omega_i$ .
- The *Cartesian product* or *product space* is defined as

$$\Omega_1 \times \cdots \times \Omega_n = \{(s_1, \dots, s_n) \mid s_i \in \Omega_i, \forall i \in \{1, \dots, n\}\} \quad (2.1.1)$$

the set of all ordered  $n$ -tuples  $(s_1, \dots, s_n)$ .

- The sample space  $\prod_{i=1}^n \Omega_i$  can be thought of as a sample space of a *composite* experiment in which all  $n$  experiments are performed.

**Identical experiment sample spaces**

- In many cases,  $\Omega_i = \Omega$  for all  $i$ , i.e. the sample space is identical for all individual experiment (e.g.  $n$  coin tosses).
- We then write  $\Omega^n = \prod_{i=1}^n \Omega$ .

**2.1.2 Events, measure and probability****Events and probability****Probability of a set**

If  $A$  is a subset of  $\Omega$ , the probability of  $A$  is a measure of the chances that the outcome of the experiment will be an element of  $A$ .

**Which sets?**

Ideally, we would like to be able to assign a probability to *every subset* of  $\Omega$ . However, for technical reasons, this is not possible.

**Example 2.1.2.** Let  $X$  be uniformly distributed on  $[0, 1]$ .

- What is the probability that  $X$  will be in  $[0, 1/4]$ ?
- What is the probability that  $X$  will be in  $[1/4, 1]$ ?
- What is the probability that  $X$  will be a rational number?



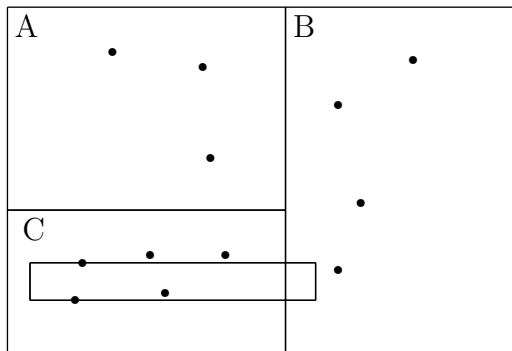


Figure 2.1: A fashionable apartment

**Measure theory primer**

Imagine that you have an apartment  $\Omega$  composed of three rooms,  $A, B, C$ . There are some coins on the floor and a 5-meter-long red carpet. We can measure various things in this apartment.

**Area**

- A:  $4 \times 5 = 20m^2$ .
- B:  $6 \times 4 = 24m^2$ .
- C:  $2 \times 5 = 10m^2$ .

**Coins on the floor**

- A: 3.
- B: 4
- C: 5.

**Length of red carpet**

- A:  $0m$
- B:  $0.5m$
- C:  $4.5m$ .

Measure the sets:  $\mathcal{F} = \{\emptyset, A, B, C, A \cup B, A \cup C, B \cup C, A \cup B \cup C\}$ . It is

easy to see that the union of any sets in  $\mathcal{F}$  is also in  $\mathcal{F}$ . In other words,  $\mathcal{F}$  is closed under union. Furthermore,  $\mathcal{F}$  contains the whole space  $\Omega$ .

Note that all those measures have an *additive property*.

### Measure and probability

As previously mentioned, the probability of  $A \subset \Omega$  is a measure of the chances that the outcome of the experiment will be an element of  $A$ . Here we give a precise definition of what we mean by measure and probability.

**Definition 2.1.1** (A field on  $\Omega$ ). *A family  $\mathcal{F}$  of sets, such that for each  $A \in \mathcal{F}$ ,  $A \subset \Omega$ , is called a field on  $\Omega$  if and only if*

1.  $\Omega \in \mathcal{F}$
2. if  $A \in \mathcal{F}$ , then  $A^c \in \mathcal{F}$ .
3. For any  $A_1, A_2, \dots, A_n$  such that  $A_i \in \mathcal{F}$ , it holds that:  $\bigcup_{i=1}^n A_i \in \mathcal{F}$ .

From the above definition, it is easy to see that  $A_i \cap A_j$  is also in the field.

**Definition 2.1.2** ( $\sigma$ -field on  $\Omega$ ). *A family  $\mathcal{F}$  of sets, such that  $\forall A \in \mathcal{F}$ ,  $A \subset \Omega$ , is called a  $\sigma$ -field on  $\Omega$  if and only if*

1.  $\Omega \in \mathcal{F}$
2. if  $A \in \mathcal{F}$ , then  $A^c \in \mathcal{F}$ .
3. For any sequence  $A_1, A_2, \dots$  such that  $A_i \in \mathcal{F}$ , it holds that:  $\bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$ .

It is easy to verify that the  $\mathcal{F}$  given in the apartment example satisfies these properties.

**Definition 2.1.3** (Measure). *A measure  $\lambda$  on  $(\Omega, \mathcal{F})$  is a function  $\lambda : \mathcal{F} \rightarrow \mathbb{R}^+$  such that*

1.  $\lambda(\emptyset) = 0$ .
2.  $\lambda(A) \geq 0$  for any  $A \in \mathcal{F}$ .
3. For any collection of subsets  $A_1, \dots, A_n$  with  $A_i \in \mathcal{F}$  and  $A_i \cap A_j = \emptyset$ .

$$\lambda\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \lambda(A_i) \quad (2.1.2)$$

It is easy to verify that the floor area, the number of coins, and the length of the red carpet are all measures. In fact, the area and length correspond to what is called a *Lebesgue measure* and the number of coins to a *counting measure*.

**Definition 2.1.4** (Probability measure). *A probability measure  $P$  on  $(\Omega, \mathcal{F})$  is a function  $P : \mathcal{F} \rightarrow [0, 1]$  such that:*

1.  $P(\Omega) = 1$
2.  $P(\emptyset) = 0$
3.  $P(A) \geq 0$  for any  $A \in \mathcal{F}$ .

4. If  $A_1, A_2, \dots$  are disjoint then

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i) \quad (\text{union})$$

$(S, \mathcal{F}, P)$  is called a probability space.

So, probability is just a special type of measure.

### The Lebesgue measure

**Definition 2.1.5** (Outer measure). Let  $(\Omega, \mathcal{F}, \lambda)$  be a measure space. The outer measure of a set  $A \subset \Omega$  is:

$$\lambda^* \triangleq \inf \left\{ \sum_k \lambda(B_k) : A \subset \bigcup_k B_k \right\}. \quad (2.1.3)$$

**Definition 2.1.6** (Inner measure). Let  $(\Omega, \mathcal{F}, \lambda)$  be a measure space. The inner measure of a set  $A \subset \Omega$  is:

$$\lambda_* \triangleq \lambda(\Omega) - \lambda(\Omega \setminus A). \quad (2.1.4)$$

**Definition 2.1.7** (Lebesgue measurable sets). A set  $A$  is (Lebesgue) measurable if the outer and inner measures are equal.

$$\lambda^*(A) = \lambda_*(A). \quad (2.1.5)$$

The common value of the inner and outer measure is called the Lebesgue measure<sup>1</sup>  $\bar{\lambda} = \lambda^*(A)$ .

### 2.1.3 Conditioning and independence

#### Independent events and conditional probability

Events correspond to sets. Thus, the probability of the event that a draw from  $\Omega$  is in  $A$  is equal to the probability measure of  $A$ ,  $P(A)$ .

**Definition 2.1.8** (Independent events). Two events  $A, B$  are independent if  $P(A \cap B) = P(A)P(B)$ . The events in a family  $\mathcal{F}$  of events are independent if for any sequence  $A_1, A_2, \dots$  of events in  $\mathcal{F}$ ,

$$P\left(\bigcap_{i=1}^n A_i\right) = \prod_{i=1}^n P(A_i) \quad (\text{independence})$$

**Definition 2.1.9** (Conditional probability). The conditional probability of  $A$  when  $B$ , s.t.  $P(B) > 0$ , is given is:

$$P(A | B) \triangleq \frac{P(A \cap B)}{P(B)}. \quad (2.1.6)$$

Of course,  $P(A \cap B) = P(A | B)P(B)$  even if  $A, B$  are not independent.

---

<sup>1</sup>It is easy to see that  $\bar{\lambda}$  is a measure.

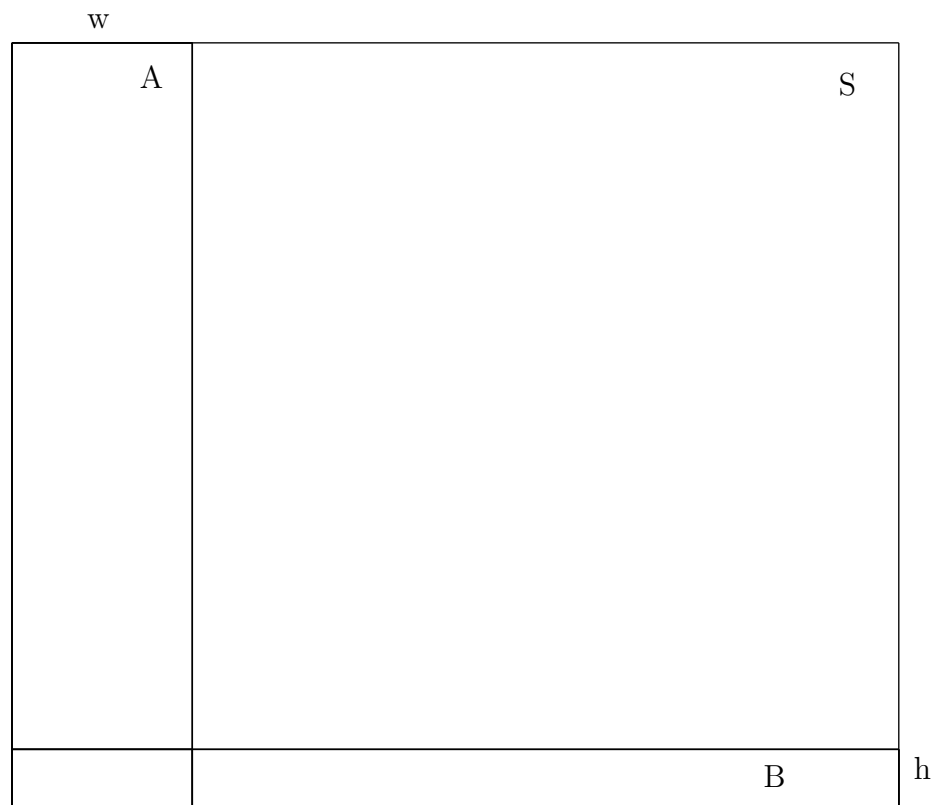


Figure 2.2: In the above case,  $S$  is a unit square and taking  $P$  to be the Lebesgue measure, we see that  $P(S) = 1 \cdot 1$ ,  $P(A) = 1 \cdot w$ ,  $P(B) = h \cdot 1$  and  $P(A \cap B) = wh$ , so  $A$  and  $B$  are independent.

**Bayes' theorem**

The following theorem trivially follows from the above discussion. However, versions of it shall be used repeatedly throughout. For this reason we present it here together with a detailed proof.

**Theorem 2.1.1** (Bayes' theorem). *Let  $A_1, A_2, \dots$  be a (possibly infinite) sequence of disjoint events such that  $\bigcup_{i=1}^{\infty} A_i = \Omega$  and  $P(A_i) > 0$  for all  $i$ . Let  $B$  be another event with  $P(B) > 0$ . Then*

$$P(A_i | B) = \frac{P(B | A_i)P(A_i)}{\sum_{j=1}^{\infty} P(B | A_j)P(A_j)} \quad (2.1.7)$$

*Proof.* From (2.1.6),  $P(A_i | B) = P(A_i \cap B)/P(B)$  and also  $P(A_i \cap B) = P(B | A_i)P(A_i)$ . Thus

$$P(A_i | B) = \frac{P(B | A_i)P(A_i)}{P(B)},$$

and we continue analyzing the denominator  $P(B)$ . First, due to  $\bigcup_{i=1}^{\infty} A_i = \Omega$  we have  $B = \bigcup_{j=1}^{\infty} (B \cap A_j)$ . Since  $A_i$  are disjoint, so are  $B \cap A_i$ . Then from the union property of probability distributions we have

$$P(B) = P\left(\bigcup_{j=1}^{\infty} (B \cap A_j)\right) = \sum_{j=1}^{\infty} P(B \cap A_j) = \sum_{j=1}^{\infty} P(B | A_j)P(A_j),$$

which finishes the proof.  $\square$

**Binomial coefficients**

Binomial coefficients appear in a lot of different distributions. They are especially useful for combinatorial problems.

$$\binom{x}{n} \triangleq \frac{\prod_{i=0}^{n-1} (x-i)}{n!}, \quad x \in \mathbb{R}, n \in \mathbb{N}, \quad (2.1.8)$$

and  $\binom{x}{0} = 1$ . It follows that

$$\binom{k}{n} = \frac{k!}{n!(k-n)!} \quad k, n \in \mathbb{N}, k \geq n. \quad (2.1.9)$$

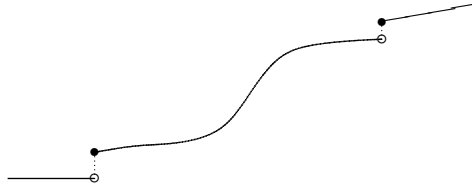
## 2.2 Random variables

**Random variables**

A random variable  $X$  is a special kind of random quantity, defined as a real function of outcomes in  $\Omega$ . Thus, it also defines a mapping from a probability measure  $P$  on  $(\Omega, \mathcal{F})$  to a probability measure  $P_X$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . More precisely, we define the following.

**Definition 2.2.1** (Measurable function). *Let  $\mathcal{F}$  on  $\Omega$  be a  $\sigma$ -field. A function  $g : \Omega \rightarrow \mathbb{R}$  is said to be measurable with respect to  $\mathcal{F}$ , or  $\mathcal{F}$ -measurable, if, for any  $x \in \mathbb{R}$ ,*

$$\{s \in \Omega \mid g(s) \leq x\} \in \mathcal{F}.$$

Figure 2.3: A distribution function  $F$ 

**Definition 2.2.2** (Random variable). *Let  $(\Omega, \mathcal{F}, P)$  be a probability space. A random variable  $X : \Omega \rightarrow \mathbb{R}$  is a real-valued,  $\mathcal{F}$ -measurable function.*

### *The distribution of $X$*

Every random variable  $X$  induces a probability measure  $P_X$  on  $\mathbb{R}$ . For any  $B \subset \mathbb{R}$  we define

$$P_X(B) = \mathbb{P}(X \in B) = P(\{s \mid X(s) \in B\}). \quad (2.2.1)$$

Thus, the probability that  $X$  is in  $B$  is equal to the  $P$ -measure of the points  $s \in \Omega$  such that  $X(s) \in B$  and also equal to the  $P_X$ -measure of  $B$ .

Here  $\mathbb{P}$  is used as a *short-hand* notation.

**Exercise 1.**  $\Omega$  is the set of 52 playing cards.  $X(s)$  is the value of each card (1, 10 for the ace and figures respectively). What is the probability of drawing a card  $s$  with  $X(s) > 7$ ?

### (Cumulative) Distribution functions

**Definition 2.2.3** ((Cumulative) Distribution function). *The distribution function of a random variable  $X$  is the function  $F : \mathbb{R} \rightarrow \mathbb{R}$ :*

$$F(t) = \mathbb{P}(X \leq t). \quad (2.2.2)$$

#### Properties

- If  $x \leq y$ , then  $F(x) \leq F(y)$ .
- $F$  is right-continuous.
- At the limit,

$$\lim_{t \rightarrow -\infty} F(t) = 0, \quad \lim_{t \rightarrow \infty} F(t) = 1.$$

### 2.2.1 Discrete and continuous random variables

#### Types of distributions

On the real line, there are two types of distributions for a random variable. Here, once more, we employ the  $\mathbb{P}$  notation as a shorthand for the probability of general events involving random variables, so that we don't have to deal with the measure notation. The two following examples should give some intuition.

##### Discrete distributions

$X : \Omega \rightarrow \{x_1, \dots, x_n\}$  takes  $n$  discrete values ( $n$  can be infinite). The probability function of  $X$  is

$$f(x) \triangleq \mathbb{P}(X = x),$$

defined for  $x \in \{x_1, \dots, x_n\}$ . For any  $B \subset \mathbb{R}$ :

$$P_X(B) = \sum_{x_i \in B} f(x_i).$$

In addition, we write  $\mathbb{P}(X \in B)$  to mean  $P_X(B)$ .

##### Continuous distributions

$X$  has a continuous distribution if there exists a *probability density function*  $f$  s.t.  $\forall B \subset \mathbb{R}$ :

$$P_X(B) = \int_B f(x) dx.$$

It is possible that  $X$  has neither a continuous, nor a discrete distribution.

### 2.2.2 Random vectors

#### Generalisation to $\mathbb{R}^m$

We can generalise to random *vectors* in a Euclidean space. Once more, there are two special cases of distributions for the random vector  $X = (X_1, \dots, X_n)$ .

##### Discrete distributions

$$\mathbb{P}(X_1 = x_1, \dots, X_m = x_m) = f(x_1, \dots, x_m)$$

##### Continuous distributions

For  $B \subset \mathbb{R}^m$

$$\mathbb{P}\{(X_1, \dots, X_m) \in B\} = \int_B f(x_1, \dots, x_m) dx_1 \cdots dx_m$$

**Measure-theoretic notation**

The previously seen special cases can be handled with a unified notation if we take advantage of the fact that probability is only a particular type of measure. As a first step, we note that summation can also be seen as integration with respect to the counting measure and that Riemann integration is integration with respect to the Lebesgue measure.

***Integral with respect to a measure  $\mu$*** 

Introduce the common notation  $\int \cdots d\mu(x)$ , where  $\mu$  is a measure. Let some real function  $g : \Omega \rightarrow \mathbb{R}$ . Then for any subset  $B \subset \Omega$  we can write

- Discrete case:  $f$  is the probability function and we choose the *counting measure* for  $\mu$ , so:

$$\sum_{x \in B} g(x)f(x) = \int_B g(x)f(x) d\mu(x)$$

Roughly speaking, the counting measure  $\mu(\Omega)$  is equal to the number of elements in  $\Omega$ .

- Continuous case:  $f$  is the probability density function and we choose the *Lebesgue measure* for  $\mu$ , so:

$$\int_B g(x)f(x) dx = \int_B g(x)f(x) d\mu(x)$$

Roughly speaking, the Lebesgue measure  $\mu(S)$  is equal to the volume of  $S$ .

In fact, since probability is a measure in itself, we do not need to complicate things by using  $f$  and  $\mu$  at the same time! This allows us to use the following notation.

**Lebesgue-Stieltjes notation**

If  $P$  is a probability measure on  $(\Omega, \mathcal{F})$  and  $B \subset \Omega$ , and  $g$  is  $\mathcal{F}$ -measurable, we write the probability that  $g(x)$  takes the value  $B$  can be written equivalently as:

$$\mathbb{P}(g \in B) = P_g(B) = \int_B g(x) dP(x) = \int_B g dP. \quad (2.2.3)$$

Intuitively,  $dP$  is related to densities in the following way. If  $P$  is a measure on  $\Omega$  and is absolutely continuous with respect to another measure  $\mu$ , then  $p \triangleq \frac{dP}{d\mu}$  is the (Radon-Nikodym) derivative of  $P$  with respect to  $\mu$ . We write the integral as  $\int gp d\mu$ . If  $\mu$  is the Lebesgue measure, then  $p$  coincides with the probability density function.



**Marginal distributions and independence**

Although this is a straightforward outcome of the set-theoretic definition of probability, we also define the marginal explicitly for random vectors.

**Marginal distribution**

The marginal distribution of  $X_1, \dots, X_k$  from a set of variables  $X_1, \dots, X_m$ , is

$$\mathbb{P}(X_1, \dots, X_k) \triangleq \int \mathbb{P}(X_1, \dots, X_k, X_{k+1} = x_{k+1}, \dots, X_m = x_m) d\mu(x_{k+1}, \dots, x_m). \quad (2.2.4)$$

In the above,  $\mathbb{P}(X_1, \dots, X_k)$  can be thought of as the probability measure for any events related to the random vector  $(X_1, \dots, X_k)$ . Thus, it defines a probability measure over  $(\mathbb{R}^k, \mathfrak{B}(\mathbb{R}^k))$ . In fact, let  $Y = (X_1, \dots, X_k)$  and  $Z = (X_{k+1}, \dots, X_m)$  for simplicity. Then define  $Q(A) \triangleq \mathbb{P}(Z \in A)$ , with  $A \subset \mathbb{R}^{m-k-1}$ . Then the above can be re-written as:

$$\mathbb{P}(Y \in B) = \int_{\mathbb{R}^{m-k-1}} \mathbb{P}(Y \in B \mid Z = z) dQ(z).$$

Similarly,  $\mathbb{P}(Y \mid Z = z)$  can be thought of as a function mapping from values of  $Z$  to probability measures. Let  $P_z(B) \triangleq \mathbb{P}(Y \in B \mid Z = z)$  be this measure corresponding to a particular value of  $z$ . Then we can write

$$\mathbb{P}(Y \in B) = \int_{\mathbb{R}^{m-k-1}} \left( \int_B dP_z(y) \right) dQ(z).$$

**Independence**

If  $X_i$  is independent of  $X_j$  for all  $i \neq j$ :

$$\mathbb{P}(X_1, \dots, X_m) = \prod_{i=1}^M \mathbb{P}(X_i), \quad f(x_1, \dots, x_m) = \prod_{i=1}^M g_i(x_i) \quad (2.2.5)$$

**2.2.3 Moments**

There are some simple properties of the random variable under consideration which are frequently of interest in statistics. Two of those properties are *expectation* and *variance*.

**Expectation**

**Definition 2.2.4.** The expectation  $\mathbb{E}(X)$  of any random variable  $X : \Omega \rightarrow R$ , where  $R$  is a vector space, with distribution  $P_X$  is defined by

$$\mathbb{E}(X) \triangleq \int_R t dP_X(t), \quad (2.2.6)$$

as long as the integral exists.

Furthermore,

$$\mathbb{E}[g(X)] = \int g(t) \, dP_X(t),$$

for any function  $g$ .

### Variance

**Definition 2.2.5.** The variance  $\mathbb{V}(X)$  of any random variable  $X : \Omega \rightarrow \mathbb{R}$  with distribution  $P_X$  is defined by

$$\begin{aligned} \mathbb{V}(X) &\triangleq \int_{-\infty}^{\infty} [t - \mathbb{E}(X)]^2 \, dP_X(t) \\ &= \mathbb{E} \left\{ [X - \mathbb{E}(X)]^2 \right\} \\ &= \mathbb{E}(X^2) - \mathbb{E}^2(X). \end{aligned} \tag{2.2.7}$$

When  $X : \Omega \rightarrow R$  with  $R$  an arbitrary vector space, the above becomes the *covariance matrix*:

$$\begin{aligned} \mathbb{V}(X) &\triangleq \int_{-\infty}^{\infty} [t - \mathbb{E}(X)] [t - \mathbb{E}(X)]^\top \, dP_X(t) \\ &= \mathbb{E} \left\{ [X - \mathbb{E}(X)] [X - \mathbb{E}(X)]^\top \right\} \\ &= \mathbb{E}(XX^\top) - \mathbb{E}(X) \mathbb{E}(X)^\top. \end{aligned} \tag{2.2.8}$$

### Divergences

One useful idea is KL-divergences on measures.

**Definition 2.2.6.** *KL-Divergence*

$$D(P \parallel Q) \triangleq \int \frac{dP}{dQ} \, dP. \tag{2.2.9}$$

### Empirical distributions

**Definition 2.2.7.** Let  $x^n = (x_1, \dots, x_n)$  drawn from a product measure  $x^n \sim P^n$  on the measurable space  $(\mathcal{X}^n, \mathfrak{F}_n)$ . Let  $\mathfrak{S}$  be any  $\sigma$ -field on  $\mathcal{X}$ . Then empirical distribution of  $x^n$  is defined as

$$\hat{P}_n(B) \triangleq \frac{1}{n} \sum_{t=1}^n \mathbb{I} \{x_t \in B\}. \tag{2.2.10}$$

## 2.3 Conclusion

### Recommended further reading

Most of this material is based on DeGroot [1970]. See Kolmogorov and Fomin [1999] for a really clear exposition of measure, starting from rectangle areas (developed from course notes in 1957). Also see Savage [1972] for a verbose, but interesting and rigorous introduction to subjective probability. More technical books, such as Ash and Doleáans-Dade [2000] are not very approachable by non-math graduates.

### Summary

- *Sample space*  $\Omega$  contains all possible *outcomes* of an experiment.
- $\sigma$ -field  $\mathcal{F}$  s.t.  $\forall A, B \in \mathcal{F}, A \subset \Omega, A \cup B \in \mathcal{F}, \Omega \in \mathcal{F}$ .
- Measurable space  $(\Omega, \mathcal{F})$ , measure space  $(\Omega, \mathcal{F}, \mu)$ .
- *Measure*  $\mu : \mathcal{F} \rightarrow \mathbb{R}$  such that  $\mu(\emptyset) = 0$ , and  $\mu(A_i) \geq 0$  for any  $A_i \in \mathcal{F}$ .  
For *disjoint*  $A_i$ ,  $\mu(\bigcup_i A_i) = \sum_i \mu(A_i)$ .
- Probability space  $(\Omega, \mathcal{F}, P)$ , with *probability measure*  $P$  such that  $P(\Omega) = 1$ .
- *Probability* that  $x \in A$ :

$$\mathbb{P}(x \in A) \triangleq P(A) = \int_A dP(t), \quad A \subset \Omega$$

- *Expectation* of  $X : \Omega \rightarrow Z$

$$\mathbb{E}(X) \triangleq \int_{\Omega} X(t) dP(t) = \int_Z u dP_X(u)$$

- *Conditional probability*

$$\mathbb{P}(A | B) = \frac{\mathbb{P}(A, B)}{\mathbb{P}(B)}, \quad P(A | B) = \frac{P(A \cap B)}{P(B)},$$

- *Marginal distribution*

$$\mathbb{P}(B) = \sum_i \mathbb{P}(B, A = i), \quad \sum_i \mathbb{P}(A = i) = 1,$$

$$P(B) = \sum_i P(B \cap A_i), \quad \bigcup_i A_i = \Omega.$$

- If  $A, B$  are *independent*

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



## Chapter 3

# Subjective probability and utility

### 3.1 Subjective probability

In order to make decisions, we need to be able to make predictions about the possible outcomes of each decision. Usually, we have *uncertainty* about what those outcomes are. This can be due to *stochasticity*, which is frequently used to model games of chance and inherently unpredictable physical phenomena. It can also be due to *partial information*, a characteristic of many natural problems. For example, it might be hard to guess at any one moment how much change you have in your wallet, whether you will be able to catch the next bus, or to remember where you left your keys.

In either case, this uncertainty can be expressed as a *subjective belief*. This does not have to correspond to reality. For example, some people believe, quite inaccurately, that if a coin comes up tails for a long time, it is quite likely to come up heads very soon. Or, you might quite happily believe your keys are in your pocket, only to realise that you left them at home as soon you arrive at the office.

In this book, we assume the view that subjective beliefs can be modelled as *probabilities*. This allows us to treat uncertainty due to stochasticity and due to partial information in a unified framework. In doing so, we shall treat each part of the problem as specifying a space of possible outcomes. What we wish to do is to find a *consistent way* of defining probabilities in the space of outcomes.

#### 3.1.1 Relative likelihood

Let us consider the simple example of guessing whether a tossed coin will come up head, or tails. Let  $\mathcal{S}$  be the sample space, and let  $A \subset \mathcal{S}$  be the set of tosses where the coin comes up heads, and  $B \subset \mathcal{S}$  be the set of tosses where it comes up tails. Here  $A \cap B = \emptyset$ , but there may be some other events such as the coin becoming lost, so it does not necessarily hold that  $A \cup B = \mathcal{S}$ . Nevertheless, we only care about whether  $A$  is more likely to occur than  $B$ . We can express that via the concept of relative likelihood:

**Definition 3.1.1.** (*The relative likelihood of two events  $A$  and  $B$* )

- If one thinks that  $A$  is more likely than  $B$ , then we write  $A \succ B$ , or equivalently  $B \prec A$ .
- If one thinks  $A$  is as likely as  $B$ , then we write  $A \approx B$ .

We also use  $\succeq$  and  $\preceq$  for at least as likely as and for no more likely than.

Let us now speak more generally about the case where we have defined an appropriate  $\sigma$ -field  $\mathcal{F}$  on  $\mathcal{S}$ . Then each element  $A_i \in \mathcal{F}$  will be a subset of  $\mathcal{S}$ . Furthermore, we have defined a relative likelihood relation for all elements  $A_i \in \mathcal{F}$ .<sup>1</sup>

As we would like to use the language of probability to talk about likelihoods, we would like to be able to define a probability measure that agrees with our given relations.

<sup>1</sup>More formally, we can define three classes:  $C_{\succ}, C_{\prec}, C_{\approx} \subset \mathcal{F}^2$  such that a pair  $(A_i, A_j) \in C_R$  if and only if it satisfies the relation  $A_i R A_j$ , where  $R \in \{\succ, \prec, \approx\}$ . It is easy to see that the three classes form a partition of  $\mathcal{F}^2$ .

A probability measure  $P : \mathcal{F} \rightarrow [0, 1]$  is said to *agree* with a relation  $A \precsim B$ , if it has the property that:  $P(A) \leq P(B)$  if and only if  $A \precsim B$ , for all  $A, B \in \mathcal{F}$ .

Of course, there are many possible measures that can agree with a given relation. It could even be that a given relational structure is incompatible with any possible probability measure. For that reason, we shall have to make some assumptions about relative likelihoods of events.

### 3.1.2 Subjective probability assumptions

Our beliefs must be *consistent*. This can be achieved if they satisfy some assumptions. First of all, it must always be possible to say whether one event is more likely than the other. Consequently, we are not allowed to claim ignorance.

**Assumption 3.1.1** (SP1). *For any pair of events  $A, B \in \mathcal{F}$ , one of the following must hold: Either  $A \succ B$ ,  $A \prec B$ , or  $A \approx B$ .*

If we can partition  $A, B$  in such a way that each part of  $A$  is less likely than its counterpart in  $B$ , then  $A$  is less likely than  $B$ . For example, let  $A_1$  be the event that it rains tomorrow morning and  $A_2$  the event that does not. Let  $B_1, B_2$  be the corresponding events for the afternoon. If it is more likely that it rains than it does not, both in the morning and the afternoon, then it is more likely that it rains sometime in the day than not. This is formalised by the following assumption.

**Assumption 3.1.2** (SP2). *Let  $A = A_1 \cup A_2$ ,  $B = B_1 \cup B_2$  with  $A_1 \cap A_2 = B_1 \cap B_2 = \emptyset$ . If  $A_i \precsim B_i$  for  $i = 1, 2$  then  $A \precsim B$ .*

We also require the simple technical assumption that any event  $A \in \mathcal{F}$  is at least as likely as the empty event  $\emptyset$ , which never happens.

**Assumption 3.1.3** (SP3). *If  $S$  is the certain event, and  $\emptyset = S^c$  the empty set, then:  $\emptyset \precsim A$  and  $\emptyset \prec S$ .*

As it turns out, these assumptions are sufficient for proving the following theorems DeGroot [1970]. The first theorem tells us that our belief must be consistent with respect to transitivity.

**Theorem 3.1.1** (Transitivity). *For all events  $A, B, D$ , if  $A \precsim B$  and  $B \precsim D$ , then  $A \precsim D$ .*

The second theorem says that if two events have a certain relation, then their negations have the converse relation.

**Theorem 3.1.2** (Complement). *For any  $A, B$ :  $A \precsim B$  iff  $A^c \succ B^c$ .*

Finally, note that if  $A \subset B$ , then it must be the case that whenever  $A$  happens,  $B$  must happen and hence  $B$  must be at least as likely as  $A$ . This is demonstrated in the following theorem.

**Theorem 3.1.3** (Fundamental property of relative likelihoods). *If  $A \subset B$  then  $A \precsim B$ . Furthermore,  $\emptyset \precsim A \precsim S$  for any event  $A$ .*

Since we are dealing with  $\sigma$ -fields, we need to introduce properties for infinite sequences of events. While these are not necessary if the field  $\mathcal{F}$  is finite, it is good to include them for generality.

**Assumption 3.1.4** (SP4). *If  $A_1 \supset A_2 \supset \dots$  is a decreasing sequence of events in  $\mathcal{F}$  and  $B \in \mathcal{F}$  is such that  $A_i \succsim B$  for all  $i$ , then  $\bigcap_{i=1}^{\infty} A_i \succsim B$ .*

As a consequence, we obtain the following dual theorem:

**Theorem 3.1.4.** *If  $A_1 \subset A_2 \subset \dots$  is an increasing sequence of events in  $\mathcal{F}$  and  $B \in \mathcal{F}$  is such that  $A_i \precsim B$  for all  $i$ , then  $\bigcup_{i=1}^{\infty} A_i \precsim B$ .*

We are now able to state a theorem for the unions of infinite sequences of disjoint events.

**Theorem 3.1.5.** *If  $(A_i)_{i=1}^{\infty}$  and  $(B_i)_{i=1}^{\infty}$  are infinite sequences of disjoint events in  $\mathcal{F}$  such that  $A_i \precsim B_i$  for all  $i$ , then  $\bigcup_{i=1}^{\infty} A_i \precsim \bigcup_{i=1}^{\infty} B_i$ .*

**Exercise 2.** *Here we prove that a probability measure  $P$  always satisfies the stipulated assumptions.*

(i) *For any events  $P(A) > P(B)$ ,  $P(A) < P(B)$  or  $P(A) = P(B)$ .*

(ii) *If  $A_i, B_i$  are partitions of  $A, B$ ,  $\forall i P(A_i) \leq P(B_i) \Rightarrow P(A) \leq P(B)$ .*

(iii) *For any  $A$ ,  $P(\emptyset) \leq P(A)$  and  $P(\emptyset) < P(\mathcal{S})$*

*Solution.* Part (i) is trivial, as  $P : \mathcal{F} \rightarrow [0, 1]$ . Part (ii) follows from  $P(A) = P(\bigcup_i A_i) = \sum_i P(A_i) \leq \sum_i P(B_i) = P(B)$ . Part (iii)  $P(\emptyset) = 0$ ,  $P(A) \geq 0$ . Also,  $P(\mathcal{S}) = 1$ .  $\square$

### 3.1.3 Assigning unique probabilities

In many cases, and particularly when  $\mathcal{F}$  is a finite field, there is a large number of probability distributions agreeing with our relative likelihoods.

How can we assign probabilities to events in an unambiguous manner?

**Example 3.1.1.** *Consider  $\mathcal{F} = \{\emptyset, A, A^c, \mathcal{S}\}$  and say  $A \succ A^c$ . Consequently,  $P(A) > 1/2$ . But this is insufficient for assigning a specific value to  $P(A)$ .*

Let  $A$  be an interval on the real line, with length  $\lambda(A)$ .

**Definition 3.1.2** (Uniform distribution).  *$x : \mathcal{S} \rightarrow [0, 1]$  has a uniform distribution on  $[0, 1]$  if, for any subintervals  $A, B$  of  $[0, 1]$ ,*

$$(x \in A) \precsim (x \in B) \quad \text{iff} \quad \lambda(A) \leq \lambda(B)$$

This means that *any* larger interval is more likely than *any* smaller interval. Now we shall connect the uniform distribution to the original sample space  $\mathcal{S}$  by assuming that there is some function with uniform distribution.

**Assumption 3.1.5** (SP5). *There exists a random variable  $x : \mathcal{S} \rightarrow [0, 1]$  with a uniform distribution in  $[0, 1]$ .*



### Constructing the probability distribution

We can now use the uniform distribution to create a unique probability measure that agrees with our likelihood relation. First, we have to map each event in  $\mathcal{S}$  to an equivalent event in  $[0, 1]$ .

**Theorem 3.1.6** (Equivalent event). *For any event  $A \in \mathcal{F}$ , there exists some  $\alpha \in [0, 1]$  such that  $A \approx (X \in [0, \alpha])$ .*

This means that we can now define the probability of an event  $A$  by matching it to a specific equivalent event on  $[0, 1]$ .

**Definition 3.1.3** (The probability of  $A$ ). *Given any event  $A$ , define  $P(A)$  to be the  $\alpha$  with  $A \approx (X \in [0, \alpha])$ .*

Hence

$$A \approx (X \in [0, P(A)]).$$

The above is sufficient to show the following theorem.

**Theorem 3.1.7** (Relative likelihood and probability). *If assumptions SP1-SP5 are satisfied, then the probability measure  $P$  defined above is unique. Furthermore, for any two events  $A, B$ ,  $A \preceq B$  iff  $P(A) \leq P(B)$ .*

### 3.1.4 Conditional likelihoods

#### Conditional likelihood

So far we have only considered the problem of forming opinions about which events are more likely *a priori*. However, we also need to have a way to incorporate evidence which may adjust our opinions. For example, while we ordinarily may think that  $A \preceq B$ , we may have additional information  $D$ , given which we think the opposite is true. We can formalise this through the notion of conditional likelihoods.

**Example 3.1.2.** *Say that  $A$  is the event that it rains in Gothenburg, Sweden tomorrow. We know that Gothenburg is quite rainy due to its oceanic climate, so we set  $A \succeq A^c$ . Now, let us try and incorporate some additional information. Let  $D$  denote the fact that good weather is forecast. I personally believe that  $(A | D) \preceq (A^c | D)$ , i.e. that good weather is more probable than rain, given the evidence of the weather forecast.*

#### Conditional likelihoods

Define  $(A | D) \preceq (B | D)$  to mean that  $B$  is at least as likely as  $A$  when it is known that  $D$  has occurred.

**Assumption 3.1.6** (CP). *For any events  $A, B, D$ ,*

$$(A | D) \preceq (B | D) \quad \text{iff} \quad A \cap D \preceq B \cap D.$$

**Theorem 3.1.8.** *If a relation  $\preceq$  satisfies assumptions SP1 to SP5 and CP, then  $P$  is the unique probability distribution such that:*

*For any  $A, B, D$  such that  $P(D) > 0$ ,*

$$(A | D) \preceq (B | D) \quad \text{iff} \quad P(A | D) \leq P(B | D)$$

**Definition 3.1.4** (Conditional probability).

$$P(A | D) = \frac{P(A \cap D)}{P(D)} \tag{3.1.1}$$

### 3.1.5 Probability elicitation

Probability elicitation is the problem of quantifying the subjective probabilities that a particular individual uses. One of the simplest, and most direct, methods, is to simply ask. However, because we cannot simply ask somebody to completely specify a probability distribution, we can ask for this distribution iteratively.

**Example 3.1.3** (Temperature prediction). *Let  $\tau$  be the temperature tomorrow at noon in Gothenburg. What are your estimates?*

*Eliciting the prior / forming the subjective probability measure  $P$*

- Select temperature  $x_0$  s.t.  $(\tau \leq x_0) \approx (\tau > x_0)$ .
- Select temperature  $x_1$  s.t.  $(\tau \leq x_1 \mid \tau \leq x_0) \approx (\tau > x_1 \mid \tau \leq x_0)$ .

Note that, necessarily,  $P(\tau \leq x_0) = P(\tau > x_0) = p_0$ . Since  $P(\tau \leq x_0) + P(\tau > x_0) = P(\tau \leq x_0 \cup \tau > x_0) = P(\tau \in \mathbb{R}) = 1$ , it follows that  $p_0 = 1/2$ . Similarly,  $P(\tau \leq x_1 \mid \tau \leq x_0) = P(\tau > x_1 \mid \tau \leq x_0) = 1/4$ .

#### Updating beliefs

Although we always start with a particular belief, this belief must be adjusted when we receive new evidence. In probabilistic inference, the updated beliefs are simply the probability of future events conditioned on observed events. This idea is captured neatly by Bayes' theorem, which links together the prior probability of events  $P(A_i)$  with their posterior probability  $P(A_i \mid B)$  given some event  $B$  and the probability  $P(B \mid A_i)$  of observing the evidence  $B$  given that hypothesis  $A_i$  is true.

**Theorem 3.1.9** (Bayes' theorem). *Let  $A_1, A_2, \dots$  be a (possibly infinite) sequence of disjoint events such that  $\bigcup_{i=1}^n A_i = \mathcal{S}$  and  $P(A_i) > 0$  for all  $i$ . Let  $B$  be another event with  $P(B) > 0$ . Then*

$$P(A_i \mid B) = \frac{P(B \mid A_i)P(A_i)}{\sum_{j=1}^n P(B \mid A_j)P(A_j)} \quad (3.1.2)$$

*Proof.* By definition,  $P(A_i \mid B) = P(A_i \cap B)/P(B)$ , and  $P(A_i \cap B) = P(B \mid A_i)P(A_i)$ , so:

$$P(A_i \mid B) = \frac{P(B \mid A_i)P(A_i)}{P(B)}, \quad (3.1.3)$$

As  $\bigcup_{i=1}^n A_i = \mathcal{S}$ , we have  $B = \bigcup_{j=1}^n (B \cap A_j)$ . Since  $A_i$  are disjoint, so are  $B \cap A_i$ . As  $P$  is a probability, the union property and an application of 3.1.3 gives

$$P(B) = P\left(\bigcup_{j=1}^n (B \cap A_j)\right) = \sum_{j=1}^n P(B \cap A_j) = \sum_{j=1}^n P(B \mid A_j)P(A_j).$$

□

**A simple exercise in updating beliefs****The area of Germany**

Form a subjective probability for the area  $a$  of Germany in  $\text{km}^2$ .

$$A_1 : a < 10^5 \text{ km}^2$$

$$A_2 : a \in [10^5, 2.5 \cdot 10^5) \text{ km}^2$$

$$A_3 : a \in [2.5 \cdot 10^5, 5 \cdot 10^5) \text{ km}^2$$

$$A_4 : a \in [5 \cdot 10^5, 10^6) \text{ km}^2$$

$$A_5 : a \geq 10^6 \text{ km}^2.$$

Choose  $P(A_i)$  for all  $i$ .

*Additional information*

- The EU's largest country is France ( $6.7 \cdot 10^5 \text{ km}^2$ ) and the smallest is Malta with  $316 \text{ km}^2$ .
- Germany is the 4th largest of the 27 EU states
- UK ( $2.4 \cdot 10^5 \text{ km}^2$ ) is the 8th largest EU state

The correct answer is  $A_3$ , since  $a = 3.57 \cdot 10^5$

**3.2 Utility theory**

While probability can be used to describe how likely an event is, utility can be used to describe how desirable it is. More concretely, our subjective probabilities are numerical representations of our beliefs and information. They can be taken to represent our “internal model” of the world. By analogy, our utilities are numerical representations of our tastes and preferences. Even if they are not directly known to us, we assume that our actions are such that we act so as to obtain maximum utility, in some sense.

**3.2.1 Rewards and preferences****Rewards**

Consider that we have to choose a *reward*  $r$  from a set  $R$  of possible rewards. While the elements of  $R$  may be arbitrary, we shall in general find that we prefer some rewards to others. In fact, some elements of  $R$  may not even be desirable. As an example,  $R$  might be a set of tickets to different musical events, or a set of financial commodities.

### Preferences

**Example 3.2.1** (Musical event tickets). *We have a set of tickets  $R$ , and we must choose the ticket  $r \in R$  we prefer best.*

- *Case 1:  $R$  are tickets to different music events at the same time, at equally good halls with equally good seats and the same price. Here preferences simply coincide with the preferences for a certain type of music or an artist.*
- *Case 2:  $R$  are tickets to different events at different times, at different quality halls with different quality seats and different prices. Here, preferences may depend on all the factors.*

**Example 3.2.2** (Route selection). *We have a set of alternate routes and must pick one.*

- *$R$  contains two routes, one short and one long, of the same quality.*
- *$R$  contains two routes, one short and one long, but the long route is more scenic.*

### Preferences among rewards

We will treat preferences in a similar manner as we have treated probabilities. That is, we will define a linear ordering among possible rewards.

Let  $a, b \in R$  be two rewards. When we *prefer*  $a$  to  $b$ , we write  $a \succ^* b$ . Conversely, when we like  $a$  *less* than  $b$  we write  $a \prec^* b$ . If we like  $a$  *as much* as  $b$ , we write  $a \approx^* b$ . We also use  $\succeq^*$  and  $\preceq^*$  for *I like at least as much as* and for *I don't like any more than*, respectively.

**Assumption 3.2.1.** *We make the following assumptions about the preference relations.*

- (i) *For any  $a, b \in R$ , one of the following holds:  $a \succ^* b$ ,  $a \prec^* b$ ,  $a \approx^* b$ .*
- (ii) *If  $a, b, c \in R$  are such that  $a \preceq^* b$  and  $b \preceq^* c$ , then  $a \preceq^* c$ .*

The first assumption means that we must always be able to decide between any two rewards. It may seem that it does not always hold in practice, since humans are frequently indecisive. This could be attributed to the difficulty of computing the value of the relation  $\preceq^*$ . Consider an algorithm that, taking  $r, r'$  as input returns  $R \in \{\succ^*, \prec^*, \approx^*\}$ . If the algorithm does not halt, then the preference relation is not defined.

The second assumption is a bit stronger. It is in fact possible to create a preference relation that will be cyclic.

**Example 3.2.3** (Counter example for transitive preferences). *Consider for example vector rewards in  $\mathcal{R} = \mathbb{R}^2$ . Given rewards  $r_i = (a_i, b_i)$ , our preference relation is:*

- $r_i \succ^* r_j$  if  $a_i \geq a_j$  and  $|b_i - b_j| < \epsilon$
- $r_i \succ^* r_j$  if  $b_i \geq b_j + \epsilon'$ .

This may correspond for example to an employer deciding to hire one of two employees,  $i, j$ , depending on their experience (a) or their school grades (b). Since grades are not very reliable, if two people have grades, then we prefer the one with the most experience. However, that may lead to a cycle. Consider a sequence of candidates  $i = 1, \dots, n$ , such that each candidate satisfies  $b_i = b_{i+1} + \delta$ , with  $\delta < \epsilon$  and  $a_i > a_{i+1}$ . Then clearly, we must always prefer  $r_i$  to  $r_{i+1}$ . However, if  $\delta n > \epsilon$ , we will prefer  $r_n$  to  $r_1$ .

### 3.2.2 Preferences among distributions

#### When we cannot select rewards directly

In most problems, we cannot choose the rewards directly. Rather, we must make some decision, and then obtain a reward depending on this decision. Since we may be uncertain about the outcome of a decision, we can specify our uncertainty regarding the rewards obtained by a decision in terms of a probability distribution.

**Example 3.2.4** (Route selection). • Each reward  $r \in R$  is the time it takes to travel from  $A$  to  $B$ .

- We prefer shorter times.
- There are two routes,  $P_1, P_2$ .
- Route  $P_1$  takes 10 minutes when the road is clear, but 30 minutes when the traffic is heavy. The probability of heavy traffic on  $P_1$  is  $q_1$ .
- Route  $P_2$  takes 15 minutes when the road is clear, but 25 minutes when the traffic is heavy. The probability of heavy traffic on  $P_2$  is  $q_2$ .

#### Preferences among probability distributions

Consequently, we have to define preferences between probability distributions, rather than rewards. We use the same notation as before. Let  $P_1, P_2$  be two distributions on  $(R, \mathcal{F}_R)$ . If we prefer  $P_1$  to  $P_2$ , we write  $P_1 \succ^* P_2$ . If we like  $P_1$  less than  $P_2$ , write  $P_1 \prec^* P_2$ . If we like  $P_1$  as much as  $P_2$ , we write  $P_1 \approx^* P_2$ . Finally, we also use  $\succsim^*$  and  $\precsim^*$  to denote strict preference relations.

### 3.2.3 Utility

The concept of utility allows us to create a unifying framework, such that given a particular set of rewards and probability distributions on them, we can define preferences among distributions automatically. The first step is to define utility as a way to define a preference relation among rewards.

**Definition 3.2.1** (Utility). The utility is a function  $U : R \rightarrow \mathbb{R}$ , such that for all  $a, b \in R$

$$a \succsim^* b \quad \text{iff} \quad U(a) \geq U(b), \quad (3.2.1)$$

The above definition is very similar to how we defined relative likelihood in terms of probability. For a given utility function, its expectation for a distribution over rewards as:

**Definition 3.2.2** (Expected utility). *The expected utility of a distribution  $P$  on  $R$  is:*

$$\mathbb{E}_P(U) = \int_R U(r) dP(r) \quad (3.2.2)$$

Finally, we make the assumption that the utility function is such that the expected utility remains consistent with the preference relations between all probability distributions we are choosing between.

**Assumption 3.2.2.** *The expected utility hypothesis The utility of  $P$  is equal to the expected utility of the reward under  $P$ . Consequently,*

$$P \succsim^* Q \quad \text{iff} \quad \mathbb{E}_P(U) \geq \mathbb{E}_Q(U). \quad (3.2.3)$$

**Example 3.2.5.** *Consider the following decision problem. You have the option of entering a lottery, for 1 CU, that gives you a prize of 10 CU. The probability of winning is 0.01. This can be formalised by making it a choice between  $P$ , where you do not enter the lottery and  $Q$ , which represents entering the lottery. Now we can calculate the expected utility for each choice. This is simply  $\mathbb{E}(U \mid$*

r	U(r)	P	Q
did not enter	0	1	0
paid 1 CU and lost	-1	0	0.99
paid 1 CU and won 10	9	0	0.01

Table 3.1: A simple gambling problem

$P) = \sum_r U(r)P(r)$  and  $\mathbb{E}(U \mid Q) = \sum_r U(r)Q(r)$  respectively.

	P	Q
$\mathbb{E}(U \mid \cdot)$	0	-0.9

Table 3.2: Expected utility for the gambling problem

### Monetary rewards

**Example 3.2.6.** *Choose between the following two gambles:*

1. *The reward is 500,000 with certainty.*
2. *The reward is 2,500,000 with probability 0.10. It is 500,000 with probability 0.89, and 0 with probability 0.01.*

**Example 3.2.7.** *Choose between the following two gambles:*

1. *The reward is 500,000 with probability 0.11, or 0 with probability 0.89.*
2. *The reward is: 2,500,000 with probability 0.1, or 0 with probability 0.9.*

**Exercise 3.** *Show that if gamble 1 is preferred in the first example, gamble 1 must also be preferred in the second example.*

**The St. Petersburg Paradox****A simple game [Bernoulli, 1713]**

- A *fair coin* is tossed until a head is obtained.
- If the first head is obtained on the  $n$ -th toss, our reward will be  $2^n$  currency units.

How much are you willing to pay, to play this game once?

- The probability to stop at round  $n$  is  $2^{-n}$ .
- Thus, the expected monetary gain of the game is

$$\sum_{n=1}^{\infty} 2^n 2^{-n} = \infty.$$

- If your utility function were linear you'd be willing to pay any amount to play.

**3.2.4 Measuring utility****Experimental measurement of utility**

**Example 3.2.8.** *We shall try and measure the utility of all monetary rewards in some interval  $[a, b]$ .*

*Let  $\langle a, b \rangle$  denote a lottery ticket that yields  $a$  or  $b$  CU with equal probability. Consider the following sequence:*

1. *Find  $x_1$  such that receiving  $x_1$  CU with certainty is equivalent to receiving  $\langle a, b \rangle$ .*
2. *Find  $x_2$  such that receiving  $x_2$  CU with certainty is equivalent to receiving  $\langle a, x_1 \rangle$ .*
3. *Find  $x_3$  such that receiving  $x_3$  CU with certainty is equivalent to receiving  $\langle x_1, b \rangle$ .*
4. *Find  $x_4$  such that receiving  $x_4$  CU with certainty is equivalent to receiving  $\langle x_2, x_3 \rangle$ .*

If  $x_1 \neq x_4$ , then your preferences do not meet the requirements of a utility function.

**Exercise 4.** 1. *Specify an amount  $a$ , then observe random value  $Y$ .*

2. *If  $Y \geq a$ , receive  $Y$ .*
3. *If  $Y < a$ , receive random reward  $X$  with known distribution (independent of  $Y$ ).*
4. *Show that you should choose  $a$  s.t.  $U(a) = \mathbb{E}[U(X)]$ .*

### 3.2.5 Convex and concave utility functions

#### Convex functions

**Definition 3.2.3.** A function  $g$  is convex on  $A$  if, for any points  $x, y \in A$ , and any  $\alpha \in [0, 1]$ :

$$\alpha g(x) + (1 - \alpha)g(y) \geq g[\alpha x + (1 - \alpha)y]$$

**Theorem 3.2.1** (Jensen's inequality). If  $g$  is convex on  $\mathcal{S}$  and  $x \in \mathcal{S}$  with measure  $P(A) = 1$  and  $\mathbb{E}(x)$  and  $\mathbb{E}[g(x)]$  exist, then:

$$\mathbb{E}[g(x)] \geq g[\mathbb{E}(x)]. \quad (3.2.4)$$

**Example 3.2.9.** If the utility function is convex, then we choose a gamble giving a random gain  $x$  rather than one giving a fixed gain  $\mathbb{E}(x)$ . Thus, a convex utility function implies risk-taking.

#### Concave functions

**Definition 3.2.4.** A function  $g$  is concave on  $\mathcal{S}$  if, for any points  $x, y \in \mathcal{S}$ , and any  $\alpha \in [0, 1]$ :

$$\alpha g(x) + (1 - \alpha)g(y) \leq g[\alpha x + (1 - \alpha)y]$$

For concave functions, an analogue of Jensen's inequality holds (in the other direction). If the utility function is concave, then we choose a gamble giving a fixed gain  $\mathbb{E}[X]$  rather than one giving a random gain  $X$ . Consequently, a concave utility function implies risk aversion.

**Example 3.2.10** (Insurance). The act of buying insurance can be related to concavity of our utility function.

Let  $x$  be the insurance cost,  $h$  our insurance cover and  $\epsilon$  the probability of needing the cover. Then we are going to buy insurance if the utility of losing  $x$  with certainty is greater than the utility of losing  $-h$  with probability  $\epsilon$ .

$$U(-x) > \epsilon U(-h) + (1 - \epsilon)U(0). \quad (3.2.5)$$

The company has a linear utility, and fixes the premium  $x$  high enough for

$$x > \epsilon h. \quad (3.2.6)$$

Consequently, we see that from (3.2.6) that  $U(-\epsilon h) \geq U(-x)$ , as  $U$  is an increasing function. Thus, from (3.2.5) we obtain  $U(-\epsilon h) > \epsilon U(-h) + (1 - \epsilon)U(0)$ . Now the  $-\epsilon h$  term is the utility of our expected monetary loss, while the right hand side is our expected utility. Consequently if the inequality holds, our utility function is (at least locally) concave.

## 3.3 Summary

### Summary

- We can subjectively indicate which events we think are more likely.
- Using relative likelihoods, we can define a *subjective probability*  $P$  for all events.



- Similarly, we can subjectively indicate *preferences for rewards*.
- We can determine a *utility function* for all rewards.
- Hypothesis: we prefer the probability distribution (over rewards) with the highest *expected utility*.
- Concave utility functions imply *risk aversion* (and convex, risk-taking).



## Chapter 4

# Decision problems

## 4.1 Introduction

In this chapter we describe how to actually formulate statistical decision problems. The simplest such problem arises when we have a choice between a number of different decisions, where each decision gives us different *rewards* with different probabilities. If these probabilities are known, then the framework of expected utility maximisation gives a solution to the problem. An example includes *gambling*, where we must choose between a number of possible lotteries, each one having different payoffs and winning probabilities.

Another classical setting is parameter estimation. Therein, we stipulate the existence of a parameterised *law of nature*, and we wish to choose a best-guess set of parameters for the law through measurements and some prior information, such as for example determining the gravitational attraction constant from observing planetary movements. These measurements are always obtained through experiments, and the automatic design of those experiment is a topic we shall consider in later chapters.

Finally, these decisions will necessarily depend on our prior information, even if have access to some additional measurements. The last section of this chapter will examine how sensitive our decisions are to the prior, and how we can choose a prior distribution so that our decisions are robust.

## 4.2 Rewards that depend on the outcome of an experiment

Consider the problem of choosing between a two different types of tickets in raffle. Each type of ticket gives you the chance to win a different prize. The first is a bicycle and the second is a tea set. As most people opt for the bicycle, your chance of actually winning it is much smaller. However, if you prefer winning a bicycle to winning the tea set, it is not clear what choice you should make in the raffle. The above is the quintessential example for problems where the reward that we obtain depends not only on our decisions, but also in the outcome of an *experiment*.

This problem can be viewed more generally, for scenarios where the reward you receive depends not only on your own choice, but also some other, unknown fact in the world. This may be something completely uncontrollable, and hence you only have to make the best possible guess.

More formally, we must make a decision  $d \in \mathcal{D}$  *before* knowing the outcome  $\omega$  of an experiment with outcomes in  $\Omega$ . After the experiment is performed, we obtain a *reward*  $r \in \mathcal{R}$  which depends on both the outcome of the experiment  $\omega$  and our decision. As discussed in the previous chapter, our preferences for some rewards over others is determined by a *utility* function  $U : \mathcal{R} \rightarrow \mathbb{R}$ , such that we prefer  $r$  to  $r'$  if and only if  $U(r) \geq U(r')$ . Now, however, we cannot choose rewards directly. Another example, which will be used throughout this section, is the following.

**Example 4.2.1** (Taking the umbrella). *We must decide whether to take an umbrella to work. Our reward is a combination of whether we get wet and the amount of objects that we carry. We would rather not get wet and not carry too many things, which can be made more precise by choosing an appropriate utility*

function. For example, we might put a value of  $-1$  for carrying the umbrella, a value of  $-10$  for getting wet. In this example, the only events of interest are whether it rains or not.

### 4.2.1 Formalisation of the problem setting

We are now ready to formulate the problem setting more precisely.

**Assumption 4.2.1** (Outcomes). *There exists a probability measure  $P$  on  $(\Omega, \mathfrak{F}_\Omega)$  such that the probability of the random outcome  $\omega$  being in  $A \subset \Omega$  is:*

$$\mathbb{P}(\omega \in A) = P(A), \quad \forall A \in \mathfrak{F}_\Omega. \quad (4.2.1)$$

This probability is completely independent of any decision that we make.

**Assumption 4.2.2** (Utilities). *Our rewards satisfy all our assumptions from Chapter 3: Preferences are transitive, all rewards are comparable and there exists a utility function  $U$ , measurable with respect to  $\mathfrak{F}_\mathcal{R}$  such that  $U(r') \geq U(r)$  iff  $r \succ^* r'$ .*

Since the random outcome  $\omega$  does not depend on our decision  $d$ , we must find a way to connect the two. This can be formalised via a reward function, so that the reward that we obtain (whether we get wet or not) depends on both our decision (to take the umbrella) and the random outcome (whether it rains).

**Definition 4.2.1** (Reward function). *A reward function  $\rho : \Omega \times \mathcal{D} \rightarrow \mathcal{R}$  defines the reward we obtain if we select  $d \in \mathcal{D}$  and the experimental outcome is  $\omega \in \Omega$ :*

$$r = \rho(\omega, d). \quad (4.2.2)$$

The decision space might be arbitrarily more complex than the one we have seen so far. For example, our decisions may be distributions over simple decisions, or functions whose value depends on future events. We shall examine those problems later in the chapter.

When we discussed the problem of choosing between distributions, in section 3.2.2, we had directly defined probability distributions on rewards. We can now formulate our problem in that setting. First, we define a set of distributions  $\{P_d \mid d \in \mathcal{D}\}$  on the reward space  $(\mathcal{R}, \mathfrak{F}_\mathcal{R})$ , such that the decision  $d$  amounts to choosing a particular distribution  $P_d$  on the rewards.

**Example 4.2.2** (Rock/Paper/Scissors). *Consider a simple game of rock-paper-scissors, where your opponent plays a move at the same time as you, so that you cannot influence his move. The opponents moves are thus  $\Omega = \{\omega_R, \omega_P, \omega_S\}$ .*

*You have studied your opponent for some time and you believe that he is most likely to play rock  $P(\omega_R) = 3/6$ , somewhat likely to play paper  $P(\omega_P) = 2/6$ , and less likely to play scissors:  $P(\omega_S) = 1/6$ . Your decision set is your own moves:  $\mathcal{D} = \{d_R, d_P, d_S\}$ , for rock, paper, scissors, respectively. The reward set is  $\mathcal{R} = \{\text{Win}, \text{Draw}, \text{Lose}\}$ .*

*What is the probability of each reward, for each decision you make? Taking the example of  $d_R$ , we see that you win if the opponent plays scissors with probability  $1/6$ , you lose if the opponent plays paper ( $2/6$ ), and you draw if he plays rock ( $3/6$ ). Consequently, we can convert the outcome probabilities to reward probabilities for every decision:*

$$\begin{array}{lll}
P_{d_R}(\text{Win}) = 1/6, & P_{d_R}(\text{Draw}) = 3/6, & P_{d_R}(\text{Lose}) = 2/6 \\
P_{d_P}(\text{Win}) = 3/6, & P_{d_P}(\text{Draw}) = 2/6, & P_{d_P}(\text{Lose}) = 1/6 \\
P_{d_S}(\text{Win}) = 2/6, & P_{d_S}(\text{Draw}) = 1/6, & P_{d_S}(\text{Lose}) = 3/6.
\end{array}$$

Of course, what we play depends on our own utility function. If we prefer winning over drawing or losing, we could for example have the utility function  $U(\text{Win}) = 1, U(\text{Draw}) = 0, U(\text{Lose}) = -1$ . Then, since  $\mathbb{E}_d U = \sum_{\omega \in \Omega} U(\omega, d) P_d(\omega)$ , we have

$$\begin{aligned}
E_{d_R} U &= -1/6 \\
E_{d_P} U &= 2/6 \\
E_{d_S} U &= -1/6
\end{aligned}$$

More generally, every decision that we make creates a corresponding probability distribution on rewards.

#### The probability measure induced by decisions

For every  $d \in \mathcal{D}$ , the function  $\rho : \Omega \times \mathcal{D} \rightarrow \mathcal{R}$  induces a probability distribution  $P_d$  on  $\mathcal{R}$ . In fact, for any  $B \in \mathfrak{F}_{\mathcal{R}}$ :

$$P_d(B) \triangleq \mathbb{P}(\rho(\omega, d) \in B) = P(\{\omega \mid \rho(\omega, d) \in B\}). \quad (4.2.3)$$

The above equation requires that the following technical assumption is satisfied:

**Assumption 4.2.3.** *The sets  $\{\omega \mid \rho(\omega, d) \in B\}$  must belong to  $\mathfrak{F}_{\Omega}$ . That is,  $\rho$  must be  $\mathfrak{F}_{\Omega}$ -measurable for any  $d$ .*

In other words, while the outcome of the experiment is independent of the decision, the distribution of rewards is effectively chosen by our decision, as before. However, this structure allows us to clearly distinguish the controllable from the random part of the rewards.

In either case, we employ the expected utility hypothesis (Assumption 3.2.2). Thus, we should choose the decision that results in the highest expected utility.

#### Expected utility

The expected utility of any decision  $d \in \mathcal{D}$  under  $P$  is: the expected utility is:

$$\mathbb{E}_{P_d}(U) = \int_{\mathcal{R}} U(r) dP_d(r) = \int_{\Omega} U[\rho(\omega, d)] dP(\omega) \quad (4.2.4)$$

From now on, we shall use the simple notation

$$U(P, d) \triangleq \mathbb{E}_{P_d} U \quad (4.2.5)$$

to denote the expected utility of  $d$  under distribution  $P$ .



Figure 4.1: Decision diagrams for the combined and separated formulation of the decision problem. Squares denote decision variables, diamonds denote utilities. All other variables are denoted by circles. Arrows denote the flow of dependency.

$\rho(\omega, d)$	$d_1$	$d_2$
$\omega_1$	dry, carrying umbrella	wet
$\omega_2$	dry, carrying umbrella	dry
$U[\rho(\omega, d)]$	$d_1$	$d_2$
$\omega_1$	0	-10
$\omega_2$	0	1
$\mathbb{E}_P(U \mid d)$	0	-1.2

Table 4.1: Rewards, utilities, expected utility for 20% probability of rain.

Instead of viewing the decision as effectively choosing a distribution over rewards (Fig. 4.1(a)) we can separate the random part of the process from the deterministic part (Fig. 4.1(b)) by considering a measure  $P$  on some space of outcomes  $\Omega$ , such that the reward depends on both  $d$  and the outcome  $\omega \in \Omega$  through a function  $\rho(\omega, d)$ . The optimal decision is of course always the  $d \in \mathcal{D}$  maximising  $\mathbb{E}(U \mid P_d)$ .

The dependency structure of this problem in either formulation can be visualised in the *decision diagram* shown in Figure 4.1(a).

**Example 4.2.3.** *You are going to work, and it might rain. The forecast said that the probability of rain ( $\omega_1$ ) was 20%. What do you do?*

- $d_1$ : Take the umbrella.
- $d_2$ : Risk it!

The reward of a given outcome and decision combination, as well as the expected utility is given in table 4.1.

### 4.2.2 Decision diagrams

Decision diagrams are also known as *decision networks* or *influence diagrams*. Like the examples shown in Figure 4.1, they are used to show dependencies between different variables. In general, these include the following types of nodes:

- Choice nodes, denoted by squares. These are nodes whose values the decision maker can directly choose. Sometimes there is more than one decision maker involved.
- Value nodes, denoted by diamonds. These are the nodes that the decision maker is interested in influencing. The utility of the decision maker is always a function of the value nodes.
- Circle nodes are used to denote all other types of variables. These include deterministic, stochastic, known or unknown variables.

The nodes are connected via directed edges. These denote the dependencies between nodes. For example, in Figure 4.1(b), the reward is a function of both  $\omega$  and  $d$ , i.e.  $r = \rho(\omega, d)$ , while  $\omega$  depends only on the probability distribution  $P$ . Typically, there must be a path from a choice node to a value node, otherwise nothing the decision maker can do will influence its utility. Nodes belonging to or observed by different players will usually be denoted by different lines or colors. In Figure 4.1(b),  $\omega$ , which is not observed, is shown in a lighter color.

### 4.2.3 Statistical estimation\*

This is especially the case in statistical problem of *parameter estimation*, such as estimating the covariance matrix of a Gaussian random variable. A simple example is estimating the distribution of votes in an election from a small sample, given below.

**Example 4.2.4 (Voting).** *Let us say for example that you wish to estimate the number of votes for different candidates in an election. The unknown parameters of the problem mainly include: the percentage of likely voters in the population, the probability that a likely voter is going to vote for each candidate. One simple way to estimate this is by polling.*

*Consider a nation with  $k$  political parties. Let  $\omega = (\omega_1, \dots, \omega_k) \in [0, 1]^k$  be the voting percentages for each party. We wish to make a guess  $d \in [0, 1]^k$ . How should we guess, given a distribution  $P(\omega)$ ? How should we select  $U$  and  $\rho$ ? This depends on what our goals is, when we make the guess.*

*If we wish to give a reasonable estimate about all the  $k$  parties votes, we can use the squared error: First, set  $\rho(\omega, d) = (\omega_1 - d_1, \dots, \omega_k - d_k)$ , our error vector  $r \in [0, 1]^k$ . Then we set  $U(r) \triangleq -\|r\|^2$ , where  $\|r\|^2 = \sum_i |x_i|^2$ .*

*If on the other hand, we just want to predict the winner of the election, then the actual percentages of all individual parties are not important. In that case, we can set  $\rho(\omega, d) = 1$  if  $\arg \max_i \omega_i = \arg \max_i d_i$  and 0 otherwise, and  $U(r) = r$ .*

- The unknown outcome of the experiment  $\omega$  is called a *parameter*.
- The set of outcomes  $\Omega$  is called the *parameter space*.

#### Losses and risks

In such problems, it is common to specify a loss instead of a utility. This is usually the negative utility:



**Definition 4.2.2** (Loss).

$$\ell(\omega, d) = -U[\rho(\omega, d)]. \quad (4.2.6)$$

Given the above, instead of the expected utility, we consider the expected loss, or risk.

**Definition 4.2.3** (Risk).

$$\sigma(P, d) = \int_{\Omega} \ell(\omega, d) dP(\omega). \quad (4.2.7)$$

Of course, the optimal decision is  $d$  minimising  $\sigma$ .

### 4.3 Bayes decisions

The decision which maximises the expected utility under a particular distribution  $P$ , is called the *Bayes-optimal* decision, or simply the *Bayes decision*. The probability distribution  $P$  is supposed to reflect all our uncertainty about the problem.

**Definition 4.3.1** (Bayes-optimal utility). *Consider an outcome (or parameter) space  $\Omega$ , decision space  $\mathcal{D}$ , and a utility function  $U : \Omega \times \mathcal{D} \rightarrow \mathbb{R}$ . For any probability distribution  $P$  on  $\Omega$ , the Bayes-optimal utility  $U^*(P)$  is defined as the smallest upper bound on  $U(P, d)$  for all decisions  $d \in \mathcal{D}$ . That is,*

$$U^*(P) = \sup_{d \in \mathcal{D}} U(P, d). \quad (4.3.1)$$

**Remark 4.3.1** (The sup notation). *A clarification of the supremum notation is in order. When we say that*

$$M = \sup_{x \in A} f(x),$$

*then: (i)  $M \geq f(x)$  for any  $x \in A$ . In other words,  $M$  is an upper bound on  $f(x)$ . (ii) for any  $M' > M$ , there exists some  $x' \in A$  s.t.  $M' < f(x')$ . In other words, there exists no smaller upper bound than  $M$ . When the function  $f$  has a maximum, then the supremum is identical to the maximum.*

As can be seen from Figure ??, for absolute loss, the optimal decision is to choose the  $d$  that is closest to the most likely  $\omega$ . However, for quadratic loss, Figure ?? appears to indicate that the optimal choice should be equal to the expected value of  $\omega$ . This is actually true in general for quadratic loss, and for  $d, \omega \in \mathbb{R}$ , as shall be seen from the following example.

**Example 4.3.1** (Quadratic loss). *Now consider  $\Omega = \mathbb{R}$  with measure  $P$  and  $\mathcal{D} = \mathbb{R}$ . For any point  $\omega \in \mathbb{R}$ , we define the utility as:*

$$U(\omega, d) = -|\omega - d|^2. \quad (4.3.2)$$

The optimal decision maximises

$$U(P, d) = - \int_{\mathbb{R}} |\omega - d|^2 dP(\omega).$$

Then, as long as  $\partial/\partial d |\omega - d|^2$  is measurable with respect to  $\mathfrak{F}_{\mathbb{R}}$

$$\frac{\partial}{\partial d} \int_{\mathbb{R}} |\omega - d|^2 dP(\omega) = \int_{\mathbb{R}} \frac{\partial}{\partial d} |\omega - d|^2 dP(\omega) \quad (4.3.3)$$

$$= 2 \int_{\mathbb{R}} (d - \omega) dP(\omega) \quad (4.3.4)$$

$$= 2 \int_{\mathbb{R}} d dP(\omega) - 2 \int_{\mathbb{R}} \omega dP(\omega) \quad (4.3.5)$$

$$= 2d - 2\mathbb{E}(\omega), \quad (4.3.6)$$

so the expected utility is maximised for  $d = \mathbb{E}(\omega)$ .

### 4.3.1 Convexity of the Bayes-optimal utility\*

We shall show now the expected utility is linear. Consequently, the Bayes-utility is convex with respect to the distribution  $P$ . This firstly implies that there is a unique “worst” distribution  $P$ , against which we cannot do very well. Secondly, we can approximate the Bayes-utility very well for all possible distributions by generalising from a small number of distributions. In order to define linearity and convexity, we must first introduce the concept of a mixture of distributions.

#### A mixture of distributions

Consider two probability measures  $P, Q$  on  $(\Omega, \mathfrak{F}_{\Omega})$ .

These define two alternative distributions for  $\omega$ . For any  $P, Q$  and  $\alpha \in [0, 1]$ , we define

$$Z_{\alpha} \triangleq \alpha P + (1 - \alpha)Q \quad (4.3.7)$$

to mean the probability measure such that  $Z_{\alpha}(A) = \alpha P(A) + (1 - \alpha)Q(A)$  for any  $A \in \mathfrak{F}_{\Omega}$ .

**Remark 4.3.2** (Linearity of the expected utility). *If  $Z_{\alpha}$  is as defined in (4.3.7), then, for any  $d \in \mathcal{D}$ :*

$$U(Z_{\alpha}, d) = \alpha U(P, d) + (1 - \alpha)U(Q, d). \quad (4.3.8)$$

*Proof.* This follows from the linearity of expectation, i.e.

$$U(Z_{\alpha}, d) = \int_{\Omega} U(\omega, d) dZ_{\alpha}(\omega) \quad (4.3.9)$$

$$= \alpha \int_{\Omega} U(\omega, d) dP(\omega) + (1 - \alpha) \int_{\Omega} U(\omega, d) dQ(\omega) \quad (4.3.10)$$

$$= \alpha U(P, d) + (1 - \alpha)U(Q, d). \quad (4.3.11)$$

□

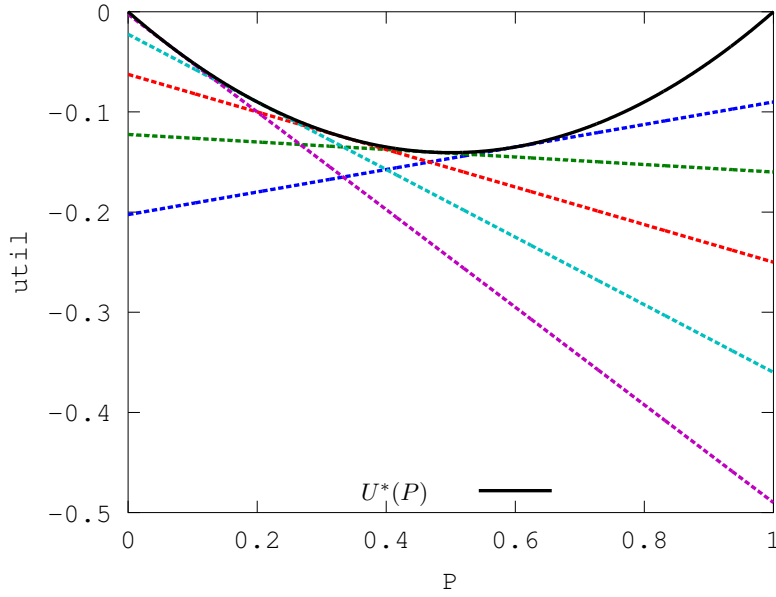


Figure 4.2: A strictly convex Bayes utility.

**Theorem 4.3.1.** For probability measures  $P, Q$  on  $\Omega$  and any  $\alpha \in [0, 1]$ ,

$$U^*[Z_\alpha] \leq \alpha U^*(P) + (1 - \alpha)U^*(Q), \quad (4.3.12)$$

where  $Z_\alpha = \alpha P + (1 - \alpha)Q$ .

*Proof.* From the definition of the expected utility (4.3.8), for any decision  $d \in \mathcal{D}$ ,

$$U(Z_\alpha, d) = \alpha U(P, d) + (1 - \alpha)U(Q, d).$$

Hence, by definition (4.3.1) of the Bayes-utility:

$$\begin{aligned} U^*(Z_\alpha) &= \sup_{d \in \mathcal{D}} U(Z_\alpha, d) \\ &= \sup_{d \in \mathcal{D}} [\alpha U(P, d) + (1 - \alpha)U(Q, d)]. \end{aligned}$$

Use  $\sup_x [f(x) + g(x)] \leq \sup_x f(x) + \sup_x g(x)$  to bound r.h.s:

$$\begin{aligned} U^*[Z_\alpha] &\leq \alpha \sup_{d \in \mathcal{D}} U(P, d) + (1 - \alpha) \sup_{d \in \mathcal{D}} U(Q, d) \\ &= \alpha U^*(P) + (1 - \alpha)U^*(Q). \end{aligned}$$

□

### Convexity of the Bayes utility

As we have proven, the expected utility is linear with respect to  $P$ . Thus, for any fixed decision  $d$  we obtain one of the lines in Fig. 4.2. Due to the theorem just proved, the Bayes risk is concave. Furthermore, the minimising decision

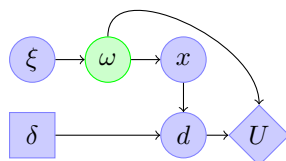


Figure 4.3: Statistical decision problem with observations

for any  $P$  is tangent to the risk at the point  $(P, U^*(P))$ . If we take a decision that is optimal with respect to some  $P$ , but the distribution is in fact  $Q \neq P$ , then we are not far from the optimal when  $P$  and  $Q$  are close and  $U^*$  is smooth. Consequently, we can trivially lower bound the Bayes utility by examining a finite set of decisions  $\hat{D}$ :

$$U^*(P) \geq \max_{d \in \hat{D}} U(P, d) \forall P$$

In addition, we can upper-bound the Bayes utility as follows. Take any two distributions  $P_1, P_2$  in the set of allowed distributions. Then, the following upper bound holds

$$U^*(\alpha P_1 + (1 - \alpha)P_2) \leq \alpha U^*(P_1) + (1 - \alpha)U^*(P_2)$$

due to convexity. The two bounds suggest an algorithm for successive approximation of the Bayes risk, by looking for the largest gap between the lower and the upper bounds.

## 4.4 Decision problems with observations

So far we have only examined problems where the outcomes were drawn from some fixed distribution. This distribution constituted our subjective belief about what the unknown parameter is. Now, we examine the case where we can obtain some observations that depend on the unknown  $\omega$  before we make our decision. These observations should give us more information about  $\omega$ , before making a decision. Intuitively, we should be able to make decisions by simply considering the posterior distribution. The following section will investigate whether this is true.

### Obtaining information

In this setting, we once more need to take some decision  $d \in \mathcal{D}$  so as to maximise expected utility. As before, we have a prior distribution  $\xi$  on some parameter  $\omega \in \Omega$ , representing what we know about  $\omega$ . Consequently, the expected utility of any fixed decision  $d$  is going to be  $\mathbb{E}_\xi(U \mid d)$ .

However, it might be possible to obtain more information about  $\omega$  before making a decision. In particular, each  $\omega$  corresponds to a *model* of the world  $\psi_\omega$ . This is expressed as a probability distribution over the observation space  $\mathcal{S}$ , such that  $\psi_\omega(X)$  is the probability that the observation is in  $X \subset \mathcal{S}$ . The set of parameters  $\Omega$  thus defines a family of models:

$$\Psi \triangleq \{\psi_\omega \mid \omega \in \Omega\}. \quad (4.4.1)$$

Now, consider the case where we take an observation  $x$  from the true model  $\psi_{\omega^*}$  before having to make a decision. We can represent the dependency of our decisions on the observations by making our decision a function of  $x$ :

**Definition 4.4.1** (Decision function). *A decision function  $\delta : \mathcal{S} \rightarrow \mathcal{D}$  maps from the set of possible observations  $\mathcal{S}$  to the set of possible decisions.*

The expected utility of a decision function  $\delta$  is:

$$U(\xi, \delta) \triangleq \mathbb{E}_{\xi} \{U[\omega, \delta(x)]\} = \int_{\Omega} \left( \int_{\mathcal{S}} U[\omega, \delta(x)] d\psi_{\omega}(x) \right) d\xi(\omega). \quad (4.4.2)$$

When the set of decision functions includes all fixed decisions, then there is a decision function  $\delta^*$  at least as good as the best fixed decision  $d^*$ . More formally:

**Remark 4.4.1.** *Let  $\mathcal{D}$  denote the set of decision functions  $\delta : \mathcal{S} \rightarrow \mathcal{D}$ . If,  $\forall d \in \mathcal{D} \exists \delta \in \mathcal{D}$  such that  $\delta(x) = d \forall x \in \mathcal{S}$ , then  $\sup_{\delta \in \mathcal{D}} \mathbb{E}_{\xi}(U | \delta) \geq \sup_{d \in \mathcal{D}} \mathbb{E}_{\xi}(U | d)$ .*

*Proof.* The proof follows by setting  $\mathcal{D}_0$  to be the set of fixed decision functions. The result follows since  $\mathcal{D}_0 \subset \mathcal{D}$ .  $\square$

This is the standard Bayesian framework for decision making. It may be slightly more intuitive in some case to use the notation  $\psi(x | \omega)$ , in order to emphasize that this is a conditional distribution. However, there is no technical difference between the two notations.

**Example 4.4.1.** *Consider the problem of deciding whether or not to go to a particular restaurant. Let  $\Omega = [0, 1]$  with  $\omega = 0$  meaning the food is in general horrible and  $\omega = 1$  meaning the restaurant is great. Let  $x_1, \dots, x_n$  be  $n$  expert opinions in  $\mathcal{S} = \{0, 1\}$  about the restaurant. Under our model, the probability of observing  $x_i = 1$  when the quality of the restaurant is  $\omega$  is given by  $\psi_{\omega}(1) = \omega$  and conversely  $\psi_{\omega}(0) = 1 - \omega$ . The probability of observing a particular<sup>1</sup> sequence  $x$  of length  $n$  is*

$$\psi_{\omega}(x) = \omega^s (1 - \omega)^{n-s}$$

with  $s = \sum_{i=1}^n x_i$ .

### Maximising utility when making observations

Statistical procedures based on the notion that a distribution can be assigned to any parameter in a statistical decision problem, as we are assuming here, are called *Bayesian statistical methods*. The scope of these methods has been the subject of much discussion in the statistical literature. See e.g. Savage [1972].

In the following, we shall look at different expressions for the expected utility. We shall overload the utility operator  $U$  for various cases: when the parameter is fixed, when the parameter is random, when the decision is fixed, and when the decision depends on the observation  $x$  and thus is random as well.

---

<sup>1</sup>We obtain a slightly different probability under the binomial model, but the end result is the same.

**Expected utility of a fixed decision  $d$  with  $\omega \sim \xi$** 

We first consider the expected utility of taking a fixed decision  $d \in \mathcal{D}$ , when  $\mathbb{P}(\omega \in A) = \xi(A)$ . This is the case we have dealt with so far.

$$U(\xi, d) \triangleq \mathbb{E}_\xi(U \mid d) = \int_{\Omega} U(\omega, d) d\xi(\omega). \quad (4.4.3)$$

**Expected utility of a decision function  $\delta$  with fixed  $\omega \in \Omega$** 

Now assume that  $\omega$  is fixed, but instead of selecting a decision directly, we select a decision that depends on the random observation  $x$ , which is distributed according to  $\psi_\omega$  on  $\mathcal{S}$ . We do this by defining a function  $\delta : \mathcal{S} \rightarrow \mathcal{D}$ .

$$U(\omega, \delta) = \int_{\mathcal{S}} U(\omega, \delta(x)) d\psi_\omega(x). \quad (4.4.4)$$

**Expected utility of a decision function  $\delta$  with  $\omega \sim \xi$** 

Now we generalise to the case where  $\omega$  is distributed with measure  $\xi$ . Note that the expectation of the previous expression (4.4.4) is by definition written as:

$$U(\xi, \delta) = \int_{\Omega} U(\omega, \delta) d\xi(\omega), \quad U^*(\xi) \triangleq \sup_{\delta} U(\xi, \delta) = U(\xi, \delta^*). \quad (4.4.5)$$

**Bayes decision rules**

We wish to construct the Bayes decision rule, that is, the decision function with maximal  $\xi$ -expected utility. However, doing so by examining all possible decision functions is hard, because (usually) there are many more decision functions than decisions. It is however, easy to find the Bayes decision for each possible observation.

**Theorem 4.4.1.** *If  $U$  is non-negative or bounded, then we can reverse the integration order of*

$$U(\xi, \delta) = \mathbb{E} \{U[\omega, \delta(x)]\} = \int_{\Omega} \int_{\mathcal{S}} U[\omega, \delta(x)] d\psi_\omega(x) d\xi(\omega),$$

*which is the normal form, to obtain the risk in extensive form.*

$$U(\xi, \delta) = \int_{\mathcal{S}} \int_{\Omega} U[\omega, \delta(x)] d\xi(\omega \mid x) df(x), \quad (4.4.6)$$

*where  $f(x) = \int_{\Omega} \psi_\omega(x) d\xi(\omega)$ .*

*Proof.* To prove this when  $U$  is non-negative, we shall use Tonelli's theorem. First we need to construct an appropriate product measure. Note that the original is written Let  $p(x | \omega) \triangleq \frac{d\psi_\omega(x)}{d\nu(x)}$  be the Radon-Nikodym derivative of  $\psi_\omega$  with respect to some dominating measure  $\nu$  on  $\mathcal{S}$ . Similarly, let  $p(\omega) \triangleq \frac{d\xi(\omega)}{d\mu(\omega)}$  be the corresponding derivative for  $\xi$ . Now, the utility can be written as:

$$U(\xi, \delta) = \int_{\Omega} \int_{\mathcal{S}} U[\omega, \delta(x)] p(x | \omega) p(\omega) d\nu(x) d\mu(\omega) \quad (4.4.7)$$

$$= \int_{\Omega} \int_{\mathcal{S}} h(\omega, x) d\nu(x) d\mu(\omega). \quad (4.4.8)$$

Clearly, if  $U$  is non-negative, then  $h$  is non-negative. Then, Tonelli's theorem applies and

$$U(\xi, \delta) = \int_{\mathcal{S}} \int_{\Omega} h(\omega, x) d\mu(\omega) d\nu(x) \quad (4.4.9)$$

$$= \int_{\mathcal{S}} \int_{\Omega} p(x | \omega) p(\omega) d\mu(\omega) d\nu(x) \quad (4.4.10)$$

$$= \int_{\mathcal{S}} \int_{\Omega} p(\omega | x) d\mu(\omega) p(x) d\nu(x) \quad (4.4.11)$$

$$= \int_{\mathcal{S}} \left[ \int_{\Omega} p(\omega | x) d\mu(\omega) \right] p(x) d\nu(x) = \int_{\mathcal{S}} \left[ \int_{\Omega} d\xi(\omega | x) \right] df(x), \quad (4.4.12)$$

where  $p(x) = df(x)/d\nu(x)$ .  $\square$

We can construct an optimal decision function  $\delta^*$  as follows. For any specific observed  $x \in \mathcal{S}$ , we set  $\delta^*(x)$  to:

$$\delta^*(x) \triangleq \arg \max_{d \in D} \mathbb{E}_{\xi}(U | x, d) = \arg \max_{d \in D} \int_{\Omega} U(\omega, d) d\xi(\omega | x).$$

So now we can plug  $\delta^*$  in the extensive form to obtain:

$$\int_{\mathcal{S}} \int_{\Omega} U[\omega, \delta^*(x)] d\xi(\omega | x) df(x) = \int_{\mathcal{S}} \left\{ \min_d \int_{\Omega} U[\omega, d] d\xi(\omega | x) \right\} df(x).$$

Consequently, there is no need to completely specify the decision function before we have seen  $x$ . In particular, this would create problems when  $\mathcal{S}$  is large.

**Definition 4.4.2** (Prior distribution). *The distribution  $\xi$  is called the prior distribution of  $\omega$ .*

**Definition 4.4.3** (Marginal distribution). *The distribution  $f$  is called the (prior) marginal distribution of  $x$ .*

**Definition 4.4.4** (Posterior distribution). *The conditional distribution  $\xi(\cdot | x)$  is called the posterior distribution of  $\omega$ .*

**Bayes decision rule.**

The *optimal decision* given  $x$ , is the optimal decision with respect to the *posterior*  $\xi(\omega | x)$ . Thus, we do not need to pre-compute the complete Bayes-optimal decision rule.

#### 4.4.1 Calculating posteriors

##### Posterior distributions for multiple observations

We now consider how we can re-write the posterior distribution over  $\Omega$  incrementally. Assume that we observe  $x^n \triangleq x_1, \dots, x_n$ . We have a prior  $\xi$  on  $\Omega$ . For the observations, we write:

**Observation probability given history  $x^{n-1}$  and parameter  $\omega$**

$$\psi_\omega(x_n | x^{n-1}) = \frac{\psi_\omega(x^n)}{\psi_\omega(x^{n-1})}$$

Now we can write the posterior as follows:

**Posterior recursion**

$$\xi(\omega | x^n) = \frac{\psi_\omega(x^n)\xi(\omega)}{f(x^n)} = \frac{\psi_\omega(x_n | x^{n-1})\xi(\omega | x^{n-1})}{f(x_n | x^{n-1})}. \quad (4.4.13)$$

Here  $f(\cdot | \cdot) = \int_\Omega \psi_\omega(\cdot | \cdot) d\xi(\omega)$  is a marginal distribution.

##### Posterior distributions for multiple independent observations

Now we consider the case where, given the parameter, the next observation does not depend on the history: If  $\psi(x_n | \omega, x^{n-1}) = \psi_\omega(x_n)$  then  $\psi_\omega(x^n) = \prod_{k=1}^n \psi_\omega(x_k)$ . Then:

**Posterior recursion with conditional independence**

$$\xi_n(\omega) \triangleq \xi_0(\omega | x^n) = \frac{\psi_\omega(x^n)\xi_0(\omega)}{f_0(x_n)} \quad (4.4.14)$$

$$= \xi_{n-1}(\omega | x_n) = \frac{\psi_\omega(x_n)\xi_{n-1}(\omega)}{f_{n-1}(x_n)}, \quad (4.4.15)$$

where we define  $\xi_t$  to be the belief at time  $t$ . Here  $f_n(\cdot | \cdot) = \int_\Omega \psi(\cdot | \cdot, \omega) d\xi_n(\omega)$  is the marginal distribution with respect to the  $n$ -th posterior.



Conditional independence allows us to write the posterior update as an identical recursion at each time  $t$ . We shall take advantage of that when we look at *conjugate prior* distributions. For such models, the recursion involves a particularly simple parameter update.

**Quick summary**

- We want to make a decision against an unknown parameter  $\omega$ .
- The risk is the negative expected utility.
- The Bayes risk is the minimum risk, and it is concave with respect to the distribution of  $\omega$ .
- Our decisions can depend on observations, via a decision function.
- We can construct a complete decision function by computing  $U(\xi, \delta)$  for all *decision functions* (normal form).
- We can instead wait until we observe  $x$  and compute  $U[\xi(\cdot | x), d]$  for all *decisions* (extensive form).
- The posterior given multiple observations can be computed recursively using independence.



## Chapter 5

# Estimation

## 5.1 Calculation of posterior distributions

In the previous unit, we have seen how to make optimal decisions with respect to a given risk function and belief. However, one important question is how belief can be calculated. It is one thing to say that we simply calculate the posterior distribution of a parameter and another thing to actually *do it*. In this unit, we shall look at cases when calculating the posterior distributions of parameters is easy. This occurs when the posterior distribution can be expressed by a function that belongs to the same family as the prior distribution, no matter what the observations are.

In the Bayesian setting, we need to calculate posterior distributions of parameters given data. The basic problem can be stated as follows. Let  $\mathcal{P} \triangleq \{P_\omega \mid \omega \in \Omega\}$  be a family of probability measures on  $(\mathcal{S}, \mathfrak{F}_\mathcal{S})$  and  $\xi$  be our prior probability measure on  $(\Omega, \mathfrak{F}_\Omega)$ . Given some data  $x \sim P_{\omega^*}$ , with  $\omega^* \in \Omega$ , how can we estimate  $\omega^*$ ? The Bayesian answer is, instead of guessing a single  $\omega^*$ , to estimate the posterior distribution  $\xi(\cdot \mid x)$ . In general, the posterior measure is a function  $\xi(\cdot \mid x) : \mathfrak{F}_\Omega \rightarrow [0, 1]$ , with:

$$\xi(B \mid x) = \frac{\int_B P_\omega(x) d\xi(\omega)}{\int_\Omega P_\omega(x) d\xi(\omega)}. \quad (5.1.1)$$

The main question is how to calculate this posterior for any value of  $x$  in practice. One imagines that if  $x$  is a complicated object, this may be a difficult job. However, there exist distribution families and priors such that this calculation is very easy. This happens when a summary of the data that does contains all necessary information can be calculated easily. Formally, this is captured via the concept of a sufficient statistic.

## 5.2 Sufficient statistics

Sometimes we want to summarise the data we have observed. This can happen when the data is a long sequence of simple observations  $x^t = (x_1, \dots, x_t)$ , with  $x_k \in \mathbb{R}$ . It may also be useful to do so when we have a single observation  $x$ , such as a high-resolution image. For some applications, it maybe sufficient to only calculate a really simple function of the data, such as the sample mean, defined below:

**Definition 5.2.1** (Sample mean). *The sample-mean  $\bar{x}_t : \mathbb{R}^t \rightarrow \mathbb{R}$  of a sequence  $x_k$  is defined as:*

$$\bar{x}_t \triangleq \frac{1}{t} \sum_{k=1}^t x_k, \quad x_k \in \mathbb{R}. \quad (5.2.1)$$

This summary, or any other function of the observations is called a *statistic*. In particular, we would be interested in using statistics with which we could completely replace all the real data in our calculations, without losing any information. Such statistics are called *sufficient*.

### 5.2.1 Formalisation of sufficient statistics

We consider the standard probabilistic setting. Let  $\mathcal{S}$  be a sample space and  $\Omega$  be a parameter space defining a family of measures on  $\mathcal{S}$ :

$$\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}.$$

In addition, we must also define an appropriate prior distribution  $\xi$  on the parameter space  $\Omega$ .

**Example 5.2.1.** Consider a sequence of random variables  $x_k$  taking values in  $\mathcal{X} = \{0, 1\}$ . Let  $\Omega = \{0, \frac{1}{2}, 1\}$  be a family of Bernoulli distributions with parameter  $\omega$ :

$$\mathbb{P}(x_k = 1 \mid \omega) = P_\omega(\{1\}) = \omega,$$

such that  $x_k \sim P_\omega$  for all  $k$ . A suitable prior distribution could be defined via  $\xi(\omega) = \theta_i$  for some vector  $\theta \in \Delta^3$  on the three-dimensional simplex.

**Definition 5.2.2.** Let  $\Xi$  be a set of prior distributions on  $\Omega$ ,  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$  be a family of distributions on  $\mathcal{S}$ . A statistic  $T : \mathcal{S} \rightarrow \mathcal{Z}$  is a sufficient statistic for  $\langle \mathcal{P}, \Xi \rangle$  if:

$$\xi(\omega \mid x) = \xi(\omega \mid x'), \quad (5.2.2)$$

for any prior  $\xi \in \Xi$  and any  $x, x' \in \mathcal{S}$  such that:

$$T(x) = T(x').$$

**Theorem 5.2.1.** A statistic  $T : \mathcal{S} \rightarrow \mathcal{Z}$  is sufficient for a family  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$  of probability functions on  $\mathcal{S}$  iff there exist functions  $u : \mathcal{S} \rightarrow (0, \infty)$ , and  $v : \mathcal{Z} \times \Omega \rightarrow [0, \infty)$  such that  $\forall x \in \mathcal{S}, \omega \in \Omega$ :

$$P_\omega(x) = u(x)v[T(x), \omega], \quad u > 0, v \geq 0. \quad (5.2.3)$$

*Proof.* The proof will be for the general case. The case when  $\Omega$  is finite is technically simpler and is left as an exercise. Assume the existence of  $u, v$ . Then for  $B \in \mathfrak{F}_\Omega$ :

$$\begin{aligned} \xi(B \mid x) &= \frac{\int_B u(x)v[T(x), \omega] d\xi(\omega)}{\int_\Omega u(x)v[T(x), \omega] d\xi(\omega)} \\ &= \frac{\int_B v[T(x), \omega] d\xi(\omega)}{\int_\Omega v[T(x), \omega] d\xi(\omega)}. \end{aligned}$$

If  $T(x) = T(x')$ , then the above is also equal to  $\xi(\omega \mid x')$ , so  $T$  is a sufficient statistic.

Conversely, let  $T$  be a sufficient statistic. Let  $\mu$  be a dominating measure on  $\mathcal{S}$  so that  $p(\omega) \triangleq \frac{d\xi(\omega)}{d\mu(\omega)}$  and take the derivative at  $B \rightarrow \omega$  so that:

$$p(\omega \mid x) \triangleq \frac{d\xi(\omega \mid x)}{d\mu(\omega)} = \frac{P_\omega(x)p(\omega)}{\int_\Omega P_\omega(x) d\xi(\omega)}$$

Consequently, we can write:

$$P_\omega(x) = \frac{p(\omega \mid x)}{p(\omega)} \int_\Omega P_\omega(x) d\xi(\omega).$$

Since  $T$  is sufficient, there is some function  $g : \mathcal{Z} \times \Omega$  such that  $p(\omega \mid x) = g[T(x), \omega]$ . Consequently, we can factorise  $P_\omega$  as:

$$P_\omega(x) = v[T(x), \omega]u(x),$$

where  $u(x) = \int_\Omega P_\omega(x) d\xi(\omega)$  and  $v[T(x), w] = g[T(x), w]/\xi(w)$ .  $\square$

With this factorisation,  $u$  is the only factor that depends directly on  $x$ .

**Example 5.2.2.** Suppose  $x^t = (x_1, \dots, x_t)$  is a random sample from a Bernoulli distribution with parameter  $\omega$ . Then the joint probability is

$$P_\omega(x^t) = \prod_{k=1}^t P_\omega(x_k) = \omega^{s_t} (1 - \omega)^{t-s_t}$$

with  $s_t = \sum_{k=1}^t x_k$  being the number of times 1 has been observed until time  $t$ . Then the statistic:

$$T(x^t) = \sum_{k=1}^t x_k.$$

satisfies (5.2.3) with  $u(x) = 1$ , while  $P_\omega(x^t)$  only depends on the data through the statistic  $s_t = T(x^t)$ .

## 5.2.2 Exponential families

Many well-known distributions such as the Gaussian, Bernoulli and Dirichlet distribution are members of the exponential family of distributions. All such distributions are factorisable in the manner shown below, while at the same time they have fixed-dimension sufficient statistics.

**Definition 5.2.3.** A distribution family  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$  with  $P_\omega$  a probability density or a probability function defined on the sample space  $\mathcal{S}$ , is said to be an exponential family if for any  $x \in \mathcal{S}$ ,  $\omega \in \Omega$ :

$$P_\omega(x) = a(\omega)b(x) \exp \left[ \sum_{i=1}^k g_i(\omega)h_i(x) \right]. \quad (5.2.4)$$

**Remark 5.2.1.** Among families of distributions satisfying certain regularity conditions, only exponential families have a fixed-dimension sufficient statistic. That is, there is some constant  $C < \infty$  such that:

$$\sup_{X \in \mathfrak{F}_\mathcal{S}} \dim(\{T(x) \mid x \in \mathcal{S}\}) \leq C. \quad (5.2.5)$$

Because of this property, exponential family distributions admit so-called *conjugate* prior distribution families. These have the property that any posterior distribution calculated will remain within the conjugate family. Frequently, because of the simplicity of the statistic used, calculation of the conjugate posterior parameters is very simple.

## 5.3 Conjugate priors

In this section, we examine some well-known conjugate families. First, we give sufficient conditions for the existence of conjugate family of priors for a given distribution family and statistic.

**Remark 5.3.1.** *If a family  $\mathcal{P}$  of distributions on  $\mathcal{S}$  has a sufficient statistic  $T : \mathcal{S} \rightarrow \mathcal{Z}$  of fixed dimension for any  $x \in \mathcal{S}$ , then there exists a conjugate family of priors  $\Xi = \{\xi_\alpha \mid \alpha \in A\}$ , where  $A$  is a set of possible parameters for the prior distribution, such that:*

1.  $P_\omega(x)$  is proportional to some  $\xi_\alpha \in \Xi$ :

$$\forall x \in \mathcal{S}, \exists \xi_\alpha \in \Xi, c > 0 : \int_B P_\omega(x) d\xi_\alpha(\omega) = c\xi_\alpha(B), \forall B \in \mathfrak{F}_\Omega$$

2. The family is closed under multiplication:

$$\forall \xi_1, \xi_2 \in \Xi, \exists \xi_\alpha \in \Xi, c > 0$$

such that:

$$\xi_\alpha = c\xi_1\xi_2.$$

### 5.3.1 Bernoulli-Beta conjugate pair

The Bernoulli-Beta conjugate pair of families is the simplest example. It is useful for problems where we wish to measure success rates of independent trials. First, we shall give details on the Bernoulli distribution. Then, we shall define the Beta distribution and describe its conjugate relation to the Bernoulli.

#### Bernoulli distribution

The Bernoulli distribution is a discrete distribution with outcomes taking values in  $\{0, 1\}$ . It is ideal for modelling the outcomes of independent random trials with fixed probability of success.

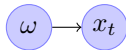


Figure 5.1: Bernoulli graphical model

**Definition 5.3.1** (Bernoulli distribution). *The Bernoulli distribution is discrete, with outcomes  $\mathcal{S} = \{0, 1\}$ , parameter  $\omega \in [0, 1]$ , and probability function:*

$$P_\omega(u) = \begin{cases} \omega, & u = 1 \\ 1 - \omega, & u = 0 \end{cases} = \omega^u(1 - \omega)^{1-u}.$$

*If  $x_t$  is distributed according to a Bernoulli distribution with parameter  $\omega$ , we write  $x_t \sim \text{Bern}(\omega)$ . The Bernoulli distribution can be extended to  $\mathcal{S} = \{0, 1\}^t$  by modelling each outcome as independent. Then  $P_\omega(x^t) = \prod_{k=1}^t P_\omega(x_k)$ . In that case, it is directly related to the Binomial distribution.*

**Definition 5.3.2** (Binomial Distribution). *Let us denote the total number of 1's observed until time  $t$  by  $s_t = \sum_{k=1}^t x_k$ . Then  $\mathbb{P}(s_t = k \mid \omega) = \binom{t}{k} \omega^k (1 - \omega)^{t-k}$  is the probability that,  $k$  out of  $t$  trials will be positive. Remember that  $\binom{x}{k} = \prod_{i=0}^{k-1} (x - i) / (1 + i)$ . If  $s_t$  is drawn from a binomial distribution with parameters  $\omega, t$ , we write  $s_t \sim \text{Binom}(\omega, t)$ .*

The difference between the two distributions, is that the Bernoulli is a distribution on a sequence of outcomes, while the binomial a distribution on the total number of positive outcomes.

**Example 5.3.1.** *A fair coin toss can be modelled as a Bernoulli distribution with  $\omega = \frac{1}{2}$ .*

Thus, the Bernoulli distribution is the simplest possible parametric model. If the  $\omega$  parameter is known, then all the observations are independent of each other. However, this is not the case when  $\omega$  is unknown. For example, let  $\Omega = \{\omega_1, \omega_2\}$ . Then  $\mathbb{P}(x^t) = \sum_{\omega \in \Omega} \mathbb{P}(x^t \mid \omega) \mathbb{P}(\omega) = \sum_{\omega} \prod_{k=1}^t \mathbb{P}(x_k \mid \omega) \mathbb{P}(\omega)$ .

### Beta distribution

The Beta distribution is a distribution on the interval  $[0, 1]$ . It has two parameters that determine the density of the observations. Because the outcomes of this distribution can be used to define the parameter of the Bernoulli distribution, we can now call the distribution's outcomes  $\omega$  and its parameter  $\alpha$ .

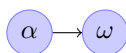


Figure 5.2: Beta graphical model

**Definition 5.3.3** (Beta distribution). *The Beta distribution has outcomes  $\omega \in \Omega = [0, 1]$  and parameters  $\alpha_0, \alpha_1 > 0$ ,  $\alpha = (\alpha_1, \alpha_0)$ . It is defined via its probability density function:*

$$f(\omega \mid \alpha) = \frac{\Gamma(\alpha_0 + \alpha_1)}{\Gamma(\alpha_0)\Gamma(\alpha_1)} \omega^{\alpha_1-1} (1 - \omega)^{\alpha_0-1}, \quad (5.3.1)$$

where  $\Gamma$  is the gamma function. If  $x_t$  is distributed according to a Beta distribution with parameters  $\alpha_1, \alpha_0$ , we write:  $x_t \sim \text{Beta}(\alpha_1, \alpha_0)$ .

The Beta family defines a family of probability measures  $\mathcal{P} = \{P_\alpha \mid \alpha_0, \alpha_1 > 0\}$ , with the probability of  $A \subset [0, 1]$ :  $\mathbb{P}(\omega \in A \mid \alpha) = P_\alpha(A) \triangleq \int_A p(\omega \mid \alpha) d\omega$ .

A Beta distribution with parameter  $\alpha$  has expectation  $\mathbb{E}(x \mid \alpha) = \alpha_1 / \|\alpha\|_1$ .

**Example 5.3.2.** *The parameter  $\omega \in [0, 1]$  of a randomly selected coin can be modelled as a Beta distribution peaking around  $1/2$ . Usually one assumes that coins are fair. However, not all coins are exactly the same. Thus, it is possible that each coin deviates slightly from fairness. We can use a Beta distribution to model how likely (we think) different values  $\omega$  of coin parameters are.*

Figure 5.3 shows the density of a Beta distribution for four different parameter vectors. When  $\alpha_0 = \alpha_1 = 1$ , the distribution is equivalent to a uniform one.



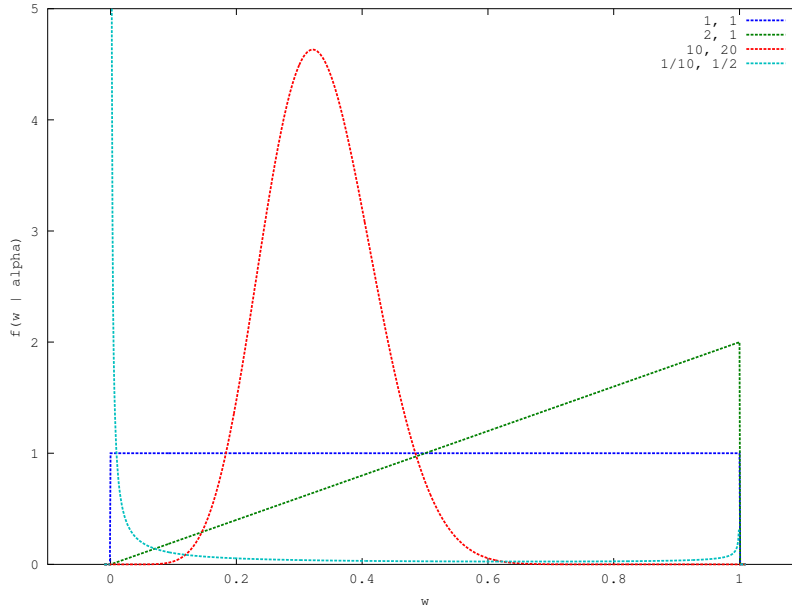


Figure 5.3: Four example Beta densities

The Beta distribution is useful for expressing probabilities of random variables in bounded intervals. In particular, since probabilities of events take values in  $[0, 1]$ , the Beta distribution is an excellent choice for expressing uncertainty about a probability.

### Beta prior for Bernoulli distributions

One simple idea is to encode our uncertainty about an unknown parameter of the Bernoulli distribution, is to use a Beta distribution. The main idea is to assume that the Bernoulli parameter  $\omega \in [0, 1]$  is unknown but fixed. We define a Beta prior distribution for  $\omega$  to represent our uncertainty. This can be summarised by a parameter  $\alpha$  and we write  $\xi_0(B) \triangleq \int p(\omega | \alpha) d\omega$  for our prior distribution  $\xi_0$ .

The posterior probability in that case is

$$p(\omega | x^t, \alpha) = \frac{\prod_{k=1}^t P_{\omega}(x_k) p(\omega | \alpha)}{\int_{\Omega} \prod_{k=1}^t P_{\omega}(x_k) p(\omega | \alpha) d\omega} \quad (5.3.2)$$

$$\propto \omega^{s_{t,1} + \alpha_1 - 1} (1 - \omega)^{s_{t,0} + \alpha_0 - 1} \quad (5.3.3)$$

and  $f(x) \propto g(x)$  means that  $f(x) = cg(x)$  for some constant  $c$ .

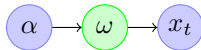


Figure 5.4: Beta-Bernoulli graphical model.

### Coin tossing

A simple illustration of the use of the Beta-Bernoulli model is to estimate the bias of a particular coin. The following figure shows a sequence of beliefs at times 0, 10, 100, 1000 respectively, from a coin with bias  $\omega = 0.6$ . To demonstrate

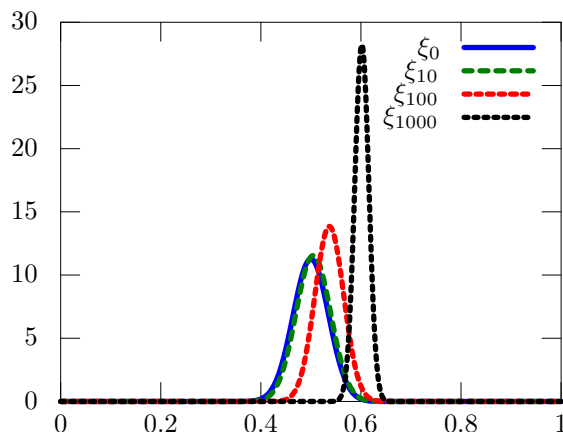


Figure 5.5: Changing beliefs as we observe tosses from a coin with probability  $\omega = 0.6$  of heads.

how belief changes, we can perform the following simple experiment. Imagine a coin such that, when it is tossed, it has a probability 0.6 of coming heads every time it is tossed, independently of previous outcomes. Thus, the distribution of outcomes is a Bernoulli distribution with parameter  $\omega = 0.6$ .

We wish to form an accurate belief about how biased the coin is, under the assumption that the outcomes are Bernoulli with parameter  $\omega$ . Our initial belief,  $\xi_0$ , is modelled as a Beta distribution on the parameter space  $\Omega = [0, 1]$ , with parameters  $\alpha_0 = \alpha_1 = 100$ . This places a strong prior on the coin being close to fair. However, we still allow for the possibility that the coin is biased.

Due to the strength of our prior, after 10 observations, the situation has not changed much and the belief  $\xi_{10}$  is very close to the starting one. However, after 100 observations, our belief has now shifted towards 0.6, the true bias of the coin. After a total of 1000 observations, our belief is centered very close to 0.6, and is now much more concentrated, reflecting the fact that we are almost certain about the value of  $\omega$ .

## 5.4 Credible intervals

### Credible interval

According to our current belief  $\xi$ , there is a certain subjective probability that the unknown parameter  $\omega$  takes a certain value. We can use this to construct intervals where we think the unknown parameter is most likely to be.

**Definition 5.4.1.** *Given some probability measure  $\xi$  on  $\Omega$  representing our*

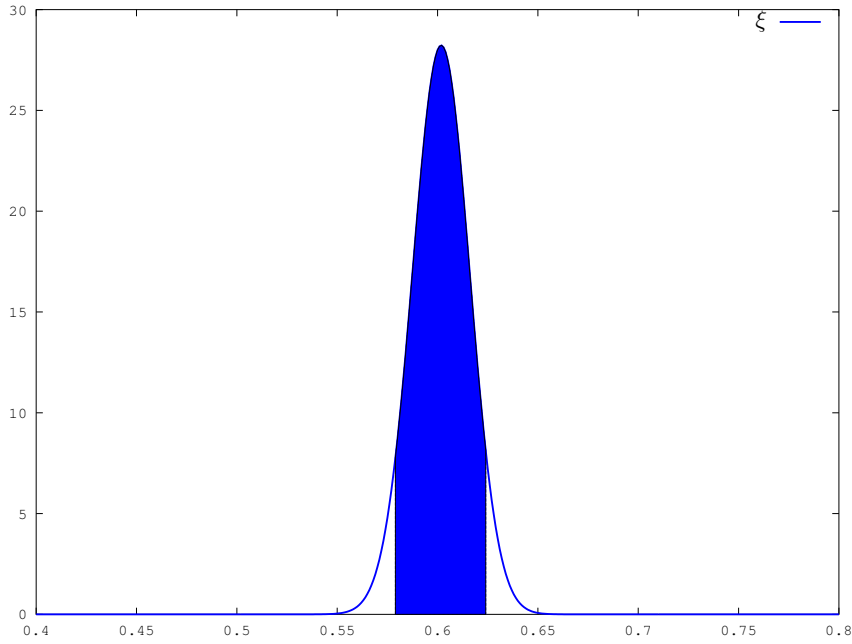


Figure 5.6: 90% credible interval after 1000 observations from a Bernoulli with  $\omega = 0.6$ .

belief and some interval  $A \subset \Omega$ ,

$$\xi(A) = \int_A d\xi = \mathbb{P}(\omega \in A \mid \xi).$$

is our subjective belief that the unknown parameter  $\omega$  is in  $A$ . If  $\xi(A) = s$ , then we say that  $A$  is an  $s$ -credible interval (or set), or an interval of size (or measure)  $s$ .

**Constructing the credible interval.**

For prior distributions on  $\mathbb{R}$ , constructing an  $s$ -credible interval is usually done by finding  $\omega_u, \omega_l \in \mathbb{R}$  such that

$$\xi([\omega_l, \omega_u]) = s.$$

However, *any* choice of  $A$  such that  $\xi(A) = s$  is valid.

Figure 5.6 shows the 90% credible interval. (The measure of  $A$  under  $\xi$  is  $\xi(A) = 0.9$ .)

What is the probability that the true value of  $\omega$  will be within a particular credible interval? This will depend on how well our prior  $\xi_0$  matches the true distribution from which the parameter  $\omega$  was drawn.

**Reliability of credible intervals**

Assume  $\phi, \xi_0$  are probability measures on the parameter set  $\Omega$ , where our prior belief is  $\xi_0$  and  $\phi$  is the actual distribution of  $\omega \in \Omega$ . Each  $\omega$  defines a measure  $P_\omega$  on the observation set  $\mathcal{S}$ . Let us construct a credible interval  $A_t \subset \Omega$  (which is a random variable  $A_t : S^t \rightarrow \mathcal{P}_\Omega$ ) such that it has measure  $s = \xi_t(A_t)$  for all  $t$ . Finally, let  $Q \triangleq \int_\omega P_\omega d\phi(\omega)$  be the marginal distribution on  $\mathcal{S}$ . Then the probability that the credible interval  $A_t$  will not include  $\omega$  is

$$Q(\{x^t \in S^t \mid \omega \notin A_t\}).$$

The main question is how this failure probability relates to  $s, t$  and  $\xi_0$ . There are some results that show that

$$Q(\{x^t \in S^t \mid \omega \in A_t\}) \leq (1-s)[1 + D(\xi_0 \parallel \phi) \mathcal{O}(t^{-1/2})].$$

Let us now design an experiment for examining how often a typical credible interval includes the parameter we are interested in. In order to do so, we will have Nature draw the parameter from some arbitrary distribution  $\phi$ , which may differ from our own assumed prior distribution  $\xi_0$ .

#### Experimental testing of a credible interval

- 1: Given a probability family  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$ .
- 2: Nature chooses distribution  $\phi$  over  $\Omega$ .
- 3: We choose another distribution  $\xi_0$  over  $\Omega$ .
- 4: **for**  $k = 1, \dots, n$  **do**
- 5:   Draw  $\omega_k \sim \phi$ .
- 6:   Draw  $x^T \mid \omega_k \sim P_{\omega_k}$ .
- 7:   **for**  $t = 1, \dots, T$  **do**
- 8:     Calculate  $\xi_t(\cdot) = \xi_0(\cdot \mid x^t)$  for all  $t$ .
- 9:     Calculate  $A_t$ , the 50% CI for all  $t$ .
- 10:    Check failure:  $\epsilon_{t,k} = \mathbb{I}\{\omega_k \notin A_t\}$
- 11:   **end for**
- 12: **end for**
- 13: Average over all  $k$ :  $\epsilon_t = \frac{1}{n} \sum_{k=1}^n \epsilon_{t,k}$ .

We performed this experiment for  $n = 1000$  trials and for  $T = 100$  observations per trial. Figure 5.7 illustrates what happens when  $\phi = \xi_0$ . We see that the credible interval is always centered around our initial mean guess and that it is quite tight. Figure 5.8 shows the average number of times the credible interval  $A_t$  around our estimated mean did not match the actual value of  $\omega_k$ . Since the measure of our interval  $A_t$  is always  $\xi_t(A_t) = 1/2$ , we expect our error probability to be  $1/2$ . This is always the case.

On the other hand, Figure 5.9 illustrates what happens when  $\phi \neq \xi_0$ . In fact in that case,  $\phi(\omega) = \delta(\omega - 0.6)$ , so that  $\omega_k = 0.6$  for all trials  $k$ . We see that the credible interval is always centered around our initial mean guess and that it is always quite tight. Figure 5.10 shows the average number of failures. We see that initially, due to the fact that our prior is different from the distribution from which the  $\omega_k$  are selected, we make many more mistakes. However, eventually, our prior is swamped by the data and our error rate converges to 50%.

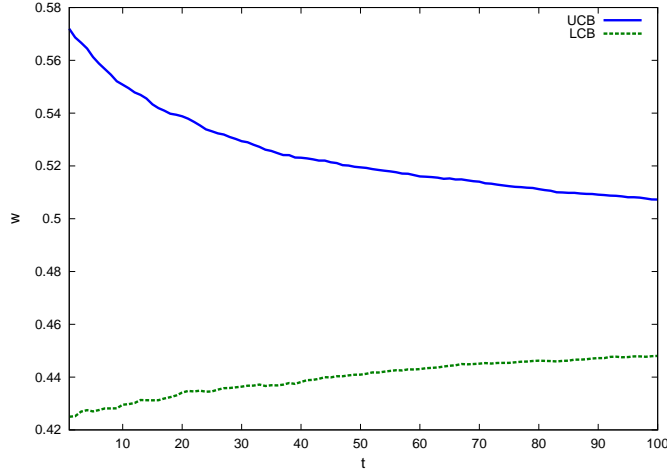


Figure 5.7: 50% credible intervals for a prior  $\text{Beta}(10,10)$ , matching the distribution of  $\omega$ .

## 5.5 Concentration inequalities

### The sample mean revisited

It is interesting to consider the case calculating a sample mean. We have seen that, for the Beta-Bernoulli conjugate prior, it is a simple enough matter to calculate a posterior distribution. From that, we can obtain a credible interval on the expected value of the unknown Bernoulli distribution. However, we would like to do the same for arbitrary distributions on  $[0, 1]$ , rather than the Bernoulli, which is defined on  $\{0, 1\}$ .

**Example 5.5.1** (Sample mean). *Let  $x^t = (x_1, \dots, x_t)$ , with  $x^t \sim P$ ,  $\mathbb{E}_P x_k = \mu \in [0, 1]$  for all  $k$ , with  $P$  an unknown distribution with support in  $[0, 1]$ . Let*

$$\bar{x}_t \triangleq \frac{1}{t} \sum_{k=1}^t x_k.$$

be the sample mean. We are going to look at a number of well-known inequalities that allow to bound  $|\hat{x}_t - \mu|$ .

### The Markov inequality

**Theorem 5.5.1** (Markov inequality). *If  $X \sim P$ , with  $P$  a distribution on  $[0, \infty)$ , then:*

$$\mathbb{P}(X \geq t) \leq \frac{\mathbb{E} X}{t}, \quad (5.5.1)$$

where  $\mathbb{P}(X \geq t) = P(\{x \mid x \geq t\})$ .

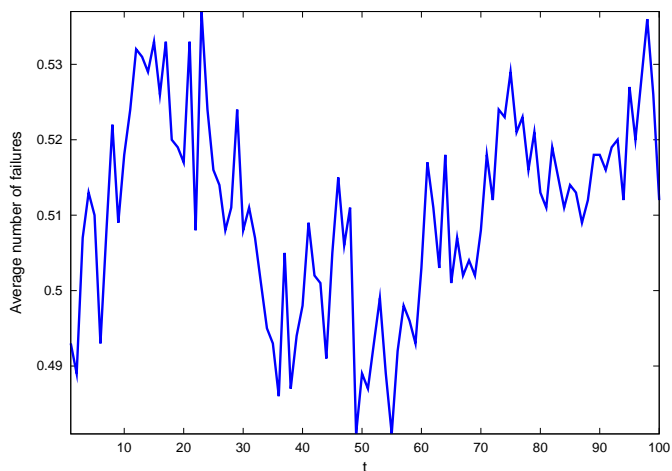


Figure 5.8: Failure rate of 50% CI for a prior  $\text{Beta}(10, 10)$ , matching the distribution of  $\omega$ .

*Proof.* The expectation of  $X$  is:

$$\mathbb{E} X = \int_0^\infty x \, dP(x) \quad (5.5.2)$$

$$= \int_0^t x \, dP(x) + \int_t^\infty x \, dP(x) \quad (5.5.3)$$

$$\geq 0 + \int_t^\infty t \, dP(x) \quad (5.5.4)$$

$$= tP(\{x \mid x \geq t\}) = t\mathbb{P}(X \geq t). \quad (5.5.5)$$

□

Consequently, we have that  $\mathbb{P}(|\bar{x}_t - p| \geq \epsilon) \leq \mathbb{E} |\bar{x}_t - p| / \epsilon$ . For  $X \in [0, 1]$ , we obtain the bound

$$\mathbb{P}(|\bar{x}_t - p| \geq \epsilon) \leq 1/\epsilon.$$

Can we do better? We might, if we take advantage of the following trick. For monotonic  $f$ ,

$$\mathbb{P}(X \geq t) = \mathbb{P}(f(X) \geq f(t)) \quad (5.5.6)$$

as  $\{x \mid x \geq t\} = \{x \mid f(x) \geq f(t)\}$ . Thus, we can apply the Markov inequality in a large number of contexts.

### The Chebyshev inequality

**Theorem 5.5.2.** *Let  $X$  be a random variable with expectation  $\mu = \mathbb{E} X$  and variance  $\sigma^2 = \mathbb{V} X$ . Then, for all  $k > 0$ :*

$$\mathbb{P}(|X - \mu| \geq k\sigma) \leq k^{-2}. \quad (5.5.7)$$

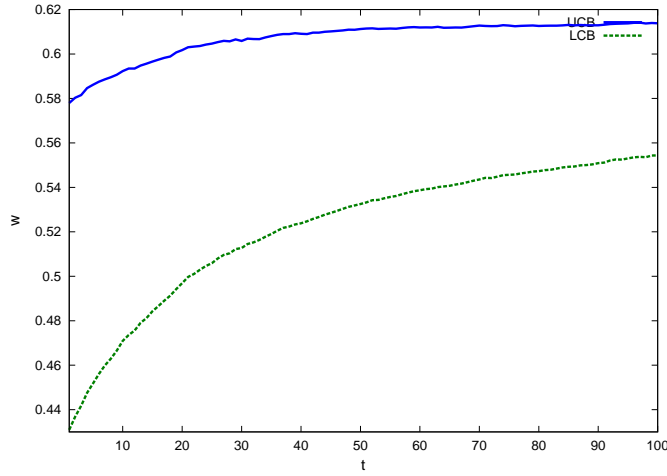


Figure 5.9: 50% credible intervals for a prior  $\text{Beta}(10, 10)$ , when  $\omega = 0.6$ .

*Proof.*

$$\begin{aligned} \mathbb{P}(|X - \mu| \geq k\sigma) &= \mathbb{P}\left(\frac{|X - \mu|}{k\sigma} \geq 1\right) = \mathbb{P}\left(\frac{|X - \mu|^2}{k^2\sigma^2} \geq 1\right) \\ &\leq \mathbb{E}\left(\frac{(X - \mu)^2}{k^2\sigma^2}\right) = \frac{\mathbb{E}(X - \mu)^2}{k^2\sigma^2} = k^{-2} \end{aligned}$$

□

**Example 5.5.2** (Application to sample mean). *It is easy to show that the sample mean has expectation  $\mu$  and variance  $\sigma_x^2/t$ , where  $\sigma_x^2 = \mathbb{V} x_k$ . Consequently:*

$$\mathbb{P}\left(|\bar{x}_t - \mu| \geq k\sigma_x/\sqrt{t}\right) \leq k^{-2}.$$

Setting  $\epsilon = k\sigma_x/\sqrt{t}$  we get  $k = \epsilon\sqrt{t}/\sigma_x$  and hence

$$\mathbb{P}(|\bar{x}_t - \mu| \geq \epsilon) \leq \frac{\sigma_x^2}{\epsilon^2 t}.$$

### Chernoff-Hoeffding bounds

The previous inequality can be quite loose. In fact, one can prove tighter bounds for the estimation of an expected value. All these bounds rest upon a different application of the Markov inequality, due to Chernoff.

#### Main idea of Chernoff bounds.

Let  $S_t = \sum_{k=1}^t X_k$ , with  $X_k \sim P$  independently, i.e.  $X^t \sim P^t$ . By definition, from Markov's inequality we obtain in turn, for any  $\alpha > 0$

$$\mathbb{P}(S_t \geq u) = \mathbb{P}(e^{\alpha S_t} \geq e^{\alpha u}) \leq e^{-\alpha u} \mathbb{E} e^{\alpha S_t} = e^{-\alpha u} \prod_k \mathbb{E} e^{\alpha X_k}. \quad (5.5.8)$$

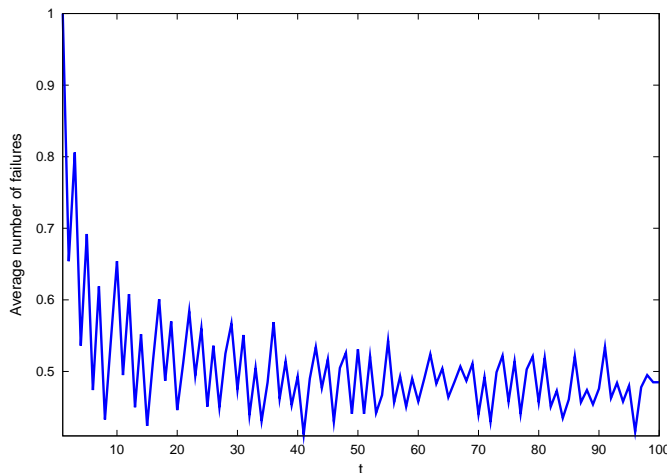


Figure 5.10: Failure rate of 50% CI for a prior  $\text{Beta}(10, 10)$ , when  $\omega = 0.6$ .

**Theorem 5.5.3.** *Hoeffding inequality (Hoeffding [1963], Theorem 2) Let  $x_k \sim P_k$  with  $x_k \in [a_k, b_k]$  with  $\mathbb{E} X_k = \mu_k$ . Then*

$$\mathbb{P}(\bar{x}_t - \mu \geq \epsilon) \leq \exp\left(-\frac{2t^2\epsilon^2}{\sum_{k=1}^t (b_k - a_k)^2}\right), \quad (5.5.9)$$

where  $\bar{x}_t = \frac{1}{t} \sum_{k=1}^t x_k$  and  $\mu = \frac{1}{t} \sum_{k=1}^t \mu_k$ .

*Proof.* Use (5.5.8), setting  $X_k = x_k - \mu_k$  so that  $S_t = t(\bar{x}_t - \mu)$  and  $u = t\epsilon$ . Then:

$$\begin{aligned} \mathbb{P}(\bar{x}_t - \mu \geq \epsilon) &= \mathbb{P}(S_t \geq u) \\ &\leq e^{-\theta u} \prod_{k=1}^t \mathbb{E} e^{\theta X_k} = e^{-\theta t \epsilon} \prod_{k=1}^t \mathbb{E} e^{\theta(x_k - \mu_k)}. \end{aligned}$$

Applying Jensen's inequality directly to the expectation does not help. However, we can use convexity in another way. Let  $f(x)$  be the linear upper bound on  $e^{\theta x}$  on the interval  $[a, b]$ , i.e.

$$f(x) = \frac{b-x}{b-a} e^{\theta a} + \frac{x-a}{b-a} e^{\theta b} \geq e^{\theta x}.$$

Then obviously  $\mathbb{E} \theta x \leq \mathbb{E} f(x)$ . Applying this to the above we get

$$\mathbb{P}(\bar{x}_t - \mu \geq \epsilon) \leq e^{-\theta t \epsilon} \prod_{k=1}^t \mathbb{E} e^{\theta(x_k - \mu_k)}.$$

Now

$$e^{\theta(x_k - \mu_k)} \leq \frac{e^{-\theta \mu_k}}{b_k - a_k} \{(b_k - \mu_k) e^{\theta a_k} + (\mu_k - a_k) e^{\theta b_k}\} = e^{F(\theta_i)},$$



where

$$F(\theta_i) = -\theta_i p_i + \ln(1 - p_i + p_i e^{\theta_i}),$$

$$\theta_i = \theta(b_i - a_i), p_i = \frac{\mu_i - a_i}{b_i - a_i}.$$

Taking derivatives and computing the Taylor expansion, we get

$$\mathbb{E} e^{\theta(x_k - \mu_k)} \leq e^{\frac{1}{8}\theta^2(b_k - a_k)^2}$$

$$\mathbb{P}(\bar{x}_t - \mu \geq \epsilon) \leq e^{-\theta t \epsilon + \frac{1}{8}\theta^2 \sum_{k=1}^t (b_k - a_k)^2}.$$

This is minimised at  $\theta = 4t\epsilon / \sum_{k=1}^t (b_k - a_k)^2$ . Plugging it into the above, we obtain the required result.  $\square$

**Example 5.5.3** (Application to sample mean). For  $x_k \in [0, 1]$ :

$$\mathbb{P}(|\bar{x}_t - \mu| \geq \epsilon) \leq 2e^{-2t\epsilon^2}$$

## 5.6 Approximate Bayesian approaches

### Monte-Carlo inference

#### Estimating expectations.

Let  $f : S \rightarrow [0, 1]$  and  $P$  a measure on  $S$ . Then

$$\mathbb{E}_P f = \int_S f(x) dP(x). \quad (5.6.1)$$

Estimating expectations is relatively easy, as long as we can generate samples from  $P$ . Then, we can our error in estimating its expectation by using the Hoeffding bound.

**Corollary 5.6.1.** Let  $\hat{f}_n = \frac{1}{n} \sum_t f(x_t)$  with  $x_t \sim P$  and  $f : S \rightarrow [0, 1]$ . Then:

$$P\left(\left\{x^n \in S^n \mid |\hat{f}_n - \mathbb{E} f| \geq \epsilon\right\}\right) \leq 2e^{-2n\epsilon^2}. \quad (5.6.2)$$

This technique is simple and fast. However, we frequently cannot sample from  $P$ , but only from some alternative distribution  $Q$ . Then it is hard to bound our error.

Another interesting application of this technique is the calculation of posterior distributions.

**Example 5.6.1** (Calculation of posterior distributions). Assume a probability family  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$  and a prior distribution  $\xi$  on  $\Omega$  such that we can draw  $\omega \sim \xi$ . The posterior distribution can be written according to (5.1.1). The nominator can be written as

$$\int_B P_\omega(x) d\xi(\omega) = \int_\Omega \mathbb{I}\{\omega \in B\} P_\omega(x) d\xi(\omega) = \mathbb{E}_\xi[\mathbb{I}\{\omega \in B\} P_\omega(x)]. \quad (5.6.3)$$

Similarly, the denominator can be written as  $\mathbb{E}_\xi[P_\omega(x)]$ . If  $P_\omega$  is bounded, then the errors can be bounded too.

An extension of this approach involves Markov chain Monte-Carlo (MCMC) methods. These are sequential sampling procedures, where data is sampled iteratively. At the  $k$ -th iteration, we obtain a sample  $x^{(k)} \sim Q_k$ , where  $Q_k$  depends on the previous sample drawn,  $x^{(k-1)}$ . Although under mild conditions  $Q_k \rightarrow P$ , there is no easy way to determine *a priori* when the procedure has converged.

### 5.6.1 Approximate Bayesian Computation

Many times, we have a more fundamental problem. The family of models that we consider An approach that has been popularised by

---

**Algorithm 1** ABC Rejection Sampling
 

---

- 1:  $\omega_k \sim \xi$ .
  - 2:  $\hat{x}^{(k)} \sim M_{\omega_k}$ .
  - 3: If  $D[T(x), T(\hat{x}^{(k)})] \leq \epsilon$  accept  $\omega_k$  as a sample from  $\xi(\omega \mid x)$ .
- 

## 5.7 Other conjugate families

### 5.7.1 Conjugates for the normal distribution

Normal distribution

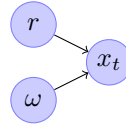


Figure 5.11: Normal graphical model

- Continuous distribution, outcomes:  $\mathcal{S} = \mathbb{R}$ .
- Parameters: mean  $\omega \in \mathbb{R}$ , precision  $r \in \mathbb{R}^+$ .
- Probability density function:  $f(x_t \mid \omega, r) = \frac{r}{\sqrt{2\pi}} \exp\left(-\frac{r}{2}(x_t - \omega)^2\right)$
- Independence:  $f(x^t \mid \omega, r) = \prod_{k=1}^t f(x_k \mid \omega, r) = \left(\frac{r}{\sqrt{2\pi}}\right)^n \exp\left(-\frac{r}{2} \sum_{k=1}^t (x_k - \omega)^2\right)$

#### Normal prior for normal distribution with known precision, unknown mean

The simplest normal estimation problem occurs when we only need to estimate the mean, but we assume that the variance, or equivalently the precision, is known. For Bayesian estimation, it is convenient to assume that the mean  $\omega$  is drawn from *another* normal distribution with known mean.

**The model**

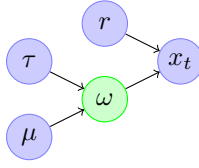


Figure 5.12: Normal with unknown mean, graphical model

- Known precision  $r$ .
- Unknown mean  $\omega$  with assumed distribution  $\omega \sim \mathcal{N}(\mu, \tau)$ .
- $\mu, \tau$  are the mean and precision of our belief about  $\omega$ , respectively.

**Theorem 5.7.1.** Let  $x^t$  be a random sample from  $\mathcal{N}(\omega, r)$ . If  $\omega \sim \xi_0 = \mathcal{N}(\mu, \tau)$ , then  $\xi_t = \mathcal{N}(\mu', \tau')$ , with

$$\mu' = \frac{\tau\mu + nr\bar{x}_t}{\tau'}, \quad \tau' = \tau + nr, \quad (5.7.1)$$

and  $\bar{x}_t \triangleq \frac{1}{t} \sum_{k=1}^t x_k$ .

### Gamma distribution

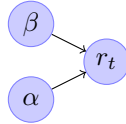


Figure 5.13: Gamma graphical model

The Gamma distribution is a distribution on the interval  $[0, \infty)$ . It has two parameters that determine the density of the observations.

- Continuous distribution, outcomes  $r \in \Omega = [0, \infty)$ .
- Parameters  $\alpha, \beta > 0$ .
- Probability density function:  $f(r \mid \alpha) = \frac{\beta^\alpha}{\Gamma(\alpha)} r^{\alpha-1} e^{-\beta r}$ , for  $r > 0$ .
- $\Gamma(\alpha) = \int_0^\infty u^{\alpha-1} e^{-u} du$ .

#### Relations to other distributions

For  $\alpha = 1$ ,  $\beta > 0$  one obtains an *exponential* distribution with parameter  $\beta$ .

For  $n \in \mathbb{N}$  and  $\alpha = n/2$ ,  $\beta = 1/2$  one obtains a  $\chi^2$  distribution with  $n$  degrees of freedom. Also, if  $x^n$  are i.i.d. standard normal, then  $\sum_{k=1}^n x_k^2$  has a  $\chi^2$  distribution with  $n$ -degrees of freedom.

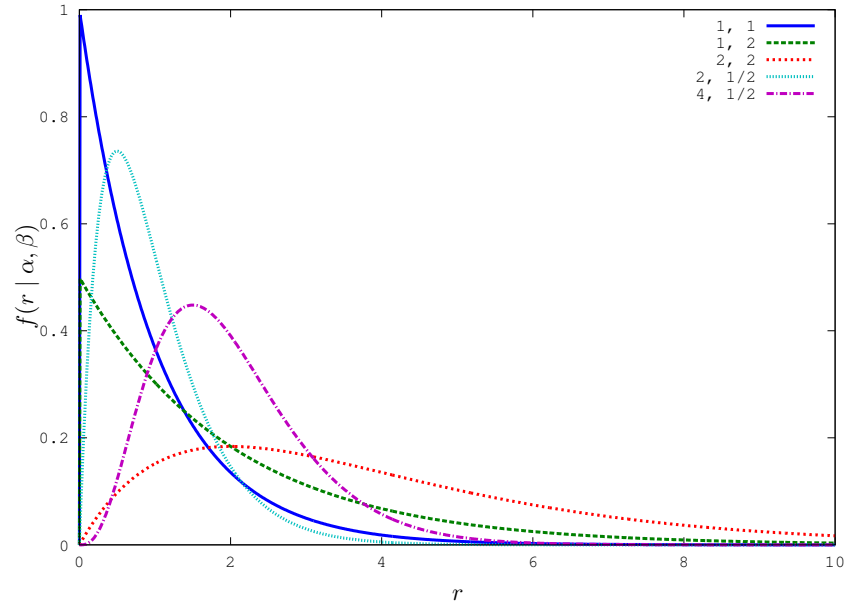


Figure 5.14: Example Gamma densities

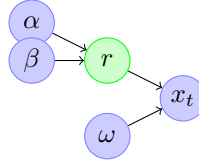


Figure 5.15: Normal with unknown precision, graphical model

### Gamma prior for the precision

**Theorem 5.7.2.** *Given a sample  $x^n$  from a normal distribution with mean  $\omega$  and unknown precision with Gamma prior  $\xi_0(r) \triangleq f(r \mid \alpha, \beta)$ , the posterior distribution is*

$$\xi_n(r) \triangleq \xi_0(r \mid x^n) = f(r \mid \alpha', \beta'), \quad (5.7.2)$$

where  $\alpha' = \alpha + \frac{n}{2}$ ,  $\beta' = \beta + \frac{1}{2} \sum_{i=1}^n (x_i - \omega)^2$ .

### Normals with unknown precision and unknown mean

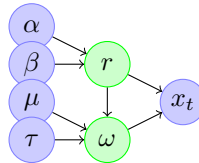


Figure 5.16: Normal with unknown mean and precision, graphical model

**Theorem 5.7.3.** *Given a sample  $x^n$  from a normal distribution with unknown mean  $\omega$  and precision  $r$ , whose prior joint distribution satisfies*

$$\omega \mid r \sim \mathcal{N}(\mu, \tau r), \quad r \sim \mathcal{Gam}(\alpha, \beta), \quad (5.7.3)$$

*the posterior distribution is*

$$\omega \mid r \sim \mathcal{N}\left(\frac{\tau\mu + n\bar{x}}{\tau + n}, (\tau + n)r\right), \quad r \sim \mathcal{Gam}\left(\alpha + \frac{n}{2}, \beta + \frac{1}{2} \sum_{i=1}^n (x_i - \bar{x})^2 + \frac{\tau n(\bar{x} - \mu)^2}{2(\tau + n)}\right). \quad (5.7.4)$$

#### *Interesting properties*

1. While  $\omega \mid r$  has normal distribution, the marginal distribution of  $\omega$  is not normal. In fact, it can be shown that it has a  $t$ -distribution.
2. There is a dependence of  $\omega$  on  $r$ . This is true for our posteriors, even if  $\omega, r$  are independent in our prior.

#### **The marginal predictive distribution**

For a normal distribution with mean  $m$ , precision  $r$ , we have

$$f(x \mid m, r) \propto r^{1/2} \exp\left(-\frac{r}{2}(m - x)^2\right).$$

For a prior  $M \mid R = r \sim \mathcal{N}(\mu, \nu r)$  and  $R \sim \mathcal{Gam}(\alpha, \beta)$ , we have the following joint distribution for the mean and precision:

$$\xi(m, r) \propto r^{1/2} e^{-(\nu r/2)(m - \mu)^2} r^{\alpha-1} e^{-\beta r}. \quad (5.7.5)$$

Now we can write the posterior marginal as

$$\xi(x) = \int f(x \mid m, r) d\xi(m, r) \quad (5.7.6)$$

$$\propto \int r e^{-\frac{r}{2}(m-x)^2} e^{-(\nu r/2)(m-\mu)^2} r^{\alpha-1} e^{-\beta r} d(m, r) \quad (5.7.7)$$

$$= \int r^\alpha e^{-\beta r} \int e^{-\frac{r}{2}(m-x)^2 - (\nu r/2)(m-\mu)^2} d(m, r) \quad (5.7.8)$$

$$= \int r^\alpha e^{-\beta r} \left( \int_{-\infty}^{\infty} e^{-\frac{r}{2}[(m-x)^2 + \nu(m-\mu)^2]} dm \right) dr \quad (5.7.9)$$

$$= \int r^\alpha e^{-\beta r} e^{-\frac{\nu r}{2(\nu+1)}(\mu-x)^2} \sqrt{\frac{2\pi}{r(1+\nu)}} dr \quad (5.7.10)$$

#### **5.7.2 Conjugates for multivariate distributions**

The binomial distribution can be extended to the multinomial and the normal distribution on the real line can be extended to the multivariate normal distribution. Fortunately, multivariate extensions exist for their corresponding conjugate priors as well.

## Multinomial-Dirichlet conjugates

### Multinomial distribution

The multinomial distribution is the extension of the binomial distribution to more than an arbitrary number of outcomes. Consider an outcome set  $S = \{1, \dots, K\}$ . This is a common model for independent random trials with a finite number of possible outcomes, such as repeated dice throws, multi-class classification problems, etc.

We now perform  $n$  trials, such that the outcome of each trial is independent of the rest. This is the an extension of a sequence of  $n$  Bernoulli trials, but with a larger set of possible outcomes in each trial.

The *multinomial* distribution gives us the probability of obtaining the  $i$ -th outcome  $n_i$  times, given that we perform a total of  $n$  trials.

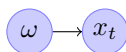


Figure 5.17: Multinomial graphical model

- Discrete distribution, outcomes  $x_t \in S = \{1, \dots, K\}$ .
- Parameter  $\omega \in \mathbb{R}_+$ ,  $\|\omega\|_1 = 1$ . Due to the second constraint, given  $\omega_1, \dots, \omega_{K-1}$ , the value of  $\omega_K$  is fully determined. However, we shall not make use of that fact.
- This is a distribution with i.i.d outcomes:  $\mathbb{P}(x_t = i \mid \omega) = \omega_i$  for all  $i$ . Note that this is the multiple outcome extension of the Bernoulli distribution, which only has two possible outcomes.
- Let us denote the number of times the  $i$ -th outcome was observed until time  $t$  by  $n_{t,i} \triangleq \sum_{k=1}^t \mathbb{I}\{x_k = i\}$ . Then:

$$\mathbb{P}(n_t \mid \omega) = \frac{t!}{\prod_{i=1}^K n_{t,i}!} \prod_{i=1}^K \omega_i^{n_{t,i}}. \quad (5.7.11)$$

### Dirichlet distribution

The Dirichlet distribution is the multivariate extension of the Beta distribution.

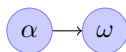


Figure 5.18: Dirichlet graphical model

The Dirichlet distribution is a distribution on the interval  $[0, 1]^K$ . It has a vector parameter that determines the density of the observations.

- Continuous distribution, outcomes  $\omega \in \Omega = \Delta^K$ , i.e.  $\|\omega\|_1 = 1$  and  $\omega_i \geq 0$ . In other words, all of the mass is on the positive  $K-1$  dimensional simplex in  $\mathbb{R}^K$ .

- Parameter vector  $\alpha \in \mathbb{R}_+^K$ .
- Probability density function:

$$f(\omega \mid \alpha) = \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod \omega_i^{\alpha_i - 1}, \quad (5.7.12)$$

$$\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du.$$

#### Dirichlet prior for multinomial distributions



Figure 5.19: Dirichlet-multinomial graphical model.

- We assume  $\omega$  is unknown, but fixed.
- We observe  $x^t = (x_1, \dots, x_t)$ .
- Our prior is determined by  $\mathcal{Dir}(\alpha)$ :  $\xi_0(\omega) \triangleq f(\omega \mid \alpha)$
- Our posterior is

$$\xi_t(\omega) \propto \prod_{i=1}^K \omega_i^{n_{t,i} + \alpha_i - 1} \quad (5.7.13)$$

$$\text{where } n_{t,i} = \sum_{k=1}^t \mathbb{I}\{x_k = i\}.$$

#### Multivariate normal conjugate families

##### Multivariate normal distribution

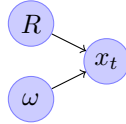


Figure 5.20: Multivariate normal graphical model

- Continuous distribution, outcomes:  $S = \mathbb{R}^K$ .
- Parameters: mean  $\omega \in \mathbb{R}^K$ , precision  $\mathbf{R} \in \mathbb{R}^{K \times K}$ , with  $x' \mathbf{R} x > 0$  for any  $x \neq 0$ .
- Probability density function:

$$f(x_t \mid w, r) = (2\pi)^{-K/2} |\mathbf{R}|^{1/2} \exp\left(-\frac{1}{2}(x_t - \omega)' \mathbf{R} (x_t - \omega)\right). \quad (5.7.14)$$

- Independence:  $f(x^t \mid \omega, \mathbf{R}) = \prod_{k=1}^t f(x_k \mid \omega, \mathbf{R})$ .

**Wishart distribution**

The Wishart distribution is a *matrix distribution* on  $\mathbb{R}^{K \times K}$  with  $n$  degrees of freedom and precision matrix  $\mathbf{T} \in \mathbb{R}^{K \times K}$ .

**Construction of the Wishart distribution**

1.  $x^n \sim \mathcal{N}(\omega, \mathbf{T})$ ,  $\omega \in \mathbb{R}^K$ ,  $\mathbf{R} \in \mathbb{R}^{K, K}$ .
2.  $\bar{x}_t \triangleq \frac{1}{n} \sum_{i=1}^n x_i$ .
3.  $\mathbf{S} = \sum_{i=1}^n (x_i - \bar{x}_t)(x_i - \bar{x}_t)'$ .

Then  $\mathbf{S}$  has a Wishart distribution with  $n - 1$  degrees of freedom and parameter matrix  $\mathbf{T}$ .

**Probability density function**

For any  $\mathbf{V} \in \mathbb{R}^{K \times K}$ , with  $\mathbf{V} > 0$ .

$$f(\mathbf{V} \mid n, \mathbf{T}) \propto |\mathbf{T}|^{n/2} |\mathbf{V}|^{(n-K-1)/2} e^{-\frac{1}{2} \text{trace}(\mathbf{T}\mathbf{V})}. \quad (5.7.15)$$

is the density for a Wishart distribution with  $n > K - 1$  degrees of freedom and precision matrix  $\mathbf{T}$ .

**Definition 5.7.1.** *The trace of a  $n \times n$  square matrix  $A$  is*

$$\text{trace}(A) \triangleq \sum_{i=1}^n a_{ii}.$$

**Normal-Wishart conjugate prior**

**Theorem 5.7.4.** *Given a sample  $x^n$  from a multivariate normal distribution in  $\mathbb{R}^K$  with unknown mean  $\omega \in \mathbb{R}^K$  and precision  $\mathbf{R} \in \mathbb{R}^{K \times K}$ , whose prior joint distribution satisfies:*

$$\omega \mid \mathbf{R} \sim \mathcal{N}(\mu, \tau \mathbf{R}), \quad \mathbf{R} \sim \text{Wish}(\alpha, \mathbf{T}), \quad (5.7.16)$$

with  $\tau > 0$ ,  $\alpha > k - 1$ ,  $\mathbf{T} > 0$ , the posterior distribution is

$$\omega \mid \mathbf{R} \sim \mathcal{N}\left(\frac{\tau \mu + n \bar{x}}{\tau + n}, (\tau + n) \mathbf{R}\right), \quad (5.7.17)$$

$$\mathbf{R} \sim \text{Wish}\left(\alpha + n, \mathbf{T} + \mathbf{S} + \frac{\tau n}{\tau + n} (\mu - \bar{x})(\mu - \bar{x})'\right), \quad (5.7.18)$$

$$\mathbf{S} = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'.$$

**5.8 Summary**

- Family of observation distributions:  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$  on  $\mathcal{S}$ .



- Family of parameter priors  $\Xi$  on  $\Omega$ .
- A *statistic*  $T$  is a function giving a *summary* of observations in  $\mathcal{S}$ .
- A *sufficient* statistic  $T$  summarises *all* relevant information for distinguishing between distributions in a family.
- $\Xi$  is *conjugate* to  $\mathcal{P}$  if  $\xi(\cdot | x) \in \Xi$  for all  $\xi \in \Xi, x \in \mathcal{S}$ .
- A credible interval/set  $A$  has a measure  $\xi(A | x)$ , which represents our belief that the true parameter  $\omega$  is in  $A$ , given observations  $x \in \mathcal{S}$  and our prior belief  $\xi \in \Xi$ .
- A confidence interval of measure  $\xi(A | x) = s$  fails with probability  $1 - s$  if  $\omega \sim \xi$ .
- The beta distribution is good for modelling parameters / observations in  $[0, 1]$ .
- The gamma distribution is good for modelling parameters / observations in  $[0, \infty)$ .
- The Dirichlet and Wishart distributions are their multivariate extensions.
- Parametric conjugate pairs: Binomial/Beta, Normal/Normal-Gamma, Multinomial/Dirichlet, Multivariate normal / Normal-Wishart, Uniform/Pareto.



## Chapter 6

# Hypothesis testing

## 6.1 Decision problems

### Hypothesis testing as a decision problem.

Hypothesis testing is a special type of decision problem, where the decision space  $\mathcal{D}$  is a set of hypotheses about the distribution that generates the observations.

#### Observations

Consider a problem where we obtain an observation  $x \in \mathcal{S}$ .  $\mathcal{S}$  can be a finite product space:  $\mathcal{S} = \mathcal{Z}^n$ , or it can be the space of all sequences of some observations  $z_t \in \mathcal{Z}$ :  $\mathcal{S} = \mathcal{Z}^* \triangleq \bigcup_{k=1}^{\infty} \mathcal{Z}^k$ .

#### Models

We have a set of probability measures on  $\mathcal{S}$ ,  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$ , indexed by  $\omega$ . We wish to choose from a set of different hypotheses about  $\omega$ . The measures  $P_\omega$  do not necessarily have to be within the same parametric family. They must however define measures on  $\mathcal{S}$ .

#### Decision space $\mathcal{D}$

Each decision  $d \in \mathcal{D}$  corresponds to a *hypothesis* about which is the correct  $\omega$ . More specifically,  $d \in \mathcal{D}$  corresponds to the hypothesis that  $\omega \in \Omega_d$ , with  $\Omega_d \subset \Omega$  for all  $d$ .

#### Loss function $\ell : \Omega \times \mathcal{D} \rightarrow \mathbb{R}$

If  $\omega^*$  is the true parameter, then the loss of decision  $d$  is  $\ell(\omega^*, d)$ .

### Decision function.

**Definition 6.1.1.** A decision function  $\delta : \mathcal{S} \rightarrow \mathcal{D}$  selects a decision  $\delta(x)$  for any observation  $x \in \mathcal{S}$ .

#### The expected loss of a decision function for a given parameter

If the true parameter is  $\omega^*$ , the expected loss of a decision function is

$$\sigma(\omega^*, \delta) \triangleq \mathbb{E}_{\omega^*}[\ell(\omega^*, \delta)] = \int_{\mathcal{S}} \ell(\omega, \delta(x)) dP_{\omega^*}(x). \quad (6.1.1)$$

This is the risk of the decision function for the case when the value of the parameter is  $\omega^* \in \Omega$ . Since  $\omega^*$  is unknown, the risk is also unknown. However, we can calculate the risk functional for all  $\omega$ .

**The (subjective) expected loss (or risk) of a decision function**  
 For *any* probability measure  $\mu$  on  $\Omega$ , our expected loss is:

$$\sigma(\mu, \delta) = \int_{\Omega} \sigma(\omega, \delta) d\mu(\omega) = \int_{\Omega} \int_{\mathcal{S}} [\ell(\omega, \delta(x))] d\mu(\omega). \quad (6.1.2)$$

*This holds both when  $\mu = \psi$ , the distribution from which  $\omega$  is actually drawn, or when  $\mu = \xi$ , our prior belief.* Thus, the meaning of the above equation depends on what  $\xi$  is. If  $\omega$  is drawn from  $\xi$  prior to our experiment, then (6.1.2) coincides with the true risk. If  $\xi$  is our subjective belief about  $\omega$ , but  $\omega$  is drawn from some other distribution  $\psi$  prior to our experiment, then the above expression is only subjective, since generally  $\rho(\xi, \delta) \neq \rho(\psi, \delta)$  if  $\psi \neq \xi$ .

### Decision problems with two points

To make the above discussion more concrete, consider the following simple problem, where  $\Omega = \{\omega_1, \omega_2\}$  and where we have two decisions, i.e.  $D = \{d_1, d_2\}$ . In *simple hypothesis testing* problems, each decision  $d_i$  corresponds to selecting  $\omega_i$  as the true parameter and we only suffer a loss when we have made the wrong choice. The loss function is shown in Table 6.1.

$\ell(\omega, d)$	$d_1$	$d_2$
$\omega_1$	0	$c_1$
$\omega_2$	$c_2$	0

Table 6.1: Cost function of a simple hypothesis testing problem

Consider that we want to decide whether  $\omega_1$  or  $\omega_2$  is true. We can do so after observing  $x$ . We construct a decision function  $\delta$ , such that for any observed value, we choose  $d_1$ , or  $d_2$ . For any selected decision function, there will be some probability that we make the wrong decision. Thus, the decision function will have an expected cost associated with it. In fact, let  $\alpha_1(\delta), \alpha_2(\delta)$  be the probabilities of making the wrong decision using  $\delta$  for the two cases,  $\omega^* = \omega_1$  and  $\omega^* = \omega_2$ . Given a probability for  $\omega$ , we can calculate the risk of our decision function. More specifically:

- Let  $\mu(\omega)$  be some (either true or subjective) prior on  $\omega$ .
- For each  $\omega \in \Omega$ , we define  $\mathbb{P}(x | \omega)$ .
- We observe  $x$  and then choose decision  $\delta(x)$ .
- Let  $\alpha_1(\delta)$  be the conditional probability that we choose  $d_2$  when  $\omega^* = \omega_1$ .

$$\alpha_1(\delta) \triangleq \mathbb{P}(\delta(x) \neq d_1 | \omega_1)$$

- Let  $\alpha_2(\delta)$  be the conditional probability that we choose  $d_1$  when  $\omega^* = \omega_2$ .

$$\alpha_2(\delta) \triangleq \mathbb{P}(\delta(x) \neq d_2 | \omega_2)$$

- Let  $b_1 \triangleq c_1\mu(\omega_1)$  and  $b_2 \triangleq c_2\mu(\omega_2)$ .

The risk of  $\delta$  is:

$$\rho(\mu, \delta) = b_1\alpha_1(\delta) + b_2\alpha_2(\delta) \quad (6.1.3)$$

**Theorem 6.1.1** (Neymann-Pearson lemma). *For any  $b_1, b_2 > 0$ , let  $\delta^*$  be a decision function such that*

$$\delta^*(x) = d_1, \quad \text{if } b_1P_{\omega_1}(x) > b_2P_{\omega_2}(x) \quad (6.1.4)$$

$$\delta^*(x) = d_2, \quad \text{if } b_1P_{\omega_1}(x) < b_2P_{\omega_2}(x), \quad (6.1.5)$$

and either of  $d_1, d_2$  otherwise. Then, for any other  $\delta$ :

$$\sigma(\mu, \delta^*) = b_1\alpha_1(\delta^*) + b_2\alpha_2(\delta^*) \leq b_1\alpha_1(\delta) + b_2\alpha_2(\delta) = \sigma(\mu, \delta)$$

#### **Interpretation.**

- When  $\mu = \psi$ , the distribution from which  $\omega^*$  is drawn,  $\rho(\mu, \delta)$  is the *actual* risk.
- When  $\mu = \xi$ , our subjective prior belief about  $\omega^*$ , then  $\rho(\mu, \delta)$  is the *subjective* risk.

### **Bayesian decisions**

#### **Construction of the decision function**

We saw in chapter 3 (Decision problems) that the optimal decision function is the one that minimises the risk relative to the current posterior. So, we can have the following procedure.

1. Given a family  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$ , decision space  $\mathcal{D}$ , loss function  $\ell : \Omega \times \mathcal{D} \rightarrow \mathbb{R}$ .
2. Select a prior  $\xi(\omega)$ .
3. Observe  $x$ , and obtain posterior  $\xi'(\omega) \triangleq \xi(\omega \mid x) \propto P_\omega(x)\xi(\omega)$ .
4. Select  $d^* = \arg \min_d \sigma(\xi(\omega), d)$ .

Thus, the optimal decision function is

$$\delta^*(x) = \arg \min_d \sigma(\xi(\omega \mid x), d).$$

#### **Open problems**

- How do we select  $\xi$ ? Ideally we would like to be able to select  $\xi$  in an objective way. However, this is usually not possible. What is done instead, is to select a  $\xi$  with certain symmetry or invariance properties, to select a  $\xi$  which is maximally pessimistic in some sense, or to actually attempt to elicit a prior  $\xi$  from experts. However, strict frequentists would say that a prior does not even exist in a single experiment.
- What is the true risk? This is a big question. The true risk depends on how fast the posterior distribution converges to the true value and how far it is from it initially. However, in general the true risk is close to the Bayes risk, also because of the convexity property of the Bayes risk.

## 6.2 Unconditional and conditional errors

### Error probabilities

Recall that the risk of decision rule  $\delta$  under  $\mu$  for the two-point problem is

$$\sigma(\mu, \delta) = b_1 \alpha_1(\delta) + b_2 \alpha_2(\delta),$$

where

$$\alpha_i(\delta) \triangleq \mathbb{P}(\delta(x) \neq d_i \mid \omega_i), \quad b_i = c_i \mu(\omega_i). \quad (6.2.1)$$

**Definition 6.2.1** (Unconditional error). *When  $c_1 = c_2 = 1$ , then  $\rho(\mu, \delta)$  is the unconditional error of decision rule  $\delta$ . In other words,*

$$\rho(\mu, \delta) = \sum_i \mu(\omega_i) \mathbb{P}(\delta(x) \neq d_i \mid \omega_i)$$

*is the a priori probability that the rule  $\delta$  will make the wrong decision, given that  $\omega$  is drawn from  $\mu$ . Once more, the above can be interpreted as a subjective quantity, if  $\mu = \xi$ , our belief about  $\omega$ , or as an actual quantity, if  $\mu = \psi$ , the true distribution of  $\omega$ .*

### Classical decisions

These types of decisions are sometimes called frequentist. However, the decisions we shall talk about here are optimal in a worst-case sense.

#### Construction of the decision function

1. Given a family  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$ , decision space  $\mathcal{D}$ , loss function  $\ell : \Omega \times \mathcal{D} \rightarrow \mathbb{R}$ .
2. We have no information on  $\psi(\omega)$ , not even a subjective  $\xi(\omega)$ .

3. Solution: minimise worst-case risk. For any  $\mu$ :

$$\sigma(\mu, \delta) = \sum_i \mu(\omega_i) c_i \mathbb{P}(\delta(x) \neq d_i \mid \omega_i). \quad (6.2.2)$$

$$\leq \sum_i \mu(\omega_i) \max_j c_j \mathbb{P}(\delta(x) \neq d_j \mid \omega_j). \quad (6.2.3)$$

$$= \max_i c_i \mathbb{P}(\delta(x) \neq d_i \mid \omega_i). \quad (6.2.4)$$

*The most powerful test for two-point problems*

- Minimise bound by equating  $c_1 \mathbb{P}(\delta(x) \neq d_1 \mid \omega_1) = c_2 \mathbb{P}(\delta(x) \neq d_2 \mid \omega_2)$ .
- Find  $X \subset \mathcal{S}$  such that  $c_1 P_{\omega_1}(S) = c_2 P_{\omega_2}(\mathcal{S} \setminus X)$ .
- $\delta(x) = d_2$  if  $x \in X$  and  $d_1$  otherwise.
- Then  $\mathbb{P}(\delta(x) \neq d_1 \mid \omega_1) = P_{\omega_1}(S)$ .

### Conditional error

Assume we observe  $x \sim w$  and let  $\mu(\omega)$  be a prior distribution on  $\Omega$ . Let us define

$$\mu(\omega \mid x) = P_\omega(x) \mu(\omega) / \mu(x),$$

to be the posterior probability of  $\omega$ .

**Definition 6.2.2** (Conditional error). *Given the above,  $\rho(\mu(\cdot \mid x), \delta)$  is the conditional error of decision rule  $\delta$  if we have observed  $x$ . In other words,*

$$\rho(\mu(\cdot \mid x), \delta(x)) = \sum_i \mu(\omega_i \mid x) \mathbb{P}(\delta(x) \neq d_i \mid \omega_i)$$

*is the posterior probability that the rule  $\delta$  will make the wrong decision, given that  $\omega$  is drawn from  $\mu$ .*

Once more, the above can be interpreted as a *subjective* quantity, if  $\mu = \xi$ , our belief about  $\omega$ , or as an *actual* quantity, if  $\mu = \psi$ , the true distribution of  $\omega$ .

### Desired properties of reported errors

Assume that we have a procedure which:

1. Given  $\mu, \delta$  and *before* seeing  $x$ , gives us an unconditional error estimate  $\alpha$ .
2. Given  $\mu, \delta$  and *after* seeing  $x$ , gives us a conditional error estimate  $\alpha(x)$ .



**Desirable properties**

$$\alpha = \sigma(\psi, \delta) \quad (6.2.5)$$

$$\alpha(x) = \sigma(\psi(\cdot | x), \delta(x)). \quad (6.2.6)$$

**6.3 Two-point test with Bernoulli trials****Example 6.3.1.** *Bernoulli trials*

1. Take a sample  $x = x_1, \dots, x_n$  from  $\mathbb{P}(x_k | w)$ .
2. Let  $s_n = \sum_{k=1}^n x_k$ .
3. From the binomial formula,

$$\mathbb{P}(s_t = k | w) = \binom{t}{k} w^k (1-w)^{t-k}, \quad \binom{x}{k} = \prod_{i=0}^{k-1} (x-i)/(1+i). \quad (6.3.1)$$

Recall that  $\binom{x}{k} = \prod_{i=0}^{k-1} (x-i)/(1+i)$ .

**The problem**

- We need to decide whether  $d_1 : \omega = \omega_1$ , or  $d_2 : \omega = \omega_2$ .
- Assume that the costs of mistakes are  $c_1, c_2 = 1$  for simplicity.
- Use rule  $\delta$  to select  $\delta(x)$  after you observe  $x$ .

**Bayesian Bernoulli trial test**

- Select a prior parameter  $\pi \in [0, 1]$ , so that  $\xi(\omega_1) = \pi = 1 - \xi(\omega_2)$ .
- Observe  $x$  and calculate:

$$\xi(\omega_1 | x) = \frac{\pi \mathbb{P}(x | \omega_1)}{\pi \mathbb{P}(x | \omega_1) + (1 - \pi) \mathbb{P}(x | \omega_2)}. \quad (6.3.2)$$

- Construct decision function

$$\delta(x) = \begin{cases} d_1, & \xi(\omega_1 | x) > \xi(\omega_2 | x) \\ d_2, & \xi(\omega_1 | x) < \xi(\omega_2 | x) \\ d_{1 \text{ w.p. } 1/2}, & \xi(\omega_1 | x) = \xi(\omega_2 | x) \end{cases}$$

In other words, this decision function contains a no-decision region of measure 0, corresponding to the case where both posterior probabilities are equal.

#### Properties

- The procedure has unconditional error  $O(e^{-t})$ .
- The conditional error is  $\max_k \xi(\omega_k) + O(e^{-t})$ .

#### Classical Bernoulli trial test

- Select error rates  $\alpha_1, \alpha_2$ .
- These can depend on the sample size: e.g.  $\alpha_1(t) = b_1 t^{-1/2}$ ,  $\alpha_2(t) = b_2 t^{-1/2}$ .
- Select a subset  $S \subset \mathcal{S}$  such that  $P_{\omega_1}(S) = \alpha_1$  and  $P_{\omega_2}(\mathcal{S} \setminus S) = \alpha_2$ , if possible.
- Observe  $x$  and use decision function

$$\delta(x) = \begin{cases} d_1, & x \in S \\ d_2, & x \notin S. \end{cases}$$

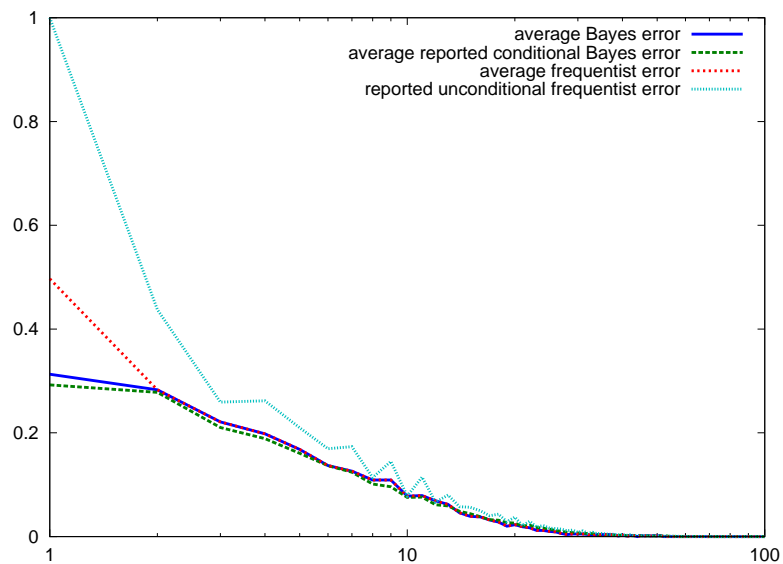
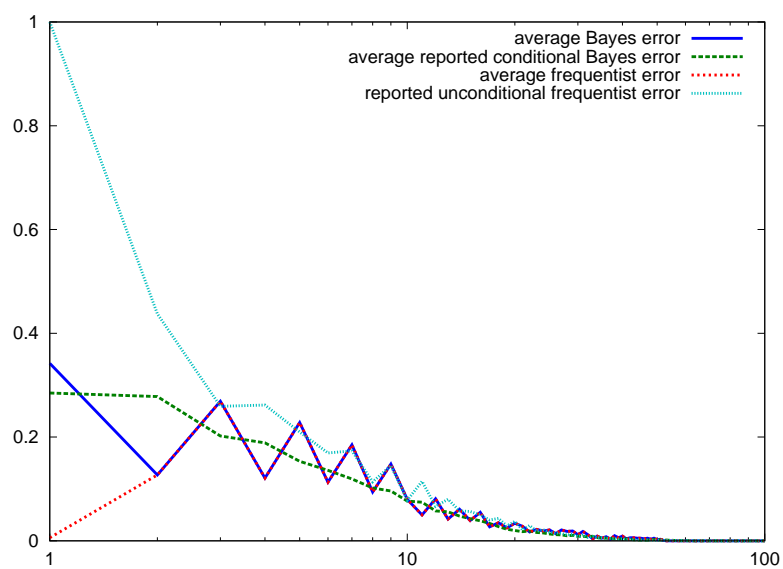
To construct  $S$ , traditionally  $\alpha_1 = \alpha_2$  and we choose a *critical value*  $c$  such that  $\mathbb{P}(s > c \mid \omega_1) = 1 - \mathbb{P}(s < c \mid \omega_2)$ , with  $s = \sum_k x_k$  being the *test statistic*. Note that  $c \approx \frac{t}{2}(\omega_1 + \omega_2)$ . However, any arbitrary  $S$  can be chosen.

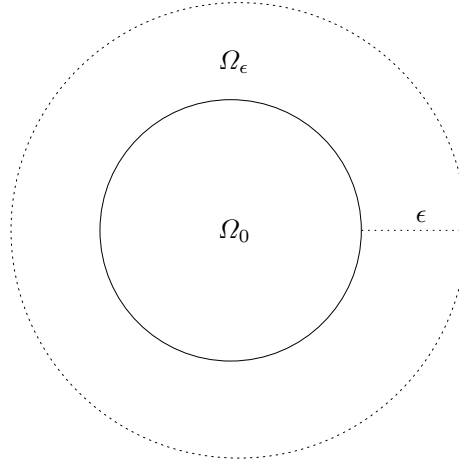
#### Properties

- The unconditional error is  $\max\{\alpha_1, \alpha_2\}$ .
- The conditional error is  $\alpha_{\delta(x)}$ . *Thus, it only depends on whether  $x \in S$ .*

#### Summary and additional remarks

- The decision rule is chosen a priori.
- The unconditional error/risk does not depend on the data.
- The conditional error/risk depends on the data.
- Bayesian procedures minimise *expected loss*.
- Classical procedures minimise *worst-case expected loss*.
- The Bayesian procedure should be consistent with classical ones.
- By selecting a *worst-case prior* we can create a Bayesian analogue of classical procedures.

Figure 6.1:  $\psi(\omega_1) = 1/2$ ,  $\omega_1 = 1/3$ ,  $\omega_2 = 3/4$ ,  $\xi(\omega_1) = 1/2$ .Figure 6.2:  $\psi(\omega_1) = 99/100$ ,  $\omega_1 = 1/3$ ,  $\omega_2 = 3/4$ ,  $\xi(\omega_1) = 1/2$ .



## 6.4 Null hypothesis tests

### The problem of null hypothesis tests

- Very frequently we wish to decide whether something is true or not.
- However, we have *no specific alternative hypothesis*.
- Consider for example we want to choose whether:  $d_0 : \omega^* \in \Omega_0$ , or  $d_1 : \omega^* \notin \Omega_0$ .
- Most of the time we can only decide  $d_0 : \omega^* \in \Omega_0$ , or  $d_1 : \omega^* \notin \Omega_\epsilon$ , with  $\Omega_\epsilon \supset \Omega_0$ .
- In the simplest case the null hypothesis is a single point:  $\Omega_0 = \{\omega_0\}$ .

### Construction of the null hypothesis test

Consider the following problem. We have a distribution in some family  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$  with unknown parameter  $\omega$ . We have two hypotheses:  $\omega = \omega_0$  and  $\omega \neq \omega_0$ . After obtaining a sample  $x = z_1, \dots, z_n$  from  $P_\omega$ , we need to decide for  $d_0 \in \mathcal{D} : \omega = \omega_0$  or  $d_1 : \omega \neq \omega_0$ . The loss of a wrong decision is 1.

#### A simple decision rule

Let us choose some  $S \subset \mathcal{S}$ , such that  $P_{\omega_0}(S) = \alpha$ . We can now define the following decision rule:

$$\delta(x) = \begin{cases} d_0, & \text{if } x \notin S, \\ d_1, & \text{if } x \in S. \end{cases}$$

*The unconditional error*

This rule has unconditional error  $\alpha$  when  $\omega = \omega_0$ . In other words, the

probability that we will make a mistake when employing this rule, and  $\omega = \omega_0$ , is  $\alpha$ . However, there is no guarantee for the error when  $\omega \neq \omega_0$ , without additional assumptions. In particular, we need to specify  $\mu$ , and  $\mathbb{P}(x \in S \mid \omega \neq \omega_0)$

$$\sigma(\mu, \delta) = \mu(\omega = \omega_0)P_{\omega_0}(S) + \mu(\omega \neq \omega_0)\mathbb{P}(x \in S \mid \omega \neq \omega_0), \quad (6.4.1)$$

where  $\mu$  is either a subjective distribution  $\xi$  or the actual distribution  $\psi$ .

**Bayesian approach**

The Bayesian approach once more involves us having a subjective belief about the nature of the problem. If, of course, we somehow know that  $\omega$  is selected according to some  $\psi$ , then we can set  $\xi = \psi$ , and we will be optimal in expectation.

**(Subjective) prior**

$$\pi_0 \triangleq \xi(\omega = \omega_0) \quad (6.4.2)$$

$$\phi(Z) \triangleq \xi(\omega \in Z \mid \omega \neq \omega_0), \quad Z \subset \Omega. \quad (6.4.3)$$

**(Subjective) posterior**

$$\pi_0(x) \triangleq \xi(\omega = \omega_0 \mid x) = \frac{\pi_0 P_\omega(x)}{\pi_0 P_\omega(x) + (1 - \pi_0)\phi(x)} \quad (6.4.4)$$

$$\phi(x) = \int_{\Omega} P_\omega(x) d\phi(\omega). \quad (6.4.5)$$

**Example 6.4.1.** Consider Bernoulli trials where we want to test whether  $\omega^* = \omega_0$ , or  $\omega^* \neq \omega_0$ . We can use a Beta density with parameters  $\beta_0, \beta_1$  for the prior distribution in the case of  $\omega^* \neq \omega_0$ :

$$\phi(Z) = \int_Z f(\omega \mid \beta),$$

with

$$f(\omega \mid \beta) = \frac{\Gamma(\beta_0 + \beta_1)}{\Gamma(\beta_0)\Gamma(\beta_1)} \omega^{\beta_1-1} (1 - \omega)^{\beta_0-1}$$

and posterior distribution also being a Beta density with parameters  $\beta_0 + t - s_t, \beta_1 + s_t$ , with  $s_t = \sum_{k=1}^t x_k$ .

**Classical approach**

Since we only have one hypothesis,  $\omega = \omega_0$ , we must *construct* an alternative hypothesis against which we can measure our other error.

**Constructing the alternative hypothesis**

1. Select  $\Omega_\epsilon \subset \Omega$ , to be the  $\epsilon$ -enlargement of  $\Omega_0$ .
2. When  $\Omega_0 = \{\omega_0\}$ ,  $\Omega_\epsilon = \{\omega \in \Omega \mid \|\omega - \omega_0\| < \epsilon\}$ .

We select  $d_1 : \omega^* \notin \Omega_\epsilon$  when  $x \in S$  and  $d_0 : \omega^* = \omega_0$  when  $x \notin S$ . Then the risk is:

$$\rho(\omega_0, \delta) = P_{\omega_0}(S),$$

and

$$\rho(\Omega \setminus \Omega_\epsilon, \delta) \leq \max \{P_\omega(x) \mid \omega \in \Omega \setminus \Omega_\epsilon\}.$$

**Example 6.4.2.** Consider Bernoulli trials where we want to test whether  $\omega^* = \omega_0$ , or  $\omega^* \neq \omega_0$ .

**The test**

- Select an error rate  $\alpha_1 > 0$ .
- Set  $s_t = \sum_{k=1}^t x_k$ .
- Decision rule, with  $c_1 < t\omega_0 < c_2 \in [0, t]$ :

$$\delta(x) = \begin{cases} d_0 & \text{if } s_t \in [c_1, c_2], \\ d_1 & \text{otherwise.} \end{cases}$$

Let  $Q_\omega^t$  be the probability measure on  $\mathbb{N}$  arising from the binomial distribution with parameters  $t, w$ .

**Unconditional error**

$$\mathbb{P}(\delta(x) \neq d_0 \mid \omega^* = \omega_0) = 1 - Q_{\omega_0}^t([c_1, c_2]) \quad (6.4.6)$$

$$\mathbb{P}(\delta(x) \neq d_1 \mid |\omega^* - \omega_0| > \epsilon) \leq \max \{\mathbb{P}(s_t \in [c_1, c_2] \mid \omega^* = w) \mid |w - \omega_0| > \epsilon\} \quad (6.4.7)$$

To calculate this, we look at two cases. The first case is:

$$\mathbb{P}(s_t \in [c_1, c_2] \mid \omega^* < \omega_0 - \epsilon) \leq \max \{Q_\omega^t([c_1, c_2]) \mid \omega \in [0, \omega_0 - \epsilon]\}$$

and assuming  $\epsilon > \omega_0 - c_1/t$

$$\mathbb{P}(s_t \in [c_1, c_2] \mid \omega^* < \omega_0 - \epsilon) \leq Q_{\omega_0 - \epsilon}^t([c_1, c_2]).$$

The second case is:

$$\begin{aligned} \mathbb{P}(s_t \in [c_1, c_2] \mid \omega^* > \omega_0 + \epsilon) &\leq \max \{Q_\omega^t([c_1, c_2]) \mid \omega \in (\omega_0 + \epsilon, 1]\} \\ &= Q_{\omega_0 + \epsilon}^t([c_1, c_2]). \end{aligned}$$

where we assumed  $\epsilon > c_2/t - \omega_0$ .

Consequently, to obtain a test with good guarantees on the  $\epsilon$ -enlargement of the null hypothesis, we need  $c_1 \geq (\omega_0 - \epsilon)t$  and  $c_2 \geq (\omega_0 + \epsilon)t$ .

## 6.5 The fallacy of $P$ -values

$P$ -values are frequently misused in applications. In fact,  $P$ -values are a side-effect of constructing a statistical test in order for it to have a certain error rate  $\alpha$ . They are a *uniformly distributed* random variable when the null-hypothesis is true. Thus, by rejecting the null-hypothesis only when the  $P$ -value is smaller than  $\alpha$ , you guarantee that the type I error is exactly  $\alpha$ . Let us discuss this point further with an example.

**$P$ -values in simple two-point tests.**

**Example 6.5.1.** • *Exponential distribution with parameter  $\omega$ :  $f_\omega(x) = \omega e^{-\omega x}$ .*

- $d_1 : \omega = 1, d_2 : \omega = 2$ .

*A classical test*

We observe  $x \in \mathbb{R}_+$ . We use the decision function

$$\delta(x) = \begin{cases} d_1, & x > c \\ d_2, & x \leq c. \end{cases}$$

*These leads to the following unconditional errors*

$$\alpha_1 \triangleq \mathbb{P}(\delta(x) \neq d_1 \mid \omega = 1) = \mathbb{P}(x \leq c \mid \omega = 1) = \int_{u=0}^c e^{-u} = 1 - e^{-c} \quad (6.5.1)$$

$$\alpha_2 \triangleq \mathbb{P}(\delta(x) \neq d_2 \mid \omega = 2) = \mathbb{P}(x > c \mid \omega = 2) = \int_{u=c}^{\infty} 2e^{-2u} = \frac{1}{2}e^{-2c}. \quad (6.5.2)$$

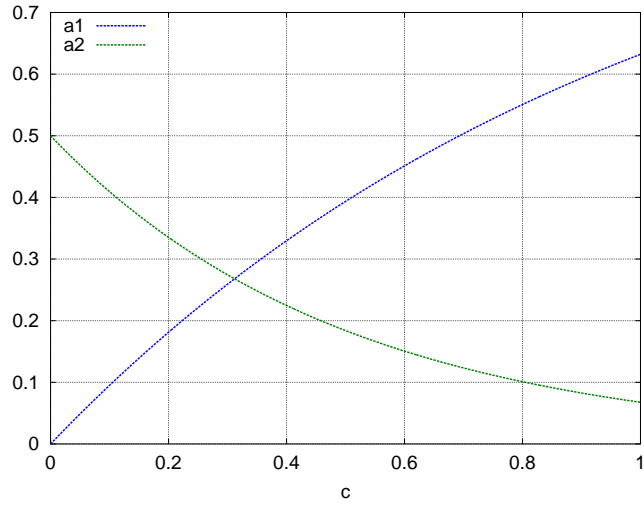


Figure 6.3: Unconditional errors for different critical values

**Selecting  $c$** 

For uniformly bounding the error, set  $\alpha_1 = \alpha_2$ :

**What is a  $P$ -value?**

**Definition 6.5.1.** Assume we want to perform a test with type I error  $\alpha_1$ . If  $p : \mathcal{S} \rightarrow [0, 1]$  is such that:

$$\mathbb{P}(p(x) \leq \alpha_1 \mid \omega_1) = \alpha_1, \quad (6.5.3)$$

then  $p(x)$  is a  $p$ -value for observations  $x$  under  $\omega_1$ .

*The uniform property of  $p$ -values.*  
 $p(x)$  is uniformly distributed in  $[0, 1]$  under  $\omega_1$

Consequently,  $p(x)$  gives *no information* about the validity of  $\omega_1$ . It also does not tell you how much the observations  $x$  might match alternative hypotheses.

**Summary**

- A test is a decision rule.
- Unconditional / conditional errors are properties of decision rules.
- Null-hypothesis tests need an  $\epsilon$ -enlargement for guarantees.
- There is always a (possibly negligible) region where no guarantees can be made.



- Classical (frequentist) methods are generally worst-case and unconditional.
- Bayesian (subjective) methods are generally expected-case and conditional.
- *P*-values are only a tool to implement classical tests and carry no information.

**Further reading**

- Bayesian methods can be extended to worst-case approaches.
- Classical methods can also report conditional errors.
- Distribution-free methods can be used for more general problems.
- Pre-validation



## Chapter 7

# Sequential sampling

## 7.1 Gains from sequential sampling

### The idea of sequential sampling

We wish to buy a large number of items in bulk. However, some portion of the items may be defective. We calculate (somehow) that if we test 100 items and 10 or more are faulty, then we should not buy. As testing is expensive, instead of testing all 100 items, we test sequentially. If, at any point, we have 10 faulty items or 91 good items, we can stop testing. This produces a decision of the same quality as testing all 100 items, but at a fraction of the cost.

### Sequential sampling

A sequential sample from some unknown distribution  $P$  is generated as follows. At time  $t$ , we have observed  $x_1, \dots, x_t$ , generated by  $P$ . These are not necessarily independently and identically distributed samples. At any time  $t$ , we can either *stop sampling* or obtain one *more* observation  $x_{t+1}$ . A sample obtained in this way is called a *sequential sample*. More formally,

**Definition 7.1.1.** *A sequential sampling procedure on a probability space  $(\mathcal{X}^*, \mathfrak{B}(\mathcal{X}^*), P)$  involves a decision function  $f : \mathcal{X}^* \rightarrow \{0, 1\}$ , such that we stop sampling at time  $t$  if and only if  $f(x^t) = 1$ , otherwise we obtain a new sample*

$$x_{t+1} \mid x^t \sim P(\cdot \mid x^t).$$

### Sampling with costs

We once more consider problems where we have some observations  $x_1, x_2, \dots$ , with  $x_t \in \mathcal{X}$ , which are drawn from some distribution with parameter  $\omega \in \Omega$ , or more precisely from a family  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$ , such that each  $(\mathcal{X}^*, \mathfrak{B}(\mathcal{X}^*), P_\omega)$  is a probability space for all  $\omega \in \Omega$ . Once more, we have a prior probability measure  $\xi$  on  $\mathfrak{B}(\Omega)$  for the unknown parameter, and we wish to take a decision  $d \in \mathcal{D}$  that maximises the expected utility according to our utility function  $U : \Omega \times \mathcal{D} \rightarrow \mathbb{R}$ .

In the classical case, we obtain a complete sample of fixed size  $n$ ,  $x^n = (x_1, \dots, x_n)$  and calculate a posterior measure  $\xi(\cdot \mid x^n)$ . We then take the decision maximising the expected utility according to our posterior.

So far, we have mainly considered decision problems where the sample size was fixed. However, frequently the sample size can also be part of the decision. Since normally larger sample sizes give us more information, in this case the decision problem is only meaningful when obtaining new samples has a cost.

#### Fixed sample size

Given  $x^n$ , find  $d \in \mathcal{D}$  maximising:

$$\mathbb{E}_\xi(U \mid d, x^n) = \int_\Omega U(w, d) d\xi(w \mid x^n). \quad (7.1.1)$$

**Samples with costs**

Consider now the case where each sample taken incurs a cost  $c$ . Then, our utility is

$$\mathbb{E}_\xi(U \mid d, x^n) = \int_{\Omega} U(w, d) d\xi(\omega \mid x^n) - cn. \quad (7.1.2)$$

**Example 7.1.1.** Consider the following decision problem.

- *Parameters:*  $\Omega = \{\omega_1, \omega_2\}$ .
- *Decisions:*  $D = \{d_1, d_2\}$ .
- *Observation distribution*  $f_i(k) = \mathbb{P}(x_t = k \mid \omega = \omega_i)$  for all  $t$  with

$$f_1(1) = 1 - \alpha, \quad f_1(2) = 0, \quad f_1(3) = \alpha, \quad (7.1.3)$$

$$f_2(1) = 0, \quad f_2(2) = 1 - \alpha, \quad f_2(3) = \alpha. \quad (7.1.4)$$

- *Utility:*  $U(\omega_i, d_j) = 0$ , for  $i = j$  and  $b < 0$  otherwise.
- *Prior:*  $\mathbb{P}(\omega = \omega_1) = \xi = 1 - \mathbb{P}(\omega = \omega_2)$ .

In this problem, it is immediately possible to distinguish  $f_1$  from  $f_2$  when you observe  $x_t = 1$  or  $x_t = 2$ . However, the values  $x_t = 3$  provide no information.

- So, the expected utility of stopping if you have only observed 3s is  $\xi b$ .
- Otherwise, it equals 0.

**7.1.1 An example sequential problem****A procedure taking  $n$  observations**

- The probability of observing  $x_t = 3$  for all  $t = 1, \dots, n$  is  $\alpha^n$ .
- The total value  $V(n)$  of the optimal procedure taking  $n$  observations is

$$V(n) = \xi b \alpha^n - cn. \quad (7.1.5)$$

We can find the optimal value more or less easily. Since  $V$  is a nice function, an approximate minimiser can be found by viewing  $n$  as a continuous variable. Taking derivatives:

$$n^* = \left\lceil \log \frac{c}{\xi b \log \alpha} \right\rceil \frac{1}{\log \alpha} \quad (7.1.6)$$

$$V(n^*) = \frac{c}{\log \alpha} \left[ 1 + \log \frac{c}{\xi b \log \alpha} \right] \quad (7.1.7)$$

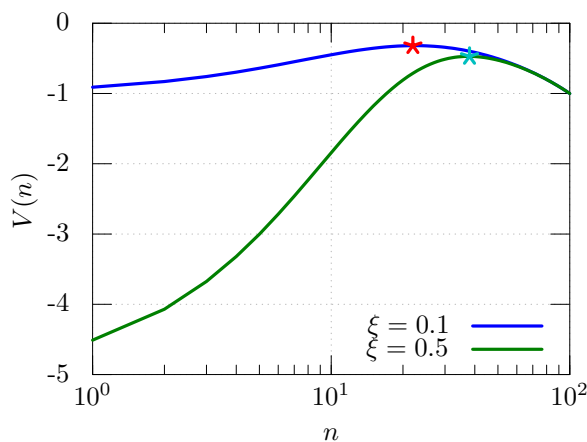


Figure 7.1: The value of taking exactly  $n$  observations under two different beliefs, for  $\alpha = 0.9$ ,  $b = -10$ ,  $c = 10^{-2}$ .

**Definition 7.1.2** (The geometric series). *The sum  $\sum_{k=0}^n x^k$  is called the geometric series and has the property*

$$\sum_{k=0}^n x^k = \frac{x^{n+1} - 1}{x - 1}. \quad (7.1.8)$$

*Taking derivatives with respect to  $x$  can result in other useful formulae.*

**A sequential procedure stopping after at most  $n^*$  steps.**

- If  $t < n^*$ , use the stopping rule  $s(x^t) = 1$  iff  $x_t = 3$ .
- In other words, stop as soon as you observe a 3, or until you reach  $t = n^*$ .
- Our posterior after stopping is, in this case, the same as in the standard procedure.
- However the number of observations  $n$  is random.

Since the probability of  $x_t = 3$  is always the same for both  $\omega_1$  and  $\omega_2$ , we have:

$$\mathbb{E}(n) = \mathbb{E}(n \mid \omega = \omega_1) = \mathbb{E}(n \mid \omega = \omega_2) < n^*$$

We can calculate the expected number of steps as follows:

$$\mathbb{E}(n) = \mathbb{E}(n \mid \omega = \omega_1) = \sum_{t=1}^{n^*} t \mathbb{P}(n = t \mid \omega = \omega_1) \quad (7.1.9)$$

$$= \sum_{t=1}^{n^*-1} t \alpha^{t-1} (1 - \alpha) + n^* \alpha^{n^*-1} = \frac{1 - \alpha^{n^*}}{1 - \alpha}, \quad (7.1.10)$$

from the *geometric series*. The value of this procedure is:

$$\bar{V} = \xi b \alpha^{n^*} - c \mathbb{E}(n) \quad (7.1.11)$$

and from the definition of  $n^*$ :

$$\bar{V} = \frac{c}{\alpha - 1} + \frac{c}{\log \alpha} \left[ 1 + \frac{c}{\xi b(1 - \alpha)} \right]. \quad (7.1.12)$$

To simplify this, note that  $e^x \geq 1 + x \Rightarrow x \geq \log(1 + x) \Rightarrow \alpha - 1 \geq \log(\alpha) \Rightarrow \frac{1}{\alpha - 1} \leq \frac{1}{\log(\alpha)}$ .

**An unbounded sequential procedure**

- Use stopping rule  $s(x^t) = 1$  iff  $x_t = 3$ .
- In other words, stop as soon as you observe a 3.

Once we observe  $x_t = 3$ , we can make a decision that has value 0. So, the value of the unbounded sequential procedure is just

$$-c \mathbb{E}(n).$$

$$\mathbb{E}(n) = \sum_{t=1}^{\infty} t \mathbb{P}(n = t \mid \omega = \omega_1) \quad (7.1.13)$$

$$= \sum_{t=1}^{\infty} t \alpha^{t-1} (1 - \alpha) = \frac{1}{1 - \alpha}. \quad (7.1.14)$$

The value of this procedure is:

$$\bar{V} = \xi b \alpha^{n^*} - c \mathbb{E}(n) \quad (7.1.15)$$

and from the definition of  $n^*$ :

$$\bar{V} = \frac{c}{\alpha - 1} + \frac{c}{\log \alpha} \left[ 1 + \frac{c}{\xi b(1 - \alpha)} \right]. \quad (7.1.16)$$

To simplify this, note that  $e^x \geq 1 + x \Rightarrow x \geq \log(1 + x) \Rightarrow \alpha - 1 \geq \log(\alpha) \Rightarrow \frac{1}{\alpha - 1} \leq \frac{1}{\log(\alpha)}$ .

**Summary**

- Bounded procedures are (in expectation) better than fixed-sampling procedures.
- Unbounded procedures are (in expectation) better than bounded procedures.
- An unbounded procedure may end up costing much more than taking a decision without observing any data.
- Essentially, with an unbounded procedure, we disregard the amount spent to time  $t$ .

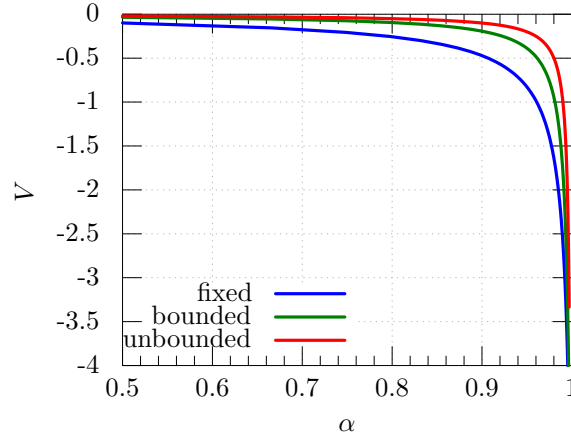


Figure 7.2: The value of three strategies for  $\xi = 1/2$ ,  $b = -10$ ,  $c = 10^{-2}$  and varying  $\alpha$ . Higher values of  $\alpha$  imply a longer time before the true  $\omega$  is known.

## 7.2 Sequential decision procedures

**Definition 7.2.1** (Sequential decision procedure). *A sequential decision procedure  $\delta = (s, d)$  is tuple composed of*

1. *A stopping rule  $s : \mathcal{X}^* \rightarrow \{0, 1\}$ .*
2. *A decision rule  $d : \mathcal{X}^* \rightarrow D$ .*

### Stopping rule

The stopping rule  $s$  specifies whether, at any given time, we should stop and make a decision in  $\mathcal{D}$  or take one more sample. That is, if

$$s(x^t) = 1,$$

stop, otherwise observe  $x_{t+1}$ .

### Decision rule

Once we have stopped (i.e.  $s(x^t) = 1$ ), we choose the decision

$$d(x^t).$$

### Deterministic stopping rules

If the stopping rule  $s$  is deterministic, then, for any  $t$ , there exists some  $B_t \subset \mathcal{X}^t$  such that

$$s(x^t) = \begin{cases} 1, & \forall x^t \in B_t \\ 0, & \forall x^t \notin B_t. \end{cases} \quad (7.2.1)$$



**Cylinder sets**

For  $r > t$ ,  $B_t$  can be also regarded as a subset of  $S_r$ : If  $x^r \in B_t$ , and  $y^r$  is such that  $y_i = x_i$  for  $i = 1, \dots, t$ , then  $y^r \in B_t$ .

### 7.3 Calculating the expected utility of a sequential decision procedure

Once more, consider a distribution family  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$  and a prior  $\xi$  over  $\mathfrak{B}(\Omega)$ . For a decision set  $\mathcal{D}$ , a utility function  $U : \Omega \times \mathcal{D} \rightarrow \mathbb{R}$ , and a sampling cost  $c$ , the risk of a sequential decision procedure is the expected **decision utility** minus the expected **sampling cost**.

$$U(\xi, \delta) = \mathbb{E}_\xi \{U[\omega, \delta(x^n)] - nc\} \quad (7.3.1)$$

$$= \sum_{n=1}^{\infty} \int_{B_n} \mathbb{E}_\xi[U(\omega, \delta(x^n)) \mid x^n] dP(x^n \mid \xi) - \sum_{n=1}^{\infty} \mathbb{P}(B_n \mid \xi) nc \quad (7.3.2)$$

$$= \sum_{n=1}^{\infty} \int_{B_n} \left\{ \int_{\Omega} U[\omega, d(x^n)] d\xi(\omega \mid x^n) \right\} dP(x^n \mid \xi) - \sum_{n=1}^{\infty} \mathbb{P}(B_n \mid \xi) nc. \quad (7.3.3)$$

Although it may seem difficult to evaluate this, it can be done by a simple dynamic programming technique called *backwards induction*.

**Definition 7.3.1** (Bounded sequential decision procedure). *A sequential decision procedure  $\delta$  is bounded if there is a positive integer  $H$  such that  $\mathbb{P}(n \leq H) = 1$ .*

- If  $\delta$  is  $H$ -bounded, then we shall take at most  $H$  samples.
- At stage  $H$ , we will have observed some sequence  $x^H$ , which gives rise to a posterior  $\xi(\omega \mid x^H)$ .
- Since we *must* stop at  $H$ , we must choose a decision  $d$  maximising

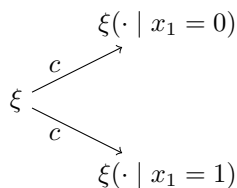
$$\mathbb{E}_\xi[U(\omega, d) \mid x^H] = \int_{\Omega} U(\omega, d) d\xi(\omega \mid x^H)$$

- The respective *value* (expected utility) is

$$V_0[\xi(\cdot \mid x^H)] \triangleq \sup_{d \in \mathcal{D}} \mathbb{E}_\xi[U(\omega, d) \mid x^H],$$

where  $V_0$  denotes the decision's expected utility.

- What about the previous step?



### A simple two-stage problem

Consider  $\mathcal{X} = \{0, 1\}$  and a prior  $\xi$  on the  $\omega$  parameter of  $\text{Bern}(w)$ . We wish to either decide now on a parameter  $\omega$ , or take one more observation, at cost  $c$ , before deciding. Thus, the problem has two stages.

- We begin with a prior  $\xi$  at the first stage. There are two possible outcomes for the **second stage**:
  1. If we have observed  $x_1 = 0$ : then our value is  $V_0[\xi(\cdot | x_1 = 0)]$ .
  2. If we have observed  $x_1 = 1$ : then our value is  $V_0[\xi(\cdot | x_1 = 1)]$ .
- At the first stage,
  1. Stop with value  $V_0(\xi)$ .
  2. Pay cost  $c$  for value:

$$V_0[\xi(\cdot | x_1)], \quad \text{with } P_\xi(x_1) = \int_{\Omega} P_\omega(x_1) d\xi(w)$$

So the expected of continuing for one more step is

$$V_1(\xi) \triangleq \int_{\mathcal{X}} V_0[\xi(\cdot | x_1)] dP_\xi(x_1).$$

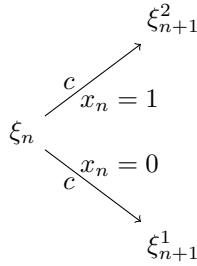
- Thus, the overall value for this problem is:

$$\max \left\{ V_0(\xi), \sum_{x_1=0}^1 V_0[\xi(\cdot | x_1)] P_\xi(x_1) - c. \right\}$$

The above is simply the maximum between the value of stopping immediately, and the value of continuing for one more step. This procedure can be applied recursively for multi-stage problems.

### Multi-stage problems

$$\begin{aligned}
 V_0(\phi) &= \sup_{d \in \mathcal{D}} \int_{\Omega} U(\omega, d) d\phi(\omega) && \text{(Immediate value)} \\
 \xi_{n+1}(\cdot) &= \xi_n(\cdot | x_n) = \xi(\cdot | x^n) && \text{(posterior)} \\
 \mathbb{E}_{\xi_n} V_0[\xi_{n+1}] &= \int_{\mathcal{X}} V_0[\xi_n(\cdot | x_n)] d\xi_n(x_n) && \text{(Next-step value)} \\
 \sigma_1(\xi_n) &= \min \{ V_0(\xi_n), \mathbb{E}_{\xi_n} V_0(\xi_{n+1}) - c. \} && (7.3.4)
 \end{aligned}$$



## 7.4 Backwards induction

The main idea expressed in the previous section is to start from the last stage of our decision problem, where the utility is known, and then move backwards. At each stage, we know the probability of reaching different points in the next stage, as well as their values. Consequently, we can compute the value of any point in the current stage as well. This notion is formalised below.

**Theorem 7.4.1** (Backwards induction). *The utility of a  $k$ -bounded optimal procedure with prior  $\xi_0$  is  $V_k(\xi_0)$  and is given by the recursion:*

$$V_{j+1}(\xi_n) = \max \{V_0(\xi_n), \mathbb{E}_{\xi_n} V_j(\xi_{n+1}) - c\}. \quad (7.4.1)$$

*After observing  $x^{k-j}$ , the value of the optimal continuation is  $V_j(\xi_{k-j})$ .*

The above theorem essentially gives a recursive calculation of the risk of the  $k$ -bounded optimal procedure. In fact, that procedure will have the following property.

**Theorem 7.4.2.** *The optimal  $k$ -bounded procedure stops at time  $t$  if*

$$V_0(\xi_t) \geq V_{k-t}(\xi_t)$$

*and chooses  $d$  maximising  $\mathbb{E}_{\xi_t} U(\omega, d)$ , otherwise takes one more sample.*

Finally, it is always useful to look ahead one more step, as shown by the following theorem.

**Theorem 7.4.3.** *For any probability measure  $\xi$  on  $\Omega$ ,*

$$V_n(\xi) \leq V_{n+1}(\xi). \quad (7.4.2)$$

## 7.5 Unbounded sequential decision procedures

Given the monotonicity of the value of bounded procedures (7.4.2), one may well ask what is the value of unbounded procedures.

**The value of unbounded sequential decision procedures**

- Let  $\delta$  be an unbounded procedure.

- The value of the procedure is

$$V(\xi, \delta) \triangleq \mathbb{E}_\xi \{V_0[\xi(\cdot \mid x^n)] - cn\}, \quad (7.5.1)$$

where  $n$  is the random number of samples taken by  $\delta$ . This is random because the observations  $x$  are random;  $\delta$  itself can be deterministic.

**Definition 7.5.1.** Let  $B_{>k} \subset \mathcal{X}^*$  be the set of sequences such that  $\delta$  takes more than  $k$  samples. Then  $\delta$  is regular if  $V(\xi, \delta) \geq V_0(\xi)$  and if, for all  $x^n \in B_{>n}$  and for all  $n \in \mathbb{N}$ ,

$$\mathbb{E}[V(\xi, \delta) \mid x^n] > V_0[\xi(\cdot \mid x^n)] - cn. \quad (7.5.2)$$

In other words, if  $\delta$  specifies that at least one observation should be taken, then the value of  $\delta$  is greater than the value of choosing a decision without any observation. Furthermore, whenever  $\delta$  specifies that another observation should be taken, the expected value of continuing must be smaller than the value of stopping. If the procedure is not regular, then there may be stages where the procedure specifies that sampling should be continued, though the value would be increased by stopping.

**Theorem 7.5.1.** If  $\delta$  is not regular, then there exists a regular  $\delta'$  such that  $V(\xi, \delta') \geq V(\xi, \delta)$ .

*Proof.* First, consider the case that  $V(\xi, \delta) \leq V_0(\xi)$ . This is equivalent to considering  $\delta'$  to be the regular procedure which chooses  $d \in D$  without any observations.

Now consider the case that  $V(\xi, \delta) < V_0(\xi)$ . Let  $\delta'$  be the procedure which stops as soon as the observed  $x^n$  does not satisfy (7.5.2).

If, for  $x^n$ ,  $\delta$  stops, then both sides of (7.5.2) are equal. Consequently,  $\delta'$  stops before  $n$ .

Finally, let

$$B_k(\delta) = \{x \in \mathcal{X}^* \mid n = k\} \quad (7.5.3)$$

be the set of observations such that exactly  $k$  samples are taken by rule  $\delta$  and

$$B_{\leq k}(\delta) = \{x \in \mathcal{X}^* \mid n \leq k\} \quad (7.5.4)$$

be the set of observations such that at most  $k$  samples are taken by rule  $\delta$ . Then

$$\begin{aligned} V(\xi, \delta) &= \sum_{k=1}^{\infty} \int_{B_k(\delta')} \{V_0[\xi(\cdot \mid x^k)] - ck\} d\xi(x^k) \\ &\geq \sum_{k=1}^{\infty} \int_{B_k(\delta')} \mathbb{E}_\xi \{V[\xi, \delta \mid x^k]\} d\xi(x^k) \\ &= \sum_{k=1}^{\infty} \mathbb{E}_\xi \{V(\xi, \delta) \mid n(\delta') = k\} \mathbb{P}_\xi(\delta' = k) = V(\xi, \delta). \end{aligned}$$

□

## 7.6 The sequential probability ratio test

**A two-point sequential decision problem.**

As an illustration, consider a sequential decision problem where we must decide for one out of two possible parameters  $\omega_1, \omega_2$ .

- Observations  $x_t \in \mathcal{X}$
- Distribution family:  $\mathcal{P} = \{P_\omega \mid \omega \in \Omega\}$ , probability spaces  $(\mathcal{X}^*, \mathfrak{B}(\mathcal{X}^*), P_\omega)$ .
- Parameter set:  $\Omega = \{\omega_1, \omega_2\}$ .
- Decision set:  $D = \{d_1, d_2\}$ .
- Prior  $\xi = \mathbb{P}(\omega = \omega_1)$ .
- Sampling cost  $c > 0$ .

$U(\omega, d)$	$d_1$	$d_2$
$\omega_1$	0	$\lambda_1$
$\omega_2$	$\lambda_2$	0

Table 7.1: The utility function, with  $\lambda_1, \lambda_2 < 0$ 

The immediate value is:

$$V_0(\xi) = \max \{\lambda_1 \xi, \lambda_2 (1 - \xi)\}. \quad (7.6.1)$$

The worst-case immediate value, i.e. the minimum, is attained when both terms are equal. Consequently, setting  $\lambda_1 \xi = \lambda_2 (1 - \xi)$ , which gives  $\xi = \lambda_2 / (\lambda_1 + \lambda_2)$ . Replacing in (7.6.1) gives:

$$V_0(\xi) \geq \frac{\lambda_1 \lambda_2}{\lambda_1 + \lambda_2}.$$

Let  $\Delta'$  denote the set of procedures  $\delta$  which take at least one observation and define:

$$V'(\xi) = \sup_{\delta \in \Delta'} V(\xi, \delta). \quad (7.6.2)$$

Then the  $\xi$ -expected utility  $V^*(\xi)$  must satisfy:

$$V^*(\xi) = \max \{V_0(\xi), V'(\xi)\}. \quad (7.6.3)$$

As we showed in Section ??,  $V'$  is a convex function of  $\xi$ . Let

$$\xi^* \triangleq \{\xi \mid \sigma_0(\xi) \leq \sigma'(\xi)\}, \quad (7.6.4)$$

be the set of priors where it is optimal to terminate sampling.

**The sequential probability ratio test**

If  $\xi \in (\xi_L, \xi_H)$ , then it is optimal to take at least one more sample. Otherwise, it is optimal to make an immediate decision with risk  $\rho_0(\xi)$ .

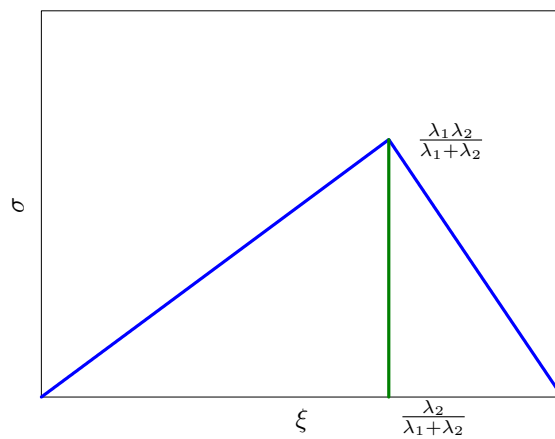
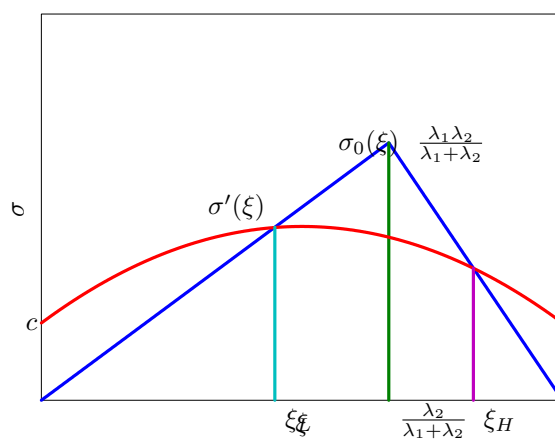


Figure 7.3: The immediate risk

Figure 7.4: The risk of the optimal continuation  $\sigma'$  versus stopping  $\sigma_0$ .

Our posterior at time  $t$  can be written as

$$\xi_t = \frac{\xi P_{\omega_1}(x^t)}{\xi P_{\omega_1}(x^t) + (1 - \xi) P_{\omega_2}(x^t)}.$$

Then, for any posterior, the optimal procedure is:

- If  $\xi_L < \xi_t < \xi_H$ , take one more sample.
- If  $x_L \geq \xi_t$ , stop and decide  $d_2$ .
- If  $x_H \leq \xi_t$ , stop and decide  $d_1$ .

If we let

$$A = \frac{\xi(1 - \xi_H)}{(1 - \xi)\xi_H}, \quad B = \frac{\xi(1 - \xi_L)}{(1 - \xi)\xi_L},$$

then the optimal procedure can be restated as:

Take another observation as long as

$$A < \frac{P_{\omega_2}(x^t)}{P_{\omega_1}(x^t)} < B.$$

If the first inequality is violated, choose  $d_1$ . If the second inequality is violated, choose  $d_2$ .

### Practical considerations

In general, determining  $\xi_L, \xi_H$ , or  $A, B$ , is *not trivial*. This is because we need to first calculate the risk of the optimal decision function that takes at least one more observation. Nevertheless, for a given pair  $A, B$ , the following frequentist property holds.

#### Frequentist property

$$\mathbb{P}(d_1 \mid \omega = \omega_2) \leq A \mathbb{P}(d_1 \mid \omega = \omega_1) \quad (7.6.5)$$

$$\mathbb{P}(d_2 \mid \omega = \omega_2) \geq B \mathbb{P}(d_2 \mid \omega = \omega_1) \quad (7.6.6)$$

Indeed, one may instead simply select  $A, B$  so as to satisfy the above property with a certain ratio. In the worst case scenario, the above are equalities. Since

$$\mathbb{P}(d_1 \mid \omega = \omega_i) + \mathbb{P}(d_2 \mid \omega = \omega_i) = 1, \quad i = 1, 2,$$

we have four equations with four unknowns which we can solve to obtain

$$\mathbb{P}(d_1 \mid \omega_1) = \frac{B - 1}{1 - A}, \quad \mathbb{P}(d_2 \mid \omega_2) = \frac{A(B - 1)}{B - A}. \quad (7.6.7)$$

Setting  $A = 1/K$ ,  $B = K$ , we get

$$\mathbb{P}(d_2 \mid \omega_1) = \mathbb{P}(d_1 \mid \omega_2) = \frac{1}{K + 1}. \quad (7.6.8)$$

Consequently, we can simply choose  $K$  so as to achieve a certain error probability. The same effect can be achieved by choosing an appropriate prior  $\xi$ , so that the SPRT not only has easily-verifiable frequentist properties, but can also be derived from Bayesian principles.

### 7.6.1 Wald's theorem

An important tool in the analysis of SPRT as well as other procedures that stop at random times is the following theorem by Wald.

**Theorem 7.6.1** (Wald's theorem). *Let  $z_1, z_2, \dots$  be a sequence of i.i.d. random variables with measure  $G$ , such that  $\mathbb{E} z_i = m$  for all  $i$ . Then for any sequential procedure with  $\mathbb{E} n < \infty$ :*

$$\mathbb{E} \sum_{i=1}^n z_i = m \mathbb{E} n. \quad (7.6.9)$$

*Proof.*

$$\begin{aligned} \mathbb{E} \sum_{i=1}^n z_i &= \sum_{k=1}^{\infty} \int_{B_k} \sum_{i=1}^k z_i dG^k(z^k) \\ &= \sum_{k=1}^{\infty} \sum_{i=1}^k \int_{B_k} z_i dG^k(z^k). \\ &= \sum_{i=1}^{\infty} \sum_{k=i}^{\infty} \int_{B_k} z_i dG^k(z^k) \\ &= \sum_{i=1}^{\infty} \int_{B_{\geq i}} z_i dG^i(z^i) \\ &= \sum_{i=1}^{\infty} \mathbb{E}(z_i) \mathbb{P}(n \geq i) = m \mathbb{E} n. \end{aligned}$$

□

We now consider application of this theorem to the SPRT. Let  $z_i = \log \frac{P_{\omega_2}(x_i)}{P_{\omega_1}(x_i)}$ . Consider the equivalent formulation of the SPRT which uses

$$a < \sum_{i=1}^n z_i < b$$

as the test. Using Wald's theorem and the previous properties and assuming  $c \approx 0$ , we obtain the following approximately optimal values for  $a, b$ :

$$a \approx \log c - \log \frac{I_1 \lambda_2 (1 - \xi)}{\xi} \quad b \approx \log \frac{1}{c} - \log \frac{I_2 \lambda_1 \xi}{1 - \xi}, \quad (7.6.10)$$

where  $I_1 = -\mathbb{E}(z \mid \omega = \omega_1)$  and  $I_2 = \mathbb{E}(z \mid \omega = \omega_2)$  is the *information*, better known as the *KL divergence*. If the cost  $c$  is very small, then the information terms vanish and we can approximate the values by  $\log c$  and  $\log \frac{1}{c}$ .

## 7.7 Martingales

Martingales are a fundamentally important concept in the analysis of stochastic processes. The main idea of a martingale is that the expectation of a random variable at time  $t + 1$  only depends on the value of another variable at time  $t$ .



An example of a martingale sequence is when  $x_t$  is the amount of money you have at a given time, and where at each time-step  $t$  you are making a gamble such that you lose or gain 1 currency unit with equal probability. Then, at any step  $t$ , it holds that  $\mathbb{E}(x_{t+1} \mid x_t) = x_t$ . This concept can be generalised as follows.

- 1.
2. A random variable  $y_n : \mathcal{S}^n \rightarrow \mathbb{R}$ .

**Definition 7.7.1.** Let  $x^n \in \mathcal{S}^n$  be a sequence of observations  $x^n \in \mathcal{S}^n$  with distribution  $P_n$ , and  $y_n : \mathcal{S}^n \rightarrow \mathbb{R}$  be a random variable. Then the sequence  $\{y_n\}$  is a martingale with respect to  $\{x_n\}$  if for all  $n$  the expectation

$$\mathbb{E}(y_n) = \int_{\mathcal{S}^n} y_n(x^n) dP_n(x^n) \quad (7.7.1)$$

exists and

$$\mathbb{E}(y_{n+1} \mid x^n) = y_n \quad (7.7.2)$$

holds with probability 1. If  $\{y_n\}$  is a martingale with respect to itself, i.e.  $y_i(x) = x$ , then we call it simply a martingale.

**Definition 7.7.2.** Similarly, sequence  $\{y_n\}$  is a super-martingale if  $\mathbb{E}(y_{n+1} \mid x^n) \leq y_n$  and a sub-martingale if  $\mathbb{E}(y_{n+1} \mid x^n) \geq y_n$ , w.p. 1.

### 7.7.1 Doob martingales

At a first glance, it might appear that martingales are not very frequently encountered, apart from some niche applications. Actually, we can always construct a martingale from any sequence of random variables as follows.

Let  $f : \mathcal{S}^m \rightarrow \mathbb{R}$  be some function that we are interested in, such that its expectation  $\mathbb{E}(f \mid x^n)$  exists. Now, let  $y_n : \mathcal{S}^n \rightarrow \mathbb{R}$  with  $n \leq m$  be :

$$y_n(x^n) = \mathbb{E}[f \mid x^n].$$

Then  $\mathbb{E}(y_{n+1} \mid x^n) = y_n$ .

### 7.7.2 The Azuma-Hoeffding inequality

**Definition 7.7.3.** A sequence  $\{y_n\}$  is a martingale difference sequence with respect to  $\{x_n\}$  if

$$\mathbb{E}(y_{n+1} \mid x^n) = 0 \quad \text{with probability 1.} \quad (7.7.3)$$

**Theorem 7.7.1.** Let  $b_k$  be a random variable depending on  $x^{k-1}$ . If  $\{y_k\}$  is a martingale difference sequence with  $y_k \in [b_k, b_k + c_k]$  w.p. 1, then setting  $s_k = \sum_{i=1}^k y_i$  it holds that

$$\mathbb{P}(s_n \geq t) \leq \exp\left(\frac{-2t^2}{\sum_{i=1}^n c_i^2}\right). \quad (7.7.4)$$

**Example 7.7.1** (Estimating a mean). Let we decide to stop if

## 7.8 Markov processes

**Definition 7.8.1** (Markov Process). *Let  $(S, \mathfrak{B}(S))$  be a measurable space and define the random sequence  $\{x_n\}$  such that*

$$\mathbb{P}(x_t \in A \mid x_{t-1}, \dots, x_1) = \mathbb{P}(x_t \in A \mid x_{t-1}). \quad (7.8.1)$$

*$x_n$  is called the state of the Markov process at time  $n$ .*

If  $\mathbb{P}(x_t \in A \mid x_{t-1} = x) = \tau(A \mid x)$  where  $\tau : \mathfrak{B}(S) \times S \rightarrow [0, 1]$ , *transition kernel*. In that case,  $\{x_n\}$  is a stationary Markov process.

**Example 7.8.1** (Markov processes). *The following are examples of Markov processes. Can you think of more?*

- *The state of a Turing machine.*
- *The state in a complete game tree.*
- *Posterior parameters.*
- *The information state defined in the backwards induction tree.*

### A Markov chain stopping problem

- Consider a stationary process with state space  $S$ .
- The transition kernel is a matrix  $\tau$ .
- At time  $t$ , we are at state  $x_t = z$  and we can either:
  - Terminate and receive reward  $b(z)$ .
  - Pay  $c(z)$  and continue to a random state  $x_{t+1}$ .
- There maybe  $S_0, S_1 \subset S$  where we may be forced to stop or continue sampling.

All the problems considered so far are stationary! This is true no matter whether the next observation is generated from the posterior state  $\xi_n$  or the true distribution  $P_\omega$ .

### Summary

- Sequential sampling is always better than a fixed sample size.
- Unbounded procedures are better than bounded procedures.
- Bounded procedures can be calculated using backwards induction.
- Unbounded procedures can be approximated as the limit of a sequence of bounded procedures.
- The sequential probability ratio test (SPRT) is a type of unbounded sequential decision procedure.

- The error probabilities of the SPRT  $A < P_{\omega_2}(x^t)/P_{\omega_1}(x^t) < B$  are approximately  $A, 1/B$ . For sample cost  $c \rightarrow 0$ , the near-optimal values are  $A = c, B = 1/c$ .
- Martingales are a special type of sequence of random variables such that  $\mathbb{E}(y_{n+1} \mid x^n) = y_n(x^n)$ .
- Concentration inequalities can be derived for martingales.
- All the above problems can be modelled as Markov processes.



## Chapter 8

# Experiment design and Markov decision processes

## 8.1 Introduction

This unit describes the very general formalism of Markov decision processes (MDPs) for dealing with various problems in sequential decision making. Thus a Markov decision process can be used to model stochastic path problems, stopping problems, reinforcement learning problems, experiment design problems, and control problems.

### 8.1.1 Experiment design: examples

The problem of experimental design originally arose in the statistical literature when considering the problem of how to best allocate treatments with unknown efficacy to patients. The problem, originally considered by Chernoff [1959, 1966], informally can be stated as follows.

We have a number of treatments of unknown efficacy. When a new patient arrives, we must choose one of them. There are two possible, slightly different, goals: either maximise the number of cured patients, or discover the best treatment. The two different problems can be formalised as follows.

**Example 8.1.1** (Adaptive treatment allocation). *Consider  $k$  treatments to be administered to  $T$  volunteers. To each volunteer only a single treatment can be assigned. At the  $t$ -th trial, we perform some experiment  $a_t \in \{1, \dots, k\}$  and obtain a reward  $r_t = 1$  if the result is successful and 0 otherwise. We wish to choose actions maximising  $\sum_t r_t$ .*

**Example 8.1.2** (Adaptive hypothesis testing). *We are given a hypothesis set  $\Omega = \{\omega_1, \omega_2\}$ , a prior  $\psi_0$  on  $\Omega$ , a decision set  $\mathcal{D} = \{d_1, d_2\}$  and a utility function  $\Omega : \mathcal{D} \times \Omega \rightarrow \mathbb{R}$ . One hypothesis  $\omega \in \Omega$  is true. We can choose from a set of  $k$  possible experiments to be performed over  $T$  trials. At the  $t$ -th trial, we choose experiment  $a_t \in \{1, \dots, k\}$  and observe outcome  $x_t \in \mathcal{X}$ , with  $x_t \sim P_\omega$  drawn from the true hypothesis. Our posterior is*

$$\xi_t(\omega) \triangleq \xi_0(\omega \mid a_1, \dots, a_t, x_1, \dots, x_t).$$

The reward is  $r_t = 0$  for  $t < T$  and

$$r_T = \max_{d \in \mathcal{D}} \mathbb{E}_{\xi_T}(U \mid d).$$

Again, we wish to maximise  $\sum_t r_t$ .

Both formalizations correspond to so-called *bandit problems* which we take a closer look at in the following section.

## 8.2 Bandit problems

The simplest bandit problem is the stochastic  $n$ -armed bandit. We are faced with  $n$  different one-armed bandit machines, such as those found in casinos. In this problem, at time  $t$ , you have to choose one *action* (i.e. machine)  $a_t \in \mathcal{A} = \{1, \dots, n\}$ . The assumption is that each time  $t$  you play a machine, you receive a bounded reward  $r_t$ , with fixed expected value  $\omega_i = \mathbb{E}(r_t \mid a_t = i)$ . Unfortunately, you do not know  $\omega_i$ . How do you then choose arms so as to maximise the total expected reward?

**The stochastic  $n$ -armed bandit problem**

- Actions  $\mathcal{A} = \{1, \dots, n\}$ .
- Expected reward  $\mathbb{E}(r_t \mid a_t = i) = \omega_i$ .
- Select actions to maximise

$$\sum_{t=0}^T \gamma^t r_t,$$

with discount factor  $\gamma \in [0, 1]$ , horizon  $T \geq 0$ .

One idea is to apply the Bayesian decision-theoretic framework we have developed earlier to maximise the reward in expectation. More specifically, given the horizon  $T \in (0, \infty]$  and the discount factor  $\gamma \in (0, 1]$ , we define our utility from time  $t$  to be:

$$U_t = \sum_{k=1}^{T-t} \gamma^k r_{t+k}. \quad (8.2.1)$$

**Decision-theoretic approach**

- Assume  $r_t \mid a_t = a \sim P_{\omega, i}$ , with  $\omega \in \Omega$ .
- Define prior  $\xi$  on  $\Omega$ .
- Select actions to maximise  $\mathbb{E}_\xi U_t = \mathbb{E}_\xi \sum_{k=1}^{T-t} \gamma^k r_{t+k}$ .

**8.2.1 Bernoulli bandits**

As a simple illustration, consider the case when the reward for choosing one of the  $n$  actions is either 0 or 1, with some fixed, yet unknown probability depending on the chosen action. This can be modelled in the standard Bayesian framework using the Beta-Bernoulli conjugate prior. More specifically, we can formalise the problem as follows.

Consider  $n$  Bernoulli distributions with unknown parameters  $\omega_i$  ( $i = 1, \dots, n$ ) such that

$$r_t \mid a_t = i \sim \text{Bern}(\omega_i), \quad \mathbb{E}(r_t \mid a_t = i) = \omega_i. \quad (8.2.2)$$

Each Bernoulli distribution thus corresponds to the distribution of rewards obtained from each bandit that we can play. In order to apply the statistical decision theoretic framework, we have to quantify our uncertainty about the parameters  $\omega$  in terms of a probability distribution. We model our belief

for each bandit's parameter  $\omega_i$  as a Beta distribution  $\mathcal{Beta}(\alpha_i, \beta_i)$ , with density  $f(\omega \mid \alpha_i, \beta_i)$  so that

$$\xi(\omega_1, \dots, \omega_n) = \prod_{i=1}^n f(\omega_i \mid \alpha_i, \beta_i).$$

Recall that the posterior of a Beta prior is also a Beta. Let

$$N_{t,i} \triangleq \sum_{k=1}^t \mathbb{I}\{a_k = i\}$$

be the number of times we played arm  $i$  and

$$\hat{r}_{t,i} \triangleq \frac{1}{N_{t,i}} \sum_{k=1}^t r_k \mathbb{I}\{a_k = i\}$$

be the *empirical reward* of arm  $i$  at time  $t$ . We can let this equal 0 when  $N_{t,i} = 0$ . Then, the posterior distribution for the parameter of arm  $i$  is

$$\xi_t = \mathcal{Beta}(\alpha_i + N_{t,i} \hat{r}_{t,i}, \beta_i + N_{t,i}(1 - \hat{r}_{t,i})).$$

Since  $r_t \in \{0, 1\}$  the possible states of our belief given some prior are  $\mathbb{N}^{2n}$ .

The state of the bandit problem is the state of our belief. A sufficient statistic for our belief is the number of times we played each bandit and the total reward from each bandit. Thus, our state at time  $t$  is entirely described by our priors  $\alpha, \beta$  (the initial state) and the vectors

$$N_t = (N_{t,1}, \dots, N_{t,i}) \tag{8.2.3}$$

$$\hat{r}_t = (\hat{r}_{t,1}, \dots, \hat{r}_{t,i}). \tag{8.2.4}$$

At any time  $t$ , we can calculate the probability of observing  $r_t = 1$  or  $r_t = 0$  if we pull arm  $i$  as:

$$\xi_t(r_t = 1 \mid a_t = i) = \frac{\alpha_i + N_{t,i} \hat{r}_{t,i}}{\alpha_i + \beta_i + N_{t,i}}$$

The next state is well-defined and depends only on the current state. For this reason, the decision-theoretic  $n$ -armed bandit problem can be formalised as a *Markov decision process*.

The number of states of this particular bandit problem is countable for  $\{0, 1\}$  rewards, though for a finite horizon  $T$  it is of order  $(2n)^T$ . In the general case, the number of states is uncountable.

## 8.2.2 Decision-theoretic bandit process

The basic bandit process can be seen in Figure 8.1. The decision-theoretic bandit process can be formalised as follows.

**Definition 8.2.1.** *Let  $\mathcal{A}$  be a set of actions, not necessarily finite. Let  $\Omega$  be a set of possible parameter values, indexing a family of probability measures  $\mathcal{P} = \{P_{\omega,a} \mid \omega \in \Omega, a \in \mathcal{A}\}$ . There is some  $\omega \in \Omega$  such that, whenever we take action  $a_t = a$ , we observe reward  $r_t$  with probability measure:*

$$P_{\omega,a}(R) \triangleq \mathbb{P}_{\omega}(r_t \in R \mid a_t = a). \tag{8.2.5}$$



Let  $\xi_1$  be a prior distribution on  $\Omega$  and let the posterior distributions be defined as:

$$\xi_{t+1}(B) \propto \int_B P_{\omega, a_t}(r_t | s_t) d\xi_t(\omega). \quad (8.2.6)$$

The next belief is random, since it depends on the random quantity  $r_t$ . In fact, the probability of next rewards if  $a_t = a$  is given by:

$$P_{\xi_t, a}(R) \triangleq \int_{\Omega} P_{\omega, a}(R) d\xi_t(\omega) \quad (8.2.7)$$

Finally, as  $\xi_{t+1}$  deterministically depends on  $\xi_t, a_t, r_t$ , the probability of obtaining a particular next belief is the same as the probability of obtaining the corresponding rewards leading to the next belief. In more detail, we can write:

$$\mathbb{P}(\xi_{t+1} = \xi | \xi_t, a_t) = \int_{\mathcal{R}} \mathbb{I}\{\xi_{t+1}(\cdot | a_t, r_t = r) = \xi\} dP_{\xi_t, a}(r). \quad (8.2.8)$$

In practice, although multiple reward sequences may lead to the same beliefs, we frequently ignore that possibility for simplicity. Then the process becomes a tree.

Since the next belief only depends on the current belief, action and reward, it satisfies the Markov property. Consequently, a decision-theoretic bandit process can be modelled as a Markov decision process.

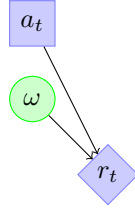


Figure 8.1: The basic bandit process. The decision maker selects  $a_t$ , while the parameter  $\omega$  of the process is hidden. It then obtains reward  $r_t$ . The process repeats for  $t = 1, \dots, T$ .

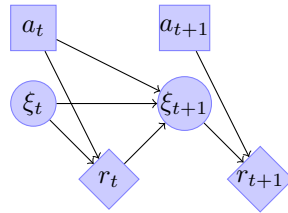


Figure 8.2: The decision-theoretic bandit process. While  $\omega$  is not known, at each time step  $t$  we maintain a belief  $\xi_t$  on  $\Omega$ . The reward distribution is then defined through our belief.

## 8.3 Markov decision processes and reinforcement learning

### Reinforcement learning

Bandit problems are one of the simplest instances of reinforcement learning problems. Informally, speaking, these are problems of learning how to act in an unknown environment, only through interaction with the environment and limited reinforcement signals.

*The reinforcement learning problem.*  
 Learning to act in an *unknown* environment, by **interaction** and **reinforcement**.

Generally, we assume that the environment that we are acting in has an underlying state  $s_t$ , which changes with time  $t$ . At the same time, the agent chooses actions  $a_t$ . We usually assume that the environment is such that its next state  $s_{t+1}$  only depends on its current state  $s_t$  and the last action taken by the agent,  $a_t$ . In addition, the agent observes a reward signal  $r_t$ , and its goal is to maximise the total reward during its lifetime.

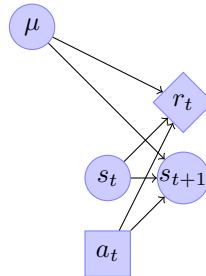
This problem is hard even in seemingly simple problems, like  $n$ -armed bandits, where the underlying state never changes. However, in many real-world applications, the state is not directly observed. Instead, we may simply have some measurements  $x_t$ , which give only partial information about the true underlying state.

Reinforcement learning problems typically fall into one of the following three groups: (1) Markov decision processes (MDPs), where the state is known; (2) Partially observable MDPs (POMDPs), where the state is hidden; and (3) potentially partially observable Markov games, where the next state also depends on the move of other agents. While all of these problem *descriptions* are different, in the Bayesian setting, they all can be reformulated as MDPs, albeit with a very large state space.

**Markov decision processes (MDP).**

At each time step  $t$ :

- We observe *state*  $s_t \in \mathcal{S}$ .
- We take *action*  $a_t \in \mathcal{A}$ .
- We receive a *reward*  $r_t \in \mathbb{R}$ .



### Markov property of the reward and state distribution

$$\mathbb{P}_\mu(s_{t+1} \mid s_1, a_1, \dots, s_t, a_t) = \mathbb{P}_\mu(s_{t+1} \in S \mid s_t, a_t) \quad (8.3.1)$$

$$\mathbb{P}_\mu(r_t \mid s_1, a_1, \dots, s_t, a_t) = \mathbb{P}_\mu(r_t \in R \mid s_t, a_t) \quad (8.3.2)$$

where  $S \subset \mathcal{S}$  and  $R \subset \mathcal{R}$ .

More formally, we have the following definition.

**Definition 8.3.1.** A Markov decision process  $\mu$  is a tuple  $\mu = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R} \rangle$ , such that  $\mathcal{P} = \{P(\cdot \mid s, a) \mid s \in \mathcal{S}, a \in \mathcal{A}\}$  is a collection of probability measures on  $\mathcal{S}$ , indexed in  $\mathcal{S} \times \mathcal{A}$  and  $\mathcal{R} = \{\rho(\cdot \mid s, a) \mid s \in \mathcal{S}, a \in \mathcal{A}\}$  is a collection of probability measures on  $\mathcal{R}$ , such that:

$$P(S \mid s, a) = \mathbb{P}_\mu(s_{t+1} \in S \mid s_t = s, a_t = a) \quad (8.3.3)$$

$$\rho(R \mid s, a) = \mathbb{P}_\mu(r_t \in R \mid s_t = s, a_t = a). \quad (8.3.4)$$

For simplicity, we shall also use

$$r_\mu(s, a) = \mathbb{E}_\mu(r_{t+1} \mid s_t = s, a_t = a), \quad (8.3.5)$$

for the expected reward.

Of course, the transition and reward distributions are different for different environments  $\mu$ . For that reason, we shall usually subscript the relevant probabilities and expectations with  $\mu$ .

### Dependencies of rewards

Sometimes it is more convenient to have rewards that depend on the next state as well, i.e.

$$r_\mu(s, a, s') = \mathbb{E}_\mu(r_{t+1} \mid s_t = s, a_t = a, s_{t+1} = s'), \quad (8.3.6)$$

though this complicates notation considerably since now the reward is obtained on the next time step. However, we can always replace this with the expected reward for a given state-action pair:

$$r_\mu(s, a) = \mathbb{E}_\mu(r_{t+1} \mid s_t = s, a_t = a) \quad (8.3.7)$$

$$= \sum_{s' \in \mathcal{S}} P_\mu(s' \mid s, a) r_\mu(s, a, s') \quad (8.3.8)$$

In fact, it is notationally more convenient to have rewards that only depend on the current state:

$$r_\mu(s) = \mathbb{E}_\mu(r_t \mid s_t = s). \quad (8.3.9)$$

Many times, for simplicity, we shall only consider the latter case.

**The agent**

The environment does not exist in isolation. The actions are taken by an agent, who is interested in obtaining high rewards. Instead of defining an algorithm for choosing actions directly, we define an algorithm for computing policies, which define distributions on actions.

**The agent's policy  $\pi$** 

$$\begin{aligned} \mathbb{P}^\pi(a_t \mid s_t, \dots, s_1, a_{t-1}, \dots, a_1) & \quad (\text{history-dependent policy}) \\ \mathbb{P}^\pi(a_t \mid s_t) & \quad (\text{Markov policy}) \end{aligned}$$

In some sense, the agent is defined by its *policy*  $\pi$ , which is a conditional distribution on actions given the history. The policy  $\pi$  is otherwise known as a *decision function* or a *strategy*. In general, the policy can be history-dependent. In certain cases, however, there are optimal policies that are Markov. This is for example the case with additive utility functions.

**Definition 8.3.2** (Utility). *Given a horizon  $T$  and a discount factor  $\gamma \in (0, 1]$ , the utility can be defined as*

$$U_t \triangleq \sum_{k=0}^{T-t} \gamma^k r_{t+k} \quad (8.3.10)$$

The agent wants to find  $\pi$  *maximising* the *expected total future reward*

$$\mathbb{E}_\mu^\pi U_t = \mathbb{E}_\mu^\pi \sum_{k=0}^{T-t} \gamma^k r_{t+k}. \quad (\text{expected utility})$$

For simplicity, we shall also use  $r_\mu^\pi(s) = \mathbb{E}_\mu^\pi(r_{t+1} \mid s_t = s)$  for the expected reward at a given state, or simply  $r(s)$  when clear from context.

**8.3.1 Value functions**

It is frequently useful to employ the following abbreviations. These will be employed in the development of both theory and algorithms and are simply the expected utility for a given policy and MDP conditioned on different states and/or actions.

**State value function**

$$V_{\mu,t}^\pi(s) \triangleq \mathbb{E}_\mu^\pi(U_t \mid s_t = s) \quad (8.3.11)$$

**State-action value function**

$$Q_{\mu,t}^\pi(s, a) \triangleq \mathbb{E}_\mu^\pi(U_t \mid s_t = s, a_t = a) \quad (8.3.12)$$

It is also useful to define the optimal policy and optimal value functions for a given MDP. In the following, a star indicated optimal quantities. The *optimal policy*  $\pi^*$

$$\pi^*(\mu) : V_{t,\mu}^{\pi^*(\mu)}(s) \geq V_{t,\mu}^\pi(s) \quad \forall \pi, t, s \quad (8.3.13)$$

dominates all other policies  $\pi$  everywhere in  $\mathcal{S}$ .

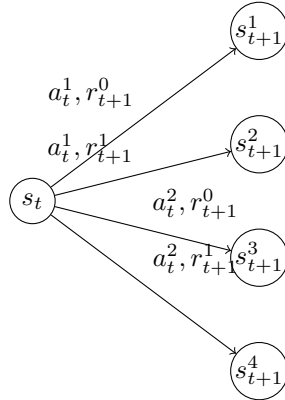
The *optimal value function*  $V^*$

$$V_{t,\mu}^*(s) \triangleq V_{t,\mu}^{\pi^*(\mu)}(s), \quad Q_{t,\mu}^*(s) \triangleq Q_{t,\mu}^{\pi^*(\mu)}(s, a). \quad (8.3.14)$$

is the value function of the optimal policy  $\pi^*$ .

#### Finding the optimal policy when $\mu$ is known

When the MDP  $\mu$  is known, the expected utility of any policy can be calculated. Therefore, one could find the optimal policy by brute force. However, there are faster methods which can be employed. First, there are iterative/offline methods where an optimal policy is found for all states of the MDP. These either try to estimate the optimal value function directly, or try to iteratively improve a policy until it is optimal. The second type of methods tries to find an optimal policy online. That is, the optimal actions are estimated only for states which can be visited in the future starting from the current state. However, the same main ideas are used in all of these algorithms.



## 8.4 Finite horizon, undiscounted problems

The conceptually simplest type of problems are finite horizon problems where  $T < \infty$  and  $\gamma = 1$ . The first thing we shall try to do is to evaluate a given policy for a given MDP.

### 8.4.1 Policy evaluation

#### Policy evaluation

An optimal policy

An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision. – Bellman.

The value function of a policy  $\pi$  (for  $\gamma = 1, T < \infty$ ) can be determined by the following recursion, noting that  $U_{t+1} = \sum_{k=1}^{T-t} r_{t+k}$ :

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (8.4.1)$$

$$= \sum_{k=0}^{T-t} \mathbb{E}_{\mu}^{\pi}(r_{t+k} \mid s_t = s) \quad (8.4.2)$$

$$= \mathbb{E}_{\mu}^{\pi}(r_t \mid s_t = s) + \mathbb{E}_{\mu}^{\pi}(U_{t+1} \mid s_t = s) \quad (8.4.3)$$

$$= \mathbb{E}_{\mu}^{\pi}(r_t \mid s_t = s) + \sum_{i \in \mathcal{S}} V_{\mu,t+1}^{\pi}(i) \mathbb{P}_{\mu}^{\pi}(s_{t+1} = i \mid s_t = s). \quad (8.4.4)$$

Note that

$$\mathbb{P}_{\mu}^{\pi}(s_{t+1} = i \mid s_t = s) = \sum_{a \in \mathcal{A}} \mathbb{P}_{\mu}(s_{t+1} = i \mid s_t = s, a_t = a) \mathbb{P}^{\pi}(a_t = a \mid s_t = s). \quad (8.4.5)$$

This derivation directly gives a number of *policy evaluation algorithms*.

Direct policy evaluation as given by Algorithm 2 is very simple. We calculate the probability of reaching any state from any other state at different times, and then add up the expected reward we would get in that state under our policy. It is easy to see that  $\hat{V}_t(s) = V_{\mu,t}^{\pi}(s)$ .

---

**Algorithm 2** Direct policy evaluation

---

```

1: for  $s \in \mathcal{S}$  do
2:   for  $t = 0, \dots, T$  do
3:
```

$$\hat{V}_t(s) = \sum_{k=t}^T \sum_{j \in \mathcal{S}} \mathbb{P}_{\mu}^{\pi}(s_k = j \mid s_t = s) \mathbb{E}_{\mu}^{\pi}(r_k \mid s_k = j).$$

```

4:   end for
5: end for
```

---

### 8.4.2 Monte-Carlo policy evaluation

Another conceptually simple algorithm is Monte-Carlo policy evaluation shown as Algorithm 3. The idea is that instead of summing over all possible states to be visited, we just draw states from the Markov chain defined jointly by the policy and the Markov decision process. Unlike direct policy evaluation the algorithm needs a parameter  $K$ , the number of trajectories to generate. Nevertheless, this is a very useful method, employed within a number of more complex algorithms.

**Algorithm 3** Monte-Carlo policy evaluation

---

```

for  $s \in \mathcal{S}$  do
  for  $k = 0, \dots, K$  do
    Choose initial state  $s_1$ .
    for  $t = 1, \dots, T$  do
       $a_t \sim \pi(a_t \mid s_t)$  // Take action
      Observe reward  $r_t$  and next state  $s_{t+1}$ .
      Set  $r_{t,k} = r_t$ .
    end for
  Save total reward:

```

$$\hat{V}_k(s) = \sum_{t=1}^T r_{t,k}.$$

```

end for

```

```

Calculate estimate:

```

$$\hat{V}(s) = \frac{1}{K} \sum_{k=1}^K \hat{V}_k(s).$$

```

end for

```

---

**Remark 8.4.1.** The estimate  $\hat{V}$  of the Monte Carlo evaluation algorithm satisfies

$$\|V - \hat{V}\|_\infty \leq \sqrt{\frac{\ln(2|\mathcal{S}|/\delta)}{2K}} \quad \text{with probability } 1 - \delta$$

*Proof.* From Hoeffding's inequality we have for any state  $s$  that

$$\mathbb{P} \left( |\hat{V}(s) - V(s)| \geq \sqrt{\frac{\ln(2|\mathcal{S}|/\delta)}{2K}} \right) \leq \delta/|\mathcal{S}|.$$

Consequently, using a union bound of the form  $P(A_1 \cup A_2 \cup \dots \cup A_n) \leq \sum_i P(A_i)$  gives the required result.  $\square$

Finally, the backwards induction algorithm shown as Algorithm 4 is similar to the backwards induction algorithm we saw for sequential sampling problems. However, here we are only evaluating a policy rather than finding the optimal policy. This algorithm is slightly less generally applicable than the Monte-Carlo method because it implicitly assumes that the policy we evaluate only takes actions depending on the current state.

**Algorithm 4** Backwards induction policy evaluation

---

```

For each state  $s \in S$ , for  $t = 1, \dots, T - 1$ :

```

$$\hat{V}_t(s) = r_\mu^\pi(s) + \sum_{j \in S} \mathbb{P}_\mu^\pi(s_{t+1} = j \mid s_t = s) \hat{V}_{t+1}(j), \quad (8.4.6)$$

```

with  $\hat{V}_T(s) = r_\mu^\pi(s)$ .

```

---

**Theorem 8.4.1.** The backwards induction algorithm gives estimates  $\hat{V}_t(s)$  satisfying

$$\hat{V}_t(s) = V_{\mu,t}^\pi(s) \quad (8.4.7)$$

*Proof.* For  $t = T - 1$ , the result is obvious. We can prove the remainder by induction. Let (8.4.7) hold for all  $t \geq n + 1$ . Now we prove that it holds for  $n$ . Note that from the recursion (8.4.6) we have:

$$\begin{aligned}\hat{V}_t(s) &= r_\mu(s) + \sum_{j \in S} \mathbb{P}_{\mu, \pi}(s_{t+1} = j \mid s_t = s) \hat{V}_{t+1}(j) \\ &= r(s) + \sum_{j \in S} \mathbb{P}_{\mu, \pi}(s_{t+1} = j \mid s_t = s) V_{\mu, t+1}^\pi(j) \\ &= r(s) + \mathbb{E}_{\mu, \pi}(U_{t+1} \mid s_t = s) \\ &= \mathbb{E}_{\mu, \pi}(U_t \mid s_t = s) = V_{\mu, t}^\pi(s),\end{aligned}$$

where the second equality is by the induction hypothesis, the third and fourth equalities are by the definition of the utility, and the last by definition of  $V_{\mu, t}^\pi$ .  $\square$

As you can see, the proof follows the same outline as the expansion of the policy value. The same theorem can be proven for history-dependent policies.

### 8.4.3 Finite horizon backwards induction

Backwards induction as given in algorithm 5 is the first non-naive algorithm for finding an optimal policy for the case where there are  $T$  stages. It is basically identical to the backwards induction algorithm we saw in Chapter 7.

---

#### Algorithm 5 Finite-horizon backwards induction

---

Input  $\mu$ , set  $\mathcal{S}_T$  of states reachable within  $T$  steps.  
 Initialise  $V_T(s) := \max_a r(s, a)$ , for all  $s \in \mathcal{S}_T$ .  
**for**  $n = T - 1, T - 2, \dots, 1$  **do**  
   **for**  $s \in \mathcal{S}_n$  **do**  
       $\pi_n(s) = \arg \max_a r(s, a) + \sum_{s' \in \mathcal{S}_{n+1}} P_\mu(s' \mid s, a) V_{n+1}(s')$   
       $V_n(s) = r(s, a) + \sum_{s' \in \mathcal{S}_{n+1}} P_\mu(s' \mid s, \pi_n(s)) V_{n+1}(s')$   
   **end for**  
**end for**  
 Return  $\pi = (\pi_n)_{n=1}^T$ .

---

**Theorem 8.4.2.** *For  $T$ -horizon problems, backwards induction is optimal, i.e.*

$$V_n(s) = V_{\mu, n}^*(s) \tag{8.4.8}$$

*Proof.* Note that the proof below also holds for  $r(s, a) = r(s)$ . First we show that  $V_t \geq V_t^*$ . For  $n = T$  we evidently have  $V_T(s) = \max_a r(s, a) = V_{\mu, T}^*(s)$ . Now assume that for  $n \geq t + 1$ , (8.4.8) holds. Then it also holds for  $n = t$ , since



for any policy  $\pi'$

$$\begin{aligned}
V_t(s) &= \max_a \left\{ r(s, a) + \sum_{j \in \mathcal{S}} P_\mu(j \mid s, a) V_{t+1}(j) \right\} \\
&\geq \max_a \left\{ r(s, a) + \sum_{j \in \mathcal{S}} P_\mu(j \mid s, a) V_{\mu, t+1}^*(j) \right\} \quad (\text{by induction assumption}) \\
&\geq \max_a \left\{ r(s, a) + \sum_{j \in \mathcal{S}} P_\mu(j \mid s, a) V_{\mu, t+1}^{\pi'}(j) \right\} \\
&\geq V_t^{\pi'}(s).
\end{aligned}$$

This holds for any policy  $\pi'$ , including  $\pi' = \pi$ , the policy returned by backwards induction. Then:

$$V_{\mu, t}^*(s) \geq V_{\mu, t}^\pi(s) = V_t(s) \geq V_{\mu, t}^*(s).$$

□

**Remark 8.4.2.** *A similar theorem can be proven for arbitrary  $\mathcal{S}$ . This requires using sup instead of max and proving the existence of a  $\pi'$  that is arbitrary-close in value to  $v$ .*

## 8.5 Infinite-horizon

When problems have no fixed horizon, they usually can be modelled as infinite horizon problems, sometimes with help of a *terminating state*, whose visit terminates the problem, or discounted rewards, which indicate that we care less about rewards further in the future. When reward discounting is exponential, these problems can be seen as undiscounted problems with random and geometrically distributed horizon. For problems with no discounting and no termination states there are some complications in the definition of optimal policy. However, we defer discussion of such problems to Chapter ??.

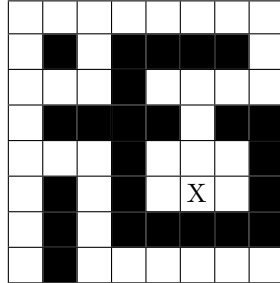
### 8.5.1 Examples

We begin with some examples.

#### Shortest-path problems

##### Deterministic shortest-path problems

Consider an agent moving in a maze, aiming to get to some terminating goal state  $X$ . That is, when reaching this state, the agent cannot move anymore, and receives a reward of 0. In general, the agent can move deterministically in the four cardinal directions, and receives a negative reward at each time step. Consequently, the optimal policy is to move to  $X$  as quickly as possible.

**Properties**

- $\gamma = 1, T \rightarrow \infty$ .
- $r_t = -1$  unless  $s_t = X$ , in which case  $r_t = 0$ .
- $\mathbb{P}_\mu(s_{t+1} = X | s_t = X) = 1$ .
- $\mathcal{A} = \{\text{North, South, East, West}\}$
- Transitions are deterministic and walls block.

Solving the shortest path problem can be done simply by looking at the distance of any point to  $X$ . Then the reward obtained by the optimal policy starting from any point, is simply the negative distance. The optimal policy simply moves to the state with the smallest distance to  $X$ .

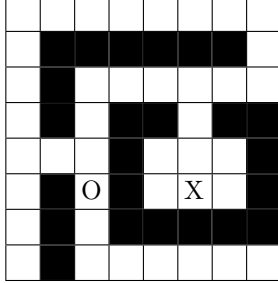
14	13	12	11	10	9	8	7
15		13					6
16	15	14		4	3	4	5
17					2		
18	19	20		2	1	2	
19		21		1	0	1	
20		22					
21		23	24	25	26	27	28

**Properties**

- $\gamma = 1, T \rightarrow \infty$ .
- $r_t = -1$  unless  $s_t = X$ , in which case  $r_t = 0$ .
- The length of the shortest path from  $s$  equals the negative value of the optimal policy.
- Also called *cost-to-go*.

**Stochastic shortest path problem with a pit**

Now assume the shortest path problem with stochastic dynamics. That is, at each time-step there is a small probability  $\omega$  that move to a random direction. In addition, there is a pit  $O$ , that is a terminating state with a reward of  $-100$ .



### Properties

- $\gamma = 1, T \rightarrow \infty$ .
- $r_t = -1$ , but  $r_t = 0$  at X and  $-100$  at O and episode ends.
- $\mathbb{P}_\mu(s_{t+1} = X | s_t = X) = 1$ .
- $\mathcal{A} = \{\text{North, South, East, West}\}$
- Moves to a random direction with probability  $\omega$ . Walls block.

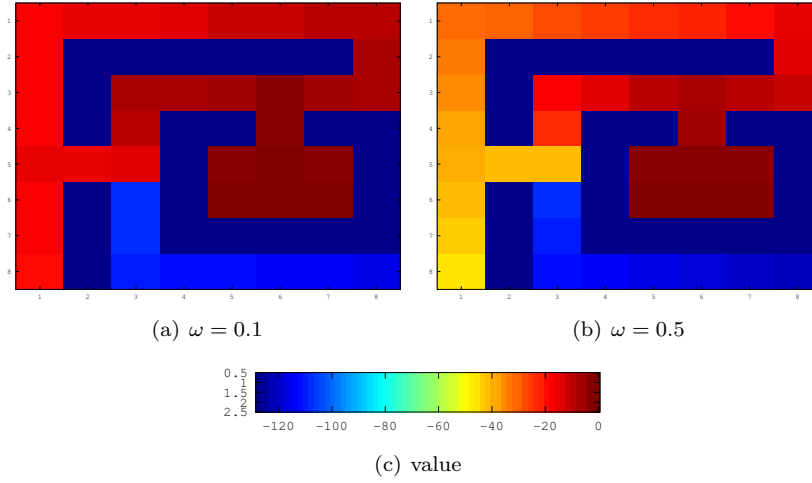


Figure 8.3: Pit maze solutions for two values of  $\omega$ .

Randomness changes the solution significantly in this environment. When  $\omega$  is relatively small, it is worthwhile (in expectation) for the agent to pass past the pit, even though there is a risk of falling in and getting a reward of  $-100$ . In the example given, even starting from the third row, the agent prefers taking the short-cut. For high enough  $\omega$ , the optimal policy avoids approaching the pit. Still, the agent prefers jumping in the pit, than being trapped at the bottom of the maze forever.

### Continuing problems

Finally, many problems have no natural terminating state, but are continuing *ad infinitum*. Frequently, we model those problems using a utility that discounts future rewards exponentially. As an example, consider the following inventory management problem. There are  $K$  storage locations, and each location  $i$  can store  $n_i$  items. At each time-step there is a probability  $\phi_i$  that a client tries to buy an item from location  $i$ , where  $\sum_i \phi_i \leq 1$ . If there is an item available,

when this happens, you gain reward 1. There are two types of actions, one for ordering a certain number  $u$  units of stock, paying  $c(u)$ . Further one may move  $u$  units of stock from one location  $i$  to another location  $j$ , paying  $\psi_{ij}(u)$ .

An easy special case is when  $K = 1$ , and we assume that deliveries happen once every  $m$  timesteps, and each time-step a client arrives with probability  $\phi$ . Then the state set  $\mathcal{S} = \{0, 1, \dots, n\}$  corresponds to the number of items we have, the action set  $\mathcal{A} = \{0, 1, \dots, n\}$  to the number of items we may order. The transition probabilities are given by  $P(s'|s, a) = \binom{m}{d} \phi^d (1 - \phi)^{m-d}$ , where  $d = s + a - s'$ , for  $s + a \leq n$ .

### 8.5.2 Markov chain theory for discounted problems

We first consider MDPs with discounted rewards, and only turn later to algorithms and bounds for the undiscounted case.

#### Discounted total reward.

Our utility in this case is:

$$U_t = \lim_{T \rightarrow \infty} \sum_{k=t}^T \gamma^k r_k, \quad \gamma \in (0, 1)$$

For simplicity, in the following we assume that rewards only depend on the current state instead of both state and action. It can easily be verified that results still hold in the latter case. We use the following notation:

- $\mathbf{v}^\pi = (\mathbb{E}^\pi(U_t \mid s_t = s))_{s \in \mathcal{S}}$  is a vector in  $\mathbb{R}^{|\mathcal{S}|}$  representing the value of policy  $\pi$ .
- $\mathbf{P}_{\mu, \pi}$  is a transition matrix in  $\mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$  for policy  $\pi$ .
- Sometimes we will use  $p(j|s, a)$  as a shorthand for  $\mathbb{P}_\mu(s_{t+1} = j \mid s_t = s, a_t = a)$ .
- $\mathbf{r}$  is a reward vector in  $\mathbb{R}^{|\mathcal{S}|}$ .
- The space of value functions  $\mathcal{V}$  is a Banach space (i.e., a complete, normed vector space) equipped with the norm

$$\|\mathbf{v}\| = \sup \{|\mathbf{v}(s)| \mid s \in \mathcal{S}\}$$

For infinite-horizon discounted MDPs, stationary policies are sufficient. This can be proven by induction, using arguments similar to other proofs given here. For a detailed set of proofs, see Puterman [1994].

**Definition 8.5.1.** A policy  $\pi$  is stationary if  $\pi(a_t \mid s_t) = \pi(a_n \mid s_n)$  for all  $n, t$ .

We now present a set of important results that link Markov decision processes to linear algebra.

**Remark 8.5.1.** We can use the Markov chain kernel  $\mathbf{P}$  to write the expected reward vector as

$$\mathbf{v}^\pi = \sum_{t=0}^{\infty} \gamma^t \mathbf{P}_{\mu, \pi}^t \mathbf{r} \quad (8.5.1)$$

*Proof.*

$$\begin{aligned} V^\pi(s) &= \mathbb{E} \left( \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right) \\ &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}(r_t \mid s_0 = s) \\ &= \sum_{t=0}^{\infty} \gamma^t \sum_{i \in \mathcal{S}} \mathbb{P}(s_t = i \mid s_0 = s) \mathbb{E}(r_t \mid s_t = i). \end{aligned}$$

Since for any distribution vector  $\mathbf{p}$  over  $\mathcal{S}$ , we have  $\mathbb{E}_{\mathbf{p}} r_t = \mathbf{p}^\top \mathbf{r}$ , the result follows.  $\square$

**Exercise 5.** Show that the expected discounted total reward of any given policy is equal to the expected undiscounted total reward with a finite, but random horizon  $T$ . In particular, let  $T$  be distributed according to a geometric distribution on  $\{1, 2, \dots\}$  with parameter  $1 - \gamma$ . Then show that:

$$\mathbb{E} \lim_{T \rightarrow \infty} \sum_{k=0}^T \gamma^k r_k = \mathbb{E} \left( \sum_{k=0}^T r_k \mid T \sim \text{Geom}(1 - \gamma) \right).$$

The value of a particular policy can be expressed as a linear equation. This is an important result, as it has led to a number of successful algorithms that employ linear theory.

**Theorem 8.5.1.** For any stationary policy  $\pi$ ,  $\mathbf{v}^\pi$  is the unique solution of

$$\mathbf{v} = \mathbf{r} + \gamma \mathbf{P}_{\mu, \pi} \mathbf{v}. \quad (8.5.2)$$

In addition, the solution is:

$$\mathbf{v}^\pi = (\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi})^{-1} \mathbf{r}. \quad (8.5.3)$$

To prove this we will need the following important theorem.

**Theorem 8.5.2.** For any bounded linear transformation  $\mathbf{A} : S \rightarrow S$  on a normed linear space  $S$  (i.e., there is  $c < \infty$  s.t.  $\|\mathbf{A}x\| := \sup_i \sum_j a_{i,j} \leq c\|x\|$  for all  $x \in S$  with spectral radius  $\sigma(\mathbf{A}) \lim_{n \rightarrow \infty} \|\mathbf{A}^n\|^{1/n} < 1$ ),  $\mathbf{A}^{-1}$  exists and is given by

$$\mathbf{A}^{-1} = \lim_{T \rightarrow \infty} \sum_{n=0}^T (\mathbf{I} - \mathbf{A})^n. \quad (8.5.4)$$

*Proof of Theorem 8.5.1.* First note that by manipulating the infinite sum in Remark 8.5.1, one obtains  $\mathbf{r} = (\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi}) \mathbf{v}^\pi$ . Since  $\|\gamma \mathbf{P}_{\mu, \pi}\| < 1 \cdot \|\mathbf{P}_{\mu, \pi}\| = 1$ , the inverse

$$(\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi})^{-1} = \lim_{n \rightarrow \infty} \sum_{t=0}^n (\gamma \mathbf{P}_{\mu, \pi})^t$$

exists by Theorem 8.5.2. It follows that

$$\mathbf{v} = (\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi})^{-1} \mathbf{r} = \sum_{t=0}^{\infty} \gamma^t \mathbf{P}_{\mu, \pi}^t \mathbf{r} = \mathbf{v}^\pi,$$

where the last step is by Remark 8.5.1 again.  $\square$

### 8.5.3 Optimality equations

Let us now look at the backwards induction algorithms in terms of operators. We introduce the operator of a policy, which is the one-step backwards induction operation for a fixed policy, and the Bellman operator, which is the equivalent operator for the optimal policy. If a value function is optimal, then it satisfies the Bellman optimality equation.

**Definition 8.5.2** (Policy and Bellman operator). *The linear operator of a policy  $\pi$  is:*

$$\mathcal{L}_\pi \mathbf{v} \triangleq \mathbf{r} + \gamma \mathbf{P}_\pi \mathbf{v} \quad (8.5.5)$$

*The (non-linear) Bellman operator in the space of value functions  $\mathcal{V}$  is defined as:*

$$\mathcal{L} \mathbf{v} \triangleq \sup_{\pi} \{\mathbf{r} + \gamma \mathbf{P}_\pi \mathbf{v}\}, \quad \mathbf{v} \in \mathcal{V} \quad (8.5.6)$$

We now show that the Bellman operator satisfies the following monotonicity properties with respect to an arbitrary value vector  $\mathbf{v}$ .

**Theorem 8.5.3.** *Let  $\mathbf{v}^* := \sup_{\pi} \mathbf{v}^\pi$ . Then for any bounded  $\mathbf{r}$ , it holds that for  $\mathbf{v} \in \mathcal{V}$ :*

- (1) *If  $\mathbf{v} \geq \mathcal{L} \mathbf{v}$ , then  $\mathbf{v} \geq \mathbf{v}^*$ .*
- (2) *If  $\mathbf{v} \leq \mathcal{L} \mathbf{v}$ , then  $\mathbf{v} \leq \mathbf{v}^*$ .*
- (3) *If  $\mathbf{v} = \mathcal{L} \mathbf{v}$ , then  $\mathbf{v}$  is unique and  $\mathbf{v} = \sup_{\pi} \mathbf{v}^\pi$ . Therefore,  $\mathbf{v} = \mathcal{L} \mathbf{v}$  is called the Bellman optimality equation.*

*Proof.* We first prove (1). A simple proof by induction over  $n$  shows that for any  $\pi$

$$\mathbf{v} \geq \mathbf{r} + \gamma \mathbf{P}_\pi \mathbf{v} \geq \sum_{k=0}^{n-1} \gamma^k \mathbf{P}_\pi^k \mathbf{r} + \gamma^n \mathbf{P}_\pi^n \mathbf{v}.$$

Since  $\mathbf{v}^\pi = \sum_{t=0}^{\infty} \gamma^t \mathbf{P}_\pi^t \mathbf{r}$  it follows that

$$\mathbf{v} - \mathbf{v}^\pi \geq \gamma^n \mathbf{P}_\pi^n \mathbf{v} - \sum_{k=n}^{\infty} \gamma^k \mathbf{P}_\pi^k \mathbf{r}.$$

The first-term on the right-hand side can be bounded by arbitrary  $\epsilon/2$  for large enough  $n$ . Also note that

$$\sum_{k=n}^{\infty} \gamma^k \mathbf{P}_{\pi}^k \mathbf{r} \geq -\frac{\gamma^n \mathbf{e}}{1-\gamma},$$

with  $\mathbf{e}$  being a unit vector, so this can be bounded by  $\epsilon/2$  as well. So for any  $\pi, \epsilon > 0$ :

$$\mathbf{v} \geq \mathbf{v}^{\pi} - \epsilon,$$

so

$$\mathbf{v} \geq \sup_{\pi} \mathbf{v}^{\pi}.$$

An equivalent argument shows that

$$\mathbf{v} \leq \mathbf{v}^{\pi} + \epsilon,$$

proving (2). Putting together (1) and (2) gives (3).  $\square$

We eventually want show that repeated application of the Bellman operator converges to the optimal value. As a preparation, we need the following theorem.

**Theorem 8.5.4** (Banach Fixed-Point theorem). *Suppose  $\mathcal{S}$  is a Banach space (i.e. a complete normed linear space) and  $T : \mathcal{S} \rightarrow \mathcal{S}$  is a contraction mapping (i.e.  $\exists \gamma \in [0, 1)$  s.t.  $\|Tu - Tv\| \leq \gamma\|u - v\|$  for all  $u, v \in \mathcal{S}$ ). Then*

- *there is a unique  $u^* \in \mathcal{S}$  s.t.  $Tu^* = u^*$ , and*
- *for any  $u^0 \in \mathcal{S}$  the sequence  $\{u^n\}$ :*

$$u^{n+1} = Tu^n = T^{n+1}u^0$$

*converges to  $u^*$ .*

*Proof.* For any  $m \geq 1$

$$\begin{aligned} \|u^{n+m} - u^n\| &\leq \sum_{k=0}^{m-1} \|u^{n+k+1} - u^{n+k}\| = \sum_{k=0}^{m-1} \|T^{n+k}u^1 - T^{n+k}u^0\| \\ &\leq \sum_{k=0}^{m-1} \gamma^{n+k} \|u^1 - u^0\| = \frac{\gamma^n(1-\gamma^m)}{1-\gamma} \|u^1 - u^0\|. \end{aligned}$$

$\square$

**Theorem 8.5.5.** *For  $\gamma \in [0, 1)$  the Bellman operator  $\mathcal{L}$  is a contraction mapping in  $\mathcal{V}$ .*

*Proof.* Let  $\mathbf{v}, \mathbf{v}' \in \mathcal{V}$ . Consider  $s \in \mathcal{S}$  such that  $\mathcal{L}\mathbf{v}(s) \geq \mathcal{L}\mathbf{v}'(s)$ , and let

$$a_s^* \in \arg \max_{a \in \mathcal{A}} \left\{ r(s) + \sum_{j \in \mathcal{S}} \gamma p_{\mu}(j \mid s, a) \mathbf{v}(j) \right\}.$$

Then

$$\begin{aligned}
0 &\leq \mathcal{L}\mathbf{v}(s) - \mathcal{L}\mathbf{v}'(s) \leq \sum_{j \in \mathcal{S}} \gamma p(j \mid s, a_s^*) \mathbf{v}(j) - \sum_{j \in \mathcal{S}} \gamma p(j \mid s, a_s^*) \mathbf{v}'(j) \\
&= \gamma \sum_{j \in \mathcal{S}} p(j \mid s, a_s^*) [\mathbf{v}(j) - \mathbf{v}'(j)] \\
&\leq \gamma \sum_{j \in \mathcal{S}} p(j \mid s, a_s^*) \|\mathbf{v} - \mathbf{v}'\| = \gamma \|\mathbf{v} - \mathbf{v}'\|.
\end{aligned}$$

Repeating the argument for  $s$  such that  $\mathcal{L}\mathbf{v}(s) \leq \mathcal{L}\mathbf{v}'(s)$ , we obtain

$$|\mathcal{L}\mathbf{r}(s) - \mathcal{L}\mathbf{r}'(s)| \leq \gamma \|\mathbf{r} - \mathbf{r}'\|.$$

Taking the supremum, we the required result follows.  $\square$

**Theorem 8.5.6.** *For discrete  $\mathcal{S}$ , bounded  $\mathbf{r}$ , and  $\gamma \in [0, 1)$*

- (i) *there is a unique  $\mathbf{v}^* \in \mathcal{V}$  such that  $\mathcal{L}\mathbf{v}^* = \mathbf{v}^*$  and such that  $\mathbf{v}^* = V_\mu^*$ ,*
- (ii) *for any stationary policy  $\pi$ , there is a unique  $\mathbf{v} \in \mathcal{V}$  such that  $\mathcal{L}_\pi \mathbf{v} = \mathbf{v}$  and  $\mathbf{v} = V_\mu^\pi$ .*

*Proof.* As the Bellman operator  $\mathcal{L}$  is a contraction by Theorem 8.5.5, application of the fixed-point Theorem 8.5.4 shows that there is a unique  $\mathbf{v}^* \in \mathcal{V}$  such that  $\mathcal{L}\mathbf{v}^* = \mathbf{v}^*$ . This is also the optimal value function due to Theorem 8.5.5. The second part of the theorem follows from the first part when considering only a single policy  $\pi$  (which then is optimal).  $\square$

### 8.5.4 MDP Algorithms

Let us now look at three basic algorithms for solving a known Markov decision process. The first, *value iteration*, is a simple extension of the backwards induction algorithm to the infinite horizon case.

#### Value iteration

In this version of the algorithm, we assume that rewards are dependent only on the state. An algorithm for the case where reward only depends on the state can be obtained by replacing  $r(s, a)$  with  $r(s)$ .

---

#### Algorithm 6 Value iteration

---

```

Input  $\mu, \mathcal{S}$ .
Initialise  $\mathbf{v}_0 \in \mathcal{V}$ .
for  $n = 1, 2, \dots$  do
  for  $s \in \mathcal{S}_n$  do
     $\pi_n(s) = \arg \max_{a \in \mathcal{A}} \{r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_\mu(s' \mid s, a) \mathbf{v}_{n-1}(s')\}$ 
     $\mathbf{v}_n(s) = r(s, \pi_n(s)) + \gamma \sum_{s' \in \mathcal{S}} P_\mu(s' \mid s, \pi_n(s)) \mathbf{v}_{n-1}(s')$ 
  end for
  break if termination-condition is met
end for
Return  $\pi_n, V_n$ .

```

---



**Theorem 8.5.7.** *The value iteration algorithm satisfies*

- $\lim_{n \rightarrow \infty} \|\mathbf{v}_n - V^*\| = 0$ .
- For each  $\epsilon > 0$  there exists  $N_\epsilon < \infty$  such that for all  $n \geq N_\epsilon$

$$\|\mathbf{v}_{n+1} - \mathbf{v}_n\| \leq \epsilon(1 - \gamma)/2\gamma. \quad (8.5.7)$$

- For  $n \geq N_\epsilon$  the policy  $\pi_\epsilon$  that takes action  $\arg \max_a r(s, a) + \gamma \sum_j p(j|s, a) \mathbf{v}_n(s')$  is  $\epsilon$ -optimal.
- $\|\mathbf{v}_{n+1} - V_\mu^*\| < \epsilon/2$  for  $n > N_\epsilon$ .

*Proof.* The first two statements follow from the fixed-point Theorem 8.5.4. Now note that

$$\|V_\mu^{\pi_\epsilon} - V_\mu^*\| \leq \|V_\mu^{\pi_\epsilon} - \mathbf{v}_n\| + \|\mathbf{v}_n - V_\mu^*\|$$

We can bound these two terms easily:

$$\begin{aligned} \|V_\mu^{\pi_\epsilon} - \mathbf{v}_{n+1}\| &= \|\mathcal{L}_{\pi_\epsilon} V_\mu^{\pi_\epsilon} - \mathbf{v}_{n+1}\| \\ &\leq \|\mathcal{L}_{\pi_\epsilon} V_\mu^{\pi_\epsilon} - \mathcal{L} \mathbf{v}_{n+1}\| + \|\mathbf{v}_{n+1} - \mathbf{v}_{n+1}\| \\ &= \|\mathcal{L}_{\pi_\epsilon} V_\mu^{\pi_\epsilon} - \mathcal{L}_{\pi_\epsilon} \mathbf{v}_{n+1}\| + \|\mathcal{L} \mathbf{v}_{n+1} - \mathcal{L} \mathbf{v}_n\| \\ &\leq \gamma \|V_\mu^{\pi_\epsilon} - \mathbf{v}_{n+1}\| + \|\mathbf{v}_{n+1} - \mathbf{v}_n\|. \end{aligned}$$

An analogous argument gives the same bound for the second term  $\|\mathbf{v}_n - V_\mu^*\|$ . Then, rearranging we obtain

$$\|V_\mu^{\pi_\epsilon} - \mathbf{v}_{n+1}\| \leq \frac{\gamma}{1 - \gamma} \|\mathbf{v}_{n+1} - \mathbf{v}_n\|, \quad \|\mathbf{v}_{n+1} - V_\mu^*\| \leq \frac{\gamma}{1 - \gamma} \|\mathbf{v}_{n+1} - \mathbf{v}_n\|,$$

and the third and fourth statements follow from the second statement.  $\square$

**Theorem 8.5.8** (Value iteration monotonicity). *Let  $\mathcal{V}$  be a value function space with Bellman operator  $\mathcal{L}$ . Then:*

1. Let  $\mathbf{v}, \mathbf{v}' \in \mathcal{V}$  with  $\mathbf{v}' \geq \mathbf{v}$ . Then  $\mathcal{L} \mathbf{v}' \geq \mathcal{L} \mathbf{v}$ .
2. Let  $\mathbf{v}_{n+1} = \mathcal{L} \mathbf{v}_n$ . If there is an  $N$  s.t.  $\mathcal{L} \mathbf{v}_N \leq \mathbf{v}_N$ , then  $\mathcal{L} \mathbf{v}_{N+k} \leq \mathbf{v}_{N+k}$  for all  $k \geq 0$  and similarly for  $\geq$ .

*Proof.* Let  $\pi \in \arg \max_\pi \mathbf{r} + \gamma \mathbf{P}_{\mu, \pi} \mathbf{v}$ . Then

$$\mathcal{L} \mathbf{v} = \mathbf{r} + \gamma \mathbf{P}_{\mu, \pi} \mathbf{v} \leq \mathbf{r} + \gamma \mathbf{P}_{\mu, \pi} \mathbf{v}' \leq \max_{\pi'} \mathbf{r} + \gamma \mathbf{P}_{\mu, \pi'} \mathbf{v}',$$

where the first inequality is due to the fact that  $\mathbf{P} \mathbf{v} \geq \mathbf{P} \mathbf{v}'$  for any  $\mathbf{P}$ .

For the second part,

$$\mathcal{L} \mathbf{v}_{N+k} = \mathbf{v}_{N+k+1} = \mathcal{L}^k \mathcal{L} \mathbf{v}_N \leq \mathcal{L}^k \mathbf{v}_N = \mathbf{v}_{N+k}.$$

since  $\mathcal{L} \mathbf{v}_N \leq \mathbf{v}_N$  by assumption and consequently  $\mathcal{L}^k \mathcal{L} \mathbf{v}_N \leq \mathcal{L}^k \mathbf{v}_N$  by part one of the theorem.  $\square$

Thus, value iteration converges monotonically to  $V_\mu^*$  if  $\mathbf{r}^0$  satisfies condition 1 of the theorem. If  $r \geq 0$ , it is sufficient to set  $\mathbf{r}^0 = \mathbf{0}$ . Then  $\mathbf{v}$  is always a lower bound on the optimal value function.

**Theorem 8.5.9.** *Value iteration converges linearly at rate  $\gamma$  and  $O(\gamma^n)$ . In addition, for  $r \in [0, 1]$  and  $\mathbf{r}^0 = \mathbf{0}$*

$$\begin{aligned}\|\mathbf{v}_n - V_\mu^*\| &\leq \frac{\gamma^n}{1 - \gamma}, \\ \|V_\mu^{\pi_n} - V_\mu^*\| &\leq \frac{2\gamma^n}{1 - \gamma}.\end{aligned}$$

*Proof.* The first part follows from the contraction property (Theorem 8.5.5):

$$\|\mathbf{v}_{n+1} - \mathbf{v}^*\| = \|\mathcal{L}\mathbf{v}_n - \mathcal{L}\mathbf{v}^*\| \leq \gamma\|\mathbf{v}_n - \mathbf{v}^*\|. \quad (8.5.8)$$

Now divide by  $\gamma^n$  to obtain the  $O(\gamma^n)$  property.  $\square$

### Policy iteration

Unlike value iteration, *policy iteration* attempts to iteratively improve a given policy. At each iteration, it calculates the value of the current policy. At the next step, it tries to obtain an improvement by calculating a policy that is greedy with respect to the previous value function. The algorithm described below can be extended to the case when the reward also depends on the action, by replacing  $\mathbf{r}$  with the policy-dependent reward vector  $\mathbf{r}_\pi$ . The following theorem can also be easily extended to this case.

---

#### Algorithm 7 Policy iteration

---

```

Input  $\mu, \mathcal{S}$ .
Initialise  $\mathbf{v}_0$ .
for  $n = 1, 2, \dots$  do
     $\pi_{n+1} = \arg \max_\pi \{\mathbf{r} + \gamma \mathbf{P}_\pi \mathbf{v}_n\}$  // policy improvement
     $\mathbf{v}_{n+1} = V_\mu^{\pi_{n+1}}$  // policy evaluation
    break if  $\pi_{n+1} = \pi_n$ .
end for
Return  $\pi_n, \mathbf{v}_n$ .
```

---

**Theorem 8.5.10.** *Let  $\mathbf{v}_n, \mathbf{v}_{n+1}$  be the value vectors generated by policy iteration. Then  $\mathbf{v}_n \leq \mathbf{v}_{n+1}$ .*

*Proof.* From the policy improvement step

$$\mathbf{r} + \gamma \mathbf{P}_{\pi_{n+1}} \mathbf{v}_n \geq \mathbf{r} + \gamma \mathbf{P}_{\pi_n} \mathbf{v}_n = \mathbf{v}_n$$

where the equality is due to the policy evaluation step for  $\pi_n$ . Rearranging, we get  $\mathbf{r} \geq (\mathbf{I} - \gamma \mathbf{P}_{\pi_{n+1}}) \mathbf{v}_n$  and hence

$$(\mathbf{I} - \gamma \mathbf{P}_{\pi_{n+1}})^{-1} \mathbf{r} \geq \mathbf{v}_n,$$

noting that the inverse is positive. Since the left side equals  $\mathbf{v}_{n+1}$  by the policy evaluation step for  $\pi_{n+1}$ , the theorem follows.  $\square$

**Corollary 8.5.1.** *If  $\mathcal{S}, \mathcal{A}$  are finite, then policy iteration terminates after a finite number of iterations and returns an optimal policy.*

*Proof.* There is only a finite number of policies, and since policies in policy iteration are monotonically improving, the algorithm must stop after finitely many iterations. Finally, the last iteration satisfies

$$\mathbf{v}_n = \max_{\pi} \mathbf{r} + \gamma \mathbf{P}_{\pi} \mathbf{v}_n. \quad (8.5.9)$$

Thus  $\mathbf{v}_n$  solves the optimality equation.  $\square$

### Modified policy iteration

The astute reader will have noticed that it may be not necessary to fully evaluate the improved policy. In fact, we can take advantage of that to speed up policy iteration. Thus, a simple variant of policy iteration involves doing only a  $k$ -step update for the policy evaluation step. For  $k = 1$ , the algorithm becomes identical to value iteration.

---

#### Algorithm 8 Modified policy iteration

---

```

Input  $\mu, \mathcal{S}$ .
Initialise  $\mathbf{v}_0$ .
for  $n = 1, 2, \dots$  do
   $\pi_n = \arg \max_{\pi} \mathbf{r} + \gamma \mathbf{P}_{\pi} \mathbf{v}_{n-1}$  // policy improvement
   $\mathbf{v}_n = \mathcal{L}_{\pi_n}^k \mathbf{v}_{n-1}$  // partial policy evaluation
  break if  $\pi_n = \pi_{n+1}$ .
end for
Return  $\pi_n, \mathbf{v}_n$ .

```

---

### Geometric view

**Definition 8.5.3.** *Difference operator* The difference operator is defined as

$$\mathcal{B}\mathbf{v} \triangleq \max_{\pi} \{\mathbf{r} + (\gamma \mathbf{P}_{\pi} - \mathbf{I})\mathbf{v}\} = \mathcal{L}\mathbf{v} - \mathbf{v}. \quad (8.5.10)$$

Thus the optimality equation can be rewritten as

$$\mathcal{B}\mathbf{v} = 0. \quad (8.5.11)$$

For any  $\mathbf{v} \in \mathcal{V}$ , define:

$$\Pi_{\mathbf{v}} \triangleq \arg \max_{\pi \in \Pi} \{\mathbf{r} + (\gamma \mathbf{P}_{\pi} - \mathbf{I})\mathbf{v}\}$$

to be the set of  $\mathbf{v}$ -improving policies.

**Theorem 8.5.11.** *For any  $\mathbf{v}, \mathbf{v}' \in \mathcal{V}$  and  $\pi \in \Pi_{\mathbf{v}}$*

$$\mathcal{B}\mathbf{v} \geq \mathcal{B}\mathbf{v} + (\gamma \mathbf{P}_{\pi} - \mathbf{I})(\mathbf{v}' - \mathbf{v}). \quad (8.5.12)$$

*Proof.* By definition,  $\mathcal{B}\mathbf{v}' \geq \mathbf{r} + (\gamma \mathbf{P}_{\pi} - \mathbf{I})\mathbf{v}'$ , while  $\mathcal{B}\mathbf{v} = \mathbf{r} + (\gamma \mathbf{P}_{\pi} - \mathbf{I})\mathbf{v}$ . Subtracting the latter from the former gives the result.  $\square$

Equation (8.5.12) is similar to the convexity equation. In fact, we can have a nice geometric view of this.

**Theorem 8.5.12.** Let  $\{\mathbf{v}_n\}$  be the sequence of value vectors obtained from policy iteration. Then for any  $\pi \in \Pi_{\mathbf{v}_n}$ ,

$$\mathbf{v}_{n+1} = \mathbf{v}_n - (\gamma \mathbf{P}_\pi - \mathbf{I})^{-1} \mathcal{B} \mathbf{v}_n. \quad (8.5.13)$$

*Proof.* By definition, we have for  $\pi \in \Pi_{\mathbf{v}_n}$

$$\begin{aligned} \mathbf{v}_{n+1} &= (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r} - \mathbf{v}_n + \mathbf{v}_n \\ &= (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} [\mathbf{r} - (\mathbf{I} - \gamma \mathbf{P}_\pi) \mathbf{v}_n] + \mathbf{v}_n. \end{aligned}$$

Since  $\mathbf{r} - (\mathbf{I} - \gamma \mathbf{P}_\pi) \mathbf{v}_n = \mathcal{B} \mathbf{v}_n$  the claim follows.  $\square$

### Temporal-Difference Policy Iteration

*Temporal-difference policy iteration* replaces the next-step value with an approximation  $\mathbf{v}_n$ . First, policy improvement is used to obtain the next policy given our approximation:

$$\mathcal{L}_{\pi_{n+1}} \mathbf{v}_n = \mathcal{L} \mathbf{v}_n. \quad (8.5.14)$$

The method uses the temporal difference error, defined as:

$$d_n(i, j) = \mathbf{v}_n(i) - [\mathbf{r}(i) + \gamma \mathbf{v}_n(j)]. \quad (\text{temporal difference error})$$

This can be seen as the difference in the estimate when we move from state  $i$  to state  $j$ . Note the similarity to the difference operator in modified policy iteration. The idea of the temporal-difference policy iteration is to use  $d_n$  as the one-stage reward for a  $\lambda$ -discounted problem. Then one can write the value vector and update for this problem as

$$\boldsymbol{\tau}_n(i) = \sum_{t=0}^{\infty} \mathbb{E}_{\pi_n, \mu} [(\gamma \lambda)^t d_n(s_t, s_{t+1}) \mid s_0 = i], \quad (8.5.15)$$

$$\mathbf{v}_{n+1} = \mathbf{v}_n + \boldsymbol{\tau}_n. \quad (8.5.16)$$

Summarizing, we obtain the following algorithm:

---

#### Algorithm 9 Temporal-Difference Policy Iteration

---

```

Input  $\mu, \mathcal{S}, \lambda$ .
Initialise  $\mathbf{v}_0$ .
for  $n = 1, 2, \dots$  do
     $\pi_n = \arg \max_{\pi} \mathbf{r} + \gamma \mathbf{P}_{\pi} \mathbf{v}_{n-1}$  // policy improvement
     $\mathbf{v}_n = \mathbf{v}_{n-1} + \boldsymbol{\tau}_k$  // temporal difference update.
    break if  $\pi_n = \pi_{n=1}$ .
end for
Return  $\pi_n, \mathbf{v}_n$ .
```

---

In fact,  $\mathbf{v}_{n+1}$  is the unique fixed point of the following equation:

$$\mathcal{D}_n \mathbf{v} \triangleq (1 - \lambda) \mathcal{L}_{\pi_{n+1}} \mathbf{v}_n + \lambda \mathcal{L}_{\pi_{n+1}} \mathbf{v}. \quad (\text{fixed point})$$

In other words, the new value vector is moved only partially towards the direction of the Bellman update. In fact, for  $\lambda = 1$ , this becomes identical to standard policy iteration. For  $\lambda = 0$ , one obtains standard value iteration.

### Linear programming

We shall now formulate the problem of finding an optimal value function for a given Markov decision process in terms of a linear program. Recall that if

$$\mathbf{v} \geq \mathcal{L}\mathbf{v}$$

then  $\mathbf{v} \geq V_\mu^*$ . In order to transform this into a linear program, we must first define a scalar function to minimise. We can do this by selecting some arbitrary distribution on the states  $\mathbf{y} \in \Delta^{|\mathcal{S}|}$ . Then we can write the following linear program.

#### Primal linear program

$$\min_{\mathbf{v}} \mathbf{y}^\top \mathbf{v},$$

such that

$$\mathbf{v}(s) - \gamma \mathbf{p}_{s,a}^\top \mathbf{v} \geq r(s, a), \quad \forall a \in \mathcal{A}, s \in \mathcal{S}.$$

Note that the inequality condition is equivalent to  $\mathbf{v} \geq \mathcal{L}\mathbf{v}$ . Consequently, the problem is to find the smallest  $\mathbf{v}$  that satisfies this inequality. When  $\mathcal{A}, \mathcal{S}$  are finite, it is easy to see that this will be the optimal value function and the Bellman equation is satisfied.

It also pays to look at the dual linear program, which is in terms of a maximisation. This time, instead of finding the minimal upper bound on the value function, we find the maximal cumulative discounted state-action visits  $x(s, a)$  that are consistent with the transition kernel of the process.

#### Dual linear program

$$\max_x \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r(s, a) x(s, a)$$

such that  $x \in \mathbb{R}_+^{|\mathcal{S} \times \mathcal{A}|}$  and

$$\sum_{a \in \mathcal{A}} x(j, a) - \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \gamma p(j | s, a) x(s, a) = y(j) \quad \forall j \in \mathcal{S}.$$

with  $\mathbf{y} \in \Delta^{|\mathcal{S}|}$ .

In this case,  $x$  can be interpreted as the discounted sum of state-action visits, as proved by the following theorem.

#### Theorem 8.5.13.

$$x_\pi(s, a) = \mathbb{E}_{\pi, \mu} \left\{ \sum \gamma^n \mathbb{I}\{s_t = s, a_t = a \mid s_0 \sim y\} \right\}$$

is a feasible solution to the dual problem. On the other hand, if  $x$  is a feasible solution to the dual problem then  $\sum_a x(s, a) > 0$ . Finally, if we define the

strategy

$$\pi(a \mid s) = \frac{x(s, a)}{\sum_{a' \in \mathcal{A}} x(s, a')}$$

then  $x_\pi = x$  is a feasible solution.

The equality condition ensures that  $x$  is consistent with the transition kernel of the Markov decision process. Consequently, the program can be seen as search among all possible cumulative state-action distributions to find the one giving the highest total reward.

### Summary

#### Markov decision processes

can represent : Shortest path problems, stopping problems, experiment design problems, multi-armed bandit problems, reinforcement learning problems.

#### Backwards induction (aka value iteration)

- In the class of dynamic programming algorithms.
- Tractable when either the state space  $\mathcal{S}$  or the horizon  $T$  are small (finite).

#### Optimal decisions and Bayesian reinforcement learning

- A known environment is represented as an MDP.
- Bandit problems can be solved by representing them as infinite-state MDPs.
- In general, an unknown environment can be represented as a distribution over MDPs.
- The decision problem can again be formulated as an infinite-state MDP.

## 8.6 Further reading

See the last chapter of [DeGroot, 1970] for further information on the MDP formulation of bandit problems in the decision theoretic setting. This was explored in more detail in Duff's PhD thesis [Duff, 2002]. When the number of (information) states in the bandit problem is finite, Gittins [1989] has proven

that it is possible to formulate simple index policies. However, this is not generally applicable. Easily computable, near-optimal heuristic strategies for bandit problems will be given in Chapter 12. The decision-theoretic solution to the unknown MDP problem will be given in Chapter 11.

Further theoretical background on Markov decision processes, including many of the theorems in Section 8.5, can be found in [Puterman, 1994]. Chapter 2 of Bertsekas and Tsitsiklis [1996] gives a quick overview of MDP theory from the operator perspective. The introductory reinforcement learning book of Sutton and Barto [1998] also explains the basic Markov decision process framework.





## Chapter 9

# Reinforcement learning and stochastic approximation

## 9.1 Introduction

In this chapter, we consider the actual problem of reinforcement learning. Up to now, we have only examined a solution method for bandit problems, which are a special case of the general reinforcement learning problem. The Bayesian decision-theoretic solution is to *reduce* the bandit problem to a *Markov decision process* which can then be solved with backwards induction.

We also have seen that Markov decision processes can be used to *describe environments* in more general reinforcement learning problems. When our knowledge of the MDP describing these problems is perfect, then we can employ a number of standard algorithms to find the optimal policy. However, in the actual reinforcement learning problem, the model of the environment is *unknown*. However, as we shall see later, ideas from both cases can be used to solve the general reinforcement learning problem.

The main idea we explore in this chapter is that we can perform approximate dynamic programming algorithms, by replacing the actual unknown dynamics of the Markov decision process with estimates. The estimates can be improved by drawing samples from the actual environment, either by actually acting within the environment or using a simulator. In both cases we end up with a number of algorithms that can be used for reinforcement learning. Although may not be performing as well as the Bayes-optimal solution, these have a low enough computational complexity that they are worth investigating in practice.

### 9.1.1 Bandit problems

#### The stochastic $n$ -armed bandit problem

Let us return to the example of bandit problems. As before, we have  $n$  actions corresponding to probability distributions  $P_i$  on the real numbers.

$$\mathcal{P} = \{P_i \mid i = 1, \dots, n\}.$$

At each time-step  $t$  we select an action  $a_t$ , obtaining a random reward distributed according to:

$$r_t \mid a_t = i \sim P_i.$$

Our objective is to find a policy  $\pi$  maximising the expected total reward.

$$\mathbb{E}_\pi U_t = \mathbb{E}_\pi \sum_{k=t}^T r_k, \quad a_t^* \triangleq \max \{\mathbb{E}(r_t \mid a_t = i) \mid i = 1, \dots, n\}.$$

Had we known the distribution, we could simply always the maximising action, as the expected reward of the  $i$ -th action can be easily calculated from  $P_i$  and the reward only depends on our current action. The situation is similar when  $\mathcal{P}$  is a parametric family unknown parameter  $\omega^*$ , outlined below.

$$\mathcal{P} = \{P_i(\cdot \mid \omega) \mid \omega \in \Omega\}, \quad r_t \mid a_t = i, \omega^* = \omega \sim P_i(r \mid \omega^*). \quad (9.1.1)$$

If in addition we have a subjective belief  $\xi$  over  $\Omega$ , we could (as explained in Sec. 8.2) in principle calculate the policy maximising the  $\xi$ -expected utility:

$$\mathbb{E}_\xi^\pi U_t = \mathbb{E}_\xi^\pi \sum_{k=t}^T r_k. \quad (9.1.2)$$

This of course will now have to be a history-dependent policy. In the remainder of this section, we shall examine algorithms which eventually converge to the optimal action, but for which we cannot always guarantee a good initial behaviour.

### 9.1.2 Estimation and Robbins-Monro approximation

The basic idea of the Robbins-Monro stochastic approximation algorithm [Robbins and Monro, 1951] is to maintain a set of *point estimates* of a parameter we want to approximate and perform *random* steps that on average move towards the solution, in a way to be made more precise later. It can in fact be seen as a generalisation of stochastic gradient descent.

---

**Algorithm 10** Robbins-Monro bandit algorithm

---

```

1: input Step-sizes  $(\alpha_t)_t$ , initial estimates  $(\mu_{i,0})_i$ , policy  $\pi$ .
2: for  $t = 1, \dots, T$  do
3:   Take action  $a_t = i$  with probability  $\pi(i \mid a_1, \dots, a_{t-1}, r_1, \dots, r_{t-1})$ .
4:   Observe reward  $r_t$ .
5:    $\mu_{t,i} = \alpha_{i,t} r_t + (1 - \alpha_{i,t}) \mu_{i,t-1}$       // estimation step
6:    $\mu_{t,i} = \mu_{j,t-1}$  for  $j \neq i$ .
7: end for
8: return  $\mu_T$ 

```

---

A simple Robbins-Monro algorithm for the  $n$ -armed bandit problem is given in Algorithm 10. The input is a particular policy  $\pi$ , that gives us a probability over the next actions given the observed history, a set of initial estimates  $\mu_{i,0}$  for the bandit means, and a sequence of step sizes  $\alpha$ .

If you examine the updates carefully, you will be able to find what the cost function you are trying to minimise is. This simple update rule can be seen as trying to minimise the expected squared error between your estimated reward, and the random reward obtained by each bandit. Consequently, the variance of the reward of each bandit plays an important role.

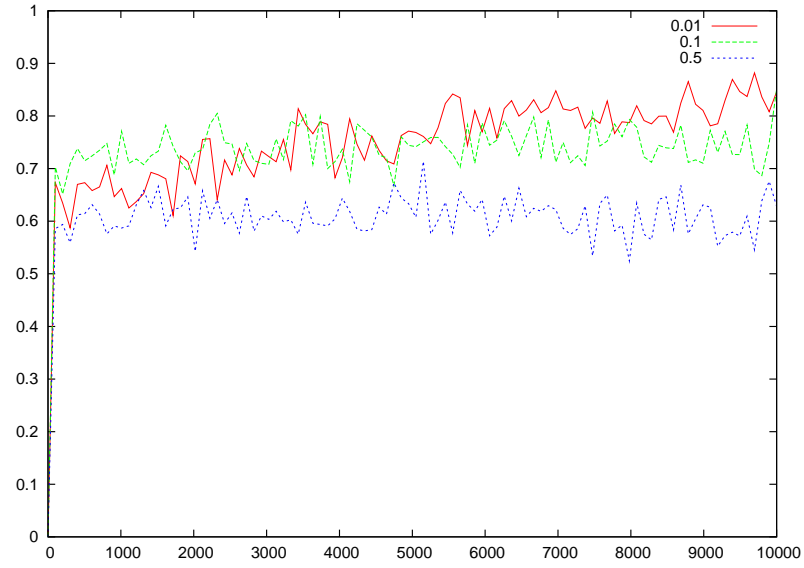
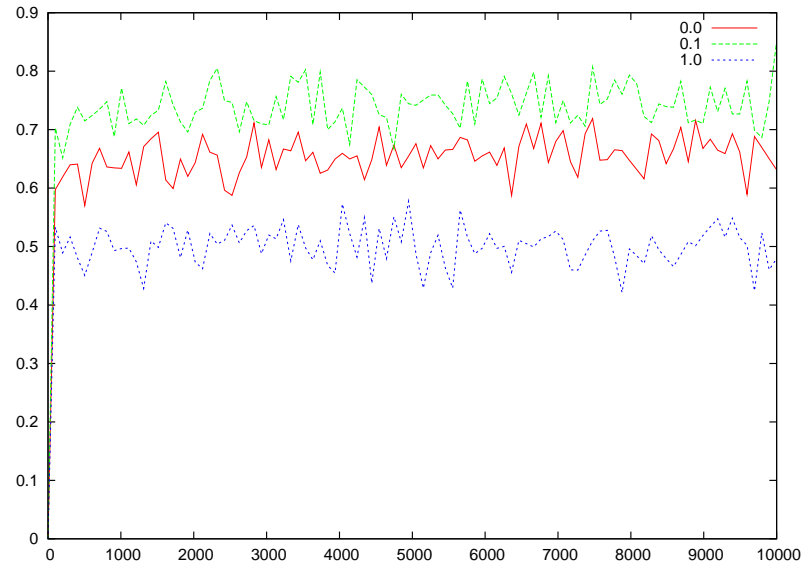
The step-sizes  $\alpha$  must obey certain constraints in order for the algorithm to work, in particular it must decay neither too slowly, nor too fast. There is one particular choice, for which our estimates are in fact the mean estimate of the expected value of the reward for each action  $i$ , which is a natural choice if the bandits are stationary.

The other question is what policy to use to take actions. We must take all actions often enough, so that we have good estimates for the expected reward of every bandit. One simple way to do it is to play the apparently best bandit most of the time, but to sometimes select bandits randomly. This is called  $\epsilon$ -greedy action selection. This ensures that all actions are tried a sufficient number of times.

**Definition 9.1.1.**  $\epsilon$ -greedy action selection

$$\hat{\pi}_\epsilon^* \triangleq (1 - \epsilon_t) \hat{\pi}_t^* + \epsilon_t \mathcal{U}nif(\mathcal{A}), \quad (9.1.3)$$

$$\hat{\pi}_t^*(i) = \mathbb{I} \left\{ i \in \hat{\mathcal{A}}_t^* \right\} / |\hat{\mathcal{A}}_t^*|, \quad \hat{\mathcal{A}}_t^* = \arg \max_{i \in \mathcal{A}} \mu_{t,i} \quad (9.1.4)$$

Figure 9.1:  $\epsilon_t = 0.1$ ,  $\alpha \in \{0.01, 0.1, 0.5\}$ .Figure 9.2:  $\epsilon_t = \epsilon$ ,  $\alpha = 0.1$ .

The main two parameters of the algorithm are randomness  $\epsilon$ -greedy action selection and the step-size. Figures 9.1 and 9.2 show the average reward obtained, if we keep the step size  $\alpha$  or the randomness  $\epsilon$  fixed, respectively. We see that there the choice of values really affects convergence.

For a fixed  $\epsilon$ , we find that larger values of  $\alpha$  tend to give a better result eventually, while smaller values have a better initial performance. This is a

natural trade-off, since large  $\alpha$  appears to “learn” fast, but it also “forgets” quickly. That is, for a large  $\alpha$ , our estimates mostly depend upon the last few rewards observed.

Things are not so clear-cut for the choice of  $\epsilon$ . We see that the choice of  $\epsilon = 0$ , is significantly worse than  $\epsilon = 0.1$ . So, that appears to suggest that there is an optimal level of exploration. How should that be determined? Ideally, we should be able to use the decision-theoretic solution seen earlier, but perhaps a good heuristic way of choosing  $\epsilon$  may be good enough.

### 9.1.3 The theory of the approximation

Here we quickly review some basic results of stochastic approximation theory. Complete proofs can be found in Bertsekas and Tsitsiklis [1996]. The main question here is Consider the algorithm

$$\mu_{t+1} = \mu_t + \alpha_t z_{t+1}. \quad (9.1.5)$$

Here  $\mu_t$  is our estimate,  $\alpha_t$  is a step-size and  $z_t$  is an observation. Let  $h_t = \{\mu_t, z_t, \alpha_t, \dots\}$  be the history of the algorithm.

There are two types of convergence conditions. The first focuses on continuity and smoothness properties, and the second on contraction properties of operators. Here, we concentrate on the first case, which also has applications in other areas.

**Assumption 9.1.1.** *Assume a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  such that:*

(i)  $f(x) \geq 0$  for all  $x \in \mathbb{R}^n$ .

(ii) (Lipschitz derivative)  $f$  is continuously differentiable (i.e. the derivative  $\nabla f$  exists and is continuous) and  $\exists L > 0$  such that:

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|, \quad \forall x, y \in \mathbb{R}^n$$

(iii) (Pseudo-gradient)  $\exists c > 0$  such that:

$$c \|\nabla f(\mu_t)\|^2 \leq -\nabla f(\mu_t)^\top \mathbb{E}(z_{t+1} \mid h_t), \quad \forall t.$$

(iv)  $\exists K_1, K_2 > 0$  such that

$$\mathbb{E}(\|z_{t+1}\|^2 \mid h_t) \leq K_1 + K_2 \|\nabla f(\mu_t)\|^2$$

Condition (ii) is a very basic condition for convergence. It basically ensures that the function is well-behaved, so that gradient-following methods can easily find the minimum. Condition (iii) combines two assumptions in one. Firstly, that expected direction of update always decreases cost, and secondly that the squared norm of the gradient is not too large relative to the size of the update. Finally, condition (iv) ensures that update is bounded in expectation relative to the gradient.

One can see how putting together the last two conditions ensures that the expected direction of our update is correct, and that its norm is bounded.

**Theorem 9.1.1.** *For the algorithm*

$$\mu_{t+1} = \mu_t + \alpha_t z_{t+1},$$

where  $\alpha_t \geq 0$  satisfy

$$\sum_{t=0}^{\infty} \alpha_t = \infty, \quad \sum_{t=0}^{\infty} \alpha_t^2 < \infty, \quad (9.1.6)$$

and under Assumption 9.1.1, with probability 1:

1. The sequence  $\{f(\mu_t)\}$  converges.
2.  $\lim_{t \rightarrow \infty} \nabla f(\mu_t) = 0$ .
3. Every limit point  $\mu^*$  of  $\mu_t$  satisfies  $\nabla f(\mu^*) = 0$ .

### A demonstration

Figure 9.3 demonstrates the convergence, or lack thereof, of our estimates  $\mu_t$  to the expected value of a random variable with mean 0.5, for three different step-size schedules, with update direction:

$$z_{t+1} = x_{t+1} - \mu_t.$$

The first one,  $\alpha_t = 1/t$ , satisfies both assumptions. The second one,  $\alpha_t = 1/\sqrt{t}$ , reduces too slowly, and the third one,  $\alpha_t = t^{-3/2}$ , approaches zero too fast.

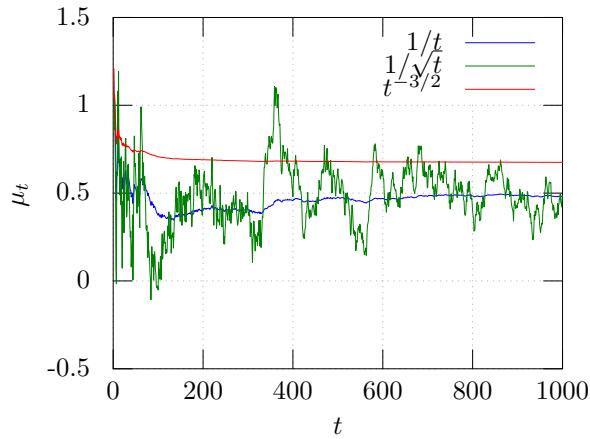


Figure 9.3: Estimation of the expectation of  $x_t \sim \mathcal{N}(0.5, 1)$  using three step-size schedules.

## 9.2 Dynamic problems

The dynamic setting presents one essential difference. Our policy now affects which sequences of states we observe. Otherwise, the algorithmic structure remains the same and is described below.

**Algorithm 11** Generic reinforcement learning algorithm

---

```

1: input Update-rule  $f : \Theta \times \mathcal{S}^2 \times \mathcal{A} \times \mathcal{R} \rightarrow \Theta$ , initial parameters  $\theta_0 \in \Theta$ ,
   policy  $\pi : \mathcal{S} \times \Theta \rightarrow \mathfrak{D}(\mathcal{A})$ .
2: for  $t = 1, \dots, T$  do
3:    $a_t \sim \pi(\cdot \mid \theta_t, s_t)$  // take action
4:   Observe reward  $r_{t+1}$ , state  $s_{t+1}$ .
5:    $\theta_{t+1} = f(\theta_t, s_t, a_t, r_{t+1}, s_{t+1})$  // update estimate
6: end for

```

---

*Questions*

- What should we estimate? For example,  $\theta_t$  could be describing a posterior distribution over MDPs, or a distribution over parameters.
- What policy should we use? For example, we could try and use the Bayes-optimal policy with respect to  $\theta$ , or some heuristic policy.

**Example 9.2.1** (The chain task). *The chain task has two actions and five states, as shown in Fig. 9.4. The reward in the leftmost state is 0.2 and 1.0 in the rightmost state, and zero otherwise. The first action (dashed, blue) takes you to the right, while the second action (solid, red) takes you to the first state. However, there is a probability 0.2 with which the actions have the opposite effects. The value function of the chain task for a discount factor  $\gamma = 0.95$  is shown in Table 9.1.*

The chain task is a very simple, but well-known task, used to test the efficacy of reinforcement learning algorithms. A variant of this task, with action-dependent rewards was used by [Dearden et al., 1998].

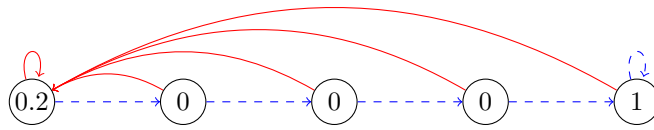


Figure 9.4: The chain task

$s$	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$
$V^*(s)$	7.6324	7.8714	8.4490	9.2090	10.209
$Q^*(s, 1)$	7.4962	7.4060	7.5504	7.7404	8.7404
$Q^*(s, 2)$	7.6324	7.8714	8.4490	9.2090	10.2090

Table 9.1: The chain task's value function for  $\gamma = 0.95$

### 9.2.1 Monte-Carlo policy evaluation and iteration

The simplest algorithm is Monte-Carlo policy evaluation. In the standard setting, we can the value function for every state by approximating the expectation with the sum of rewards obtained over multiple trajectories starting from each state. The  $k$ -th trajectory starts from some initial state  $s_0 = s$  and the next states are sampled as follows

$$a_t^{(k)} \sim \pi(a_t | h_t), r_t^{(k)} \sim \mathbb{P}_\mu(r_t | s_t^{(k)}, a_t^{(k)}) s_{t+1}^{(k)} \sim \mathbb{P}_\mu(s_{t+1} | s_t^{(k)}, a_t^{(k)}). \quad (9.2.1)$$

Then the value function satisfies

$$V_\mu^\pi(s) \triangleq \mathbb{E}_\mu^\pi(U | s_1 = s) \approx \frac{1}{K} \sum_{k=1}^K \sum_{t=1}^T r_t^{(k)},$$

where  $r_t^{(k)}$  is the sequence of rewards obtained from the  $k$ -th trajectory.

---

**Algorithm 12** Stochastic policy evaluation

---

```

1: input Initial parameters  $\mathbf{v}_0$ , Markov policy  $\pi$ .
2: for  $s \in \mathcal{S}$  do
3:    $s_1 = s$ .
4:   for  $k = 1, \dots, K$  do
5:     Run policy  $\pi$  for  $T$  steps.
6:     Observe utility  $U_k = \sum_t r_t$ .
7:     Update estimate  $\mathbf{v}_{k+1}(s) = \mathbf{v}_k(s) + \alpha_k(U_k - \mathbf{v}_k(s))$ 
8:   end for
9: end for
10: return  $\mathbf{v}_K$ 

```

---

For  $\alpha_k = 1/k$  and iterating over all  $\mathcal{S}$ , this is the same as Monte-Carlo policy evaluation.

A well-known algorithm for getting an optimal policy is policy iteration, Algorithm 7 in Section 8.5.4. This consists of estimating the value of a particular policy, and then trying to get an improved policy using this value. We can still apply the same principle for the case where we cannot exactly evaluate a policy. This is called approximate policy iteration. Unfortunately, approximate policy iteration does not necessarily converge without strong conditions on each approximation step.

---

**Algorithm 13** Approximate policy iteration

---

```

1: input Initial parameters  $\mathbf{v}_0$ , initial Markov policy  $\pi_0$ , stochastic estimator  $f$ .
2: for  $i = 1, \dots, N$  do
3:   Get estimate  $\mathbf{v}_i = f(\mathbf{v}_{i-1}, \pi_{i-1})$ .
4:   Calculate new policy  $\pi_i = \arg \max_\pi \mathcal{L} \mathbf{v}_i$ .
5: end for

```

---

#### Monte Carlo update



Note that  $s_1, \dots, s_T$  contains  $s_k, \dots, s_T$ . This suggests that we could update the value of all encountered states, as we also have the utility starting from each state. We call this algorithm every-visit Monte-Carlo.

---

**Algorithm 14** Every-visit Monte-Carlo update

---

```

1: input Initial parameters  $\mathbf{v}_k$ , trajectory  $s_1, \dots, s_T$ , rewards  $r_1, \dots, r_T$  visit counts  $n$ .
2: for  $t = 1, \dots, T$  do
3:    $U_t = \sum_{s=1}^T r_s$ .
4:    $n_t(s_t) = n_{t-1}(s_t) + 1$ 
5:    $\mathbf{v}_{t+1}(s_t) = \mathbf{v}_t(s) + \alpha_{n_t(s_t)}(s_t)(U_t - \mathbf{v}_t(s_t))$ 
6:    $n_t(s) = n_{t-1}(s)$ ,  $\mathbf{v}_t(s) = \mathbf{v}_{t-1}(s) \ \forall s \neq s_t$ .
7: end for
8: return  $\mathbf{v}_K$ 

```

---

For a proper Monte-Carlo estimate, when the environment is stationary  $\alpha_{n_t(s_t)}(s_t) = 1/n_t(s_t)$ . Nevertheless, this type of estimate can be biased, as can be seen by the following example.

**Example 9.2.2.** Consider a two-state chain with  $\mathbb{P}(s_{t+1} = 1 \mid s_t = 0) = \delta$  and  $\mathbb{P}(s_{t+1} = 1 \mid s_t = 1) = 1$ , and reward  $r(1) = 1$ ,  $r(0) = 0$ . Then the every-visit estimate is biased.

Let us consider the discounted setting. Then value of the second state is  $1/(1 - \gamma)$  and the value of the first state is  $\sum_k (\delta\gamma)^k = 1/(1 - \delta\gamma)$ . Consider the every-visit Monte-Carlo update. The update is going to be proportional to the number of steps you spend in that state. In order to avoid the bias, we must instead look at only the first visit to every state. This eliminates the dependence between states.

### Unbiased Monte-Carlo update

---

**Algorithm 15** First-visit Monte-Carlo update

---

```

1: input Initial parameters  $\mathbf{v}_1$ , trajectory  $s_1, \dots, s_T$ , rewards  $r_1, \dots, r_T$ , visit counts  $n$ .
2: for  $t = 1, \dots, T$  do
3:    $U_t = \sum_{s=1}^T r_s$ .
4:    $n_t(s_t) = n_{t-1}(s_t) + 1$ 
5:    $\mathbf{v}_{t+1}(s_t) = \mathbf{v}_t(s) + \alpha_{n_t(s_t)}(s_t)(U_t - \mathbf{v}_t(s_t))$  if  $n_t(s_t) = 1$ .
6:    $n_t(s) = n_{t-1}(s)$ ,  $\mathbf{v}_t(s) = \mathbf{v}_{t-1}(s)$  otherwise
7: end for
8: return  $\mathbf{v}_{T+1}$ 

```

---

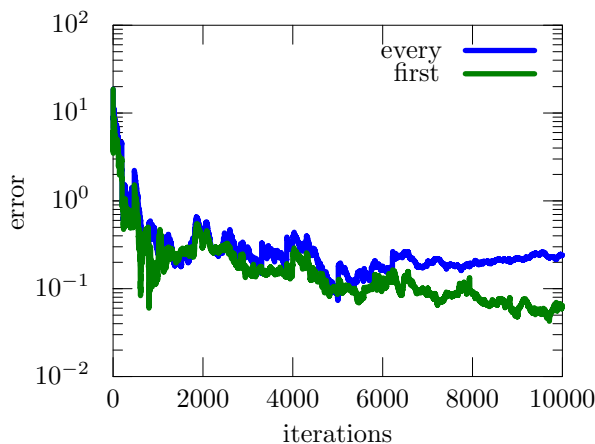


Figure 9.5: Error as the number of iterations  $n$  increases, for first and every visit Monte Carlo estimation.

### 9.2.2 Temporal difference methods

The main idea of temporal differences is to use partial samples of the utility and replace the remaining sample from time  $t$  with an estimate of the expected utility after time  $t$ . Since there maybe no particular reason to choose a specific  $t$ , frequently an exponential distribution  $t$ 's is used. Let us first look at the usual update when we have the complete utility sample  $U_k$ . The full stochastic update is of the form:

$$\mathbf{v}_{k+1}(s) = \mathbf{v}_k(s) + \alpha(U_k - \mathbf{v}_k(s)),$$

Using the *temporal difference error*  $d(s_t, s_{t+1}) = \mathbf{v}(s_t) - [\mathbf{r}(s_t) + \gamma \mathbf{v}(s_{t+1})]$ , we obtain the update:

$$\mathbf{v}_{k+1}(s) = \mathbf{v}_k(s) + \alpha \sum_t \gamma^t d_t, \quad d_t \triangleq d(s_t, s_{t+1}) \quad (9.2.2)$$

Stochastic, incremental, update:

$$\mathbf{v}_{t+1}(s) = \mathbf{v}_t(s) + \alpha \gamma^t d_t. \quad (9.2.3)$$

We have now converted the full stochastic update into an incremental update that is nevertheless equivalent to the old update. Let us see how we can generalise this to the case where we have a mixture of temporal differences.

#### Temporal difference algorithm with eligibility traces.

##### TD( $\lambda$ ).

Recall the temporal difference update when the MDP is given in analytic form.

$$\mathbf{v}_{n+1}(i) = \mathbf{v}_n(i) + \tau_n(i), \quad \tau_n(i) \triangleq \sum_{t=0}^{\infty} \mathbb{E}_{\pi_n, \mu} [(\gamma \lambda)^m d_n(s_t, s_{t+1}) \mid s_0 = i].$$

We can convert this to a stochastic update, which results in the well-known TD( $\lambda$ ) algorithm for policy evaluation.

$$\mathbf{v}_{n+1}(s_t) = \mathbf{v}_n(s_t) + \alpha \sum_{k=t}^{\infty} (\gamma\lambda)^{k-t} d_k. \quad (9.2.4)$$

Unfortunately, this algorithm is only possible to implement offline due to the fact that we are looking at future values.

This problem can be fixed by the backwards-looking Online TD( $\lambda$ ) algorithm. The main idea is to backpropagate changes in future states to previously encountered states. However, we wish to modify older states less than more recent states.

---

**Algorithm 16** Online TD( $\lambda$ )

---

```

1: input Initial parameters  $\mathbf{v}_k$ , trajectories  $(s_t, a_t, r_t)$ 
2:  $\mathbf{e}_0 = \mathbf{0}$ .
3: for  $t = 1, \dots, T$  do
4:    $d_t \triangleq d(s_t, s_{t+1})$  // temporal difference
5:    $\mathbf{e}_t(s_t) = \mathbf{e}_{t-1}(s_t) + 1$  // eligibility increase
6:   for  $s \in \mathcal{S}$  do
7:      $\mathbf{v}_{t+1}(s) = \mathbf{v}_t(s) + \alpha_t \mathbf{e}_t(s) d_t$ . // update all eligible states
8:   end for
9:    $\mathbf{e}_{t+1} = \lambda \mathbf{e}_t$ 
10: end for
11: return  $\mathbf{v}_T$ 

```

---

Figure 9.6: Eligibility traces

### 9.2.3 Stochastic value iteration methods

The main problem we had seen so far with Monte-Carlo based simulation is that we normally require a complete sequence of rewards before updating values. However, in value iteration, we can simply perform a backwards step from all the following states in order to obtain a utility estimate. This idea is explored in stochastic value iteration methods.

The standard value iteration algorithm performs a sweep over the complete state space at each iteration. However, could perform value iteration over an arbitrary sequence of states. For example, we can follow a sequence of states generated from a particular policy. This lends to the idea of *simulation-based* value iteration.

Such state sequences must satisfy various technical requirements. In particular, the policies that generate those state sequences must be *proper* for episodic problems. That is, that all policies should reach a terminating state with probability 1. For discounted non-episodic problems, this is easily achieved by using a geometric distribution for termination time. This ensures that all policies will

be proper. Alternatively, of course, we could simply select starting states with an arbitrary schedule, as long as all states are visited infinitely often in the limit.

However, value iteration also requires the Markov decision process model. The question is whether it is possible to replace the MDP model with some arbitrary estimate. This estimate can itself be obtained via simulation. This leads to a whole new family of stochastic value iteration algorithms. The most important and well-known of these is  $Q$ -learning, which uses a trivial empirical MDP model.

### Simulation-based value iteration

First, however, we shall discuss the extension of value iteration to the case where we obtain state data from simulation. This allows us to concentrate our estimates to the most useful states.

Algorithm 17 shows a generic simulation-based value iteration algorithm, with a uniform restart distribution  $\mathcal{Unif}(\mathcal{S})$  and termination probability  $\epsilon$ .

---

#### Algorithm 17 Simulation-based value iteration

---

- 1: Input  $\mu, \mathcal{S}$ .
  - 2: Initialise  $s_t \in \mathcal{S}, \mathbf{v}_0 \in \mathcal{V}$ .
  - 3: **for**  $t = 1, 2, \dots$  **do**
  - 4:  $s = s_t$ .
  - 5:  $\pi_t(s) = \arg \max_a r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbb{P}_\mu(s'|s, a) \mathbf{v}_{t-1}(s')$
  - 6:  $\mathbf{v}_t(s) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbb{P}_\mu(s'|s, \pi_t(s)) \mathbf{v}_{t-1}(s')$
  - 7:  $s_{t+1} \sim (1 - \epsilon) \cdot \mathbb{P}(s_{t+1} \mid s_t = a, \pi_t, \mu) + \epsilon \cdot \mathcal{Unif}(\mathcal{S})$ .
  - 8: **end for**
  - 9: Return  $\pi_n, V_n$ .
- 

In the following figures, we can see the error in value function estimation in the chain task when using simulation-based value iteration. It is always a better idea to use an initial value  $\mathbf{v}_0$  that is an upper bound on the optimal value function, if such a value is known. This is due to the fact that in that case, convergence is always guaranteed when using simulation-based value iteration, as long as the policy that we are using is proper.<sup>1</sup>

---

<sup>1</sup>In the case of discounted non-episodic problems, this amounts to a geometric stopping time distribution, after which the state is drawn from the initial state distribution.

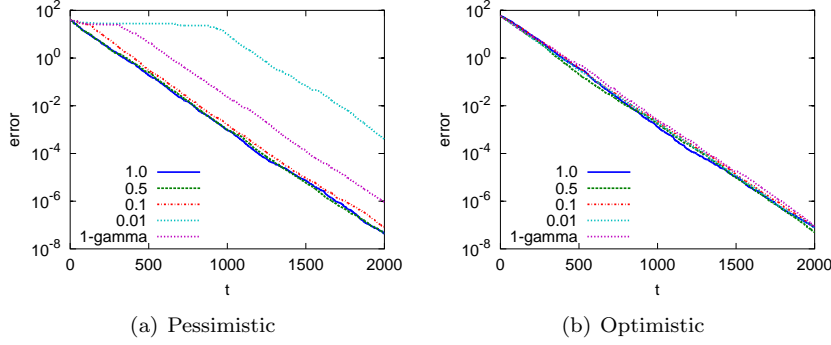


Figure 9.7: Simulation-based value iteration with pessimistic initial estimates ( $v_0 = 0$ ) and optimistic initial estimates ( $v_0 = 20 = 1/(1 - \gamma)$ ), for varying  $\epsilon$ . Errors indicate  $\|v_t - V^*\|_1$ .

As can be seen in Figure 9.7, the value function estimation error of simulation-based value iteration is highly dependent upon the initial value function estimate  $v_0$  and the exploration parameter  $\epsilon$ . It is interesting to see uniform sweeps ( $\epsilon = 1$ ) result in the lowest estimation error in terms of the value function  $L_1$  norm.

### Q-learning

Simulation-based value iteration can be suitably modified for the actual reinforcement learning problem. Instead of relying on a model of the environment, we replace arbitrary random sweeps of the state-space with the actual state sequence observed in the real environment. We also use this sequence as a simple way to estimate the transition probabilities.

---

#### Algorithm 18 Q-learning

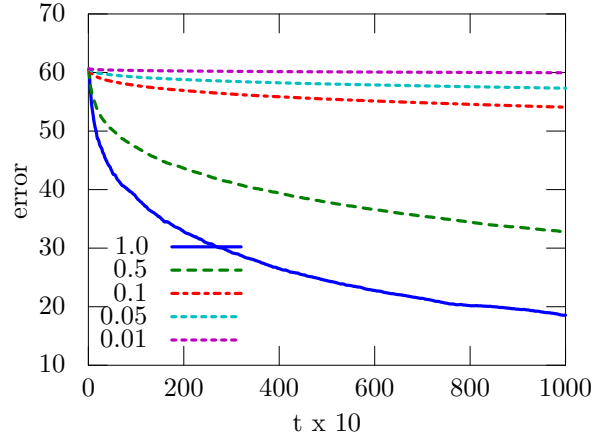
---

- 1: Input  $\mu, \mathcal{S}, \epsilon_t, \alpha_t$ .
  - 2: Initialise  $s_t \in \mathcal{S}, q_0 \in \mathcal{V}$ .
  - 3: **for**  $t = 1, 2, \dots$  **do**
  - 4:    $s = s_t$ .
  - 5:    $a_t \sim \hat{\pi}_{\epsilon_t}^*(a \mid s_t, q_t)$
  - 6:    $s_{t+1} \sim \mathbb{P}(s_{t+1} \mid s_t = a, \pi_t, \mu)$ .
  - 7:    $q_{t+1}(s_t, a_t) = (1 - \alpha_t)q_t(s_t, a_t) + \alpha_t[r(s_t) + v_t(s_{t+1})]$ , where  $v_t(s) = \max_{a \in \mathcal{A}} q_t(s, a)$ .
  - 8: **end for**
  - 9: Return  $\pi_n, V_n$ .
- 

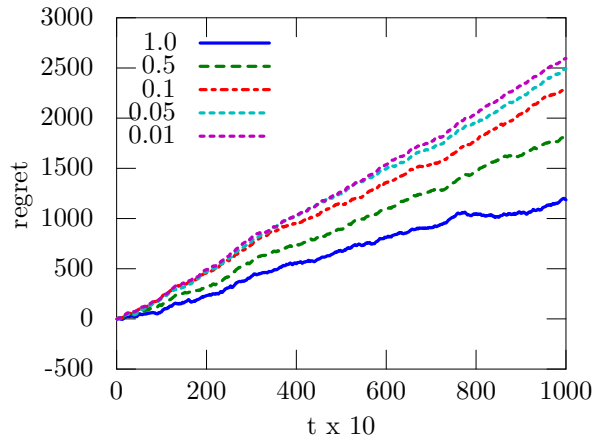
The result is Q-learning, one of the most well-known and simplest algorithms in reinforcement learning. In light of the previous theory, it can be seen as a stochastic value iteration algorithm, where at every step  $t$ , given the partial observation  $(s_t, a_t, s_{t+1})$  you have an approximate transition model for the MDP which is as follows:

$$P(s' \mid s_t, a_t) = \begin{cases} 1, & \text{if } s_{t+1} = s' \\ 0, & \text{if } s_{t+1} \neq s'. \end{cases} \quad (9.2.5)$$

Even though this model is very simplistic, it still seems to work relatively well in practice, and the algorithm is simple to implement. In addition, since we cannot arbitrarily select states in the real environment, we replace the state-exploring parameter  $\epsilon$  with a time-dependent exploration parameter  $\epsilon_t$  for the policy we employ on the real environment.



(a) Error



(b) Regret

Figure 9.8:  $Q$ -learning with  $v_0 = 1/(1 - \gamma)$ ,  $\epsilon_t = 1/n_{s_t}$ ,  $\alpha_t \in \alpha n_{s_t}^{-2/3}$ .

Figure 9.8 shows the performance of the basic  $Q$ -learning algorithm for the Chain task, in terms of value function error and regret. In this particular implementation, we used a polynomially decreasing exploration parameter and step size. Both of these depend on the number of visits to a particular state.

Of course, one could get any algorithm in between pure  $Q$ -learning and pure stochastic value iteration. In fact, variants of the  $Q$ -learning algorithm using eligibility traces (see Section 201) can be formulated in this way.

### Generalised stochastic value iteration

Finally, we can generalise the above ideas to the following algorithm. This is an online algorithm, which can be applied directly to a reinforcement learning problem and it includes simulation-based value iteration and  $Q$ -learning as special cases. There are three parameters associated with this algorithm. The first is  $\epsilon_t$ , the exploration amount performed by the policy we follow. The second is  $\alpha_t$ , the step size parameter. The third one is  $\sigma_t$ , the state-action distribution. The final parameter is the MDP estimator  $\hat{\mu}_t$ . This includes both an estimate of the transition probabilities  $\mathbb{P}_{\hat{\mu}_t}(s' | s, a)$  and of the expected reward  $r_{\hat{\mu}_t}(s, a)$ .

---

**Algorithm 19** Generalised stochastic value iteration
 

---

```

1: Input  $\hat{\mu}_0, \mathcal{S}, \epsilon_t, \alpha_t$ .
2: Initialise  $s_1 \in \mathcal{S}, \mathbf{q}_1 \in \mathcal{Q}, \mathbf{v}_0 \in \mathcal{V}$ .
3: for  $t = 1, 2, \dots$  do
4:    $a_t \sim f(\hat{\pi}_{\epsilon_t}^*(a | s_t, \mathbf{q}_t))$ 
5:   Observe  $s_{t+1}, r_{t+1}$ .
6:    $\hat{\mu}_t = \hat{\mu}_{t-1} | s_t, a_t, s_{t+1}, r_{t+1}$ .      // update MDP estimate.
7:   for  $s \in \mathcal{S}, a \in \mathcal{A}$  do
8:     With probability  $\sigma_t(s, a)$  do:


$$\mathbf{q}_{t+1}(s, a) = (1 - \alpha_t)\mathbf{q}_t(s, a) + \alpha_t \left[ r_{\hat{\mu}_t}(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbb{P}_{\hat{\mu}_t}(s' | s, a) \mathbf{v}_t(s') \right].$$


9:     otherwise  $\mathbf{q}_{t+1}(s, a) = \mathbf{q}_t(s, a)$ .
10:     $\mathbf{v}_{t+1}(s) = \max_{a \in \mathcal{A}} \mathbf{q}_{t+1}(s, a)$ ,
11:   end for
12: end for
13: Return  $\pi_n, V_n$ .
```

---

It is instructive to examine special cases for these parameters. For the case when  $\sigma_t = 1$ ,  $\alpha_t = 1$ , and when  $\hat{\mu}_t = \mu$ , we obtain standard value iteration.

For the case when  $\sigma_t(s, a) = \mathbb{I}\{s_t = s \wedge a_t = a\}$  and

$$\mathbb{P}_{\hat{\mu}_t}(s_{t+1} = s' | s_t = s, a_t = a) = \mathbb{I}\{s_{t+1} = s' | s_t = s, a_t = a\},$$

it is easy to see that we obtain  $Q$ -learning.

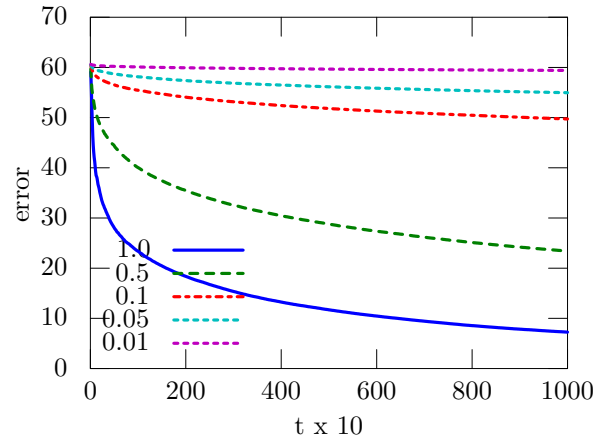
Finally, if we set  $\sigma_t(s, a) = \mathbf{e}_t(s, a)$ , then we obtain a stochastic eligibility-trace  $Q$ -learning algorithm similar to  $Q(\lambda)$ .

### Examples

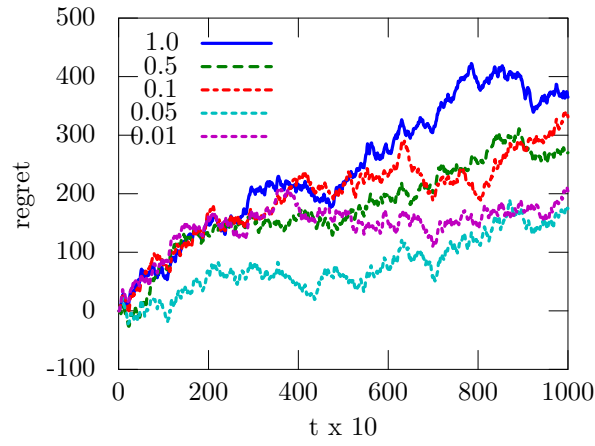
The following examples use a Dirichlet-based estimate for the transition probabilities. If  $n_t(s, a, s')$  is the number of times we have visited state  $s'$  after taking action  $a$  in state  $s$ , then the marginal probability of the next step, with a prior  $\text{Dir}(\alpha)$  is

$$\frac{n_t(s, a, s') + \alpha}{n_t(s, a) + \alpha|\mathcal{S}|}.$$

Now we can combine this with the simple  $Q$ -learning update, where we only modify the value of the current state-action pair. The error and regret are shown in figures 9.9(a) and 9.9(b)



(a) Error

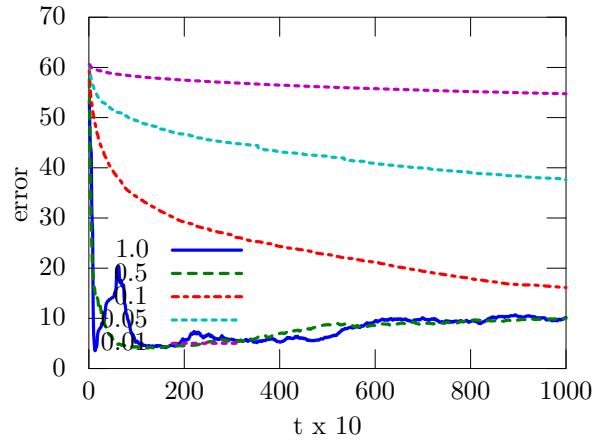


(b) Regret

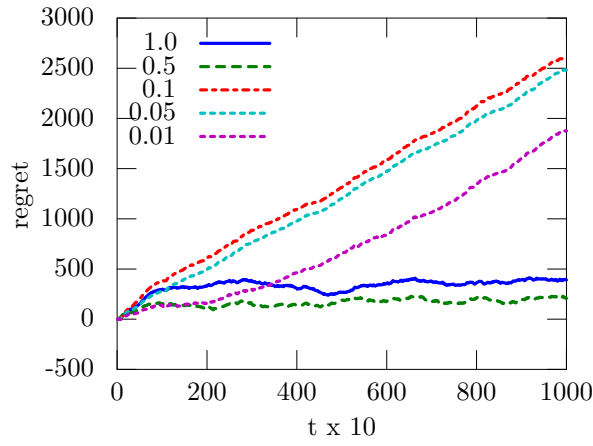
Figure 9.9: Mode-based  $Q$ -learning, i.e. GSVI with Dirichlet model and single-step updates:  $v_0 = 1/(1 - \gamma)$ ,  $\epsilon_t = 1/n_{s_t}$ ,  $\alpha_t \in \alpha n_{s_t}^{-2/3}$ .

In the following example, we again use a Dirichlet estimate but we perform a uniform sweep over the state space, i.e.  $\sigma_t = 1$ . 9.10(b) and 9.10(a)





(a) Error



(b) Regret

Figure 9.10: GSVI with Dirichlet model estimation and a uniform sweep over the state-space with  $\mathbf{v}_0 = 1/(1 - \gamma)$ ,  $\epsilon_t = 1/n_{s_t}$ ,  $\alpha_t \in \alpha n_{s_t}^{-2/3}$ .



## Chapter 10

# Approximate dynamic programming

## 10.1 Introduction

In this chapter, we consider approximate dynamic programming. This includes all methods with approximations in the maximisation step, methods where the value function used is approximate, or methods where the policy used is some approximation to the optimal policy.

We first consider the case where we have an approximate value function. Let  $\mathbf{u} \in \mathcal{V}$  be an approximate optimal value function obtained via some arbitrary method. Then we can define the greedy policy with respect to it as follows:

**Definition 10.1.1** ( $\mathbf{u}$ -greedy policy and value function).

$$\pi_{\mathbf{u}}^* \in \arg \max_{\pi} \mathcal{L}_{\pi} \mathbf{u}, \quad \mathbf{v}_{\mathbf{u}}^* = \mathcal{L} \mathbf{u}, \quad (10.1.1)$$

where  $\pi : \mathcal{S} \rightarrow \mathfrak{D}(\mathcal{A})$  maps from states to action distributions.

Although previously policies did not need to be stochastic, here we are explicitly considering stochastic policies to facilitate the approximations. Nevertheless, frequently, we cannot actually perform this maximisation if the state or action space are very large. So we define  $\phi$ , a distribution on  $\mathcal{S}$ , and parametrised sets of value functions  $\mathcal{V}_{\Theta}$  and policies  $\Pi_{\Theta}$ .

### Parameteric value function estimation

$$\mathcal{V}_{\Theta} = \{\mathbf{v}_{\theta} \mid \theta \in \Theta\}, \quad \theta^* \in \arg \min_{\theta \in \Theta} \|\mathbf{v}_{\theta} - \mathbf{u}\|_{\phi} \quad (10.1.2)$$

where  $\|\cdot\|_{\phi} \triangleq \int_{\mathcal{S}} |\cdot| d\phi$ .

In other words, we find the value function best matching the approximate value function  $\mathbf{u}$ . If  $\mathbf{u} = V^*$  then we end up getting the best possible approximation with respect to the distribution  $\phi$ .

### Parameteric policy estimation

$$\Pi_{\Theta} = \{\pi_{\theta} \mid \theta \in \Theta\}, \quad \theta^* \in \arg \min_{\theta \in \Theta} \|\pi_{\theta} - \pi_{\mathbf{u}}^*\|_{\phi} \quad (10.1.3)$$

where  $\pi_{\mathbf{u}}^* = \arg \max_{\pi \in \Pi} \mathcal{L}_{\pi} \mathbf{u}$

**Example 10.1.1.** A simple case is when  $\phi$  does not support  $\mathcal{S}$ , that is it only takes positive values for some states  $s \in \mathcal{S}$ .

#### 10.1.1 Error bounds

If the approximate value function  $\mathbf{u}$  is close to  $V^*$  then the greedy policy with respect to  $\mathbf{u}$  is close to optimal. For a finite state and action space, the following holds.

**Theorem 10.1.1.** Consider a finite MDP  $\mu$  with discount factor  $\gamma < 1$  and a vector  $\mathbf{u} \in \mathcal{V}$  such that  $\|\mathbf{u} - V_\mu^*\|_\infty = \epsilon$ . If  $\pi$  is the  $\mathbf{u}$ -greedy policy then

$$\|V_\mu^\pi - V_\mu^*\|_\infty \leq \frac{2\gamma\epsilon}{1-\gamma}.$$

In addition,  $\exists \epsilon_0 > 0$  s.t. if  $\epsilon < \epsilon_0$ , then  $\pi$  is optimal.

*Proof.* Recall that  $\mathcal{L}$  is the one-step Bellman operator and  $\mathcal{L}_\pi$  is the one-step policy operator on the value function. Then

$$\begin{aligned} \|V^\pi - V^*\|_\infty &= \|\mathcal{L}_\pi V^\pi - V^*\|_\infty \\ &\leq \|\mathcal{L}_\pi V^\pi - \mathcal{L}_\pi \mathbf{u}\|_\infty + \|\mathcal{L}_\pi \mathbf{u} - V^*\|_\infty \\ &\leq \gamma \|V^\pi - \mathbf{u}\|_\infty + \|\mathcal{L} \mathbf{u} - V^*\|_\infty \\ &\leq \gamma \|V^\pi - V^*\|_\infty + \gamma \|V^* - \mathbf{u}\|_\infty + \gamma \|\mathbf{u} - V^*\|_\infty \\ &\leq \gamma \|V^\pi - V^*\|_\infty + 2\gamma\epsilon. \end{aligned}$$

This proves the first part.

For the second part, note that the state and action sets are finite. Consequently, the set of policies is finite. Thus, there is some  $\epsilon_0 > 0$  such that the best sub-optimal policy is  $\epsilon_0$ -close to the optimal policy in value. So, if  $\epsilon < \epsilon_0$ , the obtained policy must be optimal.  $\square$

### 10.1.2 Features

Frequently, when dealing with large, or complicated spaces, it pays to project the state and action observations onto a feature space  $\mathcal{X}$ . In that way, we can make problems much more manageable. Generally speaking, a feature mapping is defined as follows.

**Feature mapping**  $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{X}$ .

For  $\mathcal{X} \subset \mathbb{R}^n$ , the feature mapping can be written in vector form:

$$f(s, a) = \begin{bmatrix} f_1(s, a) \\ \vdots \\ f_n(s, a) \end{bmatrix} \quad (10.1.4)$$

What sort of functions should we use? A common idea is to use a set of smooth functions, that are focused around a single point. One of the most usual examples are radial basis functions.

**Example 10.1.2** (Radial Basis Functions). Let  $d$  be a metric on  $\mathcal{S} \times \mathcal{A}$  and  $\{(s_i, a_i) \mid i = 1, \dots, n\}$ . Then we define each element of  $f$  as:

$$f_i(s, a) \triangleq \exp \{-d[(s, a), (s_i, a_i)]\}. \quad (10.1.5)$$

These function are sometimes called kernels.

Another common type of functions are binary functions. These effectively discretise a continuous space through either a cover or a partition.

**Definition 10.1.2.** The collection of sets  $\mathcal{G}$  is a cover of  $X$  iff  $\bigcup_{S \in \mathcal{G}} S \supset X$ .

**Definition 10.1.3.** The collection of sets  $\mathcal{G}$  is a partition of  $X$  iff

1.  $\mathcal{G}$  is a cover of  $X$
2. If  $S \neq R \in \mathcal{G}$  then  $S \cap R = \emptyset$ .
3.  $\bigcup_{S \in \mathcal{G}} S = X$ .

In reinforcement learning, these types of feature functions corresponding to partitions are usually referred to as tilings.

**Example 10.1.3 (Tilings).** Let  $\mathcal{G} = \{X_1, \dots, X_n\}$  be a partition of  $\mathcal{S} \times \mathcal{A}$  of size  $n$ . Then:

$$f_i(s, a) \triangleq \mathbb{I}\{(s, a) \in X_i\}. \quad (10.1.6)$$

Multiple tilings create a cover. These can be used without many difficulties with most discrete reinforcement learning algorithms.

## 10.2 Approximate policy iteration

### Approximate policy iteration

The main idea of approximate policy iteration is to replace the exact Bellman operator  $\mathcal{L}$  with an approximate version  $\hat{\mathcal{L}}$  and the exact value of the policy with an approximate version. In fact, in the policy improvement step, we simply try to get as close as possible to the best possible improvement, in a restricted set of policies, using an approximate operator. Similarly, in the policy evaluation step, we try to get as close as possible to the actual value of the improved policy.

---

**Algorithm 20** Generic approximate policy iteration algorithm

---

```

input Initial value function  $\mathbf{v}_0$ , approximate Bellman operator  $\hat{\mathcal{L}}$ , approxi-
mate value estimator  $\hat{V}$ .
for  $k = 1, \dots$  do
     $\pi_k = \arg \min_{\pi \in \hat{\Pi}} \left\| \hat{\mathcal{L}}_{\pi} \mathbf{v}_{k-1} - \mathcal{L} \mathbf{v}_{k-1} \right\|$  // policy improvement
     $\mathbf{v}_k = \arg \min_{\mathbf{v} \in \hat{\mathcal{V}}} \left\| \mathbf{v} - V_{\mu}^{\pi_k} \right\|$  // policy evaluation
end for

```

---

### Theoretical guarantees

**Assumption 10.2.1.** Consider a discounted problem with discount factor  $\gamma$  and iterates  $\mathbf{v}_k, \pi_k$  such that:

$$\left\| \mathbf{v}_k - V^{\pi_k} \right\|_{\infty} \leq \epsilon, \quad \forall k \quad (10.2.1)$$

$$\left\| \mathcal{L}_{\pi_{k+1}} \mathbf{v}_k - \mathcal{L} \mathbf{v}_k \right\|_{\infty} \leq \delta, \quad \forall k \quad (10.2.2)$$

**Theorem 10.2.1** (Bertsekas and Tsitsiklis [1996], proposition 6.2). Under Assumption 10.2.1

$$\limsup_{k \rightarrow \infty} \left\| V^{\pi_k} - V^* \right\|_{\infty} \leq \frac{\delta + 2\gamma\epsilon}{(1 - \gamma)^2}. \quad (10.2.3)$$

### 10.2.1 Estimation building blocks

#### Lookahead policies

Given an approximate value function  $\mathbf{u}$ , the transition model of the MDP  $P_\mu$  and expected rewards  $r_\mu$  we can always find the improving policy given in Def. 10.1.1 via the following single-step lookahead.

##### Single-step lookahead

$$\pi_{\mathbf{q}}(a \mid i) > 0 \quad \text{iff } a \in \arg \max_{a' \in \mathcal{A}} q(i, a') \quad (10.2.4)$$

$$q(i, a) \triangleq r_\mu(i, a) + \gamma \sum_{j \in \mathcal{S}} P_\mu(j \mid i, a) \mathbf{u}(j). \quad (10.2.5)$$

We are however not necessarily limited to the first-step. By looking  $T$  steps forward into the future we can improve both our value function and policy estimates.

##### $T$ -step lookahead

$$\pi(i; \mathbf{q}_T) = \arg \max_{a \in \mathcal{A}} \mathbf{q}_T(i, a), \quad (10.2.6)$$

where  $\mathbf{u}_k$  is recursively defined as:

$$q_k(i, a) = r_\mu(i, a) + \gamma \sum_{j \in \mathcal{S}} P_\mu(j \mid i, a) \mathbf{u}_{k-1}(j) \quad (10.2.7)$$

$$\mathbf{u}_k(i) = \max \{q_k(i, a) \mid a \in \mathcal{A}\} \quad (10.2.8)$$

and  $\mathbf{u}_0 = \mathbf{u}$ .

In fact, taking  $\mathbf{u} = \mathbf{0}$ , this recursion is identical to solving the  $k$ -horizon problem and at the limit we obtain solution to the original problem. In the general case, our value function estimation error is bounded by  $\gamma^k \|\mathbf{u} - V^*\|$ .

#### Rollout policies

As we have seen in Section 8.4.2 one way to obtain an the approximate value function of an arbitrary policy  $\pi$  is to use Monte Carlo estimation. That is, to simulate  $K$  sequences of state-action-reward tuples by running the policy on the MPD. More specifically, we have the following rollout estimate.

##### Rollout estimate of the $q$ -factor

$$q(i, a) = \frac{1}{K_i} \sum_{k=1}^{K_i} \sum_{t=0}^{T_k-1} r(s_{t,k}, a_{t,k}),$$

where  $s_{t,k}, a_{t,k} \sim \mathbb{P}_\mu^\pi(\cdot \mid s_0 = i, a_0 = a)$ , and  $T_k \sim \text{Geom}(1 - \gamma)$ .

This results in a set of samples of  $\mathbf{q}$ -factors. We now find a parametric policy that approximates the optimal policy with respect to our samples,  $\pi_{\mathbf{q}}^*$ . For a finite number of actions, this fitting can be seen as a classification problem. Once more, we define a distribution  $\phi$  on the states, over which we wish to perform the minimisation.

**Rollout policy estimation.**

Given a set of samples  $q(i, a)$  for  $i \in \hat{S}$ , we estimate

$$\min_{\boldsymbol{\theta}} \|\pi_{\boldsymbol{\theta}} - \pi_{\mathbf{q}}^*\|_{\phi},$$

for some  $\phi$  on  $\hat{S}$ .

### 10.2.2 The value estimation step

We can now attempt to fit a parametric approximation to a given value function  $\mathbf{v}$  or  $\mathbf{q}$ . The simplest way to do so is via a generalised linear model. A natural parametrisation for the value function is to use a generalised linear model on a set of features. Then the value function is a linear function of the features with parameters  $\boldsymbol{\theta}$ . More precisely, we can define the following model.

**Generalised linear model using features (or kernel)**

Feature mapping  $f : \mathcal{S} \rightarrow \mathbb{R}^n$ , parameters  $\boldsymbol{\theta} \in \mathbb{R}^n$ .

$$\mathbf{v}_{\boldsymbol{\theta}}(s) = \sum_{i=1}^n \theta_i f_i(s) \quad (10.2.9)$$

In order to fit a value function, we first pick a set of *representative states*  $\hat{S}$  to fit our value function  $\mathbf{v}_{\boldsymbol{\theta}}$  to  $\mathbf{v}$ . We can then estimate the optimal parameters via gradient descent.

**Fitting a value function.**

$$c(\boldsymbol{\theta}) = \sum_{s \in \hat{S}} c_s(\boldsymbol{\theta}), \quad c_s(\boldsymbol{\theta}) = \phi(s) \|\mathbf{v}_{\boldsymbol{\theta}}(s) - \mathbf{v}(s)\|_p^{\kappa}. \quad (10.2.10)$$

**Example 10.2.1.** *The case  $p = 2$ ,  $\kappa = 2$  In this case the square root and  $\kappa$  cancel out and we obtain*

$$\nabla_{\boldsymbol{\theta}} c_s = \phi(s) \sum_{a \in \mathcal{A}} \nabla_{\boldsymbol{\theta}} [\mathbf{v}_{\boldsymbol{\theta}}(s) - \mathbf{v}(s)]^2 = 2[\mathbf{v}_{\boldsymbol{\theta}}(s) - \mathbf{v}(s)] \nabla_{\boldsymbol{\theta}} \mathbf{v}_{\boldsymbol{\theta}},$$



where  $\nabla_{\theta} \mathbf{v}_{\theta}(s) = f(s)$ . Taking partial derivatives  $\partial/\partial\theta_j$ , leads to the update rule:

$$\theta'_j = \theta_j - 2\alpha\phi(s)[\mathbf{v}_{\theta}(s) - \mathbf{v}(s)]f_j(s). \quad (10.2.11)$$

### 10.2.3 Policy estimation

A natural parametrisation for the policy is to use a generalised linear model on a set of features. Then the policy can be described (up to scaling) as a linear function of the features with parameters  $\theta$ . More precisely, we can define the following model.

**Generalised linear model using features (or kernel).**

Feature mapping  $f : \mathcal{S} \rightarrow \mathbb{R}^n$ , parameters  $\theta \in \mathbb{R}^n$ .

$$\pi_{\theta}(a | s) = \frac{g(s, a)}{h(s)}, \quad g(s, a) = \sum_{i=1}^n \theta_i f_i(s, a), \quad h(s) = \sum_{b \in \mathcal{A}} g(s, b) \quad (10.2.12)$$

We are performing the intermediate step of estimating  $g$  first, because we need to make sure that the policy is a distribution over actions. An alternative method would be to directly constrain the policy parameters so the result is always a distribution, but that would require a more complex optimisation method.

In order to fit a policy, we first pick a set of representative states  $\hat{S}$  and then we find a  $\pi_{\theta}$  that approximates  $\pi$ . In order to do so, we can define an appropriate cost function and then estimate the optimal parameters via some arbitrary optimisation method.

**Fitting a policy through a cost function.**

$$c(\theta) = \sum_{s \in \hat{S}} c_s(\theta), \quad c_s(\theta) = \phi(s) \|\pi_{\theta}(\cdot | s) - \pi(\cdot | s)\|_p^{\kappa}. \quad (10.2.13)$$

The function  $\phi : \mathcal{S} \rightarrow \mathbb{R}_+$  is a weighting on the state space, such that we put more weight in more “important” states. Choosing the weights and the set of representative states  $\hat{S}$  is an interesting problem. A good choice is to relate those to the state distribution under different policies. One method to minimise the cost function is to use gradient descent. The gradient of this cost function is

$$\nabla_{\theta} c = \sum_{s \in \hat{S}} \nabla_{\theta} c_s.$$

We obtain different results for different norms. Mainly three cases are of interest:  $p = 1, p = 2, p \rightarrow \infty$ . Here, we only consider the first one.

**The case  $p = 1, \kappa = 1$ .**

The derivative can be written as:

$$\nabla_{\theta} c_s = \phi(s) \sum_{a \in \mathcal{A}} \nabla_{\theta} |\pi_{\theta}(a | s) - \pi(a | s)|,$$

$$\nabla_{\theta} |\pi_{\theta}(a | s) - \pi(a | s)| = \nabla_{\theta} \pi_{\theta}(a | s) \operatorname{sgn}[\pi_{\theta}(a | s) - \pi(a | s)]$$

The policy derivative in turn is

$$\pi_{\theta}(a | s) = \frac{h(s) \nabla_{\theta} g(s, a) - \nabla_{\theta} h(s) g(s, a)}{h(s)^2},$$

with  $\nabla_{\theta} h(s) = (\sum_{b \in \mathcal{A}} f_i(s, b))_i$  and  $\nabla_{\theta} g(s, a) = f(s, a)$ . Taking partial derivatives  $\partial/\partial\theta_j$ , leads to the update rule:

$$\theta'_j = \theta_j - \alpha \phi(s) \left( \pi_{\theta}(a | s) \sum_{b \in \mathcal{A}} f_j(s, b) - f_j(s, a) \right). \quad (10.2.14)$$

Iterating over  $(s, a)$  pairs with a decreasing step-size  $\alpha$  according to the stochastic approximation assumptions, should ensure convergence.

**Exercise 6.** Find the derivative for the two other cases, specifically:

1.  $p = 2, \kappa = 2$ .

2.  $p \rightarrow \infty, \kappa = 1$ .

**Alternative cost functions.** It is frequently a good idea to add a *penalty term* to the cost function. The purpose of this is to prevent overfitting of the parameters to a small number of observations. Frequently, this is done by constraining the parameters to be small, via a penalty term of the form  $\|\theta\|^q$ .

### 10.2.4 Rollout-based policy iteration methods

One idea for estimating the value function is to simply perform rollouts, while the policy itself is estimated in parametric form, as suggested in Bertsekas and Tsitsiklis [1996]. The first practical algorithm in this direction was Rollout Sampling Approximate Policy iteration Dimitrakakis and Lagoudakis [2008]. The main idea is to concentrate rollouts in interesting parts of the state space.

---

**Algorithm 21** Rollout Sampling Approximate Policy Iteration.

---

```

for  $k = 1, \dots$  do
  Select a set of representative states  $\hat{S}_k$ 
  for  $n = 1, \dots$  do
    Select a state  $s_n \in \hat{S}_k$  maximising  $U_n(s)$  and perform a rollout.
    If  $\hat{a}^*(s_n)$  is optimal w.p.  $1 - \delta$ , put  $s_n$  in  $\hat{S}_k(\delta)$  and remove it from  $\hat{S}_k$ .
  end for
  Calculate  $\mathbf{q}_k \approx Q^{\pi_k}$  from the rollouts.
  Train a classifier  $\pi_{\theta_{k+1}}$  on the set of states  $\hat{S}_k(\delta)$  with actions  $\hat{a}^*(s)$ .
end for

```

---

The main idea is to concentrate rollouts on promising states. We can use the empirical state distribution to select starting states. We always choose the state  $s$  with the highest upper bound  $U_n(s)$ . More specifically, we employ a Hoeffding bound to select the state with the largest gap between actions. We stop rolling out states where we are certain to have found the best action. This is done by applying the Hoeffding bound to gaps between actions.

### 10.2.5 Least Squares Methods

The main idea is to formulate the problem in linear form, using a feature space.

#### Least square value estimation

Recall that the solution of

$$\mathbf{v} = \mathbf{r} + \gamma \mathbf{P}_{\mu, \pi} \mathbf{v} \quad (10.2.15)$$

is the value function of  $\pi$  and can be obtained via

$$\mathbf{v} = (\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi})^{-1} \mathbf{r}. \quad (10.2.16)$$

However, in this setting, we do not have access to the transition matrix. In addition, when the state space is continuous (e.g.  $\mathcal{S} \subset \mathbb{R}^n$ ), the transition matrix becomes a general transition kernel. In addition, while up to now the set of value functions  $\mathcal{V}$  was a Euclidean subset, now  $\mathcal{V}$  becomes a Hilbert space.

In general, we deal with this case via projections. We project down from the infinite-dimensional Hilbert space to one with finite-dimensions. We assume that there is a projection that is complex enough for us to be able to recover the original value function.

#### Projection.

Setting  $\mathbf{v} = \Phi \boldsymbol{\theta}$  where  $\Phi$  is a feature matrix and  $\boldsymbol{\theta}$  is a parameter vector we have

$$\Phi \boldsymbol{\theta} = \mathbf{r} + \gamma \mathbf{P}_{\mu, \pi} \Phi \boldsymbol{\theta} \quad (10.2.17)$$

$$\boldsymbol{\theta} = [(\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi}) \Phi]^{-1} \mathbf{r} \quad (10.2.18)$$

Replacing the inverse with the *pseudo-inverse*, with  $\mathbf{A} = (\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi}) \Phi$

$$\tilde{\mathbf{A}}^{-1} \triangleq \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1},$$

gives us an estimate for the parameters. If the inverse exists, then it is equal to the pseudoinverse. The main idea that makes this work is to calculate everything on the empirical transition matrix, the empirical rewards and the empirical feature vectors.

**Empirical constructions.**

Given a set of data points  $\{(s_i, a_i, r_i, s'_i) \mid i = 1, \dots, n\}$ , which may not be consecutive, we define:

1.  $\mathbf{r} = (r_i)_i$ .
2.  $\Phi_i = f(s_i, a_i)$ ,  $\Phi = (\Phi_i)_i$ .
3.  $\mathbf{P}_{\mu, \pi} = \mathbf{P}_\mu \mathbf{P}_\pi$ ,  $\mathbf{P}_{\mu, \pi}(i, j) = \mathbb{I}\{j = i + 1\}$

We are now ready to define some algorithms. We begin with an algorithm that estimates an approximate value function for some policy  $\pi$  given some data  $D$  and a feature mapping  $f$ .

---

**Algorithm 22** LSTDQ - Least Squares Temporal Differences on  $\mathbf{q}$ -factors

---

**input** data  $D = \{(s_i, a_i, r_i, s'_i) \mid i = 1, \dots, n\}$ , feature mapping  $f$ , policy  $\pi$   
 $\theta = (\Phi(\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi}))^{-1} \mathbf{r}$

---

This algorithm is sufficient for performing approximate policy iteration by plugging it into the generic API algorithm to estimate a value function. Since LSTDQ returns  $\mathbf{q}$ -factors, our next policy can simply be greedy with respect to the value estimates.

---

**Algorithm 23** LSPI - Least Squares Policy Iteration

---

**input** data  $D = \{(s_i, a_i, r_i, s'_i) \mid i = 1, \dots, n\}$ , feature mapping  $f$   
 Set  $\pi_0$  arbitrarily.  
**for**  $k = 1, \dots$  **do**  
    $\theta_k = \text{LSTDQ}(D, f, \pi_{k-1})$ .  
    $\pi_k = \pi_{\Phi \theta_k}^*$ .  
**end for**

---

## 10.3 Approximate Value Iteration

Approximate algorithms can also be defined for backwards induction. The general algorithmic structure remains the same. We only need to replace the exact steps with approximations. Usually this is necessary when the value function cannot be updated everywhere exactly, possibly because our value function representations are not complex enough to capture the true value function.

### 10.3.1 Approximate backwards induction

The first algorithm is approximate backwards induction. Let us start with the basic backwards induction algorithm:

$$V_t^*(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \gamma \mathbb{E}_\mu (V_{t+1}^* \mid s_t = s, a_t = a)\} \quad (10.3.1)$$

This essentially the same both for finite and infinite-horizon problems. Now assume that the set of functions  $\mathcal{V}$  that you can use to approximate the value functions is not rich enough, so none of its members will correspond to the left side of (10.3.1). Consider then the following value function approximation.

Let our estimate at time  $t$  be  $\mathbf{v}_t \in \mathcal{V}$ , with  $\mathcal{V}$  being a set of parametrised functions. Let  $\hat{V}_t$  be our one-step update given the value function approximation at the next step,  $\mathbf{v}_{t+1}$ . Then  $\mathbf{v}_t$  will be the closest approximation in that set.

#### Iterative approximation

$$\hat{V}_t(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s'} P_\mu(s' \mid s, a) \mathbf{v}_{t+1}(s') \right\} \quad (10.3.2)$$

$$\mathbf{v}_t = \arg \min \left\{ \left\| \mathbf{v} - \hat{V}_t \right\| \mid \mathbf{v} \in \mathcal{V} \right\} \quad (10.3.3)$$

Any algorithm can be used to perform the above minimisation, including gradient descent. Now consider the case where  $\mathbf{v}$  is a parametrised function with parameters  $\boldsymbol{\theta}$ . Then it is sufficient for us to maintain the parameter  $\boldsymbol{\theta}_t$  at time  $t$ . These can be updated with a gradient scheme at every step. In the online case, our next-step estimates can be given by gradient descent:

#### Online gradient estimation

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \alpha_t \nabla_{\boldsymbol{\theta}} \left\| \mathbf{v}_t - \hat{V}_t \right\| \quad (10.3.4)$$

This gradient descent algorithm can also be made stochastic, if we sample from the probability distribution given in the iterative approximation. The next sections give some examples.

### 10.3.2 State aggregation

Partitions, or tiling of the state space, inevitably lead to what is called *state aggregation*. That is, multiple different states are seen as identical by the algorithm. Unfortunately, it is very rarely the case that aggregated states really are identical. Nevertheless, as we can see in the example below, aggregation significantly simplifies the estimation problems.

**Aggregated estimate.**

Let  $\mathcal{G} = \{S_0, S_1, \dots, S_n\}$  be a partition of  $\mathcal{S}$ , with  $S_0 = \emptyset$  and  $\theta \in \mathbb{R}^n$  and let  $f_k(s_t) = \mathbb{I}\{s_t \in S_k\}$ . Then the approximate value function is

$$v(s) = \theta(k), \quad \text{if } s \in S_k, k \neq 0. \quad (10.3.5)$$

That is, the value of every state corresponds to the value of the  $k$ -th set in the partition. Of course, this is only a very rough approximation if the sets  $S_k$  are very large. However, this is a very nice approach to use for gradient descent updates, as only one parameter needs to be updated at every step.

**Online gradient estimate.**

Consider the case  $\|\cdot\| = \|\cdot\|_2^2$ . For  $s_t \in S_k$ :

$$\theta_{t+1}(k) = (1 - \alpha)\theta_t(k) + \alpha \max_{a \in \mathcal{A}} r(s_t, a) + \gamma \sum_j P(j | s_t, a) v_t(s) \quad (10.3.6)$$

For  $s_t \notin S_k$ :

$$\theta_{t+1}(k) = \theta_t(k). \quad (10.3.7)$$

Of course, whenever we perform the estimation online, we are limited to estimation on the sequence of states  $s_t$  that we visit. Consequently, estimation on other states may not be very good. It is indeed possible that we will suffer from oscillation problems.

**10.3.3 Representative states**

A rather different idea is to choose only some representative states on which to perform the approximation. The main assumption is that the value of all other states can be represented as a convex combination of the value of the representative states.

**Representative states approximation.**

Let  $\hat{S}$  be a set of  $n$  representative states and  $\theta \in \mathbb{R}^n$  and a feature mapping  $f$ :

$$\sum_{i=1}^n f_i(s) = 1, \quad \forall s \in \mathcal{S}.$$

The feature mapping is used to perform the convex combination. For any given state  $s$ , it has higher value for representative states  $i$  which are “closer” to it. In general, the feature mapping is fixed, and we just want to find a set of parameters for the values of the representative states.

We focus here on the online estimate. At time  $t$ , for each representative state  $i$ , we obtain a new estimate of its value function and plug it back in.

**Representative state update.**

For  $i \in \hat{S}$ :

$$\boldsymbol{\theta}_{t+1}(i) = \max_{a \in \mathcal{A}} \left\{ r(i, a) + \gamma \int \mathbf{v}_t(s) dP(s | i, a) \right\} \quad (10.3.8)$$

with

$$\mathbf{v}_t(s) = \sum_{i=1}^n f_i(s) \boldsymbol{\theta}_t(i). \quad (10.3.9)$$

When the summation is not possible, we may instead approximate the expectation with a Monte-Carlo method. One particular problem with this method arises when the transition kernel is very sparse. Then we are basing our estimates on approximate values of other states, which may be very far from any other representative state.

**Bellman error methods**

The problems with the representative state update can be alleviated through Bellman error minimisation. The idea here is to obtain as a *consistent* value function as possible. The basic Bellman error minimisation is as follows:

$$\min_{\boldsymbol{\theta}} \|\mathbf{v}_{\boldsymbol{\theta}} - \mathcal{L}\mathbf{v}_{\boldsymbol{\theta}}\| \quad (10.3.10)$$

This is different from the approximate backwards induction algorithm we saw previously, since the same parameter  $\boldsymbol{\theta}$  appears in both sides of the equality. Furthermore, if the norm has support in all of the state space and the approximate value function space contains the actual set of value functions then the minimum is 0 and we obtain the optimal value function.

**Gradient update.**

When the norm is

$$\|\mathbf{v}_{\boldsymbol{\theta}} - \mathcal{L}\mathbf{v}_{\boldsymbol{\theta}}\| = \sum_{s \in \hat{S}} D_{\boldsymbol{\theta}}(s)^2, \quad D_{\boldsymbol{\theta}}(s) = \mathbf{v}_{\boldsymbol{\theta}}(s) - \max_{a \in \mathcal{A}} \int_S \mathbf{v}_{\boldsymbol{\theta}}(j) dP(j | s, a). \quad (10.3.11)$$

then the gradient update becomes

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \alpha D_{\boldsymbol{\theta}_t}(s_t) \nabla_{\boldsymbol{\theta}} D_{\boldsymbol{\theta}_t}(s_t) \quad (10.3.12)$$

$$\nabla_{\boldsymbol{\theta}} D_{\boldsymbol{\theta}_t}(s_t) = \nabla_{\boldsymbol{\theta}} \mathbf{v}_{\boldsymbol{\theta}_t}(s_t) - \int_S \nabla_{\boldsymbol{\theta}} \mathbf{v}_{\boldsymbol{\theta}_t}(j) dP(j | s_t, a_t^*) \quad (10.3.13)$$

$$a_t^* = \arg \max_{a \in \mathcal{A}} \left\{ r(s_t, a) + \gamma \int_S \mathbf{v}_{\boldsymbol{\theta}}(j) dP(j | s_t, a) \right\} \quad (10.3.14)$$

We can also construct a  $Q$ -factor approximation for the case where no model is available. This is going to be simply done by replacing  $P$  with the empirical transition observed at time  $t$ .

#### A litany of approximation algorithms

- Fitted  $Q$ -iteration Antos et al. [2008b].
- Fitted value iteration Munos and Szepesvári [2008].
- Rollout sampling policy iteration Dimitrakakis and Lagoudakis [2008]
- State aggregation Singh et al. [1995], Bernstein [2007]
- Bellman error minimisation Antos et al. [2008a], Dimitrakakis [2013], Ghavamzadeh and Engel [2006]
- Least-squares methods Bradtke and Barto [1996], Boyan [2002], Lagoudakis and Parr [2003].



## Chapter 11

# Bayesian reinforcement learning

## 11.1 Introduction

Bayesian reinforcement learning connects all elements previously seen in the course. Firstly, how to express uncertainty and preferences via probabilities and utilities. Secondly, how to make decisions under uncertainty, including how to make decisions maximising the utility in different conditions. Thirdly, how to adjust our subjective belief in the face of new evidence. Fourthly, optimal experiment design: how to make decisions in problems where our decisions can affect the evidence we obtain. These problems can be modelled as Markov decision processes. We also consider the problem of finding optimal policies for Markov decision processes.

In the previous two chapters, we have considered stochastic algorithms for acting within Markov decision processes. These stochastic analogues of deterministic MDP algorithms can also be used in the context of “learning” the optimal policy while acting in the MDP itself, even if the MDP parameters are not known. In the case where the MDP is very large, or the state/action spaces are continuous, it is necessary to approximate it. These can be used in conjunction with stochastic approximations.

Now, however, we come full circle to the setting of subjective probability and utility. We shall try and solve the reinforcement learning problem directly. Here, we are acting in an MDP which is not known, but we have a subjective belief about what the MDP is.

## 11.2 Bayesian reinforcement learning

The reinforcement learning problem can be formulated as the problem of learning to act in an unknown environment, only by interaction and reinforcement. All of those elements of the definition are important. Firstly and foremostly it is a *learning* problem. Consequently, we have only partial prior knowledge about the environment we are acting in. This knowledge is arrived at via *interaction* with the environment. We do not have a fixed set of data to work with, but we must actively explore the environment to understand how it works. Finally, there is an intrinsic *reinforcement* that punishes some behaviours and rewards others. We can formulate some of these problems as Markov decision processes.

### Markov decision processes (MDP) as an environment.

We are in some *environment*  $\mu$ , where at each time, we: step  $t$ :

- Observe *state*  $s_t \in \mathcal{S}$ .
- Take *action*  $a_t \in \mathcal{A}$ .
- Receive *reward*  $r_t \in \mathbb{R}$ .

In these types of problems, the environment state and our action fully determines the distribution of the immediate reward, as well as that of the next state, as described in Definition 8.3.1. When  $\mu$  is unknown, the probability of the immediate reward is given by  $P_\mu(r_t \mid s_t, a_t)$  and that of next state by  $P_\mu(s_{t+1} \mid s_t, a_t)$ .

However, we now assume that we do not know  $\mu$ . The structure of the unknown MDP process is shown in Figure 11.1 below,

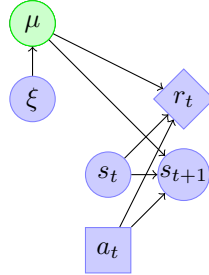


Figure 11.1: The unknown Markov decision process

#### The optimal policy for a given $\mu$

When  $\mu$  is known, we wish to find a *policy*  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  maximising the *utility* in expectation. This requires us to solve the maximisation problem  $\max_{\pi} \mathbb{E}_{\mu}^{\pi} U$ , where the utility is an additive function of rewards,  $U = \sum_{t=1}^T r_t$ . When  $\mu$  is *known*, we can use standard algorithms, such as value or policy iteration. However, knowing  $\mu$  is contrary to the problem definition.

In Chapter 9 we have seen a number of stochastic approximation algorithms which allow us to learn the optimal policy for a given MDP eventually. However, these generally give few guarantees on the performance of the policy while learning. How can we create an algorithm for optimal learning MDPs? This should trade off exploring the environment to obtain further knowledge, and simultaneously exploiting its knowledge.

The solution is rather simple, conceptually. Within the subjective probabilistic framework, we only need to define a prior belief  $\xi$  on the set of MDPs  $\mathcal{M}$ , and then find the policy that maximises the expected utility with respect to the prior. The value of information is automatically taken into account in this model.

This should not be too surprising, as we have previously seen it in two Bayes-optimal construction. The first was the simple optimal stopping procedure in Section 7.4, which introduced the backwards induction algorithm. The second was the optimal experiment design problem, which resulted in the bandit Markov decision process of Section 8.2. Let us now formulate the reinforcement learning problem as a Bayesian maximisation problem.

Let  $\xi$  be a prior over  $\mathcal{M}$  and  $\Pi$  be a set of policies. Then the expected utility of the optimal policy is:

$$U_{\xi}^* \triangleq \max_{\pi \in \Pi} \mathbb{E}(U \mid \pi, \xi) = \max_{\pi \in \Pi} \int_{\mathcal{M}} \mathbb{E}(U \mid \pi, \mu) d\xi(\mu) \quad (11.2.1)$$

Finding the optimal policy is not easy as in general the optimal policy  $\pi$  must now map from *complete histories* to actions: *Planning* must take into account *future learning*.

### Policy types

We use  $\Pi$  to denote the set of all policies. We use  $\Pi_k$  to denote the set of  $k$ -order Markov policies. Important special cases are the set of *blind* policies  $\Pi_0$  and the set of *memoryless* policies  $\Pi_1$ . A policy in  $\pi \in \bar{\Pi}_k \subset \Pi_k$  is *stationary*, when  $\pi(A \mid s_{t-k+1}^t, a_{t-k+1}^{t-1}) = \pi(A \mid s^k, a^{k-1})$  for all  $t$ .

Generally speaking, the Bayes-optimal policies are history-dependent, as shown by the following counterexample.

**Example 11.2.1.** Consider two MDPs,  $\mu_1, \mu_2$  with states  $\mathcal{S} = \{1\}$  and actions  $\mathcal{A} = \{1, 2\}$ . In the  $i$ -th MDP, whenever you take action  $a_t = i$ , you obtain reward  $r_t = 1$ , otherwise you obtain reward 0. The expected utility of a memoryless policy taking action  $i$  with probability  $\pi(i)$  would be

$$\mathbb{E}_\xi^\pi U = T \sum_i \xi(\mu_i) \pi(i),$$

for horizon  $T$ . Consequently, if your prior is not uniform, you select the action corresponding to the MDP with the highest prior probability. Then, the maximal expected utility is:

$$\max_{\pi \in \Pi_1} \mathbb{E}_\xi^\pi U = T \max_i \xi(\mu_i).$$

In this case, we are certain which one is the right MDP as soon as we take one action. We can then follow the policy which selects the apparently best action at first, and then switches to the best action for the MDP we have seen. Then, our utility is simply  $\max_i \xi(\mu_i) + (T - 1)$ .

Given the above general remarks, let us now discuss how the optimal policies can be constructed. Firstly, we must examine how to update the belief. Given that, we shall examine methods for obtaining near-optimal policies.

### 11.2.1 Updating the belief

Strictly speaking, in order to update our belief, we must condition the prior distribution on all the information. This includes the sequence of observations up to this at point in time, including the states  $s^t$ , actions  $a^{t-1}$ , and rewards  $r^{t-1}$ , as well the policy  $\pi$  that we followed. Let  $D_t = \langle s^t, a^{t-1}, r^{t-1} \rangle$  be the observed data to time  $t$ . Then

$$\xi(B \mid D_t, \pi) = \frac{\int_B \mathbb{P}_\mu^\pi(D_t) d\xi(\mu)}{\int_{\mathcal{M}} \mathbb{P}_\mu^\pi(D_t) d\xi(\mu)}. \quad (11.2.2)$$

However, as we shall see in the following remark, we can usually ignore the policy in the calculations.

**Remark 11.2.1.** The dependence on the policy can be removed, since the posterior is the same for all policies that put non-zero mass on the observed data: Let  $D_t \sim \mathbb{P}_\mu^\pi$ . Then it is easy to see that  $\forall \pi' \neq \pi$  such that  $\mathbb{P}_\mu^{\pi'}(D_t) > 0$ ,

$$\xi(B \mid D_t, \pi) = \xi(B \mid D_t, \pi').$$

Finally, since we are dealing with MDPs, the posterior calculation is easy to perform incrementally. This also more clearly demonstrates why there is no dependence on the policy. Let  $\xi_t$  be the (random) posterior at time  $t$ . Then, the next-step belief is going to be:

$$\xi_{t+1}(B) \triangleq \xi(B \mid D_{t+1}) = \frac{\int_B \mathbb{P}_\mu^\pi(D_t) d\xi(\mu)}{\int_{\mathcal{M}} \mathbb{P}_\mu^\pi(D_t) d\xi(\mu)} \quad (11.2.3)$$

$$= \frac{\int_B \mathbb{P}_\mu(s_{t+1}, r_t \mid s_t, a_t) \pi(a_t \mid s^t, a^{t-1}, r^{t-1}) d\xi(\mu \mid D_t)}{\int_{\mathcal{M}} \mathbb{P}_\mu(s_{t+1}, r_t \mid s_t, a_t) \pi(a_t \mid s^t, a^{t-1}, r^{t-1}) d\xi(\mu \mid D_t)} \quad (11.2.4)$$

$$= \frac{\int_B \mathbb{P}_\mu(s_{t+1}, r_t \mid s_t, a_t) d\xi_t(\mu)}{\int_{\mathcal{M}} \mathbb{P}_\mu(s_{t+1}, r_t \mid s_t, a_t) d\xi_t(\mu)} \quad (11.2.5)$$

The above calculation is easy to perform for arbitrarily complex MDPs when the set  $\mathcal{M}$  is finite. The posterior calculation is also simple under certain conjugate priors.

**Exercise 7.** A practical case is when we have an independent belief over the transition probabilities of each state-action pair. Consider the case where we have  $n$  states and  $k$  actions. Similar to the product-prior in the bandit case in Section 8.2, we assign a probability (density)  $\xi_{s,a}$  to the probability vector  $\theta_{(s,a)} \in \Delta^n$ . We can then define our joint belief on the  $(nk) \times n$  matrix  $\Theta$  to be

$$\xi(\Theta) = \prod_{s \in \mathcal{S}, a \in \mathcal{A}} \xi_{s,a}(\theta_{(s,a)}).$$

- (i) Derive the updates for a product-Dirichlet prior on transitions.
- (ii) Derive the updates for and a product-Normal-Gamma prior on rewards.
- (iii) What would be the meaning of using a Normal-Wishart prior on rewards?

## 11.3 Finding Bayes-optimal policies

The problem of policy optimisation in the Bayesian case is much harder than in the known-MDP case. This is simply because of the history dependence, which has two effects. Firstly, it makes the policy space much larger, as we need to consider history dependent policies. However, even we consider only memoryless policies, it does not make dynamic programming easier.

In this section, we first consider two simple heuristics for finding optimal policies. Then we examine policies which try and construct upper and lower bounds on the expected utility. Finally, we consider finite-lookahead backwards induction, that uses the same upper and lower bounds to perform efficient tree search.

### 11.3.1 The expected MDP heuristic

One simple heuristic is to simply calculate the expected MDP for a given belief  $\xi$ :

$$\bar{\mu}_\xi \triangleq \mathbb{E}_\xi \mu.$$

Then, we simply calculate the optimal memoryless policy for  $\bar{\mu}_\xi$ :

$$\pi^*(\bar{\mu}_\xi) \in \arg \max_{\pi \in \Pi_1} V_{\bar{\mu}_\xi}^\pi,$$

where  $\Pi_1 = \{\pi \in \Pi \mid \mathbb{P}_\pi(a_t \mid s^t, a^{t-1}) = \mathbb{P}_\pi(a_t \mid s_t)\}$ . Finally, we execute  $\pi^*(\bar{\mu}_\xi)$  on the real MDP. The algorithm can be written as follows. Unfortunately, this

---

**Algorithm 24** The expected MDP heuristic

---

```

for  $k = 1, \dots$  do
   $\mu_k \triangleq \mathbb{E}_{\xi_{t_k}} \mu$ .
   $\pi_k \approx \arg \max_{\pi} \mathbb{E}_{\mu_k}^\pi U$ .
  for  $t = 1 + T_{k-1}, \dots, T_k$  do
    Observe  $s_t$ .
    Update belief  $\xi_t(\cdot) = \xi_{t-1}(\cdot \mid s_t, a_{t-1}, r_{t-1}, s_{t-1})$ .
    Take action  $a_t \sim \pi_k(a_t \mid s_t)$ .
    Observe reward  $r_t$ .
  end for
end for

```

---

approach may be far from the optimal policy in  $\Pi_1$ , as shown by the following counterexample.

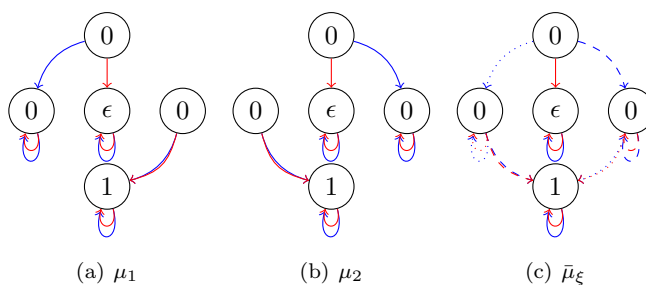


Figure 11.2: The two MDPs and the expected MDP from example ??

**Example 11.3.1** (Counterexample<sup>1</sup>). *In this example, illustrated in Figure 11.2,  $\mathcal{M} = \{\mu_1, \mu_2\}$  is the set of MDPs, and the belief is  $\xi(\mu_1) = \theta$ ,  $\xi(\mu_2) = 1 - \theta$ . All transitions are deterministic, and there are two actions, the blue and the red action. When we calculate the expected MDP, we see that now the state with reward 1 is reachable.*

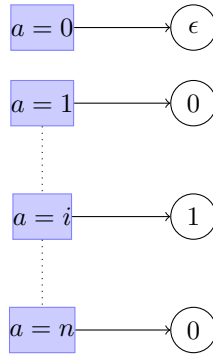
*Consequently, when  $T \rightarrow \infty$ , the  $\bar{\mu}_\xi$ -optimal policy is not optimal in  $\Pi_1$  if:*

$$\epsilon < \frac{\gamma\theta(1-\theta)}{1-\gamma} \left( \frac{1}{1-\gamma\theta} + \frac{1}{1-\gamma(1-\theta)} \right)$$

*In this example,  $\bar{\mu}_\xi \notin \mathcal{M}$ .*

---

<sup>1</sup>Based on one by Remi Munos

Figure 11.3: The MDP  $\mu_i$  from example 11.3.2

### 11.3.2 The maximum MDP heuristic

An alternative idea is to simply pick the maximum-probability MDP, as shown in Algorithm 25. This at least guarantees that the MDP that you are acting optimally for is actually within the set of MDPs. However, it may still be the case that the resulting policy is sub-optimal, as shown by the following counterexample.

---

**Algorithm 25** The maximum MDP heuristic

---

```

for  $k = 1, \dots$  do
   $\mu_k \triangleq \arg \max_{\mu} \xi_{t_k}(\mu)$ .
   $\pi_k \approx \arg \max_{\pi} \mathbb{E}_{\mu_k}^{\pi} U$ .
  for  $t = 1 + T_{k-1}, \dots, T_k$  do
    Observe  $s_t$ .
    Update belief  $\xi_t(\cdot) = \xi_{t-1}(\cdot \mid s_t, a_{t-1}, r_{t-1}, s_{t-1})$ .
    Take action  $a_t \sim \pi_k(a_t \mid s_t)$ .
    Observe reward  $r_t$ .
  end for
end for

```

---

**Example 11.3.2** (Counterexample for  $\hat{\mu}_{\xi}^* \triangleq \arg \max_{\mu} \xi(\mu)$ ). Let the MDP set be  $\mathcal{M} = \{\mu_i \mid i = 1, \dots, n\}$  with  $\mathcal{A} = \{0, \dots, n\}$ . In all MDPs,  $a_0$  gives a reward of  $\epsilon$  and the MDP terminates. In the  $i$ -th MDP, all other actions give you a reward of 0 apart from the  $i$ -th action which gives you a reward of 1. Then the MDP terminates. The MDP is visualised in Figure 11.3.

For this problem, the  $\xi$ -optimal policy takes action  $i$  iff  $\xi(\mu_i) \geq \epsilon$ , otherwise takes action 0. On the other hand, the  $\hat{\mu}_{\xi}^*$ -optimal policy takes  $a = \arg \max_i \xi(\mu_i)$ . Thus, this policy is sub-optimal if  $\max_i \xi(\mu_i) < \epsilon$ .

For smooth beliefs,  $\bar{\mu}_{\xi}$  is close to  $\hat{\mu}_{\xi}^*$ , and in this case, those heuristics might be reasonable. However, they can be shown to be sub-optimal even for very simple stopping problems.

### 11.3.3 Bounds on the expected utility

Given that these heuristics are incorrect, what can we actually do? The first thing to try is to calculate the expected utility of some arbitrary policy. As it turns out, this operation is relatively simple in the Bayesian case, even when the set of MDPs is infinite.

Policy evaluation is particularly simple in Bayesian MDP problems. We simply apply the basic utility theory definitions. We first define the Bayes-value function of a policy  $\pi$  to be the expected utility under that policy and our belief  $\xi$ .

**Expected utility of a policy  $\pi$  for a belief  $\xi$**

$$V_\xi^\pi(s) \triangleq \mathbb{E}_\xi^\pi(U \mid s_t = s) \quad (11.3.1)$$

$$= \int_{\mathcal{M}} \mathbb{E}_\mu^\pi(U \mid s_t = s) d\xi(\mu) \quad (11.3.2)$$

$$= \int_{\mathcal{M}} V_\mu^\pi(s) d\xi(\mu) \quad (11.3.3)$$

---

**Algorithm 26** Bayesian Monte-Carlo policy evaluation

---

**input** policy  $\pi$ , belief  $\xi$   
**for**  $k = 1, \dots, K$  **do**  
     $\mu_k \sim \xi$ .  
     $v_k = V_{\mu_k}^\pi$   
**end for**  
 $u = \frac{1}{K} \sum_{k=1}^K v_k$ .  
**return**  $u$ .

---

The value of any policy gives us a natural lower bound on the Bayes-optimal value function.

We can also get the following upper bounds:

$$V_\xi^* \triangleq \sup_\pi \mathbb{E}_\xi^\pi(U) = \sup_\pi \int_{\mathcal{M}} \mathbb{E}_\mu^\pi(U) d\xi(\mu) \quad (11.3.4)$$

$$\leq \int_{\mathcal{M}} \sup_\pi V_\mu^\pi d\xi(\mu) = \int_{\mathcal{M}} V_\mu^* d\xi(\mu) \triangleq V_\xi^+ \quad (11.3.5)$$

**Bounds on  $V_\xi^* \triangleq \max_\pi \mathbb{E}(U \mid \pi, \xi)$**

Given the previous development, it is easy to see that the following inequalities always hold:

$$V_\xi^\pi \leq V_\xi^* \leq V_\xi^+, \quad \forall \pi. \quad (11.3.6)$$

These bounds are geometrically demonstrated in Fig. 11.4. They are entirely analogous to the Bayes risk bounds of Sec. ??.



**Algorithm 27** Bayesian Monte-Carlo upper bound

---

```

input policy  $\pi$ , belief  $\xi$ 
for  $k = 1, \dots, K$  do
     $\mu_k \sim \xi$ .
     $\mathbf{v}_k = V_{\mu_k}^*$ 
end for
 $\mathbf{u}^* = \frac{1}{K} \sum_{k=1}^K \mathbf{v}_k$ .
return  $\mathbf{u}^*$ .

```

---

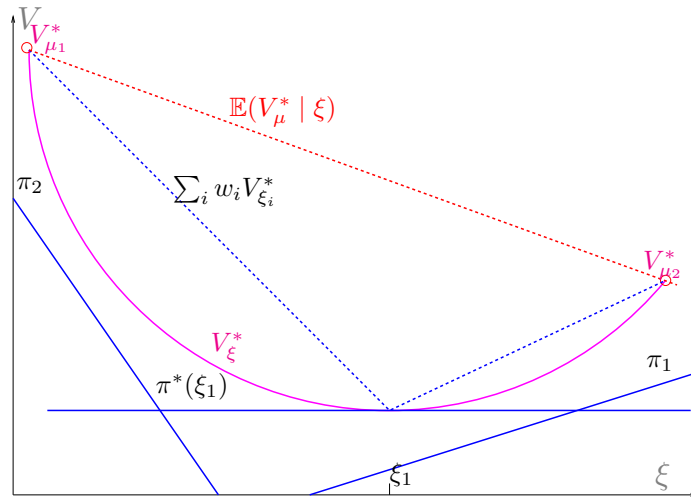
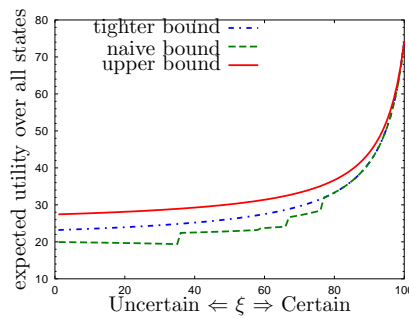


Figure 11.4: A geometric view of the bounds

**11.3.4 Tighter lower bounds**

One idea to get a better lower bound is to simply find better policies. This idea was explored in Dimitrakakis [2011].



The main idea was to maximise try and find the best memoryless policies. This can be done approximately by assuming that the belief is nearly constant over time, and performing backwards induction on  $n$  MDPs simultaneously. While this greedy procedure might not find the optimal memoryless policy, it still improves the lower bounds considerably.

The multi-MDP backwards induction procedure simply involves calculating the expected utility of a particular policy over all MDPs.

$$Q_{\xi,t}^{\pi}(s, a) \triangleq \int_{\mathcal{M}} \left\{ \bar{R}_{\mu}(s, a) + \gamma \int_{\mathcal{S}} V_{\mu,t+1}^{\pi}(s') d\mathcal{T}_{\mu}^{s,a}(s') \right\} d\xi(\mu) \quad (11.3.7)$$

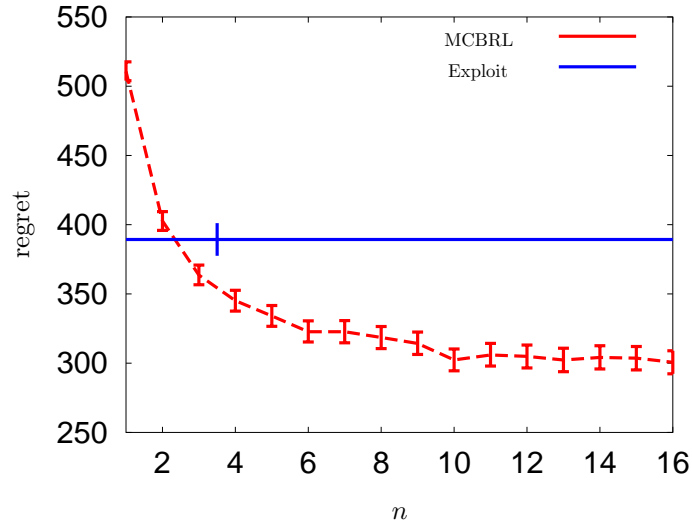
The algorithm greedily performs backwards induction as shown in Algorithm 28. However, this is not an optimal procedure, since the belief at any time-step  $t$  is not constant. Indeed, as the policy is memoryless,  $\xi(\mu \mid s_t, \pi) \neq \xi(\mu \mid s_t, \pi')$ . This is because the probability of being at a particular state is different under different policies and at different time-steps (e.g. if you consider periodic MDPs). For the same reason, this type of backwards induction may not converge in the manner of value iteration.

---

**Algorithm 28** Multi-MDP backwards induction
 

---

- 1:  $\text{MMBLM}, \xi, \gamma, T$
  - 2: Set  $V_{\mu,T+1}(s) = 0$  for all  $s \in \mathcal{S}$ .
  - 3: **for**  $t = T, T-1, \dots, 0$  **do**
  - 4:   **for**  $s \in \mathcal{S}, a \in \mathcal{A}$  **do**
  - 5:     Calculate  $Q_{\xi,t}(s, a)$  from (11.3.7) using  $\{V_{\mu,t+1}\}$ .
  - 6:   **end for**
  - 7:   **for**  $s \in \mathcal{S}$  **do**
  - 8:      $a_{\xi,t}^*(s) \in \arg \max_{a \in \mathcal{A}} Q_{\xi,t}(s, a)$ .
  - 9:     **for**  $\mu \in \mathcal{M}$  **do**
  - 10:        $V_{\mu,t}(s) = Q_{\mu,t}(s, a_{\xi,t}^*(s))$ .
  - 11:     **end for**
  - 12:   **end for**
  - 13: **end for**
- 



**MCBRL: Application to Bayesian RL**

1. For  $i = 1, \dots$
2. At time  $t_i$ , sample  $n$  MDPs from  $\xi_{t_i}$ .
3. Calculate best memoryless policy  $\pi_i$  wrt the sample.
4. Execute  $\pi_i$  until  $t = t_{i+1}$ .

***Relation to other work***

- For  $n = 1$ , this is equivalent to the Thompson sampling used by Strens [2000].
- Unlike BOSS Asmuth et al. [2009] it does not become more optimistic as  $n$  increases. BOSS takes multiple samples and constructs the most optimistic MDP possible in this set.
- BEETLE Poupart et al. [2006], Poupart and Vlassis [2008] is a belief-sampling approach. It examines a set of possible future beliefs and approximates the value of each belief with a lower bound. In essence, it then creates the set of policies which are optimal with respect to these bounds.
- Furnstion and Barber Furnstion and Barber [2010] use approximate inference to estimate policies. These use the expectation-maximisation view Toussaint et al. [2006] of reinforcement learning.

**Generalisations**

- Policy search for improving lower bounds.
- Search enlarged class of policies
- Examine all history-based policies.

**The augmented MDP**

We are given an initial belief  $\xi_0$  on a set of MDPs  $\mathcal{M}$ . Each  $\mu \in \mathcal{M}$  is a tuple  $(\mathcal{S}, \mathcal{A}, P_\mu, \mathbf{r})$ , with state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , transition kernel  $P_\mu$  and reward vector  $\mathbf{r}$ . We now construct the following augmented Markov decision process:  $(\mathcal{S} \times \Xi, \mathcal{A}, P, \mathbf{r})$ , with factorised transition probabilities:

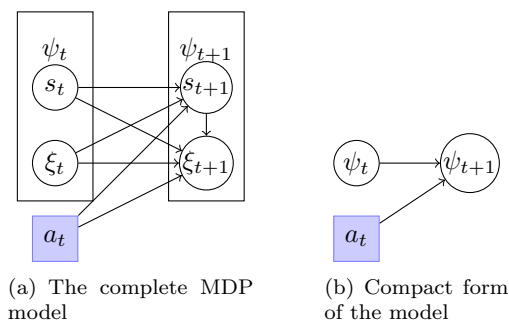


Figure 11.5: Belief-augmented MDP

The optimal policy for the augmented MDP is the  $\xi$ -optimal for the original problem.

$$P(s_{t+1} \in S \mid \xi_t, s_t, a_t) \triangleq \int_S P_\mu(s_{t+1} \in S \mid s_t, a_t) d\xi_t(\mu) \quad (11.3.8)$$

$$\xi_{t+1}(\cdot) = \xi_t(\cdot \mid s_{t+1}, s_t, a_t) \quad (11.3.9)$$

and reward  $r_t = \rho(s_t, a_t)$ . In the above,

- $\xi_t$  our belief over MDPs  $\mu \in \mathcal{M}$  at time  $t$ .
- $s_t$  is the observed state of the unknown MDP at time  $t$ .
- $P_\mu$  is the transition kernel of the MDP  $\mu$ .
- $a_t$  is our action at time  $t$ .
- For simplicity, we assume that  $r_t$  be known.

One of the first treatments of this idea was due to Bellman [1957]. Although the idea was well-known in the statistical community DeGroot [1970], the popularisation of the idea in reinforcement learning was achieved with Duff's thesis Duff [2002].

### Belief-augmented MDP tree structure

Given a belief over MDPs, we can create an *augmented* MDP with state space  $\Psi = \mathcal{S} \times \Xi$ . This has a pseudo-tree structure (since belief states might repeat). Consider an MDP family  $\mathcal{M}$  with  $\mathcal{A} = \{a^1, a^2\}$ ,  $\mathcal{S} = \{s^1, s^2\}$ .

$$\psi_t = (s_t, \xi_t)$$

### 11.3.5 Stochastic branch and bound

#### Branch and bound

Branch and bound is a general technique for solving large problems. It can be applied in all cases where upper and lower bounds on the value of solution sets can be found

##### Value bounds

Let upper and lower bounds  $q^+$  and  $q^-$  such that:

$$q^+(\psi, a) \geq Q^*(\psi, a) \geq q^-(\psi, a) \quad (11.3.10)$$

$$v^+(\psi) = \max_{a \in \mathcal{A}} Q^+(\psi, a), \quad v^-(\psi) = \max_{a \in \mathcal{A}} Q^-(\psi, a). \quad (11.3.11)$$

$$q^+(\psi, a) = \sum_{\psi'} p(\psi' | \psi, a) [r(\psi, a, \psi') + V^+(\psi')] \quad (11.3.12)$$

$$q^-(\psi, a) = \sum_{\psi'} p(\psi' | \psi, a) [r(\psi, a, \psi') + V^-(\psi')] \quad (11.3.13)$$

**Remark 11.3.1.** If  $q^-(\psi, a) \geq q^+(\psi, b)$  then  $b$  is sub-optimal at  $\psi$ .

##### Stochastic branch and bound for belief tree search Dimitrakakis [2010, 2008]

- (Stochastic) Upper and lower bounds on the values of nodes (via Monte-Carlo sampling)
- Use upper bounds to expand tree, lower bounds to select final policy.
- Sub-optimal branches are quickly discarded.

## 11.4 Partially observable Markov decision processes

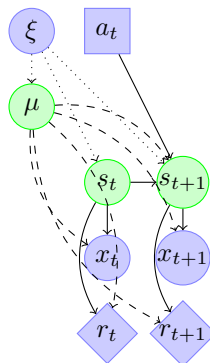
In most real applications, the state of the system is not observed.

##### Partially observable Markov decision processes (POMDP)

When acting in  $\mu$ , each time step  $t$ :

- The system state  $s_t \in \mathcal{S}$  is not observed.
- We receive an observation  $x_t \in \mathcal{X}$  and a reward  $r_t \in \mathcal{R}$ .
- We take action  $a_t \in \mathcal{A}$ .

- The system transits to state  $s_{t+1}$ .



**Definition 11.4.1.** *Partially observable Markov decision process (POMDP)* A POMDP  $\mu \in \mathcal{M}_P$  is a tuple  $(\mathcal{X}, \mathcal{S}, \mathcal{A}, P)$  where  $\mathcal{X}$  is an observation space,  $\mathcal{S}$  is a state space,  $\mathcal{A}$  is an action space, and  $P$  is a conditional distribution on observations, states and rewards. The following Markov property holds:

$$\mathbb{P}_\mu(s_{t+1}, r_t, x_t \mid s_t, a_t, \dots) = P(s_{t+1} \mid s_t, a_t)P(x_t \mid s_t)P(r_t \mid s_t) \quad (11.4.1)$$

#### Belief state in POMDPs when $\mu$ is known

In POMDPs, we can similarly define a *belief state* summarising our knowledge. This takes the form of a probability distribution on the hidden state variable  $s_t$ . If  $\mu$  defines starting state probabilities, then the belief is not subjective

##### Belief $\xi$

For any distribution  $\xi$  on  $\mathcal{S}$ , we define:

$$\xi(s_{t+1} \mid a_t, \mu) \triangleq \int_{\mathcal{S}} P_\mu(s_{t+1} \mid s_t a_t) d\xi(s_t) \quad (11.4.2)$$

When there is no ambiguity, we shall use  $\xi$  to denote arbitrary marginal distributions on states and state sequence given the belief  $\xi$ .

##### Belief update

$$\xi_t(s_{t+1} \mid x_{t+1}, r_{t+1}, a_t, \mu) = \frac{P_\mu(x_{t+1}, r_{t+1} \mid s_{t+1}) \xi_t(s_{t+1} \mid a_t, \mu)}{\xi_t(x_{t+1} \mid a_t, \mu)} \quad (11.4.3)$$

$$\xi_t(s_{t+1} \mid a_t, \mu) = \int_{\mathcal{S}} P_\mu(s_{t+1} \mid s_t, a_t, \mu) d\xi_t(s_t) \quad (11.4.4)$$

$$\xi_t(x_{t+1} \mid a_t, \mu) = \int_{\mathcal{S}} P_\mu(x_{t+1} \mid s_{t+1}) d\xi_t(s_{t+1} \mid a_t, \mu) \quad (11.4.5)$$

**Example 11.4.1.** If  $\mathcal{S}, \mathcal{A}, \mathcal{X}$  are finite, and then we can define

- $\partial_t(j) = P(x_t \mid s_t = j)$
- $\mathbf{A}_t(i, j) = P(s_{t+1} = j \mid s_t = i, a_t)$ .
- $\mathbf{b}_t(i) = \xi_t(s_t = i)$

We can then use Bayes theorem:

$$\mathbf{b}_{t+1} = \frac{\text{diag}(\mathbf{p}_{t+1}) \mathbf{A}_t \mathbf{b}_t}{\mathbf{p}_{t+1}^\top \mathbf{A}_t \mathbf{b}_t}, \quad (11.4.6)$$

**When the POMDP  $\mu$  is unknown**

$$\xi(\mu, s^t \mid x^t, a^t) \propto P_\mu(x^t \mid s^t, a^t) P_\mu(s^t \mid a^t) \xi(\mu) \quad (11.4.7)$$

#### Cases

- Finite  $\mathcal{M}$ .
- Finite  $\mathcal{S}$
- General case

#### Strategies for POMDPs

- Bayesian RL on POMDPs? *EXP inference and planning*
- Approximations and stochastic methods.
- Policy search methods.





## Chapter 12

# Distribution-free reinforcement learning

## 12.1 Introduction

The Bayesian framework requires specifying a prior distribution  $\xi_0$ . For many reasons, we may frequently be unable to specify such a prior distribution. In addition, as we have seen, the Bayes-optimal solution is frequently intractable. Here we shall take a look at a number of *heuristic* algorithms that do not require specifying a prior distribution. Instead, they employ the heuristic of “optimism under uncertainty” to select policies. This idea is very similar to heuristic search algorithms, such as  $A^*$ . All these algorithms select the policy with the highest optimistic value, i.e. with the highest upper bound in its value. The upper bound can be interpreted as the maximum utility we could obtain by playing this policy now, even if we have to switch policies later. In general we want the upper bound to

1. Be as tight as possible
2. Hold with high probability.

We begin with an introduction to these ideas in bandit problems, when the objective is to maximise total reward. We then expand this discussion to structured bandit problems, which have many applications in optimisation. Finally, we look at the case of maximising total reward in unknown MDPs. The same main ideas can be used, but the very definition of an optimal MDP policy is not trivial when we wish to maximise total reward. For this reason, we shall go over the various optimality criteria we can use. We then briefly discuss a nearly-optimal reinforcement learning algorithm.

## 12.2 Bandit problems

First of all, let us remind the reader of the stochastic bandit problem. We have a choice between a set of  $K$  bandits, corresponding to an action set  $\mathcal{A} = \{1, \dots, K\}$ . The expected reward we get when we play the  $i$ -th bandit is  $\mu_i \triangleq \mathbb{E}(r_t \mid a_t = i)$ . What we wish to do is to maximise the total reward during our life-time  $\sum_{t=1}^T r_t$ , where  $T$  may be random. What is a good heuristic strategy?

Let  $\pi_t$  define a probability distribution over the arms at time  $t$ . Let  $\mu^* \triangleq \max_i \mu_i$  be the highest average reward we can achieve and let  $\pi^*$  be the policy that always plays the arm with the highest average reward. It is important to recognise that maximising total reward is equivalent to minimising total regret with respect to that policy

**Definition 12.2.1** (Total regret). *The (total) regret of a policy  $\pi$  relative to the optimal fixed policy  $\pi^*$  is:*

$$L_T(\pi) \triangleq \sum_{t=1}^T r_t^* - r_t^\pi. \quad (12.2.1)$$

where  $r_t^* \triangleq r_t^{\pi^*}$  is the reward obtained by  $\pi^*$  and  $r_t^\pi$  are the rewards generated by the policy  $\pi$

Here we are comparing with the best fixed policy. This can be seen as the policy that would be played by oracle that knows which is the best bandit on average. We could also compare with the oracle that knows exactly how much reward each bandit would give at different times, but it is hard to prove anything meaningful about that, as the rewards at each time-step are by definition unpredictable. Taking expectations, we have that the expected (total) regret is

$$\mathbb{E} L_T(\pi) \triangleq T\mu^* - \mathbb{E}_\pi \sum_{t=1}^T r_t. \quad (12.2.2)$$

Now we consider a simple algorithm that uses the empirical average rewards obtained by each bandit.

**Empirical average**

$$\hat{\mu}_{t,i} \triangleq \frac{1}{n_{t,i}} \sum_{k=1}^t r_{k,i} \mathbb{I}\{a_k = i\}, \quad n_{t,i} \triangleq \sum_{k=1}^t \mathbb{I}\{a_k = i\}.$$

Using the empirical averages directly is not a very good idea, because you might get stuck with a sub-optimal bandit. A better idea is to play bandits optimistically. That is, as long as a particular bandit has a significant chance of being the best, you play it. One way to implement this is through the following algorithm, which assumes that the rewards are bounded, i.e. that  $r_t \in \mathcal{R} \subset \mathbb{R}$ .

---

**Algorithm 29** Optimistic initial values

---

**Input**  $\mathcal{A}, \mathcal{R}$   
 $r_{\max} \triangleq \max \mathcal{R}$   
**for**  $t = 1, \dots$  **do**  
 $u_{t,i} = \frac{n_{t-1,i} \hat{\mu}_{t-1,i} + r_{\max}}{n_{t-1,i} + 1}$   
 $a_t = \arg \max_{i \in \mathcal{A}} u_{t,i}$   
**end for**

---

If you analyse this decision rule carefully, you can see that the algorithm chooses the arm with maximal  $\hat{\mu}_i + O(1/n_i)$ . That is, it adds a small bonus value to each arm, depending on how many times the arm has been played.

**A simple analysis in the deterministic case**

When there is no randomness, the algorithm is easy to analyse. Firstly, it must hold that  $r_{t,i} = \mu_{t,i}$  for all bandits. Secondly, note that  $u_{t,i} \geq \mu_i$  for all  $t, i$ . At time  $t$ , we will play arm  $i$  only if  $u_{t,i} \geq u_{t,j}$  for all  $j$ . However,  $u_{t,j} \geq \mu_j$  and so a necessary condition for us to play arm  $i$  is that it  $u_{t,i} > \mu_j$  for all arms, including the optimal arm. From this, we obtain that we play  $i$  at most

$$n_{t,i} \leq \frac{r_{\max}}{\Delta_i}$$

times, where  $\Delta_i \triangleq \mu^* - \mu_i$ . This is easy to see since, for us to play  $i$  it must hold that

$$\begin{aligned} \frac{n_{t,i}\hat{\mu}_{t,i} + r_{\max}}{n_{t,i} + 1} &\geq \mu^* \\ n_{t,i}\hat{\mu}_{t,i} + r_{\max} &\geq \mu^*(n_{t,i} + 1) \\ r_{\max} - \mu^* &\geq \Delta_i n_{t,i} \end{aligned}$$

Since every time we play  $i$  we lose  $\Delta_i$ , the regret is

$$L_T \leq \sum_{i \neq j} \Delta_i \frac{r_{\max} - \mu^*}{\Delta_i} = (K - 1)(r_{\max} - \mu^*).$$

Unfortunately this algorithm does not have very good properties in the stochastic case. However, an other algorithm, with a choice of actions based upon concentration inequalities, does.

### 12.2.1 UCB

The idea of UCB is to

---

#### Algorithm 30 UCB1

---

**Input**  $\mathcal{A}, \mathcal{R}$   
 $\hat{\mu}_{0,i} = r_{\max}, \forall i.$   
**for**  $t = 1, \dots$  **do**  
     $u_{t,i} = \hat{\mu}_{t-1,i} + \sqrt{2 \frac{\ln t}{n_{t-1,i}}}.$   
     $a_t = \arg \max_{i \in \mathcal{A}} u_{t,i}$   
**end for**

---

Note that the decision rule in this case uses the heuristic bound of the form  $\hat{\mu}_i + O(\sqrt{\ln t / n_i})$  to select the action.

**Theorem 12.2.1** (Auer et al Auer et al. [2002]). *The expected regret of UCB1 after  $T$  rounds is at most*

$$c_1 \sum_{i: \mu_i < \mu^*} \left( \frac{\ln T}{\Delta_i} \right) + c_2 \sum_{j=1}^K \Delta_j$$

*Proof.* First we prove that

$$\mathbb{E} n_{t,i} \leq O \left( \frac{\ln T}{\Delta_i^2} \right)$$

Then we note that the expected regret can be written as

$$\sum_{i: \mu_i < \mu^*} \Delta_i \mathbb{E} n_{t,i}$$

due to Wald's identity. □

Let  $B_{t,s} = \sqrt{(2 \ln t)/s}$ . Then we can prove  $\forall c \in \mathbb{Z}$ :

$$\begin{aligned}
n_{T,i} &= 1 + \sum_{t=K+1}^T \mathbb{I}\{a_t = i\} \\
&\leq c + \sum_{t=K+1}^T \mathbb{I}\{a_t = i \wedge n_{t-1,i} \geq c\} \\
&\leq c + \sum_{t=K+1}^T \mathbb{I}\left\{\hat{\mu}_{n_{t-1}^*}^* + B_{t-1,n_{t-1}^*} \leq \max \hat{\mu}_{n_i(t-1),i} + B_{t-1,n_i(t-1)}\right\} \\
&\leq c + \sum_{t=K+1}^T \mathbb{I}\left\{\min_{0 < s < t} \hat{\mu}_s^* + B_{t-1,s} \leq \max_{c \leq s_i < t} \hat{\mu}_{s_i,i} + B_{t-1,s_i}\right\} \\
&\leq c + \sum_{t=1}^{\infty} \sum_{s=1}^{t-1} \sum_{s_i=c}^{t-1} \mathbb{I}\{\hat{\mu}_s^* + B_{t-1,s} \leq \hat{\mu}_{s_i,i} + B_{t-1,s_i}\}
\end{aligned}$$

When the indicator function is true one of the following holds:

$$\hat{\mu}_s^* \leq \mu^* - B_{t,s} \quad (12.2.3)$$

$$\hat{\mu}_{s_i,i} \geq \mu_i + B_{t,s_i} \quad (12.2.4)$$

$$\mu^* < \mu_i + 2B_{t,s_i} \quad (12.2.5)$$

*Proof idea*

- Bound the probability of the first two events.
- Choose  $c$  to bound the last term.

From Hoeffding bound:

$$\mathbb{P}(\hat{\mu}_s^* \leq \mu^* - B_{t,s}) \leq e^{-4 \ln t} = t^{-4} \quad (12.2.6)$$

$$\mathbb{P}(\hat{\mu}_{s_i,i} \geq \mu_i + B_{t,s_i}) \leq e^{-4 \ln t} = t^{-4} \quad (12.2.7)$$

Setting  $c = \lceil (8 \ln n)/\Delta_i^2 \rceil$  makes the last event false as  $s_i \geq c$ .

$$\mu^* - \mu_i - 2B_{t,s_i} = \mu^* - \mu_i - 2\sqrt{(2 \ln t)/s_i} \geq \mu^* - \mu_i - \Delta_i = 0.$$

Summing up all the terms completes the proof.

## 12.3 Structured bandit problems

### Bandits and optimisation

- Continuous stochastic functions Kocsis and Szepesvári [2006], Auer et al. [2007], Bubeck et al. [2011]
- Constrained deterministic distributed functions Ottens et al. [2012]

**First idea** Auer et al. [2007]

**Solve a sequence of discrete bandit problems.**

At epoch  $i$ , we have some interval  $A_i$

- Split the interval  $A_i$  in  $k$  regions  $A_{i,j}$
- Run UCB on the  $k$ -armed bandit problem.
- When a region is sub-optimal with high probability, remove it!

**Tree bandits** Bubeck et al. [2011]

Create a tree of coverings, with  $(h, i)$  being the  $i$ -th node at depth  $h$ .  $\mathcal{D}$  are the descendants and  $\mathcal{C}$  the children of a node.

At time  $t$  we pick node  $H_t, I_t$ . Each node is picked at most once.

$$n_{h,i}(T) \triangleq \sum_{t=1}^T \mathbb{I}\{(H_t, I_t) \in \mathcal{D}(h, i)\} \quad (\text{visits of } (h, i))$$

$$\hat{\mu}_{h,i}(T) \triangleq \frac{1}{n_{h,i}(T)} \sum_{t=1}^T r_t \mathbb{I}\{(H_t, I_t) \in \mathcal{C}(h, i)\} \quad (\text{reward from } (h, i))$$

$$C_{h,i}(T) \triangleq \hat{\mu}_{h,i}(T) + \sqrt{\frac{2 \ln T}{n_{h,i}(T)}} + nu_1 \rho^h \quad (\text{confidence bound})$$

$$B_{h,i}(T) \triangleq \min \left\{ C_{h,i}(T), \max_{(h+1,j) \in \mathcal{C}(h,i)} B_{h+1,j} \right\} \quad (\text{child bound})$$

## 12.4 Reinforcement learning problems

### 12.4.1 Optimality Criteria

In all previous cases, we assumed a specific discount rate, or horizon for our problem. Now we shall examine different choices and how they affect the existence of an optimal policy. As mentioned previously, the following two views of discounted reward processes are equivalent.

**Infinite horizon, discounted**

Discount factor  $\gamma$  such that

$$U_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k} \quad \Rightarrow \quad \mathbb{E} U_t = \sum_{k=0}^{\infty} \gamma^k \mathbb{E} r_{t+k} \quad (12.4.1)$$

**Geometric horizon, undiscounted**

At each step  $t$ , the process terminates with probability  $1 - \gamma$ :

$$U_t^T = \sum_{k=0}^{T-t} r_{t+k}, \quad T \sim \text{Geom}(1 - \gamma) \Rightarrow \mathbb{E} U_t = \sum_{k=0}^{\infty} \gamma^k \mathbb{E} r_{t+k} \quad (12.4.2)$$

$$V_{\gamma}^{\pi}(s) \triangleq \mathbb{E}(U_t \mid s_t = s)$$

**The expected total reward criterion**

$$V_t^{\pi, T} \triangleq \mathbb{E}_{\pi} U_t^T, \quad V^{\pi} \triangleq \lim_{T \rightarrow \infty} V^{\pi, T} \quad (12.4.3)$$

**Dealing with the limit**

- Consider  $\mu$  s.t. the limit exists  $\forall \pi$ .

$$V_+^{\pi}(s) \triangleq \mathbb{E}_{\pi} \left( \sum_{t=1}^{\infty} r_t^+ \mid s_t = s \right), \quad V_-^{\pi}(s) \triangleq \mathbb{E}_{\pi} \left( \sum_{t=1}^{\infty} r_t^- \mid s_t = s \right) \quad (12.4.4)$$

$$r_t^+ \triangleq \max\{-r, 0\}, \quad r_t^- \triangleq \max\{r, 0\}. \quad (12.4.5)$$

- Consider  $\mu$  s.t.  $\exists \pi^*$  for which  $V^{\pi^*}$  exists and

$$\lim_{T \rightarrow \infty} V^{\pi^*, T} = V^{\pi^*} \geq \limsup_{T \rightarrow \infty} V^{\pi, T}.$$

- Use optimality criteria sensitive to the divergence rate.

**The average reward (gain) criterion****The gain  $g$** 

$$g^{\pi}(s) \triangleq \lim_{T \rightarrow \infty} \frac{1}{T} V^{\pi, T}(s) \quad (12.4.6)$$

$$g_+^{\pi}(s) \triangleq \limsup_{T \rightarrow \infty} \frac{1}{T} V^{\pi, T}(s), \quad g_-^{\pi}(s) \triangleq \liminf_{T \rightarrow \infty} \frac{1}{T} V^{\pi, T}(s) \quad (12.4.7)$$

If  $\lim_{T \rightarrow \infty} \mathbb{E}(r_T \mid s_0 = s)$  exists then it equals  $g^\pi(s)$ .

Let  $\Pi$  be the set of all history-dependent, randomised policies.

Using our overloaded symbols, we have that  $\pi^*$  is *total reward optimal* if

$$V^{\pi^*}(s) \geq V^\pi(s) \quad \forall s \in \mathcal{S}, \pi \in \Pi.$$

$\pi^*$  is *discount optimal* for  $\gamma \in [0, 1)$  if

$$V_\gamma^{\pi^*}(s) \geq V_\gamma^\pi(s) \quad \forall s \in \mathcal{S}, \pi \in \Pi.$$

$\pi^*$  is *gain optimal* if

$$g^{\pi^*}(s) \geq g^\pi(s) \quad \forall s \in \mathcal{S}, \pi \in \Pi.$$

### Overtaking optimality

$\pi^*$  is *overtaking optimal* if

$$\liminf_{T \rightarrow \infty} [V^{\pi^*, T}(s) - V^{\pi, T}(s)] \geq 0 \quad \forall s \in \mathcal{S}, \pi \in \Pi.$$

However, no overtaking optimal policy may exist.

$\pi^*$  is *average-overtaking optimal* if

$$\liminf_{T \rightarrow \infty} \frac{1}{T} [V^{\pi^*, T}(s) - V_+^\pi(s)] \geq 0 \quad \forall s \in \mathcal{S}, \pi \in \Pi.$$

### Sensitive discount optimality

$\pi^*$  is *n-discount optimal* for  $n \in \{-1, 0, 1, \dots\}$  if

$$\liminf_{\gamma \uparrow 1} (1 - \gamma)^{-n} [V_\gamma^{\pi^*}(s) - V_\gamma^\pi(s)] \geq 0 \quad \forall s \in \mathcal{S}, \pi \in \Pi.$$

A policy is *Blackwell optimal* if  $\forall s, \exists \gamma^*(s)$  such that

$$V_{\gamma^*}^{\pi^*}(s) - V_{\gamma^*}^\pi(s) \geq 0, \quad \forall \pi \in \Pi, \gamma^*(s)\gamma < 1.$$

**Lemma 12.4.1.** *If a policy is m-discount optimal then it is n-discount optimal for all  $n \leq m$ .*

**Lemma 12.4.2.** *Gain optimality is equivalent to  $-1$ -discount optimality.*

## 12.4.2 UCRL

### An upper-confidence bound algorithm

Confidence region  $M_t$  such that

$$\mathbb{P}(\mu \notin M_t) < \delta \tag{12.4.8}$$

Optimistic value for policy  $\pi$ :

$$V_+^\pi(M_t) \triangleq \max \{V_\mu^\pi \mid \mu \in M_t\} \tag{12.4.9}$$



**UCRL Jacksh et al. [2010] outline**

- At round  $k$ , start time  $t_k$ , calculate  $M_{t_k}$ .
- Choose  $\pi_k \in \arg \max_{\pi} V_+^{\pi}(M_{t_k})$ .
- Execute  $\pi_k$ , observe rewards and update model until  $t_{k+1}$ .

**The confidence region**

Let  $M_t$  be a set of plausible MDPs for time  $t$  with transitions  $\tau$  s.t.:

$$\left\| \mathbf{P}(\cdot | s, a) - \hat{\mathbf{P}}_t(\cdot | s, a) \right\|_1 \leq \sqrt{\frac{n \ln T}{N_t(s, a)}}, \quad \forall s \in \mathcal{S}, a \in \mathcal{A}, \quad (12.4.10)$$

where  $\hat{\mathbf{P}}_t(\cdot | s, a)$  is the empirical transition probability.

Then  $\mathbb{P}(\mu \in M_t) > 1 - nkT^{-2}$ , via a bound due to Weissman Weissman et al. [2003].

***Changing set of plausible MDPs***

- This implies that we may have to switch policies.
- We do so when  $N_t(s, a)$  doubles for some  $s, a$ .

**Calculating the upper bound**

In effect, create an *augmented MDP*

$$Q_t(s, a) = r(s, a) + \max \left\{ \sum_{s' \in \mathcal{S}} \mathbf{P}(s' | s, a) V_{t+1}(s') \mid \|\mathbf{P} - \hat{\mathbf{P}}\|_1 \leq \epsilon \right\} \quad (12.4.11)$$

$$V_t(s) = \max_{a \in \mathcal{A}} Q_t(s, a) \quad (12.4.12)$$

**Comparison with Bayesian upper bound****High-probability value function bound**

$$V_+^* = \max \{ V_{\mu}^* \mid \mu \in M_t \}, \quad \mathbb{P}(\mu^* \in M_t) \geq 1 - \delta.$$

**Highly credible value function bound**

$$V_+^* = \max \{ V_\mu^* \mid \mu \in M_t \}, \quad \xi_t(M_t) \geq 1 - \delta.$$

**Bayesian value function bound (e.g. Dimitrakakis [2011])**

$$V_+^* = \int_{\mathcal{M}} V_\mu^* d\xi_t(\mu) \quad \xi_t = \xi_0(\cdot \mid s_t, r_t, \dots)$$

**12.4.3 Bibliographical remarks**

Different optimality criteria are treated in detail in Puterman [1994] Chapter 5.

**.1 Symbols**

## **.2 Index**

# Index

Bayes' theorem, 34

concave function, 40

convex function, 40

gamma function, 64

Jensen inequality, 40

likelihood

- conditional, 33
- relative, 30

linear programming, 141

Markov decision process, 118, 120, 121,  
**123**, 126, 143

martingale, 113

policy iteration

- modified, 139
- temporal-difference, 140

preference, 36

probability

- subjective, 30

reward, 35

utility, 37

Utility theory, 35

value iteration, 136

Wald's theorem, 112

**.3 Glossary**

# Bibliography

- A. Antos, C. Szepesvári, and R. Munos. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71(1):89–129, 2008a.
- Andr  Antos, R mi Munos, and Csaba Szepesvari. Fitted q-iteration in continuous action-space mdps. In J.C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*. MIT Press, Cambridge, MA, 2008b.
- Robert B. Ash and Catherine A. Dole ans-Dade. *Probability & Measure Theory*. Academic Press, 2000.
- J. Asmuth, L. Li, M. L. Littman, A. Nouri, and D. Wingate. A Bayesian sampling approach to exploration in reinforcement learning. In *UAI 2009*, 2009.
- P. Auer, R. Ortner, and C. Szepesvari. Improved Rates for the Stochastic Continuum-Armed Bandit Problem. In *Learning Theory*, volume 4539 of *Lecture Notes in Computer Science*, page 454. Springer, 2007.
- Peter Auer, Nicol  Cesa-Bianchi, and Paul Fischer. Finite time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2/3):235–256, 2002.
- Richard Ernest Bellman. A problem in the sequential design of experiments. *Sankhya*, 16:221–229, 1957.
- A. Bernstein. Adaptive state aggregation for reinforcement learning. Master’s thesis, Technion – Israel Institute of Technology, 2007.
- Dimitri P. Bertsekas and John N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1996.
- J.A. Boyan. Technical update: Least-squares temporal difference learning. *Machine Learning*, 49(2):233–246, 2002.
- S.J. Bradtke and A.G. Barto. Linear least-squares algorithms for temporal difference learning. *Machine Learning*, 22(1):33–57, 1996.
- S bastien Bubeck, R mi Munos, Gilles Stoltz, and Csaba Szepes v ri. X-armed bandits. *Journal of Machine Learning Research*, 12:1655–1695, 2011.
- Herman Chernoff. Sequential design of experiments. *Annals of Mathematical Statistics*, 30(3):755–770, 1959.

- Herman Chernoff. Sequential Models for Clinical Trials. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Vol.4*, pages 805–812. Univ. of Calif Press, 1966.
- Richard Dearden, Nir Friedman, and Stuart J. Russell. Bayesian Q-learning. In *AAAI/IAAI*, pages 761–768, 1998. URL [citeseer.ist.psu.edu/dearden98bayesian.html](http://citeseer.ist.psu.edu/dearden98bayesian.html).
- Morris H. DeGroot. *Optimal Statistical Decisions*. John Wiley & Sons, 1970.
- Christos Dimitrakakis. Tree exploration for Bayesian RL exploration. In *Computational Intelligence for Modelling, Control and Automation, International Conference on*, pages 1029–1034, Wien, Austria, 2008. IEEE Computer Society. ISBN 978-0-7695-3514-2. doi: <http://doi.ieeecomputersociety.org/10.1109/CIMCA.2008.32>.
- Christos Dimitrakakis. Complexity of stochastic branch and bound methods for belief tree search in Bayesian reinforcement learning. In *2nd international conference on agents and artificial intelligence (ICAART 2010)*, pages 259–264, Valencia, Spain, 2010. ISNTICC, Springer.
- Christos Dimitrakakis. Robust bayesian reinforcement learning through tight lower bounds. In *European Workshop on Reinforcement Learning (EWRL 2011)*, number 7188 in LNCS, pages 177–188, 2011.
- Christos Dimitrakakis. Monte-carlo utility estimates for bayesian reinforcement learning. In *IEEE 52nd Annual Conference on Decision and Control (CDC 2013)*, 2013. arXiv:1303.2506.
- Christos Dimitrakakis and Michail G. Lagoudakis. Rollout sampling approximate policy iteration. *Machine Learning*, 72(3):157–171, September 2008. doi: 10.1007/s10994-008-5069-3. Presented at ECML’08.
- Michael O’Gordon Duff. *Optimal Learning Computational Procedures for Bayes-adaptive Markov Decision Processes*. PhD thesis, University of Massachusetts at Amherst, 2002.
- Milton Friedman and Leonard J. Savage. The expected-utility hypothesis and the measurability of utility. *The Journal of Political Economy*, 60(6):463, 1952.
- Thomas Furnston and David Barber. Variational methods for reinforcement learning. In Yee Whye Teh and Mike Titterton, editors, *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 9 of *JMLR : W&CP*, pages 241–248, Chia Laguna Resort, Sardinia, Italy, 2010.
- Mohammad Ghavamzadeh and Yaakov Engel. Bayesian policy gradient algorithms. In *NIPS 2006*, 2006.
- C. J. Gittins. *Multi-armed Bandit Allocation Indices*. John Wiley & Sons, New Jersey, US, 1989.



- Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, March 1963.
- Thomas Jacksh, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.
- Levente Kocsis and Csaba Szepesvári. Bandit based Monte-Carlo planning. In *Proceedings of ECML-2006*, 2006.
- AN Kolmogorov and SV Fomin. *Elements of the theory of functions and functional analysis*. Dover Publications, 1999.
- M.G. Lagoudakis and R. Parr. Least-squares policy iteration. *The Journal of Machine Learning Research*, 4:1107–1149, 2003.
- R. Munos and C. Szepesvári. Finite-time bounds for fitted value iteration. *The Journal of Machine Learning Research*, 9:815–857, 2008.
- Brammert Ottens, Christos Dimitrakakis, and Boi Faltings. DUCT: An upper confidence bound approach to distributed constraint optimization problems. In *AAAI 2012*, 2012.
- P. Poupart, N. Vlassis, J. Hoey, and K. Regan. An analytic solution to discrete Bayesian reinforcement learning. In *ICML 2006*, pages 697–704. ACM Press New York, NY, USA, 2006.
- Pascal Poupart and Nikos Vlassis. Model-based Bayesian reinforcement learning in partially observable domains. In *International Symposium on Artificial Intelligence and Mathematics (ISAIM)*, 2008.
- Marting L. Puterman. *Markov Decision Processes : Discrete Stochastic Dynamic Programming*. John Wiley & Sons, New Jersey, US, 1994.
- H. Robbins and S. Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407, 1951.
- Leonard J. Savage. *The Foundations of Statistics*. Dover Publications, 1972.
- S. Singh, T. Jaakkola, and M.I. Jordan. Reinforcement learning with soft state aggregation. *Advances in neural information processing systems*, pages 361–368, 1995.
- Malcolm Strens. A Bayesian framework for reinforcement learning. In *ICML 2000*, pages 943–950, 2000.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- Marc Toussaint, Stefan Harmelign, and Amos Storkey. Probabilistic inference for solving (PO)MDPs, 2006.
- T. Weissman, E. Ordentlich, G. Seroussi, S. Verdu, and M.J. Weinberger. Inequalities for the  $L_1$  deviation of the empirical distribution. *Hewlett-Packard Labs, Tech. Rep*, 2003.