## Machine learning in science and society

From automated science to beneficial artificial intelligence

Christos Dimitrakakis

January 24, 2023

## **Preface**

This book is a technical introduction to important aspects of machine learning in science and society. It covers the most fundamental concepts in privacy, fairness, causality, reproducibility and experiment design. While I have tried to make the book as rigorous as possible, in the interest of brevity some technical details are omitted.

However, even though the book is not addressed to a technical audience, every concept is necessarily connected to precise technical definition. This includes the bare minimum amount of theoretical development for a thorough understanding of the main concepts. Throughout the book, we use Bayesian decision theory and graphical models to offer a unifying perspective on those issues.

Throughout the book, we will use the following blocks of text for particular types of material.

#### **♦** An example.

Such blocks contain illustrations and examples.

#### **∧** An important note.

Important caveats or properties of algorithms will be given in such blocks.

#### Theoretical note.

Discussions of a more theoretical nature, or notes that require some additional thinking will be placed in such blocks.

#### **Exercise**.

These exercises are meant to be done individually, either while reading the text, or in class.

#### **Group** activity.

These activities are meant to be done in class.

The first chapter offers an introduction to the field of machine learning in general and to fairness, privacy and reproducibility in particular. Chapter 2 discusses privacy. Chapter 3 talks about fairness.

I whole-heartedly recommend the book "The Ethical Algorithm" as a non-technical companion to this book and "The foundations of differential privacy" as a more thorough overview of differential privacy.

# Contents

1	Intr	oduction 7
	1.1	Introduction to machine learning
	1.2	Data analysis, learning and planning
		1.2.1 Experiment design
		1.2.2 Inference
	1.3	Book overview
<b>2</b>	Priv	racy 19
	2.1	Things we do with data
	2.2	Statistical disclosure
	2.3	Algorithmic privacy
	2.4	Simple anonymisation and $k$ -anonymity
	2.5	Differential privacy
	2.6	The local and central privacy models
		2.6.1 Properties of differential privacy
	2.7	The Laplace mechanism
	2.8	Interactive data access
		2.8.1 Utility of queries
		2.8.2 Differential privacy as a hypothesis testing game
	2.9	Advanced topics
		2.9.1 Relaxations of differential privacy
	2.10	Discussion
		Exercises
	2.11	22.00.0000
3	Fair	ness. 47
	3.1	Introduction
	3.2	Group fairness
		3.2.1 Measuring group fairness
		3.2.2 Trading off utility and fairness
	3.3	Individual fairness
	3.4	Exercises
4	Ren	roducibility 59
	4.1	Reproducibility
		4.1.1 The human as an algorithm
		4.1.2 Algorithmic sensitivity
		4.1.3 Beyond the data you have: simulation and replication

6 CONTENTS

5	Cau	ısality	69
	5.1	Introduction	69
		5.1.1 Common structural assumptions	72
	5.2	Interventions	76
	5.3	Policy evaluation and optimisation	78
		5.3.1 Monte-Carlo estimation	79
	5.4	Individual effects and counterfactuals	81
		5.4.1 Disturbances and structural equation models	81
		5.4.2 Example: Learning instrumental variables	82
		5.4.3 Discussion	83
		5.4.4 Exercises	84
A	Dec	sision and Learning Theory	85
	A.1	Hierarchies of decision making problems	85
		A.1.1 Simple decision problems	85
		A.1.2 Decision rules	88
		A.1.3 Statistical testing*	89
	A.2	Formalising classification problems	97
	A.3	Beliefs and probabilities	99
		A.3.1 Probability and Bayesian inference	104
	A.4	Classification with stochastic gradient descent*	108
		A.4.1 Neural network models	110
	A.5	Nearest neighbours	113
	A.6	Naive Bayes classifiers*	118
В	Em	pirical estimates and sampling	<b>121</b>
$\mathbf{C}$	Gra	phical models	123
	C.1	•	127
	C.2		129
	C.3		130
	C.4	9	132

## Chapter 1

## Introduction

### 1.1 Introduction to machine learning

Problems in machine learning are similar to problems in science. Scientists must plan experiments intelligently and collect data. The must be able to use the data to verify a different hypothesis. More generally, they must be able to make decisions under uncertainty (Without uncertainty, there would be no need to gather more data). Similar problems appear in more mundane tasks, like learning to drive a car.

For that reason, science is a very natural application area for machine learning. We can model the effects of climate change and how to mitigate it; discover structure in social networks; map the existence of dark matter in the universe by intelligently shifting through weak gravitational lens data, and not only study the mechanisms of protein folding, but discover methods to synthesize new drugs.

We must be careful, however. In many cases we need to be able to interpret what our model tells us. We also must make sure that the any results we obtain are reproducible.

While machine learning models in science are typically carefully handcrafted by scientists and experts in machine learning and statistics, this is not typically the case in everyday applications. Nevertheless, well-known or home-grown machine learning models are being deployed across the application spectrum. This involve home assistants that try and get you want, web advertising, which tries to find new things for you to want, lending, which tries to optimally lend you money so that you buy what you didn't need before. We also have autonomous vehicles, which take you were you want to go, and ridesharing services, which do the same thing, but use humans instead. Finally, there are many applications in public policy, such as crime prevention, justice, and disease control which use machine learning. In all those cases, we have to worry about a great many things that are outside the scope of the machine learning problems itself. These are (a) privacy: you don't want your data used in ways that you have not consented to (b) fairness: you don't want minorities to be disadvantaged and (c) safety: you don't want your car to crash.

## 1.2 Data analysis, learning and planning

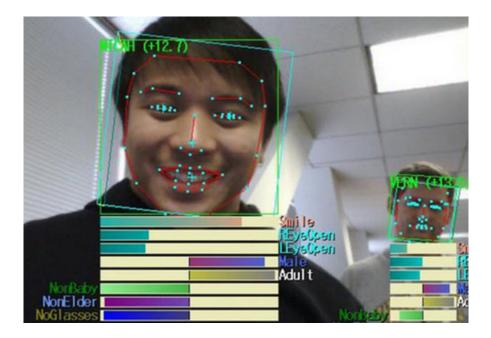
To make the above more concrete, let's have a look at a number of problems in machine learning. These involve learning from and analysing data, including inferring decision rules, and constructing complex plans using the evidence gleaned from the data. Machine learning problems are commonly separated in three different types: supervised, unsupervised and reinforcement

learning. Typical supervised learning problems include classification and regression, while unsupervised problems include compression, clustering and topic modelling. Reinforcement learning, on the other hand, is concerned with artificially intelligent agents more generally, with examples including game playing and adaptive control. Their main differences are two. Firstly, the type of feedback we have about learning performance. Secondly, and perhaps more importantly, whether or not the problem involves active data collection. In this course, we will try and take a global view of these problems in the context of decision theory.

#### Can machines learn from data?



An unsupervised learning problem: topic modelling

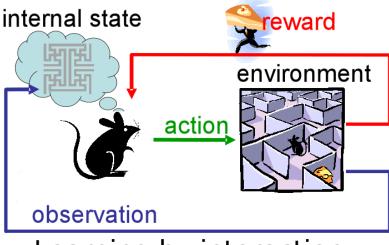


A supervised learning problem: object recognition

You can use machine learning just to analyse, or find structure in the data. This is generally called unsupervised learning. One such example is topic modelling, where you let the algorithm find topics from a corpus of text. These days machines are used to learn from in many applications. These include speech recognition, facial authentication, weather prediction, etc. In general, in these problems we are given a *labelled* dataset with, say, example images from each class. Unfortunately this does not scale very well, because obtaining labels is expensive.

This is partially how science works, because what we need to do is to find a general rule of nature from data. Starting from some hypothesis and some data, we reach a conclusion. However, many times we may need to actively experiment to obtain more data, perhaps because we found that our model is wrong.

Can machines learn from their mistakes?



Learning by interaction

So, what happens when we make a mistake? Can we somehow recognise it? Humans and other animals can actually learn from their mistakes. Consider the proverbial rat in the maze. At some intervals, the experimenter places some cheese in there, and the rat must do a series of actions to obtain it, such as navigating the maze and pulling some levers. It doesn't know how to get to the cheese easily, but it slowly learns the layout of the maze through observation, and in the end, through trial-and-error it is able to get to the cheese very efficiently.

We can formalise this as a reinforcement learning problem, where the rat takes a series of actions; at each step it also obtains a reward, let's say equal to 0 when it has no cheese, and 1 when it eats cheese. Then we can declare that the rat's utility is the sum of all rewards over time, i.e. the total amount of cheese it can eat before it dies. The rat needs to explore the environment in order to be able to get to the cheese.

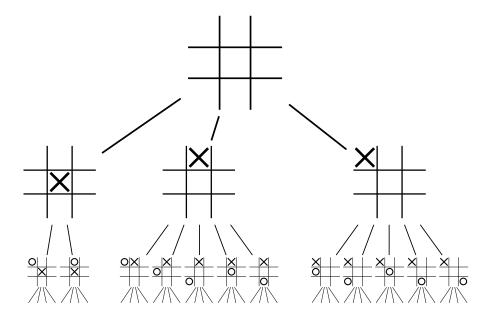
An example in robotics is trying to teach a robot to flip pancakes. One easy thing we can try is to show the robot how to do it, and then let it just copy the demonstrated movement. However, this doesn't work! The robot needs to explore variations of the movement, until it manages to successfully flip pancakes. Again, we can formulate this as a reinforcement learning problem, with a reward that is high whenever the pancake's position is flipped, and on the pan; and low everywhere else. Then the robot can learn to perform this behaviour through trial and error. It's important to note that in this example, merely demonstration is not enough. Neither is reinforcement learning enough. The same thing is true for the recent success of AlphaGo in beating a master human: apart from planning, they used both demonstration data and self-play, so that it could learn through trial and error.

Can machines make complex plans?



I suppose the first question is whether machines can plan ahead. Indeed, even for large problems, such as Go, machines can now perform at least as well as top-rated humans. How is this achieved?

Machines can make complex plans!



The basic construction is the planning tree. This is an enumeration of all possible future events. If a complete enumeration is impossible, a partial tree is constructed. However this requires evaluating non-terminal game positions. In the old times, this was done with heuristics, but now this is data-driven, both through the use of expert databases, and through self-play and reinforcement learning.

#### 1.2.1 Experiment design

An example that typifies trial and error learning are bandit problems. Imagine that you are in a Casino and you wish to maximise the amount of money you make during the night. There are a lot of machines to play. If you knew which one was the best, then you'd just play it all night long. However, you must also spend time trying out different machines, in order to get an estimate of how much money each one gives out. The trade off between trying out different machines and playing the one you currently think is best is called the exploration-exploitation trade-off and it appears in many problems of experiment design for science.

Let's say we want to build a robot scientist and tell it to discover a cure for cancer. What does the scientist do and how can the robot replicate it??

Simplifying the problem a bit, consider that you have a large number of drug candidates for cancer and you wish to discover those that are active against it. The ideas is that you select some of them, then screen them, to sort them into active and inactive. However, there are too many drugs to screen, so the process is interactive. At each cycle, we select some drugs to screen, classify them, and then use this information to select more drugs to screen. This cycle, consequently has two parts: 1. Selecting some drugs given our current knowledge. 2. Updating our knowledge given new evidence.

#### Drawing conclusions from results

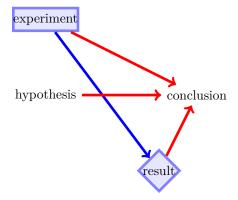


Figure 1.1: Dependence diagram between selection of an experiment, formulation of a hypothesis, and drawing of a conclusion. The result depends only on the experiment. However, the conclusion depends on the experiment, hypothesis and the obtained result. The red lines indicate computational dependencies, while the blue lines indicate physical dependencies.

In general, we would like to have some method which can draw conclusions from results. This involves starting with a hypothesis, performing an experiment to verify or refute it, obtain some experimental result; and then concluding for or against the hypothesis. Here the arrows show dependencies between these variables. So what do we mean by "hypothesis" in this case?

#### 1.2.2 Inference.

Let's take the example of planetary orbits. Here Tycho famously spent 20 years experimentally measuring the location of Mars. He had a hypothesis: that planetary orbits were circular, but he didn't know which were the right orbits. When he tried to fit his data to this hypothesis, he concluded a specific circular orbit for Mars ...around Earth.

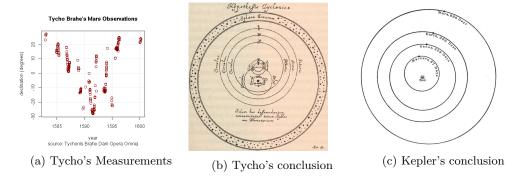
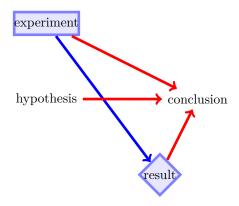


Figure 1.2: How given the same data, one can reach different conclusion depending on one's modelling assumptions.

Kepler had a more general hypothesis: that orbits could be circular or elliptic, and he actually accepted that the planets orbited the sun. This led him to the broadly correct model of all planets being in elliptical orbits around the sun. However, the actual verification that all things do not revolve around earth, requires different experiments.

Later on, Gauss collected even more experimental data to calculate the orbit of Ceres. He



did this using one of the first formal statistical methods; this allowed him to avoid cheating (like Kepler did, to accentuate his finding that orbits were elliptical).

It is quite easy to draw the wrong conclusions from applying machine learning / statistics to your data. For example, it was fashionable to perform fMRI studies in humans to see whether some neurons have a particular functional role. There were even articles saying that "we found the neurons encoding for Angelina Jolie". So some scientists tried to replicate those results. They took a dead salmon, and put it an fMRI scanner. They checked its brain activity when it was shown images of happy or sad people. Perhaps surprisingly, they found an area of the brain that was correlated with the pictures - so it seemed, as though the dead salmon could distinguish photos of happy people from sad ones. However, this was all due to a misapplication of statistics. In this course, we will try and teach you to avoid such mistakes.

#### A simple simulation study

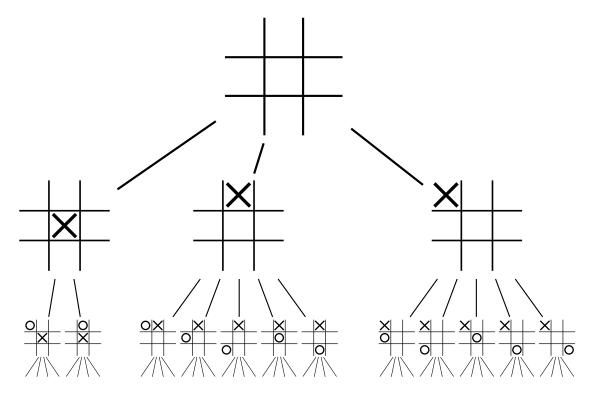
Sometimes we want to use a simple simulation study to understand how well our methods work. The following code is an example of how to do this. Here we are doing a simplified fMRI analysis, but the general idea is not significaey, antly different from what people actually do in the field. src/reproducibility/mri\_analysis.ipynb

#### Planning future experiments

So far we have focused only on the problem of data analysis. However, we also need to think about the problem of planning for experiments. This is called *experiment design*. Experiment designs are usually fixed. In that case, we can use assumptions about the data and how much accuracy we need to design the experiment. However, it is also possible to have an adaptive experiment design where our future experiments depend on our current conclusions.

In general, optimal experiment design is indeed difficult, especially in setting such as drug discovery where the number of experiments is huge. However, conceptually, there is a simple and elegant solution to this problem.

#### Planning experiments is like Tic-Tac-Toe



The basic idea is to think of experiment design as a game between the scientist and Nature. At every step, the scientist plays an X to denote an experiment. Then Nature responds with an Observation. The main difference from a game is that Nature is (probably) not adversarial. We can also generalise this idea to problems in robotics, etc.

These kinds of techniques, coming from the reinforcement learning literature have been successfully used at the university of Manchester to create a robot, called Eve, that recently (re)-discovered a malaria drug.

#### 1.3 Book overview

While in the past machine learning was mainly a field of research, with only some applications, but now machine learning is pervasive: Our phones, cars, watches, doorbells, kettles are connected to the internet and send a continuous stream of data to companies. In addition, many companies and government actors use machine learning algorithms to make or support decisions. This creates a number of problems in privacy, fairness and safety.

#### The view from statistics

While in machine learning people generally discuss mainly supervised, unsupervised or reinforcement learning, these are actually three rather broad, but still limited categories of problems. There are many other problems, such as semi-supervised learning, active learning, imitation/apprenticeship learning, inverse reinforcement learning, preference elicitation, adaptive control, and possibly many others to come. The statistical view is that there basically three types of problems:

**Inference.** Given what we know, what can we say about how the world works, the current state of the world or events in the past?

**Prediction.** Can we predict specific evens in the future? This frequently is done through some type of inference.

**Decision making.** Given what we know, and what we want to achieve, what is the best decision we can make? This typically requires some ability to predict the effect of our actions. We will encounter many specific inference, prediction and decision making tasks during this course.

#### Technical topics

The book covers a number of technical topics. Sections marked with an asterisk (\*) may be skipped safely upon a first reading without compromising understanding of the remaining material.

#### Machine learning problems

- Unsupervised learning. Loosely speaking, this is simply the problem of estimating some structure from data. In statistical terms, it is usually the problem of estimating some joint distribution of random variables under some model assumptions. Problems in unsupervised learning include clustering, anomaly detection, compression.
- Supervised learning. In this setting data can be split in two groups of variables. One group that is always available, and another group that must be predicted. A special case of the problem is when we wish to estimate some function  $f: \mathcal{X} \to \mathcal{Y}$  from data. Classical problems in this setting are classification and regression.
- Reinforcement learning. This is a very general sequential decision problem, where an agent must learn how to behave optimally in an unknown environment only by limited feedback and reinforcement. The standard setting involves the agent trying to maximise its (expected) cumulative reward over time.

#### Algorithms and models

- Bayesian inference and decision theory. We will cover the basic ideas of probabilistic graphical models and decision diagrams more generally. These allow us to formalise statistical inference and decision problems, as well as conditional dependence relationships between variables. It is also a way for us to represent model architectures, as well as some fairness notions.
- Differential privacy. This book focuses on an information-theoretic notion of privacy, differential privacy. This quantifies the amount of information leakage about any one person due to the public release of aggregate statistics.
- Stochastic optimisation and neural networks. Optimisation is at the core of modern machine learning techniques. We will cover the basic of stochastic gradient descent. We will also examine neural networks from both a probabilistic and optimisation perspective.

#### Further reading

- $\bullet$  Bennett et al. <sup>2</sup> describe how the usual uncorrected analysis of fMRI data leads to the conclusion that the dead salmon can reason about human images.
- Bennett et al. 1 discuss how to perform analyses of medical images in a principled way. They also introduce the use of simulations in order to test how well a particular method is going to perform.

#### Resources

- Course code and notes: https://github.com/olethrosdc/ml-society-science
- $\bullet \ \ Book: \ https://github.com/olethrosdc/ml-society-science/book.pdf$

## Chapter 2

# **Privacy**

One interpretation of privacy is simply the ability to maintain a personal secret. Is it possible to maintain perfect secrecy? Can the thoughts in our head be perfectly safe, or can they be revealed indirectly through our actions?

While we certainly *can* securely store and transmit data using cryptography, the problems we will discuss in this chapter are *not* solvable solely through cryptography. We are interested in scenarios where we must make a public decision or release some summary statistics that depend on private data, in such a way as to minimise harm to the individuals contributing their data.

If we never have to reveal any of our computations, cryptography is what is needed. For example, through homomorphic computation, an untrusted party can even perform some computations on encrypted data, returning an encrypted result to us, while learning nothing about the original secret. As long as the data, and the results of any computation on it, are kept under lock and key, our personal information cannot be revealed.

However, sometimes we must make public decisions or release public information based on this data. Then it can be revealed indirectly. For example, you can trust your doctor to maintain confidentiality, but when you go to the pharmacy to get the prescribed medicine, somebody can infer the medical condition you suffer from.

It is even possible to learn personal information from aggregate statistics. Let us say your doctor publishes a list of cases of different diseases every week, together with some other information such as the approximate patient age. Even though your own data is mixed with that of all other patients, it is possible to infer your diagnosis, especially with some additional side-information: If somebody knows you were the only person in that age group visiting the doctor that week, they will learn your diagnosis.

From that point of view, it is not the data itself, but the complete process of data collection, treatment and public release that can be characterised as private or non-private. For that reason, we will emphasise an *algorithmic* view of privacy: participants entrusts their data to an algorithm, which produces a useful output in return. The algorithmic process is typically not fully automated, as it also depends on some human input: One example is a medical study examining different treatments for a disease. While humans select and administer the treatments, they will typically rely on a randomised strategy for assigning treatments to individuals, and use a statistical method to report their results. Given this mixture of ad-hoc decisions and formal algorithmic methods, is it possible to guarantee privacy in any sense? What kind of guarantees can we make?

Generally speaking, an algorithm has good privacy properties, if the amount of information that can be revealed through the algorithm's output about any individual contributing data



Just because they're the problem,

is bounded. In particular, we are interested in how much an adversary can learn about any individual's input to the algorithm from the algorithm's output.

This does not preclude learning general facts about individuals from the output. For example, a study about the use of steroids in sports may show that 90% of sprinters with times under 10 seconds are using steroids, while only 50% of slower sprinters do so. Any sprinter with a time under 10 seconds is thus suspected of using steroids by association. However, it does not matter if their data has been used in the study. The publication of the result does not impact their privacy, but the amount of harm it does to them does not depend in their participation: the same statistical result would have been obtained with or without them.

In this chapter, we will look at two formal concepts of privacy protection: k-anonymity and differential privacy. The first is a simple method for anonymising databases. However, it provides only limited resistance to identification. The latter is a more general concept, which provides full information-theoretic protection to individuals. A major problem with any privacy definition and method, however is correct interpretation of the privacy concept used, and correct implementation of the algorithm used.

## 2.1 Things we do with data

How do we use data in the first place? Sometimes we simply publish the data, perhaps after some initial processing. Publication of datasets is useful for researchers that want to do perform further analysis on the data. Usually though, we collect data in order to calculate specific statistics. For example the census collects data about the number of people in different households, wages, etc, and then publishes tables detailing the average age and number of people per household in different areas of the country. This demographic information is useful for policy makers, urban planning to organise voting centers and the distribution of police and fire stations, as well as other public services.

Fundamentally, privacy in statistics is an *issue of trust*. The analyst, be it a human, or an automated service, will use your data to make decisions. You must also decide who to trust and how much. Do we trust the data analyst? How much privacy are we willing to sacrifice to the analyst? How much to the public at large? What you want out of the service. Is the service important enough to sacrifice significant amounts of privacy? What is an acceptable trade-off between utility and privacy? These are difficult questions and are hard to quantify, hence we assume that we have already decided how to answer them, and we simply want to find an appropriate methodology for achieving a good result.

Consider a researcher wishing to collect data for a statistical analysis. If the analysis is eventually published, this creates two possible scenarios for that may lead to privacy violations.

- Publication of "anonymised" data. Sometimes we may collect data in order to publish the dataset itself for other researchers to use. This is common practice in machine learning, with image classification datasets being a good example. However, there is always some privacy risk even if identifying information such as names are removed.
- Public data analysis. In this setting, we only publish summary statistics or models about the data. A prime example is a national census analysis, which may provide detailed demographic information for all towns and regions in a nation. Although only aggregate data is published, it is theoretically possible to infer personal information. For that reason, the US Census Bureau conducted its analysis in 2020 using differential private algorithms.

#### 

Cryptography provides secure communication and computation, authentication and verification. These are used to establish secure channels with somebody that we trust. However, the privacy violations we are concerned with relate to *publicly released* outputs of algorithms. It is not important whether or not all the data and computation are encrypted: as long as the algorithm generates a public output, it is theoretically possible for somebody to learn something about the algorithm's input.

#### 2.2 Statistical disclosure

"Data is everywhere", said the statistician, "data, data!".

In the past, statistical analysis was performed with laboriously collected and annotated datasets. Even as recently as in the early 21st century, databases for machine learning were limited to a few thousand entries at most. At the time of writing, not only has the size of datasets become extremely large, but the sources of data are much more diverse. Data are collected and commercialised whenever we visit a website and even as we walk around with our phone. To a limited extent, there is a tradeoff between what we can get out of a service and what we pay into it. Many free services such as navigation software rely on collecting user data to perform better: If you can tell that there is a traffic jam in Central Avenue, you can after all try and take another route. As long as informed consent exists, use of private data is generally regarded as unproblematic.<sup>2</sup>

However, even apparently benign data collected with appropriate consent can lead to serious and unexpected privacy violations. There are three famous examples of this: Firstly, the identification of people in supposedly anonymous health data in the 1990s in the state of Massachussets, which we will go over in detail in this chapter. Secondly, the identifications of users through anonymised movie ratings in the Netflix dataset. Finally, the ability to discover if any given individual's data is contained in a pooled genomic study.

#### Anonymisation

 $<sup>^{1}</sup>$ If somebody knows that the analysis is being conducted, however, they could still learn something private from the fact that the analysis has *not* been published.

<sup>&</sup>lt;sup>2</sup>This is actually underscored by the GDPR legislation, which focuses on consent and data use methods.

Data is collected for many reasons. Any typical service you might want to use will require a minimal amount of data. For example, a dating service will at a minimum require your age and location, as shown in the example below.

Birthday	Name	Height	Weight	Age	Postcode	Profession
06/07	Li Pu	190	80	60-70	1001	Politician
06/14	Sara Lee	185	110	70+	1001	Rentier
01/01	A. B. Student	170	70	40-60	6732	Time Traveller

EXAMPLE 1 (Typical relational database in a dating website.). If somebody is a user in a dating website, you expect them to give some minimal personal information to be stored in the site's database. This might include their birthday, location and profession, for example.

When you submit your data to a service, you expect it to be used responsibly. For the dating service, you expect it to use the information to find good matches, based on your preferences and location. Whenever somebody uses the service, they obtain some information about you, at least indirectly. For example, if they make a query for similar singles in their neighbourhood, and they see your profile, then they have learned that you live nearby.

If we wish to publish a database, frequently we need to protect the identities of the people involved. A simple idea is to erase directly identifying information. However, this does not really work most of the time, especially since attackers can have side-information that can reveal the identities of individuals in the original data, when combined with the published dataset.

The simple act of hiding or using random identifiers is called anonymisation. However this is generally insufficient as other identifying information may be used to re-identify individuals in the data. In particular, even if somebody is unable to infer the information of any individual from the published, anonymised, dataset, they may be able to do so via some side-information. All that is needed is another dataset with some columns in common with the dataset they want to target.

#### Record linkage

As an example, in the 1990s a the governor of Massachussets, decided to publish anonymised information about the health records of individual state employs. They were careful to hide all obviously identifying information such as their name, and thus claimed that there are no potential privacy violations. However, they left in some information that they thought would be useful for researchers: the postcode, the birthdate, and the sex of each individual.

This allowed a PhD student, McSweeney, to perform the following linkage attack. She order a floppy disk containing database of voting records in the state. This contained names and addresses, as well as the postcode, birth date and sex of every voter. In that way, she was able to cross-reference this data, and so obtain the identities of individuals in that database. The first record she obtained was that of the governor himself, Bill Weld. Later, she estimated that approximately 87% of Americans are uniquely identifiable through those three attributes.

Clearly, anonymisation is not enough. So, is there a way to formally guarantee privacy? In the next section we will go over the solution of McSweeny, which provides us with an algorithmic definition of anonymity. While this provides some degree of protection against certain linkage attacks, it is sadly insufficient to protect privacy in general. Later, we will introduce the notion of differential privacy, which protects individuals against statistical disclosure in an information-theoretic manner.

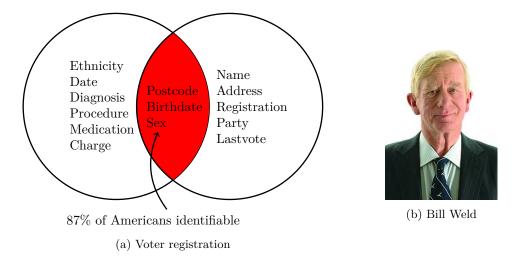


Figure 2.1: Linkage attack on a anonymised patient information database. A misguided attempt to release medical information as public data (left circle in figure (a)) by then-governor Bill Weld resulted in a simple linkage attack. This was done by accessing a party registration database (right circle in figure (b)) which contained three common fields with the medical database.

### 2.3 Algorithmic privacy.

The above discussion should help emphasise that a dataset cannot be characterised as private or non-private. We should actually start from the observation that any data contributed by an individual should be considered private, because it could be combined with other datasets to obtain sensitive personal information.

For that reason, we wish to ensure that, whenever individuals contribute data for processing, the result of the processing cannot be used to reveal whatever information they had contributed. Thus, the idea of privacy only applies on *algorithms* that are used on personal data. Generally speaking, we wish for algorithms to both be useful and have good privacy properties. Unfortunately, as we shall see in the sequel, there is a distinct trade-off between privacy and utility. However, let us specify more precisely what we mean by an algorithm.

#### What is an algorithm?

Any functional process is an algorithm. The information given to an algorithm for processing is called an *input* to the algorithm. The result of the process is called the algorithm's *output*. In this book, identify inputs with observations or features  $\mathcal{X}$ , and outputs with actions  $\mathcal{A}$  by the algorithm: i.e. we view the algorithm as taking actions that have an observable effect on its external environment. Algorithms are usually deterministic: given the same input, they produce the same output. However, it is also useful to consider *stochastic* algorithms, whose output randomly changes, even when the input is the same.

**Definition 2.3.1** (Stochastic algorithm  $\pi$ ). A stochastic algorithm  $\pi$  with input domain  $\mathcal{X}$  and output domain  $\mathcal{A}$  is a mapping  $\pi: \mathcal{X} \to \Delta(\mathcal{A})$  from observations  $\mathcal{X}$  to *distributions* over outputs  $\Delta(\mathcal{A})$ . When  $\mathcal{A}$  is finite, we will write  $\pi(a \mid x)$  to denote the conditional probability that the algorithm outputs a given input x.

So, one way of seeing a stochastic algorithm is as a *collection* of probability distributions  $\pi(a \mid x)$  over  $\mathcal{A}$ , one for each value of x. This is exactly equivalent to the definition of a

deterministic algorithm  $f: \mathcal{X} \to \mathcal{A}$ , which would be a simple function. This would define one specific value in  $\mathcal{A}$  (instead of a distribution) for each value of x.

The above definition is fine when the output domain  $\mathcal{A}$  is finite. Then outputs  $a \in \mathcal{A}$  cannot have infinitesimally small probabilities. However, in many problems  $\mathcal{A}$  is a Euclidean subset, e.g. the interval [0,1]. In that case, we have two choices: either treat  $\pi$  as a probability density, or as a probability measure. The latter is more general, and applies no matter what  $\mathcal{A}$  is. Formally, we can write this as follows.

#### The general case.

It is best to think of  $\pi(\cdot|x)$  as a different probability measure over  $\mathcal{A}$  for every possible value of x. Then we write

$$\pi(A \mid x) \triangleq \mathbb{P}_{\pi}(a \in A \mid x), \qquad A \subset \mathcal{A},$$

for the conditional probability that algorithm's output is in some set  $A \subset \mathcal{A}$  given input x. This allows us to treat the case when  $\mathcal{A}$  is continuous or discrete with the same notation.

How can we actually obtain a stochastic algorithm? The usual method is to define an algorithm as usual, but with a source of randomness as part of its inputs, fancifully called *random coins*. By using those random coins as part of the calculations, we obtain a random result, even if every step of the algorithm is deerministic.

#### Algorithmic randomness.

We can construct a random algorithm through access to a random coin  $\omega$  taking values in  $\Omega$ . We can then define the output through a deterministic function  $a_{\pi}(\omega, x)$ . Since  $\omega$  is random, the output of the function is also random. If P is the probability distribution of  $\omega$ , then:

$$\pi(A \mid x) = P(\{\omega \in \Omega \mid a_{\pi}(\omega, x) \in A\}),$$

i.e. the probability that the algorithm's output is in A is equal to the measure of the values of  $\omega$  for which the  $a_{\pi}(\omega, x) \in A$ .

#### What is private?

The overall framework we will use in this book is that everything pertaining to an individual is by definition private. While in reality there might be some attributes that individuals might want to remain secret and some which are not considered important, it is generally impossible to predict potential harm from the disclosure of any individual information. Thus, this chapter takes an expansive view of privacy: any information about an individual, which includes even the very fact that they might be part of a database is considered private.

We assume that there is a *private input*, which is processed by an *algorithm*, which then generates a *public output*. So, our main focus is how the algorithm generates the output, because the algorithm defines how the output and the input are linked. Our goal is to obtain algorithms that can be generally applied and that can have good privacy-preserving properties.

## 2.4 Simple anonymisation and k-anonymity

One of the basic ideas for protecting individual information is to simply remove identifying information. However, this is easier said than done, as potentially any information can be used identify an individual in a dataset. Indeed, simply removing the names of people from a database

is a very weak method, as there almost always exist enough information in the database to reidentify individuals. However, a slightly stronger, though not perfect, method for preventing re-identification is given by the framework of k-anonymity.

#### k-anonymity







(b) Sweeney

The concept of k-anonymity was introduced by Samarati and Sweeney <sup>15</sup> and provides some guarantees against inferring personal information from a single database. This requires the analyst to first determine the variables of interest (which should not be modified), and then determine which variables are *quasi-identifiers*, i.e. they could be potentially used to identify somebody in the database. *It's the analyst's job to define quasi-identifiers*. However, in general all variables should be considered quasi-identifiers.

**Definition 2.4.1** (k-anonymity). A database provides k-anonymity if for every person in the database is indistinguishable from k-1 persons with respect to quasi-identifiers.

This hope is that, if the database satisfies k-anonymity it can be safely released, without revealing any private information directly. As you can see, the definition of k-anonymity relates to the algorithm output, and not the algorithm itself. Because of this, it is not possible to give formal guarantees that hold generally for a k-anonymous database, as the result of the process does not tell us anything about how much information it conveys about the original input.

But first, let us walk through an extended example to explain the concept.

#### k-anonymity example

In particular, let us say that the analyst simply wants to calculate some statistics about how different professions correlate with age, weight, height and where people live. Some areas of the country might produce more politicians, for example. And taller people may be more successful in politics. The initial data collected might look like the table below. It was obvious to the analyst, that even if he did remove all the names, somebody knowing where A. B. Student lived and saw the table would have little trouble recognising them.

Birthday	Name	Height	Weight	Age	Postcode	Profession
06/07	Li Pu	190	80	65	1001	Politician
06/14	Sara Lee	185	110	67	1001	Rentier
06/12	Nikos Karavas	180	82	72+	1243	Politician
01/01	A. B. Student	170	70	52	6732	Time Traveller
05/08	Li Yang	175	72	35	6910	Politician

Table 2.1: 1-anonymity.

After thinking about it for a bit, the analyst decides to remove the birthday and name, and broadly categorise people according to their height in increments of 10cm, the weight in increments of 20cm, and keep just the first digit of the postcode. Now that looked much more reasonable. Still, somebody that knows that Li Yang and A. B. Student are in the table, as well as their postcodes, and they also know that Li Yang is a politician, can infer that A. B. Student is a Time Traveller.

Height	Weight	Age	Postcode	Profession
180-190	80+	60+	1*	Politician
180-190	80+	60+	1*	Rentier
180-190	80+	60+	1*	Politician
170-180	60-80	20-60	6*	Time Traveller
170-180	60-80	20-60	6*	Politician

Table 2.2: 2-anonymity: the database can be partitioned in sets of at least 2 records

However, with enough information, somebody may still be able to infer something about the individuals. In the example above, it remains true that if somebody knows that both A. B. Student and Li Yang are in the database, as well as their postcodes, as well as that Li Yang is a politician, they can infer A. B. Student's profession. Fortunately, there is a way to protect individual information from adversaries with arbitrary side-information. This is given by differential privacy.

## 2.5 Differential privacy

This section introduces one of the main tools for giving formal guarantees about the privacy of any algorithm ran on a dataset, differential privacy. This will provide individual-level privacy, in the sense that an algorithm that is differentially private guarantees that no adversary can significantly increase their knowledge about any particular individual by observing the algorithm's output. This is independent of the adversary's existing knowledge, or computational power.

While k-anonymity can protect against specific re-identification attacks when used with care, it says little about what to do when the adversary has a lot of knowledge. For example, if the adversary knows the data of everybody that has participated in the database, it is trivial for them to infer what our own data is. For some particularly sensitive datasets, we may want for the adversary to be unable to tell whether or not your data was part of the base. Differential privacy offers protection against adversaries with unlimited side-information or computational power. Intuitively, an algorithmic computation is differentially-private if an adversary cannot distinguish two "neighbouring" databases based on the result of the computation. Informally, two databases are neighbours when they are identical apart from the data of one person. A

differentially private algorithm, because of its randomness, makes it impossible for somebody to tell from the algorithm's output whether any specific individual's data was in the database.

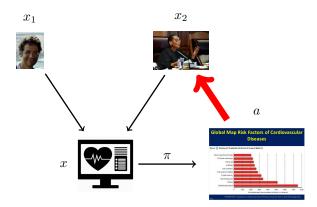


Figure 2.3: If two people contribute their data  $x = (x_1, x_2)$  to a medical database, and an algorithm  $\pi$  computes some public output a from x, then it should be hard infer anything about the data from the public output.

Consider the example given in Figure 2.3, where two people contribute their data to a medical database. The *i*-th individual contributes data  $x_i$ , and the complete dataset is  $x = (x_1, x_2)$ . The algorithm  $\pi$  defines a distribution over the set of possible outputs  $a \in \mathcal{A}$ , with  $\pi(a|x)$  being the probability that the algorithm outputs a if the data is x. If the algorithm was deterministic, then it might be possible for an adversary to invert the computation and obtain x from a. But even if that is not possible, maybe they can learn *something* about the data from the output. In the section below, we will formalise this notion.

#### Privacy desiderata

Consider a scenario where n persons give their data  $x_1, \ldots, x_n$  to an analyst. This analyst then performs some calculation on the data and publishes the result through a randomised algorithm  $\pi$ , where for any output a, and any dataset x,  $\pi(a|x)$  is the probability that the algorithm generates a. The following properties are desirable from a general standpoint.

**Anonymity.** Individual participation in the study remains a secret. From the release of the calculations results, nobody can significantly increase their probability of identifying an individual in the database.

**Secrecy.** The data of individuals is not revealed. The release does not significantly increase the probability of inferring individual's information  $x_i$ .

**Side-information.** Even if an adversary has arbitrary side-information, he cannot use that to amplify the amount of knowledge he would have obtained from the release.

**Utility.** The released result has, with high probability, only a small error relative to a calculation that does not attempt to safeguard privacy.

#### The prevalence of drugs in sport

Let's say you need to perform a statistical analysis of the drug-use habits of athletes. Obviously, even if you promise the athlete not to reveal their information, you still might not convince them. Yet, you'd like them to be truthful. The trick is to allow them to randomly change their answers, so that you can't be *sure* if they take drugs, no matter what they answer.

#### Algorithm for randomising responses about drug use

- 1. Flip a coin.
- 2. If it comes heads, respond truthfully.
- 3. Otherwise, flip another coin and respond yes if it comes heads and no otherwise.

#### **T** Calculating the true rate of responses.

Assume that the observed rate of positive responses in a sample is p, that everybody follows the protocol, and the coin is fair. Then, what is the true rate q of drug use in the population?

The problem with this approach, of course, is that we are effectively throwing away half of our data sources. In particular, if we repeated the experiment with a coin that came heads at a rate  $\epsilon$ , then our error bounds would scale as  $O(1/\sqrt{\epsilon n})$  for n data points.

This algorithm is very specific: it assumes binary responses, and it uses a fair coin, which introduces a lot of noise. Since the coin flips make the responses noisy, we may want to have some way of controlling it. In general, we want to consider an algorithm that takes data  $x_1, \ldots, x_n$  from n users transforms it randomly to  $a_1, \ldots, a_n$  using the following mapping.

#### The binary randomised response mechanism

**Definition 2.5.1** (Randomised response). The *i*-th user, whose data is  $x_i \in \{0,1\}$ , responds with  $a_i \in \{0,1\}$  with probability

$$\pi(a_i = j \mid x_i = k) = p, \qquad \pi(a_i = k \mid x_i = k) = 1 - p,$$

where  $j \neq k$ .

Given the complete data x, the algorithm's output is  $a = (a_1, \ldots, a_n)$ . Since the algorithm independently calculates a new value for each data entry, the output probability is

$$\pi(a \mid x) = \prod_{i} \pi(a_i \mid x_i)$$

This mechanism satisfies the formal notion of  $\epsilon$ -differential privacy, which will be given in Definition 2.5.2. In a more general setting, we may have multiple possible responses. While the algorithm can be trivially generalised to n-ary outputs, the special case of when the outputs are integers is deferred until later.

#### $\mathbf{\hat{c}}_{s}^{s}$ The original randomised response mechanism $^{16}$

As first proposed by Warner, the mechanism distributes spinners to people. The spinner has a probability p of landing on A and 1-p of landing on B. This can be easily arranged by having the corresponding areas to have the appropriate proportions. The interesting thing about this mechanism is that the responders must merely say whether or not the spinner landed on the group they identify with. They do not have to reveal a group. This makes the mechanism feel like you are revealing less, even if it is just a special case of a general randomised response mechanism.

#### What can we learn from the output?

In the Bayesian setting, we can think of the adversary as having some prior belief  $\beta(x_i)$ , expressed as a probability distribution over the secret value of each individual.

After the adversary observes the output, they can form a posterior belief  $\beta(x_i \mid a_i)$ , representing the information they collected. The following example quantifies the amount of information gained by the adversary.

EXAMPLE 2. For simplicity, consider only one individual, and let the adversary have a prior  $\beta(x=0) = 1 - \beta(x=1)$  over the values of the true response of an individual. We use the randomised response mechanism with parameter p and the adversary observes the randomised data a=1 for that individual, then what is  $\mathbb{P}^{\pi}_{\beta}(x=1 \mid a=1)$ , assuming the adversary knows the mechanism?

*Proof.* Bayes's theorem states that<sup>3</sup>

$$\mathbb{P}^{\pi}_{\beta}(x \mid a) = \pi(a \mid x)\beta(x)/\mathbb{P}^{\pi}_{\beta}(a),$$

where

$$\mathbb{P}^{\pi}_{\beta}(a) = \sum_{x'} \pi(a \mid x') \beta(x').$$

Let  $q = \beta(x = 1)$ . Then we have:

$$\mathbb{P}^{\pi}_{\beta}(x=1 \mid a=1) = pq/[pq + (1-p)(1-q)].$$

It is particularly interesting to consider the case where q=1/2. Then

$$\mathbb{P}^{\pi}_{\beta}(x=1\mid a=1)=p,$$

so we only have limited evidence for whether x = 1. Now consider the case where we have some arbitrary prior q, and p = 1/2. This means that the output of the algorithm is completely random. Consequently:

$$\beta(x = 1 \mid a = 1) = q = \beta(x = 1).$$

So, in this scenario we learn nothing from the algorithm's output.

#### Differential privacy.

Now let us take a look at a way to characterise the the inherent privacy properties of algorithms. This is called differential privacy, and it can be seen as a bound on the information an adversary with arbitrary power or side-information could extract from the result of a computation  $\pi$  on the data. For reasons that will be made clear later, this computation has to be stochastic.

<sup>&</sup>lt;sup>3</sup>If there are too many symbols for you, you can write Bayes theorem simply with P(x|a) = P(a|x)P(x)/P(a).

#### $\stackrel{\text{\tiny iii}}{\rhd}$ $\epsilon$ -Differential Privacy

**Definition 2.5.2.** A stochastic algorithm  $\pi: \mathcal{X} \to \Delta(\mathcal{A})$ , where  $\mathcal{X}$  is endowed with a neighbourhood relation N, is said to be  $\epsilon$ -differentially private if

$$\left| \ln \frac{\pi(A \mid x)}{\pi(A \mid x')} \right| \le \epsilon, \qquad \forall x N x', \quad A \subset \mathcal{A}.$$
 (2.5.1)

This form is related to standard notions of statistical distance such as the KL divergence  $\sum_{a} \ln \frac{\pi(a|x)}{\pi(a|x')} \pi(a \mid x)$ . However, the above inequality can be equivalently rewritten as

$$\pi(A \mid x) \le e^{\epsilon} \pi(A \mid x').$$

Frequently (and particularly when A is finite) we can also use  $\pi(a \mid x)$  without any technical difficulties.

Typically, algorithms are applied to datasets  $x = (x_1, ..., x_n)$  composed of the data of n individuals. Thus, all privacy guarantees relate to the data contributed by these individuals.

#### Neighbourhoods

Differential privacy guarantees that it is hard to distinguish neighbouring datasets. Hence, the definition of neighbourhood we use reflects what we want to protect. It makes sense to define neighbourhoods in terms of changes in one individual's data, because then an adversary cannot learn about the values of any particular individual.

In this book we will use two definitions of neighbourhoods. The first is constructed so that the adversary cannot distinguish whether or not any particular individual's information is in the dataset. If not, then they cannot infer the individual's presence from the output of the algorithm.

**Definition 2.5.3** (Participation neighbourhood). If two datasets x, x' are neighbours, then we write xNx', and it holds that

$$x = (x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n), \qquad x' = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n),$$

for some i, i.e. if one dataset is missing an element.

Under this neighbourhood relation, two datasets are neighbours if one contains the data of one individual, and the other does not, but they are otherwise the same. Then, if the algorithm is differentially private with respect to this neighbourhood, it is hard to tell if any single person's data has been used in the calculation. This is an important concept if the participation in the dataset is itself sensitive. For example, if somebody is enlisted in study for experimental treatment of a rare disease, then the mere fact that they are part of the study is strong evidence that they have the disease.

The second definition is slightly weaker. Here, two datasets are neighbours if they are identical, apart from the data of one individual, which is changed. It determines whether or not we can distinguish between different *values* of individual data. If not, then they cannot infer whether the data submitted by the individual has a particular value.

**Definition 2.5.4** (Edit neighbourhood). We say that two datasets x, x' are neighbours, and write  $xN_ex'$  if

$$x = (x_1, \dots, x_i, \dots, x_n), \qquad x' = (x_1, \dots, x'_i, \dots, x_n), \qquad x_i \neq x'_i.$$

i.e. if one dataset has an altered element.

If x, x' are 1-neighbours under the second definition, then they are 2-neighbours under the first definition. To see this, create a new dataset  $\hat{x}$  from x, without the data of person i. Then  $xN\hat{x}$  and  $\hat{x}Nx'$ , since we can change the data of one person by removing the original data  $x_i$  and then re-inserting the altered data  $x_i'$ . This is illustrated in the example below.

#### Neihgbourhood example

In this example, we have three datasets, x, x' and hatx. In the second dataset,  $\hat{x}$ , the highlighted row in x is missing. In the third dataset, x', another row with the same name is added, but the height and weight are changed.

Birthday	Name	Height	Weight
06/07	Li Pu	190	80
06/14	Sara Lee	185	110
06/12	$John \ Smith$	170	82
01/01	A. B. Student	170	70
05/08	Li Yang	175	72

Table 2.3: Data x

In fact,  $\hat{x}$  is obtained through a deletion from x. Of course, the operation can be reversed: x can be obtained from an addition to  $\hat{x}$ .

Birthday	Name	Height	Weight
06/07	Li Pu	190	80
06/14	Sara Lee	185	110
01/01	A. B. Student	170	70
05/08	Li Yang	175	72

Table 2.4:  $\hat{x}$ , 1-Neighbour of x.

We can now instead add another row to  $\hat{x}$ , which will be similar to the row we had removed. This will have the effect of altering the original data in x.

Birthday	Name	Height	Weight
06/07	Li Pu	190	80
06/14	Sara Lee	185	110
06/12	John Smith	180	80
01/01	A. B. Student	170	70
05/08	Li Yang	175	72

Table 2.5: x', 2-Neighbour of x.

Since x, x' only differ in the contents of a single row, they are edit-neighbours.

As we hinted earlier, the randomised response mechanism satisfies  $\epsilon$ -DP. The  $\epsilon$  parameter is dependent on p, with higher values giving a smaller  $\epsilon$ , and thus better privacy protection.

Remark 2.5.1. The randomised response mechanism with  $p \leq 1/2$  is  $(\ln \frac{1-p}{p})$ -DP with respect to the edit neighbourhood  $N_e$ .

*Proof.* Consider  $x = (x_1, \ldots, x_j, \ldots, x_n), x' = (x_1, \ldots, x_j', \ldots, x_n)$ . Then

$$\pi(a \mid x) = \prod_{i} \pi(a_i \mid x_i)$$

$$= \pi(a_j \mid x_j) \prod_{i \neq j} \pi(a_i \mid x_i)$$

$$\leq \frac{1 - p}{p} \pi(a_j \mid x'_j) \prod_{i \neq j} \pi(a_i \mid x_i)$$

$$= \frac{1 - p}{p} \pi(a \mid x')$$

 $\pi(a_j = k \mid x_j = k) = 1 - p$  so the ratio is  $\max\{(1 - p)/p, p/(1 - p)\} \le (1 - p)/p$  for  $p \le 1/2$ .

#### Moving to a new neighbourhood

Is the randomised-response mechanism for it to be differentially private with respect to the insertion-neighbourhood definition? If not, is it possible to modify it in order to satisfy that privacy definition? *Hint: The mechanism must be able to hide the participation of a single individual in the database.* 

Finally, it may be convenient to the look at neighbourhoods in terms of a distance between datasets. Let  $\mathbb{N}^{|\mathcal{X}|}$  be the set of all possible dataset histograms, i.e. counts of different possible rows in each dataset. Then the distance between two datasets is simply the total difference in counts:

**Definition 2.5.5** (Hamming distance between datasets). Consider two datasets, x, x' with n and n' elements respectively. Then, we define their Hamming distance as:

$$\|\boldsymbol{x} - \boldsymbol{x}'\|_{1} = \sum_{j} |\sum_{i=1}^{n} \mathbb{I}\{x_{i} = j\} - \sum_{i=1}^{n'} \mathbb{I}\{x_{i} = j\}| = \sum_{j \in \mathcal{X}} |n_{j}(\boldsymbol{x}) - n_{j}(\boldsymbol{x}')|,$$
 (2.5.2)

where  $n_j(\mathbf{x})$  is the number of elements in  $\mathbf{x}$  equal to j.

Let us see how the Hamming distance relates to the two neighbourhoods we defined. In particular,  $\|\boldsymbol{x} - \boldsymbol{x}'\|_1 = 1$  if and only if  $\boldsymbol{x} N \boldsymbol{x}'$ . To see this, notice that if  $\boldsymbol{x}'$  has one row less than  $\boldsymbol{x}$ , but is otherwise identical, then there is exactly one value j for which  $n_j(\boldsymbol{x}) = n_j(\boldsymbol{x}') + 1$ , with  $n_j(\boldsymbol{x}) = n_j(\boldsymbol{x}')$  otherwise. Consequently,  $\|\boldsymbol{x} - \boldsymbol{x}'\|_1 = 1$ . The reverse direction follows by contradiction: if the Hamming distance is one, then the two databases must differ in at least one element. If they differ in more, then their distance has to be larger than one.

On the other hand, for the second neighbourhood definition, things are slightly different. There, we assume that  $x_i \neq x_i'$  for one i. Without loss of generality, let us say that  $x_i = j$  and  $x_i' = k$ , with  $j \neq k$ . Then  $n_j(\mathbf{x}) = n_j(\mathbf{x}') + 1$  and  $n_k(\mathbf{x}') = n_k(\mathbf{x}) + 1$ . Consequently  $\|\mathbf{x} - \mathbf{x}'\|_1 = 2$ .

### 2.6 The local and central privacy models

So far, we have only seen the concept of differential privacy applied as a method to generate a privatised version of a dataset. This can then be used to perform arbitrary computations.

Typically, however, we want to perform some specific calculations on the data. Is it possible to compute things privately using the original data, but so that the result of the computation does not leak any information? Clearly this should depend somehow on the computation you wish to perform.

For concreteness, let us assume that you have defined a function  $f: \mathcal{X} \to \mathcal{Y}$ , which you wish to compute on arbitrary data  $x \in \mathcal{X}$ . If you wish to preserve privacy, however, you need the computation to satisfy some formal guarantee like differential privacy. The simplest way to do this is by simply generating a dataset a with a DP algorithm and then processing the data with the function we have already specified. This corresponds to the local privacy model.

#### The local privacy model.

Given a function  $f: \mathcal{X} \to \mathcal{Y}$ , and a set of individual data  $\mathbf{x} = (x_i)_{i=1}^n$ , use a differentially private algorithm  $\pi(\mathbf{a}|\mathbf{x})$  to generate  $\mathbf{a} \in \mathcal{X}$ . Then output  $f(\mathbf{a})$ .

The advantage of the local model is that the calculation f does not need to be modified. The individuals do not have to trust anybody, as they can modify their data locally. However, they must take care that the DP calculation is performed correctly.

Typically, in the local privacy model, the i-th individual's data  $x_i$  is used to generate a private response  $a_i$ . This means that individual data is only available to a local private mechanism. This model allows us to publish a complete dataset of private responses. In the central privacy model, the data x is collected and the result a is published by a trusted curator.

#### The central privacy model.

A trusted curator obtains the data  $\boldsymbol{x}$  of all individuals and selects a DP policy  $\pi(\boldsymbol{a} \mid \boldsymbol{x})$  to generate the output a, so that  $\boldsymbol{a} \approx f(\boldsymbol{x})$ .

The main advantage of the centralised model is that approximating f with a differentially-private version can be much more accurate than simply using the original function with noisy data. However, obtaining such an approximation is not always easy.

The dependency diagrams in Figures 2.4 shows how the output depends on the individual data, the non-private and private algorithm.

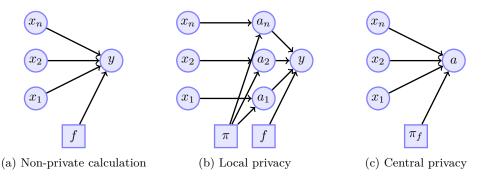


Figure 2.4: Non-privacy, local privacy, and centralised privacy models. In the first case, the value of y depends directly on the function f to be calculated, as well as the data x. In the local privacy model, noise is added to the individual data with a differentially private algorithm  $\pi$ , and the function is calculated on the output. In the central model, a stochastic version  $\pi_f$  of the original function f is calculated on the data.

#### Trust models.

Google uses a local privacy model to collect data from Android phone users. How does this compare to Google gathering data in the clear and then performing private computations? Who do the users have to trust? What are the privacy risks in either case?

#### 2.6.1 Properties of differential privacy.

Remark 2.6.1. Any differentially private algorithm must be stochastic.

To prove that this is necessary, consider the example of counting how many people take drugs in a competition. If the adversary only doesn't know whether you in particular take drugs, but knows whether everybody else takes drugs, it's trivial to discover your own drug habits by looking at the total. This is because in this case,  $f(x) = \sum_i x_i$  and the adversary knows  $x_i$  for all  $i \neq j$ . Then, by observing f(x), he can recover  $x_j = f(x) - \sum_{i \neq j} x_i$ . Consequently, it is not possible to protect against adversaries with arbitrary side information without stochasticity.

Randomness is also necessary in cryptography. In that setting, Alice wants to communicate a message x to Bob. In a secret key cryptographic scheme, information theoretic security is achieved by making sure that the encrypted message a comes from a uniform distribution  $\pi(a|x)$ . This is achieved by uniformly selecting a key and selecting an appropriate hash function

Depending on the DP scheme, each query answered may leak privacy. In particular, if we always respond with an  $\epsilon$ -DP mechanism, after T queries our privacy guarantee is  $T\epsilon$ . There exist mechanisms that do not respond to each query independently, which can reduce the total privacy loss, but those are outside the scope of this chapter.

**Definition 2.6.1** (*T*-fold composition). In this privacy model, we compose a mechanism out of *T* differential private mechanisms  $\pi_1, \ldots, \pi_T$ . This composition is fixed ahead of time.

**Theorem 2.6.1.** For any  $\epsilon > 0$ , the class of  $\epsilon$ -differentially private mechanism satisfy  $T\epsilon$ -differential privacy under T-fold composition. More generally, if each mechanism is  $\epsilon_i$ -DP, the composed mechanism is  $\sum_{i=1}^{T} \epsilon_i$ -DP.

**Theorem 2.6.2** (Post-processing). Let mechanism  $\pi(a \mid x)$  be  $\epsilon$ -DP. Applying any transformation  $f: A \to Y$  to the output of the mechanism to obtain y = f(a), results in another  $\epsilon$ -DP mechanism.

The composition theorem is a very useful tool, as it allows us to create new mechanisms that are composed of simpler parts. As a first example, we can look at how we can use it to create a randomised response algorithm when the respondents provide multiple attributes.

#### Randomised response for multiple attributes.

Up to now we have been discussing the case where each individual only has one attribute. However, in general each individual t contributes multiple data  $x_{t.i}$ , which can be considered as a row  $x_t$  in a database. Then the mechanism can release each  $a_{t.i}$  independently.

For n users and k attributes, if the release of each attribute i is  $\epsilon$ -DP then the data release is  $k\epsilon$ -DP. Thus to get  $\epsilon$ -DP overall, we need  $\epsilon/k$ -DP per attribute.

The result follows immediately from the composition theorem. We can see each attribute release as the result of an individual query. More generally, if each attribute i is released with an  $\epsilon_i$ -DP mechanism, the overall mechanism is  $\sum_i \epsilon_i$ -DP.

#### Differential privacy from a Bayesian viewpoint.

Bayesian inference offers a simple way to explain the meaning of differential privaacy. Intuitively, using the Bayesian formalism, we can show that, no matter what prior knowledge the adversary has, they cannot infer a lot from the private release. We are specifically interested in adversary knowledge about the dataset, before and after the mechanism's release. In particular, assume that the adversary knows that the data is either  $\boldsymbol{x}$  or  $\boldsymbol{x}'$ . For concreteness, assume the data is either

$$\boldsymbol{x} = (x_1, \dots, x_i = 0, \dots, x_n)$$

where  $x_i$  is the data of person i, or the alternative dataset:

$$\mathbf{x}' = (x_1, \dots, x_i' = 1, \dots, x_n).$$

In other words, the adversary knows the data of all people apart from one, the j-th person. They only need to dinstiguish one possible value from another. Without loss of generality, we can model the adversary as having some arbitrary prior belief

$$\beta(\boldsymbol{x}) = 1 - \beta(\boldsymbol{x}')$$

for the two cases. Assume the adversary knows the output a of a mechanism  $\pi$ . What can we say about the posterior distribution of the adversary  $\beta(\boldsymbol{x} \mid a, \pi)$  after having seen the output, if  $\pi$  is  $\epsilon$ -DP? How does it depend on  $\epsilon$ ?

Solution. We can write the adversary posterior as follows.

$$\mathbb{P}_{\beta}^{\pi}(\boldsymbol{x} \mid a) = \frac{\pi(a \mid \boldsymbol{x})\beta(\boldsymbol{x})}{\pi(a \mid \boldsymbol{x})\beta(\boldsymbol{x}) + \pi(a \mid \boldsymbol{x}')\beta(\boldsymbol{x}')} 
\geq \frac{\pi(a \mid \boldsymbol{x})\beta(\boldsymbol{x})}{\pi(a \mid \boldsymbol{x})\beta(\boldsymbol{x}) + \pi(a \mid \boldsymbol{x})e^{\epsilon}\beta(\boldsymbol{x}')}$$
(from DP definition)
$$= \frac{\beta(\boldsymbol{x})}{\beta(\boldsymbol{x}) + e^{\epsilon}\beta(\boldsymbol{x}')}.$$
(2.6.2)

Note that  $\pi(a \mid \boldsymbol{x}) \leq e^{\epsilon} \pi(a \mid \boldsymbol{x}')$  and conversely  $\pi(a \mid \boldsymbol{x}') \leq e^{\epsilon} \pi(a \mid \boldsymbol{x})$ . We can also then bound the quantity from above:

$$\mathbb{P}^\pi_\beta(\boldsymbol{x}\mid a) \leq \frac{\beta(\boldsymbol{x})}{\beta(\boldsymbol{x}) + e^{-\epsilon}\beta(\boldsymbol{x}')}.$$

Consequently,  $\lim_{\epsilon \to 0} \mathbb{P}^{\pi}_{\beta}(\boldsymbol{x} \mid a) = \beta(\boldsymbol{x})$ , hence the information gained by the adversary is bounded by the prior and the information loss  $\epsilon$ .

### 2.7 The Laplace mechanism

Many times we already have a function  $f: \mathcal{X} \to \mathbb{R}$  we want to calculate on data, and we would like to make the function preserve privacy. This clearly falls within the *central privacy* model: The analyst has access to the data and function, and wishes for the algorithm generating the final output to be private. One solution, that can satisfy differential privacy, is to add noise to the function's output. We start by first calculating the value of the function for the data we

have, f(x), and then we add some random noise  $\omega$ , hence our calculation is random:

$$a = f(x) + \omega$$
.

The amount and type of noise added, together with the smoothness of the function f, determine the amount of privacy we have. One of the simplest noise-adding mechanisms is to add zero-mean Laplace noise. Then we write  $\omega \sim \text{Laplace}(\lambda)$ .<sup>4</sup> The mechanism is defined below.

**Definition 2.7.1** (The Laplace mechanism). For any function  $f: \mathcal{X} \to \mathbb{R}$ , the output a of the mechanism is sampled from a Laplace distribution with mean f(x) and scaling parameter  $\lambda$ , i.e.

$$\pi(a \mid x) = Laplace(f(x), \lambda). \tag{2.7.1}$$

The probability density function of the Laplace distribution with mean  $\mu$  and scaling  $\lambda$  is given by:

$$p(\omega \mid \mu, \lambda) = \frac{1}{2\lambda} \exp\left(-\frac{|\omega - \mu|}{\lambda}\right).$$

and has mean  $\mu$  and variance  $2\lambda^2$ .

#### The Laplace mechanism for averages

Here we have n individuals for which we wish to calculate the average salary.

- The *i*-th person receives salary  $x_i$
- We wish to calculate the average salary in a private manner.

We can do this in two ways. By using a DP mechanism on each individual salary and then calculating the average, or by first calculating the average and then applying a DP mechanism to the result. In particular, we can add try adding Laplace noise in both cases.

The local model. In this case,  $\pi(a \mid x)$  is obtained by independent Laplace noise  $\omega_i$  for each individual:

- Obtain  $y_i = x_i + \omega_i$ , where  $\omega_i \sim \text{Laplace}(\lambda)$ .
- Return  $a = n^{-1} \sum_{i=1}^{n} y_i$ .

The centralised model. In this case,  $\pi(a \mid x)$  is obtained by averaging first and adding noise later.

- Calcualte  $y = n^{-1} \sum_{i=1}^{n} x_i$ .
- Return  $a = y + \omega$ , where  $\omega \sim \text{Laplace}(\lambda')$ .

We must tune  $\lambda, \lambda'$  appropriately in to obtain the needed  $\epsilon$ -DP guarantee.

In the centralised privacy model, the non-private calculation directly outputs y = f(x). We can approximate this with a DP calculation with distribution  $\pi(a \mid x)$ . The Laplace mechanism does so by first calculating y = f(x) and then generating  $\pi(a \mid y)$  so that  $\mathbb{E}_{\pi}[a \mid x] = f(x)$ .

Let us now talk about how the Laplace mechanism that can be used both in the centralised and local model when  $\mathcal{X} \subset \mathbb{R}^n$  and  $\mathcal{Y} \subset \mathbb{R}$ .

<sup>&</sup>lt;sup>4</sup>When unspecified, the mean parameter is assumed to be zero.

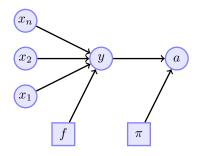


Figure 2.5: Laplace mechanism

#### DP properties of the Laplace mechanism

To use the Laplace mechanism in the centralised setting, we must relate it to a specific function f that we wish to compute. In particular, the mechanism works by adding noise to the output of the function f in a carefully calibrated manner so that we achieve exactly  $\epsilon$ -differential privacy.

**Definition 2.7.2** (Sensitivity). The sensitivity of a function f is

$$\mathbb{L}\left(f\right) \triangleq \sup_{xNx'} \left| f(x) - f(x') \right|$$

If we define a metric d, so that d(x, x') = 1 for xNx', then:

$$|f(x) - f(x')| \le \mathbb{L}(f) d(x, x'),$$

i.e. f is  $\mathbb{L}(f)$ -Lipschitz with respect to d.

EXAMPLE 3. If  $f: \mathcal{X} \to [0, B]$ , e.g.  $\mathcal{X} = \mathbb{R}$  and  $f(x) = \min\{B, \max\{0, x\}\}$ , then  $\mathbb{L}(f) = B$ .

EXAMPLE 4. If  $f:[0,B]^n \to [0,B]$  is  $f = \frac{1}{n} \sum_{t=1}^n x_t$ , then  $\mathbb{L}(f) = B/n$ .

*Proof.* Consider two neighbouring datasets x, x' differing in example j. Then

$$f(x) - f(x') = \frac{1}{n} [f(x_j) - f(x'_j)] \le \frac{1}{n} [B - 0]$$

**Theorem 2.7.1.** The Laplace mechanism on a function f with sensitivity  $\mathbb{L}(f)$ , ran with Laplace( $\lambda$ ) is  $\mathbb{L}(f)/\lambda$ -DP. Consequently, if the Laplace mechanism is ran with  $\lambda = \mathbb{L}(f)/\epsilon$ , then it is  $\epsilon$ -DP.

Proof.

$$\frac{\pi(a\mid x)}{\pi(a\mid x')} = \frac{e^{|a-f(x')|/\lambda}}{e^{|a-f(x)|/\lambda}} \le \frac{e^{|a-f(x)|/\lambda + \mathbb{L}(f)/\lambda}}{e^{|a-f(x)|/\lambda}} = e^{\mathbb{L}(f)/\lambda}$$

The first step follows from the definition of the Laplace mechanism. The inequality follows from the fact that for xNx', we have  $|f(x) - f(x')| \le \mathbb{L}(f)$ . In particular,

(a) If  $f(x') - a \ge 0$  then

$$|f(x') - a| = f(x') - a \le \mathbb{L}(f) + f(x) - a \le L + |a - f(x)|,$$

as  $f(x') \leq f(x) + \mathbb{L}(f)$  from the Lipschitz property.

(b) If f(x') - a < 0 then

$$|f(x') - a| = a - f(x') \le \mathbb{L}(f) - f(x) + a \le L + |a - f(x)|,$$

as 
$$-f(x') \leq \mathbb{L}(f) - f(x)$$
.

Replacing into the exponential gives us the required inequality. The final result is obtained with elementary algebra.  $\Box$ 

# DP properties of the average in the local and central model.

What is the effect of applying the Laplace mechanism in the local versus centralised model? Let us continue the average example. Here let us assume  $x_i \in [0, B]$  for all i.

# The Laplace mechanism in the local privacy model

The sensitivity of the individual data is B, so to obtain  $\epsilon$ -DP we need to use  $\lambda = B/\epsilon$ . The variance of each component is  $2(B/\epsilon)^2$ , so the total variance is  $2B^2/\epsilon^2n$ .

# The Laplace mechanism in the centralised privacy model

The sensitivity of f is B/n, so we only need to use  $\lambda = \frac{B}{n\epsilon}$ . The variance of a is  $2(B/\epsilon n)^2$ .

Thus the two models have a significant difference in the variance of the estimates obtained, for the same amount of privacy. While the central mechanism has variance  $O(n^{-2})$ , the local one is  $O(n^{-1})$  and so our estimates will need much more data to be accurate under this mechanism. In particular, we need square the amount of data in the local model as we need in the central model. Nevertheless, the local model may be the only possible route if we have no specific use for the data.

# 2.8 Interactive data access.

In a lot of applications, we cannot pre-define the computation that we want to perform. After we have collected the data, we need to be perform an adaptive data analysis. This requires performing a sequence of computations on the data, where the result of one computation determines what computation we will do next.

We can think of this as performing queries to the database. You may be familiar with database access languages such as SQL, where you ask a question such as "what is the sum of the attribute age in table students?" and obtain the exact sum back. It is possible to answer such queries through a differentially private mechanism. However, the more queries you answer, the more you reveal about the original data. In addition, since an adversary may be cleverly adjusting the queries to more efficiently discover a secret, we need the concept of adaptive composition.

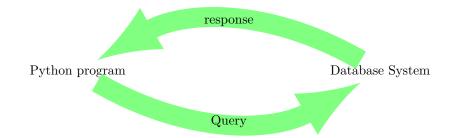


Figure 2.6: Database access model

# 2.8.1 Utility of queries

Rather than saying that we wish to calculate a private version of some specific function f, sometimes it is more useful to consider the problem from the perspective of the utility of different answers to queries. More precisely, imagine the interaction between a database system and a user:

## Interactive queries

- System has data x.
- At time t, user asks query  $q_t = q$ .
- System responds with  $a_t = a$ .
- There is a common utility function  $U: \mathcal{X}, \mathcal{A}, \mathcal{Q} \to \mathbb{R}$ .

We wish to give the most useful answer, by return a that maximises U, i.e.  $a = \arg\max_{a'} U(x, a', q)$ , but are constrained by the fact that we also want to preserve privacy.

The utility U(x, a, q) describes how appropriate each answer a given by the system for a query q is given the data x. It can be seen as how useful the response is  $^5$  It allows us to quantify exactly how much we would gain by replying correctly. The exponential mechanism, described below is a simple differentially private mechanism for responding to queries while trying to maximise utility for any possible utility function.

# The Exponential Mechanism.

Here we assume that we can answer queries q, whereby each possible answer a to the query has a different utility to the DM: U(q, a, x). The idea is that the best answer to the query should have the highest utility. The further away from the correct answer the response is, the lower the utility should be. This allows us to quantify how much we value correct answers to a query.

As an example, if the optimal response to a query is given by a real-valued function f(q, x), then one possible utility function is  $U(q, a, x) = -[f(q, x) - a]^2$ . This results in the correct response having a utility of zero, and responses whose value is further away will have negative values.

<sup>&</sup>lt;sup>5</sup>This is essentially the utility to the user that asks the query, but it could be the utility to the person that answers. In either case, the motivation does not matter the action should maximise it, but is constrained by privacy.

The exponential mechanism allows us to answer queries in a "soft" manner, by using this utility function. It assigns lower probability to answers with lower utility, hence the better responses have a higher likelihood of being selected.

**Definition 2.8.1** (The Exponential mechanism). For any utility function  $U : \mathcal{Q} \times \mathcal{A} \times \mathcal{X} \to \mathbb{R}$ , define the policy, which implicitly depends on U and  $\epsilon$ , as:

$$\pi(a \mid x, q) \triangleq \frac{e^{\epsilon U(q, a, x)/2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}},$$
(2.8.1)

where  $\mathbb{L}(U(q)) \triangleq \max_{a} \sup_{x \in Nx'} |U(q, a, x) - U(q, a, x')|$  denotes the sensitivity of a query.

Clearly, when  $\epsilon \to 0$ , this mechanism is uniformly random. When  $\epsilon \to \infty$  the action maximising U(q,a,x) is always chosen (see Exercise 4). Although the exponential mechanism can be used to describe some known DP mechanisms (see Exercise 5), its best use is in settings where there is a natural utility function.

**Theorem 2.8.1.** The exponential mechanism is  $\epsilon$ -differentially private.

*Proof.* We only need to look at the ratio between two different distributions of answers to queries. Then we have:

$$\begin{split} \frac{\pi(a \mid x, q)}{\pi(a \mid x', q)} &= \frac{e^{\epsilon U(q, a, x)/2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}} \times \frac{\sum_{a'} e^{\epsilon U(q, a', x')/\mathbb{L}(U(q))}}{e^{\epsilon U(q, a, x')/2\mathbb{L}(U(q))}} \\ &= e^{\epsilon U(q, a, x)/2\mathbb{L}(U(q)) - \epsilon U(q, a, x')/2\mathbb{L}(U(q))} \times \frac{\sum_{a'} e^{\epsilon U(q, a', x')/2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}} \\ &= e^{\epsilon [U(q, a, x) - U(q, a, x')]/2\mathbb{L}(U(q))} \times \frac{\sum_{a'} e^{\epsilon U(q, a', x')/2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}} \\ &\leq e^{\epsilon/2} \times \frac{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}} \\ &\leq e^{\epsilon/2} \times \frac{\sum_{a'} e^{\epsilon [(U(q, a', x) + \mathbb{L}(U(q)))/2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}} \leq e^{\epsilon} \end{split}$$

# $\stackrel{\text{\tiny 11}}{\hookrightarrow}$ Interactive differential privacy.

So far we only defined differential privacy in terms of a fixed mechanism. However, in general the mechanism must respond to sequential queries  $q_1, \ldots, q_t$ , with answers  $a_1, \ldots, a_t$ . The answer can depend arbitrarily on the sequence of responses and queries. So, privacy must hold no matter what the previous sequence of questions and answers. This gives rise to the following definition.

**Definition 2.8.2** (Interactive DP). An algorithm  $\pi$  satisfies interactive  $\{\epsilon_1, \ldots, \epsilon_t\}$  differential privacy with respect to N if

$$\ln \frac{\pi(a_t|x, q_1, \dots, q_t, a_1, \dots, a_{t-1})}{\pi(a_t|x', q_1, \dots, q_t, a_1, \dots, a_{t-1})} \le \epsilon_t, \quad \forall x N x'.$$
(2.8.2)

Since the queries are not defined in advance, this means that  $q_t$  may also depend on all the previous answers  $a_1, \ldots, a_{t-1}$ .

Answering queries independently in the interactive setting is sufficient for achieving differential privacy  $\epsilon_t$  at each step t. The total privacy loss is at most  $\sum_{i=1}^{t} \epsilon_i$ . However, it is possible to bound the privacy loss more tightly with advanced composition theorems.

# 2.8.2 Differential privacy as a hypothesis testing game.

Another way of viewing differential privacy is as a simple game between an adversary interacting with an honest data analyst. At time t, the adversary chooses two arbitrary neighbouring databases x, x'. The analyst randomly chooses between  $x_t$  to be one of the two databases. It then generates  $a_t \sim \pi(a \mid x_t)$ , and shows it to the adversary. The adversary must now decide if  $x_t = x$  or  $x_t = x'$ , i.e. which dataset was used by the analyst. More precisely, the adversary must choose between hypothesis 0 ( $x_t = x$ ) and hypothesis 1, ( $x_t = x'$ ).

Since the adversary can only base their decision on  $a_t$ , they must use a decision rule of the form  $f: \mathcal{A} \to \{0,1\}$  so as to decide between the two hypotheses. That means that they must partition  $\mathcal{A}$  in two subsets,  $A_0, A_1$  so that whenever  $a_t \in A_0$  they decide for  $\boldsymbol{x}$  and whenever  $a_t \in A_1$  they decide for  $\boldsymbol{x}'$ . Because we know the that the probability of the response given the dataset depends only on the mechanism  $\pi$ , we can write the probability of a true and a false positive as:

$$p_{\text{TP}} = \pi(A_0 \mid \boldsymbol{x}), \qquad p_{\text{FP}} = \pi(A_1 \mid \boldsymbol{x}).$$

This is simply because the adversary answers correctly whenever  $a_t \in A_0$ , and incorrectly otherwise. Similarly, the probability of a true and false negative as

$$p_{\text{TN}} = \pi(A_0 \mid \boldsymbol{x}), \qquad p_{\text{FN}} = \pi(A_0 \mid \boldsymbol{x}').$$

The probability with which the mechanism selects either dataset does not factor into these quantities, are these are simply the conditional probability of the incorrect decision region for each possible dataset.

**Theorem 2.8.2.** If the mechanism  $\pi$  is  $(\epsilon, \delta)$ -differentially private then the false positive rate (FPR) and false negative rate (FNR) of the adversary are linked as follows:

$$p_{TP} \leq p_{FP}e^{\epsilon} + \delta$$
,

i.e. the true positive rate cannot be much higher than the false positive rate. and conversely:

$$p_{TN} \leq p_{FN}e^{\epsilon} + \delta$$
,

i.e. the true negative rate caanot be much higher than the false negative rate.

*Proof.* The proof is by direct substitution and the definition of  $(\epsilon, \delta)$ -differential privacy:

$$p_{\text{TP}} = \pi(A_1|\mathbf{x}') < \pi(A_1|\mathbf{x})e^{\epsilon} + \delta = p_{\text{FP}}e^{\epsilon} + \delta.$$

Similarly

$$p_{\text{TN}} = \pi(A_0|\mathbf{x}) < \pi(A_0|\mathbf{x}')e^{\epsilon} + \delta = p_{\text{FN}}e^{\epsilon} + \delta.$$

# 2.9 Advanced topics.

In this section, we give a brief overview of some more advanced topics on privacy. The interested reader is urged to look into the given references.

# 2.9.1 Relaxations of differential privacy.

In practice, satisfying differential privacy with a small  $\epsilon$  results in mechanisms with low utility. It is possible to relax the definition of differential privacy to include an additive term, as shown below, which greatly improves the privacy-utility trade-off we can achieve.

**Definition 2.9.1.** Approximate differential privacy. A stochastic algorithm  $\pi: \mathcal{X} \to \Delta(\mathcal{A})$ , where  $\mathcal{X}$  is endowed with a neighbourhood relation N, is said to be  $(\epsilon, \delta)$  differentially private if

$$\pi(A \mid x) \le \pi(A \mid x')\epsilon + \delta, \quad \forall xNx', \quad A \subset \mathcal{A}.$$
 (2.9.1)

One simple interpretation of this definition is that the mechanism is  $(\epsilon, 0)$ -DP with probability  $1 - \delta$ . Conversely, this implies that with a very small probability  $\delta$ , the complete dataset might be revealed. For that reason, the parameter  $\delta$  must be as small as possible.

Approximate differential privacy is satisfied by the *Gaussian* mechanism. This is similar to the Laplace mechanism, but scales noise to the  $\ell_2$  sensitivity of the function.

# The Gaussian mechanism

Given a function  $f: \mathcal{X} \to \mathbb{R}^k$ , the Gaussian mechanism with parameter  $\sigma$  outputs  $a \sim \mathcal{N}(f(x), \sigma)$ .

Let  $\ell_2$  sensitivity of f be

$$\mathbb{L}(f)_2 = \max\{\|f(x) - f(y)\|_2 \mid xNy\}.$$

For  $\sigma \geq 2\ln(1.25/\delta)\frac{\Delta_2(f)}{\epsilon}$ , then the Gaussian mechanism is  $(\epsilon, \delta)$ -differentially private.

# Example Renyi differential privacy.

The careful reader will have noticed that differential privacy is essentially a bound on the distributional distance of mechanism outputs for similar inputs. In particular, the distance used for  $\epsilon$ -DP is  $D_{\infty}(P\|Q) \triangleq \sup_A |\ln P(A)/Q(A)|$ . We can use any other divergence between distributions, such as the KL divergence, as long as this becomes unbounded for distributions with unequal support.

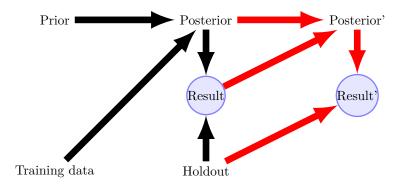
#### Reproducibility

Training and testing, overfitting on the test set.

# The unfortunate practice of adaptive analysis

In the ideal data analysis, we start from some prior hypothesis, then obtain some data, which we split into training and holdout. We then examine the training data and obtain a posterior that corresponds to our conclusions. We can then measure the quality of these conclusions in the independent holdout set.

2.10. DISCUSSION 43



However, this is not what happens in general. Analysts typically use the same holdout repeatedly, in order to improve the performance of their algorithms. This can be seen as indirectly using the holdout data to obtain a new posterior, and so it is possible that you can overfit on the holdout data, even if you never directly see it. It turns out we can solve this problem if we use differential privacy, so that the analyst only sees a differentially private version of queries.

# The reusable holdout 96

One idea to solve this problem is to only allow the analyst to see a private version of the result. In particular, the analyst will only see whether or not the holdout result is  $\tau$ -close to the training result.

# Algorithm parameters

- Performance measure f.
- Threshold  $\tau$ . How close do we want f to be on the training versus holdout set?
- Noise  $\sigma$ . How much noise should we add?
- Budget B. How much are we allowed to learn about the holdout set?

# Algorithm idea

Run algorithm  $\lambda$  on data  $D_T$  and get e.g. classifier parameters  $\theta$ . Run a DP version of the function  $f(\theta, D_H) = \mathbb{I}\{U(\theta, D_T) \ge \tau U(\theta, D_H)\}$ .

So instead of reporting the holdout performance at all, you just see if you are much worse than the training performance, i.e. if you're overfitting. The fact that the mechanism is DP also makes it difficult to learn the holdout set. See the thresholdout link for more details.

# 2.10 Discussion

# The definition of differential privacy

• First rigorous mathematical definition of privacy.

 $<sup>^6\</sup>mathrm{Also}$  see https://ai.googleblog.com/2015/08/the-reusable-holdout-preserving.html

- Relaxations and generalisations possible.
- Connection to learning theory and reproducibility.

#### Current uses

- Apple. DP is used internally in the company to "protect user privacy". It is not clear exactly what they are doing but their efforts seem to be going in the right direction.
- Google. The company has a DP API available based on randomised response, RAP-POR
- Uber. Elastic sensitivity for SQL queries, which is available as open source. This is a good thing, because it is easy to get things wrong with privacy.
- US 2020 Census. It uses differential privacy to protect the condidentiality of responders' information while maintaining data that are suitable for their intended uses.

# Open problems

- Complexity of differential privacy.
- Verification of implementations and queries.

# Available privacy toolboxes

## Differential privacy

- Open DP <a href="https://opendp.org/">https://opendp.org/</a>
- https://github.com/bmcmenamin/thresholdOut-explorationsThreshold out
- https://github.com/steven7woo/Accuracy-First-Differential-PrivacyAccuracy-constrained DP
- https://github.com/menisadi/pydpVarious DP algorithms
- https://github.com/haiphanNJIT/PrivateDeepLearning Deep learning and DP

#### k-anonymity

• https://github.com/qiyuangong/Mondrian Mondrian k-anonymity

# Learning outcomes

# Understanding

- Linkage attacks and k-anonymity.
- Inferring data from summary statistics.
- The local versus centralised differential privacy model.

2.11. EXERCISES 45

• False discovery rates.

#### Skills

- Make a dataset satisfy k-anonymity with respect to identifying attributes.
- Apply the randomised response and Laplace mechanism to data.
- Apply the exponential mechanism to simple decision problems.
- Use differential privacy to improve reproducibility.

# Reflection

- How can potentially identifying attributes be chosen to achieve k-anonymity?
- How should the parameters of the two ideas,  $\epsilon$ -DP and k-anonymity be chosen?
- Does having more data available make it easier to achieve privacy?

### Further reading

- *k*-anonymity <sup>15</sup>
- Randomness, privacy, and the US 2020 census  $^{11}$
- The paper introducing differential privacy<sup>8</sup>
- Differential privacy book <sup>7</sup>
- Bayesian inference and privacy<sup>5</sup>
- Local differential privacy and statistics <sup>6</sup>
- Local differential privacy and applications <sup>17</sup>
- The exponential mechanism <sup>14</sup>

# 2.11 Exercises

EXERCISE 1. Show that the randomised response mechanism, as originally defined, is not differentially private with respect to the addition/deletion neighbourhood.

EXERCISE 2. Define a variant of the binary randomised response mechanism that satisfies differential privacy with respect to the insertion and deletion neighbourhood. More precisely, it should hold that for any dataset pair  $\mathbf{x} = (x_1, \dots, x_n)$ , and  $\mathbf{x}' = (x_1, \dots, x_n)$ ,  $\mathbf{x} = (x_1, \dots, x_n, x_{n+1})$ , the mechanism satisfies

$$\pi(A|\mathbf{x}) \le e^{\epsilon} \pi(A|\mathbf{x}'), \quad \forall A \subset \mathcal{A}.$$

EXERCISE 3. The mean estimator  $\hat{\theta}$  for the randomised response mechanism with flipping probability p, where the true data generation process is Bernoulli with parameter  $\theta$ , is defined as

$$\hat{\theta}(a) = \frac{\bar{a} - p}{1 - 2p},$$

where  $\bar{a} = \frac{1}{T} \sum_{t=1}^{T} a_t$  is the mean of the observed answers. Prove that this estimator is unbiased, i.e. that

$$\mathbb{E}[\hat{\theta}] = \theta.$$

where the expectation is taken over the unobserved data randomness and the randomness of the mechanism.

EXERCISE 4. Prove that the exponential mechanism is uniform when  $\epsilon \to 0$  and deterministically returns the maximising action when  $\epsilon \to \infty$ .

EXERCISE 5. Prove that the exponential mechanism, when we used to calculate a noisy version of the function q(x) with  $q: \mathcal{X} \to \mathbb{R}$  and utility  $U(a,q,x) = -|q(x)-a|^p$ , results in the Laplace mechanism for p=1 and the Gaussian mechanism for p=2.

EXERCISE 6. Consider the following relaxation of differential privacy to KL divergences. A mechanism is  $\epsilon$ -KL-private if

$$\sum_{a} \ln \frac{\pi(a|x)}{\pi(a|x')} \pi(a|x) \le \epsilon$$

Prove that two-folded composition of such a mechanism is  $2\epsilon$ -KL-private.

# Chapter 3

# Fairness.

When machine learning algorithms are applied at scale, it can be difficult to imagine what their effects might be. In this part of the book, we consider notions of fairness as seen through the prism of conditional independence and meritocracy. The first notion requires that we look deeper into directed graphical models.

The problem of fairness in machine learning and artificial intelligence has only recently been widely recognised. When any algorithm is implemented at scale, no matter the original objective and whether it is satisfied, it has significant societal effects. In particular, even when considering the narrow objective of the algorithm, even if it improves it overall, it may increase inequality.

In this course we will look at two aspects of fairness. The first has to do with disadvantaged populations that form distinct social classes due to a shared income stratum, race or gender. The second has to do with meritocratic notions of fairness.

This chapter requires some knowledge of decision theory (See Appendix A) to understand the principles of expected utility maximisation and graphical models (See Appendix C) to understand the idea of conditional independence.

# 3.1 Introduction.

Fairness is a concept that has received much attention recently when applied to large-scale algorithmic decision making. However, the very concept of fairness is not well-defined and encompasses many different ideas. Some of those relate to fair treatment of individuals: *Meritocracy* is the idea that people should receive rewards according to their merit. *Equal treatment* is the related notion that similar people should be treated similarly under similar circumstances. Some concepts are related more to the treatment of different groups: *Proportional representation* is the idea that proportions of different groups in society should be reflected in every facet of society. Finally *non-discrimination* captures the notion of not treating people differently depending on sensitive characteristics.

#### Meritocracy

Meritocracy embodies the principle that merit should be rewarded. A common example are admissions to universities. Some type of summary, typically a grade obtained from high school, is used to represent the underlying merit of individuals.

EXAMPLE 5 (College admissions). In this example, we have two students. In terms of grades, student B is clearly better. If we can only accept one of them, and given no other information, it seems like the

natural choice is student B.

- Student A has a grade 4/5 from Gota Highschool.
- Student B has a grade 5/5 from Vasa Highschool.

Grades, by themselves, are typically insufficient information. It might be that grades from some high-schools are inflated and do not represent the quality of individuals accurately. So, let us suppose we now consider the information.

EXAMPLE 6 (Additional information). In particular, let us suppose that we have statistics on how well students from different high school do, depending on their high school grade.

- 70% of admitted Gota graduates with 4+ get their degree.
- 50% of admitted Vasa graduates with 5 get their degree.

All other thing being equal, it is now more likely that student A will graduate. So perhaps we should take in A and not B.

We still don't know how a specific student will do!

I must emphasise that these are only statistics, and not necessarily predictive of the students' ability. Ideally, we would like to admit the students that we expect to do well, given the information that we have. However, this information is typically not enough for us to make reliable predictions. In addition, we might want to also make sure that everybody has a chance to obtain a good education. In order to achieve this, we might want to promote ethnic or gender equality through university admissions. Unfortunately, there is no ideal solution and we must always balance the benefit of individual students with that of specific societal groups as well as society as a whole.

#### Solutions

These solution methods are not completely exclusive, and can be implemented simultaneously to some extent.

- Admit everybody? This suggests that everybody is admitted to at least one university, perhaps even their university of choice. However, it requires that there is enough teaching capacity for all students in the first year. Subsequently, we expect the students who were not qualified to drop out. Of course, this is unfair to the qualified students, as it drains resources that could have been used for them.
- Admit *randomly*? Completely random decisions are not considered fair, because they do not take into account any information. However, randomisation can also be used in conjunction with grades to ensure that everybody has a shot.
- Use *prediction* of individual academic performance? The more information we have, the better we can predict academic performance. A grade from high school is one indicator, but more data can be used to obtain better predictions. Of course, no prediction is perfect.
- Should we take into account *group membership* or other population information? For many reasons, students in some groups can perform differently in standardised tests, even though their innate talents may be no different than students not in the group. The classical example of this is high school teachers discouraging girls from mathematics.

EXAMPLE 7 (Hiring decisions.). As a further example, consider gender balance in hiring decisions. Typically, received applications are screened, so that some applicants undergo through an interview process.

At the end, some of the interviewed applicants will be hired. There are two decision points here, with most people being cut off at the first point: the screening. To automate this process, Amazon worked on resume-screen program. <sup>1</sup> However, this was scrapped after it was discovered that it predominantly favoured men. The reason is not entirely clear, but it was probably due to the fact that they trained the system on their own screening decisions, and given that the tech industry predominantly hires men in the first place, women were likely rejected in the screening phase.

# 3.2 Group fairness.

Let us now take a look at concepts of fairness related to *group membership*. This includes concepts such as equal treatment, equality of opportunity and generally lack of discrimination. The general idea is that we would like for members of society to follow trajectories through life that do not strongly depend on their membership in sensitive groups. For example, gender should not play a role in academic achievement. Ethnic heritage should have no influence on annual income. Unfortunately, the underlying societal dynamics create situations where group membership becomes important.

In this section, and throughout this chapter, we will imagine that an individual is interacting with a system, which makes decisions about the individual, such as whether or not to give them a loan. These result in a certain outcome, such as the individual using the borrowed money to invest in a business and then having a particular annual income. Crucially, these outcomes can be correlated with group membership, because of societal dynamics. The system designer should make sure that not only the system does what it intends, such as giving loans to people that are expected to repay them, but also that it does not create inequalities between different groups. As we will see later, depending on our fairness definition, and on the societal dynamics, this is not always easy.

#### Bail decisions

For a more detailed example, let us consider bail decisions in the US court system. When a defendant is charged, the judge has the option to either place them in jail pending trial, or set them free, under the condition that the defendant pays some amount of 'bail'. The amount of bail (if any) is set to deter flight or a relapse.

This process sometimes includes the use of a software tool called COMPAS, which gives risk scores for the possibility of flight, recidivism or violent behaviour. These scores are taken into account by judges when making decisions. In some cases, it appears as though automating this procedure might lead to better outcomes. But is that generally true?

In this setting, a defendant t appears before a judge with observable features  $x_t \in \mathcal{X}$ , and a sensitive group variable  $z_t \in \mathcal{Z}$ . The judge employs a specific policy  $\pi$  to select a decision  $a_t \in \mathcal{A}$ , with

$$\pi(a_t \mid x_t, z_t) \tag{3.2.1}$$

denoting the probability of action  $a_t$  given the individual's features, as well as the sensitive variable. We can assume that the policy is fixed ahead of time, and thereafter decisions are made according to the policy. After the judge makes their decision, they observe an outcome  $y_t$  sampled from some potentially unknown distribution with parameter  $\theta$ :

$$P_{\theta}(y_t \mid x_t, z_t), \qquad P_{\theta}(y_t \mid x_t, z_t, a_t), \tag{3.2.2}$$

 $<sup>^{1}</sup>$ https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MKO8G

denoting the probability of  $y_t$  given the individual's features and the action taken. Here, there are two possibilities. Either the outcome  $y_t$  depends only on the observed features, or the action as well. The correct formulation depends on the meaning of the action.

Actions that affect the outcome. Lt us first consider a simple set of actions  $\mathcal{A} = \{0, 1\}$ , where a defendant is granted  $(a_t = 1)$  or denied  $(a_t = 0)$  bail. If a defendant is not given bail, they must remain in jail, and will always be at the trial  $(y_t = 0)$ . This costs both the government and the defendant, so the judge prefers not to deny bail too often. On the other hand, if the defendant is granted bail, there is a chance they will re-offend, or will fail to attend trial  $(y_t = 1)$ . Consequently, we can define the following utility function<sup>2</sup> for the judge, that roughly reflects those preferences:

$$U(a,y) = a - y. (3.2.3)$$

While this weighs the judge's two main concerns equally, we can also imagine different formulations: e.g. if the judge thinks it is much more important to keep people out of jail than to prevent re-offences until trial, they might select a function such as : U(a, y) = 10a - y.

Actions that do not affect the outcome. In this scenario, the only way for the actions to not affect the outcome is if everybody is released, with the judge's action only being a confidential note on the case file. Then their action cannot affect the outcome, since it is only revealed to the judge. The utility function can be based on the accuracy of predictions, so we can simply set it to U(a, y) = a - y.

The optimisation problem. Now, given the parameter  $\theta$  of the unknown distribution, the judge can find a decision rule  $\pi$  maximising utility in expectation

$$\mathbb{E}_{\theta}^{\pi}(U) = \sum_{x,z} P_{\theta}(x,z) \sum_{a} \pi(a|x,z) \sum_{y} P_{\theta}(y|a,x,z) U(a,y). \tag{3.2.4}$$

In practice, of course,  $\theta$  is not known but we have some training data  $D = (x_t, z_t, a_t, y_t)_{t=1}^T$ , collected with some historical policy  $\pi_0$ , and we have to resort to one of the following solutions. Firstly, estimating some  $\hat{\theta}$  from the data and replacing that in the expectation. Secondly, calculating a posterior distribution  $\beta(\theta|D)$  and maximising  $\int_{\Theta} \mathbb{E}_{\theta}^{\pi}(U) \, \mathrm{d}\beta(\theta|D)$ . Thirdly, maximising an unbiased estimate of expected utility, as shown below:

$$\begin{split} \mathbb{E}^{\pi}_{\theta}(U) &= \sum_{x,z} P_{\theta}(x,z) \sum_{a} \pi(a|x,z) \sum_{y} P_{\theta}(y|a,x,z) U(a,y) \\ &\approx \frac{1}{T} \sum_{t} \sum_{a} \pi(a|x_{t},z_{t}) \sum_{y} P_{\theta}(y|a,x_{t},z_{t}) U(a,y) \approx \frac{1}{T} \sum_{t} \frac{\pi(a_{t}|x_{t},z_{t})}{\pi_{0}(a_{t}|x_{t},z_{t})} U(a_{t},y_{t}). \end{split}$$

The first approximation simply replaces the feature distribution with its empirical approximation. The second approximation is a bit more delicate, and is a type of importance sampling, as explained in more detail in Section 5.3.1. This is because, in the original data we have only some specific  $a_t, y_t$  for every individual t, where  $a_t$  was selected by some historical policy  $\pi_0$ . For that reason, we must use importance sampling in order to estimate the utility of any policy  $\pi$ . Unfortunately, the policy  $\pi_0$  is also unknown and must be estimated from the data. So, we are left with deciding between modelling  $P_{\theta}(y|a,x,z)$  or  $\pi_0(a|x,z)$ .

However, doing so may have unintended effects, as treatment of different groups may be unequal. For that reason, we may have to add explicit fairness constraints to our optimisation problem.

<sup>&</sup>lt;sup>2</sup>See Appendix A

#### Demographic parity and equality of opportunity.

Let us take a look at how many individuals are given different scores. Figure ?? shows the proportions of individuals obtaining different risk scores as a function of their ethnic group. While the general population, and Caucasians in particular, have a distribution sharply concentrated in low scores, Black defendants obtain much higher risk scores: they are almost as likely to be ranked a 10 as the are to be scored a 1.

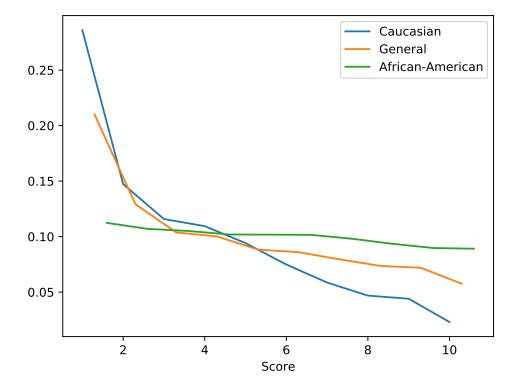


Figure 3.1: Apparent bias in risk scores towards black versus white defendants.

Typically, we would like to demographic parity across groups, either with respect to the decisions, or with respect to the outcomes. Decision parity is called *equality of opportunity* and satisfies the following condition

$$\mathbb{P}_{\theta}^{\pi}(a_t|z_t) = \mathbb{P}_{\theta}^{\pi}(a_t). \tag{3.2.5}$$

In other words, the action (risk score, in this case) is independent of the group. As we see in Figure ??, these distributions are quite different. The curve for "General" corresponds to  $\mathbb{P}^{\pi}_{\theta}(a_t)$ , while the other two curves correspond to  $\mathbb{P}^{\pi}_{\theta}(a_t|z_t)$  for  $z_t$  being Caucasian and African-American respectively.

*Demographic parity* instead relates to outcomes. If this condition is satisfied, then the probability distribution of outcomes is independent of group membership.

$$\mathbb{P}_{\theta}^{\pi}(y_t|z_t) = \mathbb{P}_{\theta}^{\pi}(y_t). \tag{3.2.6}$$

#### Calibration.

On the other hand, the scores generated by the software seemed to be very predictive on whether or not defendants would re-offend, independently of their race. Figure 3.2 shows that, if an individual obtains a high score then they are very likely to re-offend, and conversely, they are unlikely to re-offend when they have a low score.

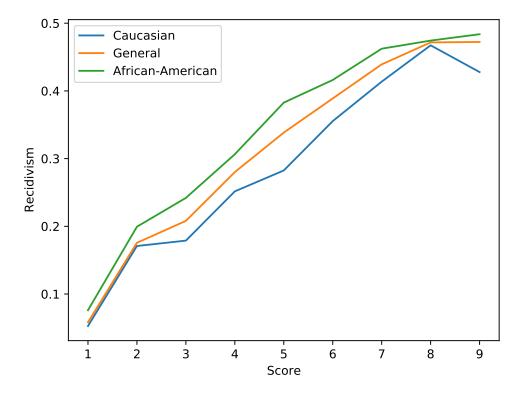


Figure 3.2: Recidivism rates by risk score.

This concept can be quantified in terms of the conditional distribution of outcomes given the score:

$$\mathbb{P}_{\theta}^{\pi}(y_t|a_t, z_t) = \mathbb{P}_{\theta}^{\pi}(y_t|a_t). \tag{3.2.7}$$

This means that  $y_t$  is conditionally independent of  $z_t$  given the score  $a_t$ . So, in some sense, the score is sufficient for us to predict the outcome, and knowing the race does not help us predict any better.

### Balance

While the system's predictions seem to be calibrated against the chance of recidivism, this does not mean that race plays no role. Figure 3.3 breaks down the population in people that re-offended and those that did. For each sub-population, we then plot the proportion of people receiving different scores by race. While people generally have a small probability of obtaining a high risk score, we see that Black defendants obtain much higher scores.

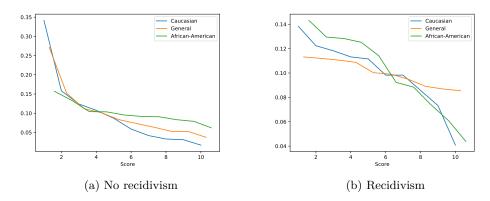


Figure 3.3: Score breakdown based on recidivism rates.

Balance can also be interpreted as a probabilistic condition, of the following form:

$$\mathbb{P}_{\theta}^{\pi}(a_t|y_t, z_t) = \mathbb{P}_{\theta}^{\pi}(a_t|y_t). \tag{3.2.8}$$

Here we have that  $a_t$  is conditionally independent of  $z_t$  given the outcome  $y_t$ . So, if we partition the population according to their outcome, we should find that given the outcome, the distribution of scores is the same no matter what their race. However, this is not what the example shows: there is a strong dependence on race.

How can we explain this discrepancy? We can show that in fact, each one of these different measures of bias in our decision rules can be seen as a notion of conditional independence. Had  $z_t$  been an independent variable, i.e. so that  $P_{\theta}(x,z) = P_{\theta}(x)P_{\theta}(z)$  and  $P_{\theta}(y|x,a,z) = P_{\theta}(y|x,a)$ , then it would have no role. Unfortunately, though, it does influence the distributions of the other variables.

Example 8 (Classification). One setting where the above discussion applies directly are classification problems. There, the outcomes y are simply labels. the decision space  $\mathcal A$  is to just predicting a label, and the utility function is the classification accuracy, so that  $U(a,y) = \mathbb{I}\{a=y\}$ . In this setting, y does not directly depend on a, since a is just a prediction of the latent label. More precisely,  $y \perp \!\!\! \perp a \mid x$ . This allows us to write:

$$\mathbb{P}_{\theta}^{\pi}(y_t \mid a_t, x_t, z_t) = \mathbb{P}_{\theta}^{\pi}(y_t \mid x_t, z_t).$$

EXAMPLE 9 (Regression). We can also apply these ideas to regression problems. We are again asked to simply predict a latent variable y. For concreteness, assume  $y \in \mathbb{R}$ . Then, the decision space  $\mathcal{A} = \mathbb{R}$  as well. A standard utility function is the negative squared prediction error, so that  $U(a, y) = -(a - y)^2$ .

In this setting we can also relax our framework to work with expectations instead of probabilities. For example, calibration and balance can be written as the conditions

$$\mathbb{E}_{\theta}^{\pi}(y_t|a_t, z_t) = \mathbb{E}_{\theta}^{\pi}(y_t|a_t), \qquad \mathbb{E}_{\theta}^{\pi}(a_t|y_t, z_t) = \mathbb{E}_{\theta}^{\pi}(a_t|y_t),$$

respectively.

# 3.2.1 Measuring group fairness.

The above discussion focused on the setting where we have a precise model parameter  $\theta$  and policy  $\pi$ . However, we frequently do not know at least one of those things. For example, when we only have some observational data generated by some historical policy  $\pi$ , need to estimate  $\mathbb{P}^{\pi}_{\theta}$  somehow. Moreover, these conditions will not be perfectly satisfied. How can we measure the deviation from these conditions?

Empirical expected utility. Let us start with expect utility, and write it in full:

$$\mathbb{E}_{\theta}^{\pi}(U) = \sum_{x,z,y} \mathbb{P}_{\theta}^{\pi}(y, a, x, z) U(a, y).$$

The  $y_t, a_t, x_t, z_t$  is already sampled from  $\mathbb{P}^{\pi}_{\theta}$ , so we can approximate the expected utility with

$$\hat{E}_n(U) = \sum_{t=1}^n U(a_t, y_t), \qquad a_t, y_t \sim \mathbb{P}_{\theta}^{\pi}.$$

Thus, estimating the expected utility of a particular policy is always possible, as long as we can sample actions and outcomes from  $\mathbb{P}^{\pi}_{\theta}$ . Of course, this requires having some specific model for sampling outcomes for any given action, as the policy in historical data might have taken different actions from the ones we are taking.

**Deviation from balance.** How about measuring whether balance is satisfied? Unfortunately none of these fairness conditions will be satisfied in practice. So what are we going to do? A simple idea is to simply look at how far away the distribution is from independence. In particular, let us look at the difference between those distributions

$$\mathbb{P}^{\pi}_{\theta}(a|y,z), \qquad \mathbb{P}^{\pi}_{\theta}(a|y)$$

for all values of y, z. Let us first look at the total variation distance:

$$\|\mathbb{P}_{\theta}^{\pi}(a|y,z) - \mathbb{P}_{\theta}^{\pi}(a|y)\|_{1} = \sum_{a} \|\mathbb{P}_{\theta}^{\pi}(a|y,z) - \mathbb{P}_{\theta}^{\pi}(a|y)\|_{1}.$$

We can then sum over all possible values of y, z to obtain a single measure of how much we deviate from balance. To obtain an empirical estimate, we can simply sum over all observations to obtain the empirical models:

$$\hat{P}_n(a|y,z) = \sum_t \frac{\mathbb{I}\left\{a_t = a \land y_t = y \land z_t = i\right\}}{\mathbb{I}\left\{y_t = y \land z_t = i\right\}}$$

We can then plug those into our original measure:

$$\|\mathbb{P}_{\theta}^{\pi}(a|y,z) - \mathbb{P}_{\theta}^{\pi}(a|y)\|_{1} \approx \sum_{a} \|\hat{P}_{n}(a|y,z) - \hat{P}_{n}(a|y)\|_{1}.$$

However, we can use something else as well, such as the KL divergence:

$$D(\mathbb{P}^{\pi}_{\theta}(a|y,z)||\mathbb{P}^{\pi}_{\theta}(a|y)) = \sum_{a} \ln \frac{\mathbb{P}^{\pi}_{\theta}(a|y,z)}{\mathbb{P}^{\pi}_{\theta}(a|y)} \mathbb{P}^{\pi}_{\theta}(a|y,z)$$

# 3.2.2 Trading off utility and fairness.

We have now selected a notion of utility and of fairness. We have also seen how to measure utility and fairness in practice. How can we choose a policy  $\pi$  that optimises for both? The general idea is to formalise the problem in terms of a constrained or unconstrained optimisation.

## Formalisation of the problem

Let us use  $F(\pi, \theta)$  for the measure of deviation from fairness and  $U(\pi, \theta)$  for utility.

We can deal with the problem in two ways. The first is to formulate a combined objective that mixes the utility and fairness objectives. Then we can formalise the problem as an unconstrained optimisation. The second is to optimise one objective and place constrains on the other. For example, we can place a limit on how much unfairness we are willing to tolerate and maximise utility. Or, we can place a limit on how far away we want to be from the solution that optimises for utility only, and optimises for fairness. Which method to use is application-dependent, but also has some computational impact.

## Unconstrained optimisation.

The simplest idea is to treat the utility and fairness symmetrically. Then we can simply optimise for their mixture. Let  $\lambda \in [0, 1]$ :

$$\max_{\pi}(1-\lambda)U(\pi,\theta) - \lambda F(\pi,\theta)$$

In this setting, we have an unconstrained optimisation problem. These are typically easier to solve, and we can use simple gradient methods. In particular, note that the gradient with respect to the policy of the unconstrained objective is

$$\nabla_{\pi}[(1-\lambda)U(\pi,\theta) - \lambda F(\pi,\theta)] = (1-\lambda)\nabla_{\pi}U(\pi,\theta) - \lambda\nabla_{\pi}F(\pi,\theta).$$

Consequently, we only need to minimise the sum of these gradients.

#### Constrained fairness.

In this setting, we have an upper bound  $\epsilon$  on the amount of fairness deviation we can allow. Let  $\epsilon \geq 0$ :

$$\max_{\pi} U(\pi, \theta)$$
s.t. $F(\pi, \theta) \le \epsilon$ .

# Constrained utility.

In this setting, we want to limit the amount of sub-optimality  $\epsilon$  we have relative to the optimal policy that ignores fairness. Let  $\epsilon \geq 0$ :

$$\in_{\pi} F(\pi, \theta)$$
 s.t. $U(\pi, \theta) \ge \max_{\pi'} U(\pi', \theta) - \epsilon$ 

# 3.3 Individual fairness

An altogether different type of fairness concept is individual fairness. This relates to the notion of meritocracy. This means that "better" people should have "better" outcomes. But what do we mean by better? This depends strongly on the context.

## Meritocracy.

#### Inherent worth.

- a: Action taken by the decision maker. For example, a could specify which applicants are admitted to university, with  $a_i$  specifying the university to which an applicant is admitted
- $x_i \in \mathbb{R}^d$ : Individual features. That's what we observe about the individual.
- $w_i \in \mathbb{R}$ : Individual worth. This can be taken to mean how well the individual is expected to perform in university, i.e. their expected grade. However, this quantity is not directly observed.
- $u_i(a) \in \mathbb{R}$ : Individual utility. This is how much the individual gains by the decision. For simplicity, we can imagine that i is only impacted by the decision regarding them, rather than anybody else, so we can just define a function  $u_i(a_i)$  instead.

A decision a is fair if, for all i, j it holds that: if  $w_i > w_j$  then  $u_i(a) > u_j(a)$ .

The problem in this setting, is that it is not clear what an individual's worth is. It is also not clear what their utility is. In a sense, if i and j are both admitted to the same university, is it worth the same to them?

#### Smoothness.

# Treating similar individuals similarly.

- $d(x_i, x_j)$ : Distance between individuals i, j.
- $D[\pi(a_i|x_i), \pi(a_i|x_i)]$ : Distance between policy decisions.
- Utility function  $U(\pi)$ .

Our goal is to find a decision rule  $\pi$  maximising  $U(\pi)$ , so that it is Lipschitz-smooth with respect to D,d:

$$D[\pi(a_i|x_i), \pi(a_i|x_i)] \le d(x_i, x_i)$$

When  $D(P,Q) = \max_{a \in A} |\log[P(a)/Q(a)]|$  then  $\pi$  satisfies differential privacy with  $d(x,y) = \epsilon ||x-y||_1$ .

#### Smoothness and parity.

To define groups, we select  $S, T \subset \mathcal{X}$ .

### **Definition 3.3.1.** $\epsilon$ -parity

$$\|\pi(a|x \in S) - \pi(a|x \in T)\|_1 \le \epsilon.$$

**Definition 3.3.2.** Bias for some action a:

$$\max_{\pi} \pi(a|x \in S) - \pi(a|x \in T) \qquad \text{s.t. } \pi \text{ is Lipschitz}$$

The main problem here is that if we are required to be smooth, but one set has most of its members in one subset.

3.4. EXERCISES 57

# 3.4 Exercises

EXERCISE 7. Consider the following problem, where z is Bernoulli(1/2) and  $x\mathbb{R}^2$  is a normally distributed variable with mean (0, z) and identity covariance matrix. There are two actions and outcomes, with the outcome being  $y_t = 1$  with probability  $1/(1 + e^{x_a})$ . [Add more details]

# Chapter 4

# Reproducibility

# 4.1 Reproducibility

One of the main problems in science is reproducibility: when we are trying to draw conclusions from one specific data set, it is easy to make a mistake. For that reason, the scientific process requires us to use our conclusions to make testable predictions, and then test those predictions with new experiments. These new experiments should bear out the results of the previous experiments. In more detail, reproducibility can be thought of as two different aspects of answering the question "can this research be replicated?"

# Computational reproducibility: Can the study be repeated?

Can we, from the available information and data, exactly reproduce the reported methods and results?

This is something that is useful to be able to even to the original authors of a study. The standard method for achieving this is using version control tools so that the exact version of algorithms, text and data used to write up the study is appropriately labelled. Ideally, any other researcher should be able to run a single script to reproduce all of the study and its computations. The following tools are going to be used in this course:

- jupyter notebooks for interactive note taking.
- svn, git or mercurial version control systems for tracking versions, changes and collaborating with others.

# Scientific reproducibility: Is the conclusion correct?

Can we, from the available information and a *new* set of data, reproduce the conclusions of the original study?

Here followup research may involve using exactly the same methods. In AI research would mean for example testing whether an algorithm is really performing as well as it is claimed, by testing it in new problems. This can involve a re-implementation. In more general scientific research, it might be the case that the methodology proposed by an original study is flawed, and so a better method should be used to draw better conclusions. Or it might simply be that more data is needed.

When publishing results about a *new method*, computational reproducibility is essential for scientific reproducibility.



A simple example is the 2016 presidential election in the United States, where polling forecasts were not able to predict the outcome. In general, while we can make models about people's opinions regarding candidates in order to predict voting totals, the test of these models comes in the actual election. Unfortunately the only way we have of tuning our models is on previous elections, which are not that frequent, and on the results of previous polls. In addition, predicting the winner of an election is slightly different from predicting how many people are prepared to vote for them across the country. This, together with other factors such as shifting opinions, motivation and how close the sampling distribution is to the voting distribution have

a significant effect on accuracy.

The same thing can be done in when dealing purely with data, by making sure we use some of the data as input to the algorithm, and other data to measure the quality of the algorithm itself. In the following, we assume we have some algorithm  $\lambda: \mathcal{D} \to \Pi$ , where  $\mathcal{D}$  is the universe of possible input data and  $\Pi$  the possible outputs, e.g. all possible classification policies. We also assume the existence of some quality measure U. How should we measure the quality of our algorithmic choices?

Take classification as an example. For a given training set, simply memorising all the labels of each example gives us perfect performance on the training set. Intuitively, this is not a good measure of performance, as we'd probably get poor results on a freshly sampled set. We can think of the training data as input to an algorithm, and the resulting classifier as the algorithm output. The evaluation function also requires some data in order to measure the performance of the policy. This can be expressed into the following principle.

# 

Data used for estimation cannot be used for evaluation.

This applies both to computer-implemented and human-implemented algorithms.

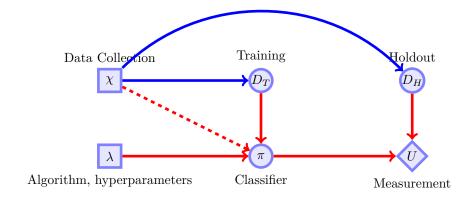


Figure 4.1: The decision process in classification.

One can think of the decision process in classification as follows. First, we decide to collect some data according to some experimental protocol  $\chi$ . We also decide to use some algorithm (with associated hyperparameters)  $\lambda$  together with data  $D_T$  we will obtain from our data collection in order to obtain a classification policy  $\pi$ . Typically, we need to measure the quality of a policy according to how well it classifies on unknown data. This is because our policy has been generated using  $D_T$ , and so any measurement of its quality using the same data  $D_T$ , is going to be biased.

Classification accuracy. For classification problems, there is a natural metric U to measure: The classification accuracy of the classifier. If the classification decisions are stochastic, then the classifier assigns probability  $\pi(a \mid x)$  to each possible label a, and our utility is simply the identity function  $U(a, y) \triangleq \mathbb{I}\{a = y\}$ .

#### Classification accuracy

For any fixed classifier  $\pi$ , the classification accuracy under the data distribution  $\chi$  is:

$$\mathbb{E}_{\chi}[U(\pi)] = \sum_{x,y} \underbrace{\mathbb{P}_{\chi}(x,y)}_{\text{Data probability}} \overbrace{\pi(a=y\mid x)}^{\text{Decision probability}}$$

The classification accuracy of policy  $\pi$  under  $\chi$  is the expected number of times the policy decides  $\pi$  chooses the correct class. However, when approximating  $\chi$  with a sample  $D_H$ , we instead obtain the empirical estimate:

$$\mathbb{E}_{D_H} U(\pi) = \sum_{(x,y) \in D_H} \pi(a = y \mid x) / |D_H|.$$

Of course, there is no reason to limit ourselves to the identity function. The utility could very well be such that some errors are penalised more than other errors. Consider for example an intrusion detection scenario: it is probably more important to correctly classify intrusions. This can be taken into account by weighing correct decisions about intrusions more relative to correct decisions about normal data.

# 4.1.1 The human as an algorithm

#### The human as an algorithm.

The same way with which an algorithm creates a model from some prior assumptions and data, so can a human select an algorithm and associated hyperparameters by executing an algorithm herself. This involves trying different algorithms and hyperparameters on the same training data  $D_T$  and then measuring their performance in the holdout set  $D_H$ .

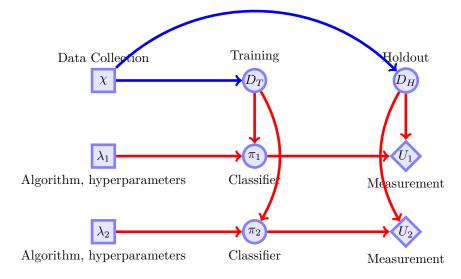


Figure 4.2: Selecting algorithms and hyperparameters through holdouts

#### Holdout sets

To summarise, holdout sets are used in order to be able to evaluate the performance of specific algorithms, or hyparameter selection.

- Original data D, e.g.  $D = (x_1, \ldots, x_T)$ .
- Training data  $D_T \subset D$ , e.g.  $D_T = x_1, \ldots, x_n, n < T$ .
- Holdout data  $D_H = D \setminus D_T$ , used to measure the quality of the result.
- Algorithm  $\lambda$  with hyperparametrs  $\phi$ .
- Get algorithm output  $\pi = \lambda(D_T, \phi)$ .
- Calculate quality of output  $U(\pi, D_H)$

We start with some original data D, e.g.  $D = (x_1, \ldots, x_T)$ . We then split this into a training data set  $D_T \subset D$ , e.g.  $D_T = x_1, \ldots, x_n$ , n < T and holdout dataset  $D_H = D \setminus D_T$ . This is used to measure the quality of selected algorithms  $\lambda$  and hyperparameters  $\phi$ . We run an algorithm/hyperparameter combination on the training data and obtain a result  $\pi = \lambda(D_T, \phi)$ . We then calculate the quality of the output  $U(\pi, D_H)$  on the holdout set. Unfortunately, the combination that appears the best due to the holdout result may look inferior in a fresh sample. Following the principle of "data used for evaluation cannot be used for estimation", we must measure performance on another sample. This ensures that we are not biased in our decision about what is the best algorithm.

#### Holdout and test sets for unbiased algorithm comparison

Consider the problem of comparing a number of different algorithms in  $\Lambda$ . Each algorithm  $\lambda$  has a different set of hyperparameters  $\Phi_{\lambda}$ . The problem is to choose the best parameters for each algorithm, and then to test them independently. A simple meta-algorithm for doing this is based on the use of a *holdout* set for choosing hyperparameters for each algorithm, and a *test* set to measure algorithmic performance.

# Algorithm 1 Unbiased adaptive evaluation through data partitioning

```
Partition data into D_T, D_H, D^*.

for \lambda \in \Lambda do

for \phi \in \Phi_{\lambda} do

\pi_{\phi,\lambda} = \lambda(D_T, \phi).

end for

Get \pi_{\lambda}^* maximising U(\pi_{\phi,\lambda}, D_H).

u_{\lambda} = U(\pi_{\lambda}^*, D^*).

end for

\lambda^* = \arg \max_{\lambda} u_{\lambda}.
```

$$\lambda(D_T) = \arg\max_{y} U(y, D_T)$$

we typically obtain a biased estimate, which depends both on the algorithm itself and the training data. For K-NN in particular, when we measure accuracy on the training data, we can nearly always obtain near-perfect accuracy, but not always perfect. Can you explain why?

<sup>&</sup>lt;sup>1</sup>As typically algorithms are maximising the quality metric on the training data,

## Final performance measurement

When comparing many algorithms, where we must select a hyperparameter for each one, then we can use one dataset as input to the algorithms, and another for selecting hyperparameters. That means that we must use another dataset to measure performance. This is called the testing set. Figure 4.3 illustrates this.

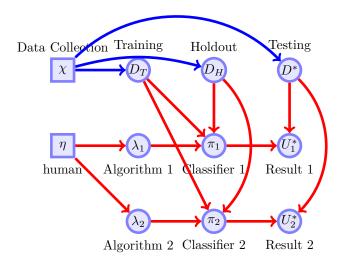


Figure 4.3: Simplified dependency graph for selecting hyperparameters for different algorithms, and comparing them on an independent test set. For the i-th algorithm, the classifier model is

# 4.1.2 Algorithmic sensitivity

The algorithm's output does have a dependence on its input, obviously. So, how sensitive is the algorithm to the input?

#### Independent data sets

One simple idea is to just collect independent datasets and see how the output of the algorithm changes when the data changes. However, this is quite expensive, as it not might be easy to collect data in the first place.

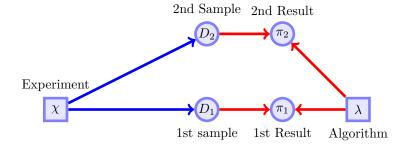


Figure 4.4: Multiple samples

65

## Bootstrap samples

A more efficient idea is to only collect one dataset, but then use it to generate more datasets. The simplest way to do that is by sampling with replacement from the original dataset, new datasets of the same size as the original. Then the original dataset is sufficiently large, this is approximately the same as sampling independent datasets. As usual, we can evaluate our algorithm on an independent data set.

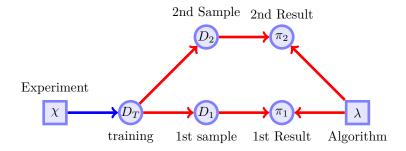


Figure 4.5: Bootstrap replicates of a single sample

#### **Bootstrapping**

Bootstrapping is a general technique that can be used to:

- Estimate the sensitivity of  $\lambda$  to the data x.
- Obtain a distribution of estimates  $\pi$  from  $\lambda$  and the data x.
- When estimating the performance of an algorithm on a small dataset  $D^*$ , use bootstrap samples of  $D^*$ . This allows us to take into account the inherent uncertainty in measured performance. It is very useful to use bootstrapping with pairwise comparisons.

# **Bootstrapping**

- 1. **Input** Training data D, number of samples k.
- 2. **For** i = 1, ..., k
- 3.  $D^{(i)} = \text{Bootstrap}(D)$
- 4. **return**  $\{D^{(i)} \mid i = 1, ..., k\}$ . where Bootstrap(D) samples with replacement |D| points from  $D_T$ .

In more detail, remember that even though the test score is an *independent* measurement of an algorithm's performance, it is *not* the actual expected performance. At best, it's an unbiased estimate of performance. Hence, we'd like to have some way to calculate a likely performance range from the test data. Bootstrapping can help: by taking multiple samples of the test set and calculating performance on each one, we obtain an empirical distribution of scores.

Secondly, we can use it to tell us something about the sensitivity of our algorithm. In particular, by taking multiple samples from the training data, we can end up with multiple models. If the models are only slightly different, then the algorithm is more stable and we can be more confident in its predictions.

Finally, bagging also allows us to generate probabilistic predictions from deterministic classification algorithms, by simply averaging predictions from multiple bootstrapped predictors. This is called  $bagging\ predictors^3$ .

# Cross-validation

While we typically use a single training, hold-out and test set, it might be useful to do this multiple times in order to obtain more robust performance estimates. In the simplest case, cross-validation can be used to obtain multiple training and hold-out sets from a single dataset. This works by simply partitioning the data in k folds and then using one of the folds as a holdout and the remaining k-1 as training data. This is repeated k times. When k is the same size as the original training data, then the method is called leave-one-out cross-validation.

#### k-fold Cross-Validation

- 1. Input Training data  $D_T$ , number of folds k, algorithm  $\lambda$ , measurement function U
- 2. Create the partition  $D^{(1)} \dots, D^{(k)}$  so that  $\bigcup_{i=1}^k D^{(k)} = D$ .
- 3. Define  $D_T^{(i)} = D \setminus D^{(i)}$
- 4.  $\pi_i = \lambda(D_T^{(i)})$
- 5. **For**  $i = 1, \ldots, k$ :
- 6.  $\pi_i = \lambda(D^{(i)})$
- 7.  $u_i = U(\pi_i)$
- 8. **return**  $\{y_1, \ldots, y_i\}$ .

## Online evaluation

We can get around this problem if we consider online evaluation of learning algorithms. This means that the learning algorithm is always evaluated on previously unseen data. However, when new data is seen, it can be used by the algorithm to learn.

Example 10. Online prediction accuracy

- Adaptive decision rule  $\pi$
- At time t
  - 1.  $\pi$  predicts  $a_t$
  - 2. The true data  $x_t$  is observed and we see whether  $a_t = y_t$ .
  - 3.  $\pi$  adapts to the new data  $x_t$

For this example, you can consider the decision rule  $\pi$  as being conditioned on the previous data, i.e.

$$\pi(a_t \mid x_1, \ldots, x_t)$$

# 4.1.3 Beyond the data you have: simulation and replication

In the end, however, you are always limited by the data you actually have. The more you tweak your models and algorithms to improve performance with your current dataset, the more you are simply engineering a solution that is adapted to the specific data. This may not generalise

well, even if you are using cross-validation or bootstrapping every step of the way. How can you then make sure that your methodology is robust?

The first method is simply to simulate data from an artificial process, where the ground truth (e.g. the labels) is known, and where the dataset size and dimensionality is similar to the one you have. You can use this to see whether the overall process is robust, and this can give you confidence that you will get reasonable results when using real data. The second method requires actually collecting new data, and repeating the study. If the results can be replicated, the original study was not a fluke.

#### Simulation

Simulation can be extremely useful. It allows you to examine the performance of various methods as you change aspect of the data-generating process without ever having to look at the data in detail. Since the data is synthetically generated, you always know the ground truth, so you know precisely how good your methods are going to be. This is useful in particular when you want to perform a null hypothesis test, and want to see under which conditions you actually accept or reject a null hypothesis.

A good example of the use of simulation to validate a method is in the article by Bennett et al. where they discuss the use of corrections for multiple comparison tests. This followed their study of uncorrected methods for fMRI analysis<sup>2</sup>, where they found that commonly used such methods would detect meaningful brain activity in a dead salmon. They use simulation to select a correction method that would be neither too conservative (i.e. not detecting any significant brain activity) nor too sensitive (i.e. detecting activity where there is none).

#### Steps for a simulation pre-study

- 1. Criteria: Define your quality criteria, e.g. a utility function U.
- 2. Algorithms: Define the class of algorithms, models or to consider.
- 3. **Simulation**: Create a simulation that allows you to collect data similar to the real one
- 4. **Protocol**: This defines the processing and model pipeline, as you would do it with the actual data.
- 5. **Pre-study**: Collect data from the simulation and analyse it with the created according to your protocl.
- 6. **Adjustment:** If the results are not as expected, alter the simulation, protocol etc as needed. The aim is to find a *protocol* for the pre-study that would lead to (a) reproducible results (b) allows you to analyse the behaviour of different algorithms under different assumptions.
- 7. **Use:** Apply the protocol to actual data, whether it has been already collected, or will be collected according to the protocol.

# Independent replication

The gold standard for reproducibility is independent replication. Simply have another team try and reproduce the results you obtained, using completely new data. If the replication is successful, then you can be pretty sure there was no flaw in your original analysis. In typical

scientific publishing, the replication study is done by a different set of authors than those of the original study.

#### Replication study

- 1. Reinterpret the original hypothesis and experiment.
- 2. Collect data according to the original protocol, unless flawed. It is possible that the original experimental protocol had flaws. Then the new study should try and address this through an improved data collection process. For example, the original study might not have been double-blind. The new study can replicate the results in a double-blind regime.
- 3. Run the analysis again, *unless flawed*. It is possible that the original analysis had flaws. For example, possible correlations may not have been taken into account.
- 4. See if the conclusions are in agreement.

#### Learning outcomes

#### Understanding

- What is a hold-out set, cross-validation and bootstrapping.
- The idea of not reusing data input to an algorithm to evaluate it.
- The fact that algorithms can be implemented by both humans and machines.

## Skills

- Use git and notebooks to document your work.
- Use hold-out sets or cross-validation to compare parameters/algorithms in Python.
- Use bootstrapping to get estimates of uncertainty in Python.

# Reflection

- What is a good use case for cross-validation over hold-out sets?
- When is it a good idea to use bootstrapping?
- How can we use the above techniques to avoid the false discovery problem?
- Can these techniques fully replace independent replication?

# Chapter 5

# Causality

# 5.1 Introduction

## Headaches and aspirins

Causal questions do not just deal with statistical relationships. The meaning of these questions is slightly different depending on whether we are talking about the population at large, or a specific individual. For populations, the main question is whether or not our actions have a causal effect. In observational data, we also need to consider the *direction of causation*.

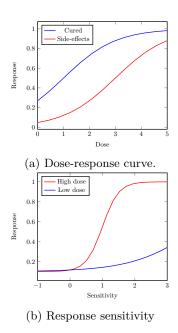


Figure 5.1: Investigating the response of the population to various doses of the drug.

EXAMPLE 11 (Population effects). We can ask ourselves two different questions about the effect of population effect aspirin on headaches.

• Is aspirin an effective cure for headaches?

• Does having a headache lead to aspirin-taking?

To examine the effect of aspirin on the population, we can look at the dose-response curve in Figure 5.1a. We first have to define what it means to be 'cured'. We also need to define possible side-effects. We can then measure how increased doses of aspirin lead to different outcomes. If the experiment has been properly conducted, this dose-response curve suggests that aspirin does have an effect in curing headaches, though it also has some side-effects. However, not all individuals are the same. Let us assume each person has a different *sensitivity* to aspirin, so that some individuals responds better to treatment, especially at high doses, and some do not respond very much at all, as shown in Figure 5.1b

For individuals, the first question is, what is the possible effect of our actions? This is called the *effect of causes*. The second question is, what was the reason for something happening? That is called the *cause of effects?* 



5.1. INTRODUCTION 71

EXAMPLE 12 (Individual effects). We can ask ourselves two different questions about the individual effect of aspirin on headaches.

- Effects of Causes: Will my headache pass  $if\ I\ take$  an aspirin?
- Causes of Effects: Would my headache have passed if I had not taken an aspirin?

In order to be able to meaningfully talk about effects and causes we must also introduce decisions. Formally, there is nothing different in the decisions in this section and those introduced in Section A.1. However, in this case we will try and use decisions to model outside interventions in a "natural" system, whereby a null decision means that we do not intervene.

#### Overview

#### Inferring causal models

We can distinguish different *models* from observational or experimental data.

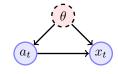
#### Inferring individual effects

The effect of possible intervention on an individual is not generally determinable. We usually require strong assumptions.

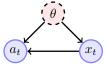
#### Decision-theoretic view

There are many competing approaches to causality. We will remain within the decision-theoretic framework, which allows us to crisply define both our knowledge and assumptions.

# What causes what?



(a) Independence of  $a_t$ .



(b) Independence of  $x_t$ .

Example 13. Suppose we have data  $x_t, a_t$  where

- $x_t$ : lung cancer
- $a_t$ : smoking

Does smoking cause lung cancer or does lung cancer make people smoke? Can we compare the two models above to determine it?

The answer is no. Let us consider two different parametrisations of the distribution. One in which  $a_t$  generates  $x_t$ , and the converse, for any given parameter value  $\theta$ , as given below:

$$P_{\theta}(D) = \prod_{t} P_{\theta}(x_{t}, a_{t}) = \prod_{t} P_{\theta'}(x_{t} \mid a_{t}) P_{\theta'}(a_{t}) = \prod_{t} P_{\theta''}(a_{t} \mid x_{t}) P_{\theta''}(x_{t}).$$

In particular, for any parametrisation  $\theta$  of the joint distribution, there exists a  $\theta'$  and  $\theta''$  giving the same probabilities for all  $x_t, a_t$ . For the example above, we can look at Bernoulli distributions for both variables so that  $P_{\theta}(x_t = x, a_t = a) = \theta_{x,a}$ . Then

$$\theta'_{a} = \sum_{x} \theta_{x,a}, \qquad \theta'_{x|a} = \theta_{x,a}/\theta'_{a}$$

$$\theta''_{x} = \sum_{a} \theta_{x,a}, \qquad \theta'_{a|x} = \theta_{x,a}/\theta''_{x}.$$

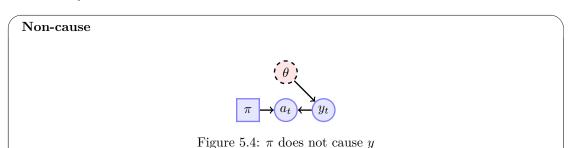
This means we can define prior distributions  $\beta, \beta', \beta''$  in these three parameter spaces that give exactly the same results, e.g. by modelling each parameter as an independent Beta distribution. So, clearly simply looking at a simple graphical model does not let us answer this question.

# 5.1.1 Common structural assumptions

In order to be clear about what constitutes an observation by the experimenter and what is a decision, we must clearly separate random variables from decision variables. The individual actions may be random variables, but they will depend on decisions taken. As we will see later, this is useful for modelling interventions.

# Basic causal structures

Directed graphical models are not sufficient to determine causality by themselves, as they only determine correlations between random variables. If we have decision variables, however, we can always determine whether or not our decisions influence outcomes.



In the diagram above, we see that  $y_t \perp \!\!\! \perp \pi$ .

EXAMPLE 14. Consider the model

$$y_t \sim \mathcal{N}(0, 1)$$

$$a_t \mid y_t, \pi \sim \mathcal{N}(y_t + \pi, 1)$$

5.1. INTRODUCTION 73

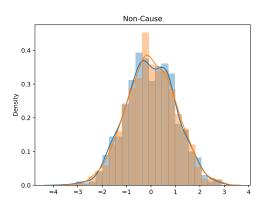


Figure 5.5:  $\mathbb{P}^{\pi}(y_t)$  for  $\pi \in \{-1, 1\}$  when  $\pi$  is not a cause for  $y_t$ 

In this example, we see the  $y_t$  is independent of the policy  $\pi$ . However,  $y_t$  is not independent of the action taken, as the action depends on  $y_t$  directly. The correlation between y, a is shown in Figure 5.8a.

## No confounding

Confounding is a term that indicates the existence of latent variables that create dependencies between  $y_t, \pi, a_t$ . We are sure that there is no confounding whenever  $y_t \perp \!\!\! \perp \pi \mid a_t$ , as captured by the diagram in Figure 5.6. In this case  $\pi$  is a direct cause for  $y_t$  through  $a_t$ .

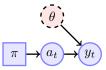


Figure 5.6: No confounding:  $\pi$  causes  $y_t$ 

EXAMPLE 15. Consider the model

$$a_t \sim \mathcal{N}(\pi, 1)$$
$$y_t \mid a_t, \pi \sim \mathcal{N}(a_t, 1)$$

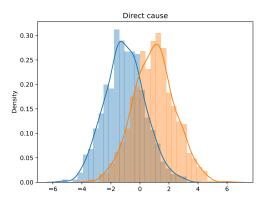


Figure 5.7:  $\mathbb{P}^{\pi}(y_t)$  for  $\pi \in \{-1,1\}$  when  $\pi$  is a direct cause for  $y_t$ 

We can see how the distribution of  $y_t$  changes when  $\pi$  changes in Figure 5.7. In this case there is also a correlation between  $a_t, y_t$  as seen in Figure 5.8.

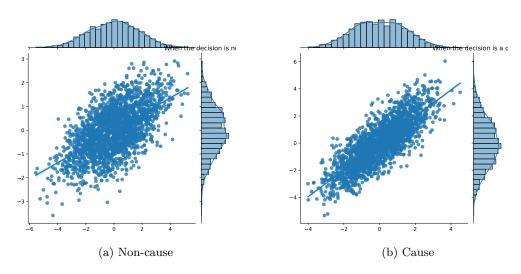


Figure 5.8: Correlation between  $a_t$  and  $y_t$ 

## Covariates

## Sufficient covariate

Sometimes the variable of interest is not conditionally independent of the treatment, unless there exists a *sufficient covariate*  $x_t$  such that  $y_t \perp \!\!\! \perp \pi \mid a_t, x_t$ . If  $x_t$  is not observed, then it is sometimes called a confounder.

5.1. INTRODUCTION 75

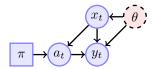


Figure 5.9: Sufficient covariate  $x_t$ 

Example 16. Consider the model

$$x_t \sim \mathcal{N}(0, 1)$$

$$a_t \sim \mathcal{N}(x_t + \pi, 1)$$

$$y_t \mid a_t, \pi \sim \mathcal{N}(x_t + a_t, 1),$$

Here  $x_t$  influences the outcome  $y_t$ , but also directly influences  $a_t$  through the policy  $\pi$ . As we can see in Figure 5.7, the policy then has an influence on  $y_t$ 

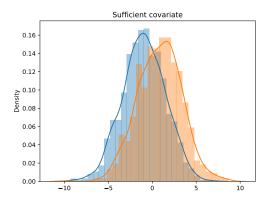


Figure 5.10:  $\mathbb{P}^{\pi}(y_t)$  for  $\pi \in \{-1,1\}$  when  $\pi$  is a direct cause for  $y_t$ 

**Definition 5.1.1** (non-confounder). A latent variable  $x_t$  is not a confounder if either  $x_t \perp \!\!\! \perp a_t \mid \pi$  or  $x_t \perp \!\!\! \perp y_t \mid a_t$ 

#### Instrumental variables and confounders

If the sufficient covariate  $x_t$  is not observed, we may still have another variable available,  $z_t$ , on the basis of which we make our decisions. This is called an *instrumental variable*. More formally,  $z_t$  is an instrumental variable with respect to  $y_t$  since the latent variables on which  $y_t$  depends satisfy  $x_t$ ,  $\theta \perp \!\!\! \perp \pi$  and the outcome variable is independent of the policy given the action and the two latent variables:  $y \perp \!\!\! \perp \pi \mid a_t, x_t, \theta$ .

In this case  $z_t$  and  $x_t$  are dependent, but the effect of the treatment depends on  $x_t$  directly. As  $x_t$  is a latent covariate, it can be called a *confounder*.

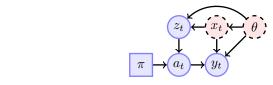


Figure 5.11: Instrumental variable  $z_t$ , confounder  $x_t$ 

EXAMPLE 17. Consider the model

$$x_t \sim \mathcal{N}(0, 1)$$

$$z_t \sim \mathcal{N}(x_t, 1)$$

$$a_t \sim \mathcal{N}(z_t + \pi, 1)$$

$$y_t \mid a_t, \pi \sim \mathcal{N}(x_t + a_t, 1)$$

In this scenario,  $x_t$  directly influences the outcome  $y_t$ , but is not observed.

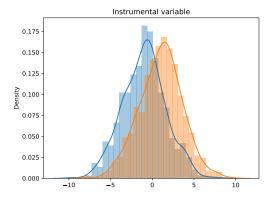


Figure 5.12:  $\mathbb{P}^{\pi}(y_t)$  for  $\pi \in \{-1,1\}$  when  $\pi$  is a direct cause for  $y_t$ 

Finally, variables which are neither confounders, nor instrumental, are called nuisance variables. Typically these include the latent variables of the problem, and any other variable that is marginalised out.

## 5.2 Interventions

Interventions are of primary interest when we have a set of observational data, collected under a *null* or *default* policy  $\pi_0$ . We then wish to intervene with some policy  $\pi$  in order to maximise our utility function, or to simply try and estimate the exact relationships between variables.

## Modelling interventions

• Observational data D. This represents data we have collected from some previous regime. In order to be able to model the problem precisely, we must posit the existence of some default policy  $\pi_0$  under which the data was collected.

<sup>&</sup>lt;sup>1</sup>Hence, it can be called a confounder.

5.2. INTERVENTIONS

• Policy space  $\Pi$ . This must include  $\pi_0$ , as well as any other policies that the decision maker may be able to choose in the future.

77

#### Default policy

The space of policies  $\Pi$  includes a default policy  $\pi_0$ , under which the data was collected. The policy  $\pi_0$  might already be known, if for example the data was collected with a specific algorithm. However, frequently  $\pi_0$  is not given, and must also be inferred from the data. In that case, it can be seen as an additional parameter, complementary to  $\theta$ .

### Intervention policies

Except  $\pi_0$ , policies  $\pi \in \Pi$  represent different interventions specifying a distribution  $\pi(a_t \mid x_t)$ .

- Direct interventions. The simplest scenario is when we are able to choose a  $\pi$  for which we know  $\pi(a_t \mid x_t)$ . This counts as a direct intervention, as we can specify any distribution of actions allowed in  $\Pi$ . If  $\Pi$  includes all conditional distributions, we can select an arbitrary action for every individual. This assumption is plausible for algorithmic decision making such as recommendation systems, but implausible when the actions are taken by another agent, such as a human.
- Indirect interventions and non-compliance. In this scenario we assume that, while we are free to choose policies from  $\Pi$ , we do not know what distribution  $\pi(a_t \mid x_t)$  each policy specifies. In algorithmic decision making, this occurs whenever  $\pi$  represents hyperparameters and algorithms for which we have insufficient information. Then policies must be evaluated through some type of black-box (e.g. A/B) testing. When the actions are taken by (human) agents, the policy implies making a recommendation, which may not be followed by the agent. If we denote the recommendation by  $v_t$ , then we can decompose the actual distribution of actions as follows

$$\pi(a_t \mid x_t) = \sum_{z_t} \pi(a_t \mid z_t, x_t) \pi(z_t \mid x_t).$$

In this scenario we can freely specify the recommendation policy  $\pi(v_t \mid x_t)$ , but the response to our recommendations  $\pi(a_t \mid v_t, x_t)$  must be estimated. For that reason, it is usually simpler to simply marginalise out  $a_t$ . Another approach is to consider non-compliance as a confounder  $x_t$ , and our recommendation  $v_t$  as an instrumental variable.

EXAMPLE 18 (Weight loss). Consider weight loss. We can collect observational data from a population of overweight adults over a year. We can imagine that x represents the weight and vital statistics of an individual and y their change in weight after a year. We may also observe their individual actions a, such as whether or not they are following a particular diet or exercise regime. Under the default policy  $\pi_0$ , their actions are determined only the individuals. Consider an alternative policy  $\pi$ , which prescribes diet and exercise regimes. Due to non-compliance, actual actions taken by individuals may differ from prescribed actions. In addition, actions might not be observed.

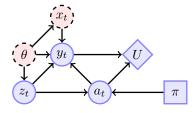


Figure 5.13: Model of non-compliance as a confounder.

## 5.3 Policy evaluation and optimisation

#### The value of an observed policy

In this section, we will focus on the general model of Figure C.7. If we have data  $D = \{(x_t, a_t, y_t) \mid t \in [T]\}$  generated from some policy  $\pi_0$ , we can always infer the average quality of each action a under that policy.

$$\hat{\mathbb{E}}_D(U \mid a) \triangleq \frac{1}{\mid \{t \mid a_t = a\} \mid \sum_{t: a_t = a} U(a_t, y_t)}$$
(5.3.1)

$$\approx \mathbb{E}_{\theta}^{\pi_0}(U \mid a) \qquad (a_t, y_t) \sim \mathbb{P}_{\theta}^{\pi_0}. \qquad (5.3.2)$$

Can we calculate the value of another policy? As we have seen from Simpson's paradox, it is folly to simply select

$$\hat{a}_D^* \in \arg\max_{a} \hat{\mathbb{E}}_D(U \mid a),$$

as the action also depends on the observations x through the policy. To clarify this, let us look again at the model shown in Figure C.7.

$$x_t \mid \theta \sim P_{\theta}(x)$$

$$y_t \mid \theta, x_t, a_t \sim P_{\theta}(y \mid x_t, a_t)$$

$$a_t \mid x_t, \pi \sim \pi(a \mid x_t).$$

Assume that  $x \in \mathcal{X}$ , a continuous space, but  $y \in \mathcal{Y}$  is discrete. In this scenario, then the value of an action under a policy  $\pi$  is nonsensical, as it does not really depend on the policy itself:

$$\mathbb{E}_{\theta}^{\pi}(U \mid a) = \int_{\mathcal{X}} dP_{\theta}(x) \sum_{y \in \mathcal{Y}} P_{\theta}(y \mid x, a) U(a, y).$$

We see that there is a clear dependence on the distribution of x, and there is no dependence on the policy any more. In fact, equation above only tells us the expected utility we'd get if we always chose the same action a. But what is the optimal policy? First, we have to define the value of a policy.

The value of a policy

$$\mathbb{E}_{\theta}^{\pi}(U) = \int_{\mathcal{X}} dP_{\theta}(x) \sum_{a \in \mathcal{A}} \pi(a \mid x) \sum_{y \in \mathcal{Y}} P_{\theta}(y \mid x, a) U(a, y)$$

The optimal policy under a known parameter  $\theta$  is given simply by

$$\max_{\pi \in \Pi} \mathbb{E}^{\pi}_{\theta}(U),$$

where  $\Pi$  is the set of allowed policies.

#### 5.3.1 Monte-Carlo estimation

#### Monte-Carlo estimation

The simplest method to estimate the value of an alternative policy is to use Monte-Carlo estimation and importance sampling. However, this estimate can have a very high variance if the alternative policy is very different from the original policy.

#### Importance sampling

We can obtain an unbiased estimate of the utility in a model-free manner through importance sampling:

$$\mathbb{E}_{\theta}^{\pi}(U) = \int_{\mathcal{X}} dP_{\theta}(x) \sum_{a} \mathbb{E}_{\theta}(U \mid a, x) \pi(a \mid x)$$

$$\approx \frac{1}{T} \sum_{a} \sum_{t} \mathbb{E}_{\theta}(U \mid a, x_{t}) \pi(a \mid x_{t}), \qquad x_{t} \sim P_{\theta}(x)$$

$$= \frac{1}{T} \sum_{t} \sum_{a} \mathbb{E}_{\theta}(U \mid a, x_{t}) \frac{\pi(a \mid x_{t})}{\pi_{0}(a \mid x_{t})} \pi_{0}(a \mid x_{t})$$

$$\approx \frac{1}{T} \sum_{t=1}^{T} U_{t} \frac{\pi(a_{t} \mid x_{t})}{\pi_{0}(a_{t} \mid x_{t})}, \qquad a_{t} \mid x_{t} \sim \pi_{0}, \quad U_{t} \mid x_{t}, a_{t} \sim P_{\theta}(U_{t} \mid x_{t}, a_{t})$$

where the last step involves a single sample of  $a_t, U_t$ .

## Bayesian estimation

Unfortunately this method has high variance. If we  $\pi_0$  is given, we can calculate the utility of any policy to whatever degree of accuracy we wish. We begin with a prior  $\beta$  on  $\Theta$  and obtain the following, assuming the policy  $\pi_0$  is stationary.

$$\beta(\theta \mid D, \pi_0) \propto \prod_t \mathbb{P}_{\theta}^{\pi_0}(x_t, y_t, a_t)$$

$$\mathbb{E}_{\beta}^{\pi}(U \mid D) = \int_{\Theta} \mathbb{E}_{\theta}^{\pi}(U) \, \mathrm{d}\beta(\theta \mid D)$$

$$= \int_{\Theta} \int_{\mathcal{X}} \, \mathrm{d}P_{\theta}(x) \sum_{t=1}^{T} \sum_{a} \mathbb{E}_{\theta}(U \mid a, x) \pi(a \mid x) \, \mathrm{d}\beta(\theta \mid D).$$

## Causal inference and policy optimisation

Causal inference requires building a complete model for the effect of both the model parameter  $\theta$ , representing nature, and the policy  $\pi$ , representing the decision maker. This means that we have to be explicit about the dependencies of random variables on the model and the policy.

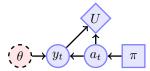


Figure 5.14: Simple decision problem.

EXAMPLE 19. Let  $a_t, y_t \in \{0, 1\}, \theta \in [0, 1]^2$  and

$$y_t \mid a_t = a \sim \mathcal{B}ernoulli(\theta_a)$$

Then, by estimating  $\theta$ , we can predict the effect of any action. How can we estimate  $\theta$  from historical data? We simply have to select the right parameter value. Simply put, each choice of a corresponds to one part of the parameter vector. This means that the maximum likelihood estimate

$$\hat{\theta}_a \triangleq \frac{1}{|\{t \mid a_t = a\}|} \sum_{\{t \mid a_t = a\}} y_t$$

is valid. We can also consider a product-Beta prior  $\mathcal{B}eta(\alpha_a^0, \beta_b^0)$  for each one of the Bernoulli parameters, so that the posterior after t observations is

$$\alpha_a^t = \alpha_a^0 + \sum_{\{t \mid a_t = a\}} y_t, \qquad \beta_a^t = \beta_a^0 + \sum_{\{t \mid a_t = a\}} (1 - y_t).$$

How can we optimise the policy? Let us parametrise our policy with  $\pi(a_t = 1) = w$ .

For a fixed  $\theta$ , the value of the policy is

$$\mathbb{E}_{\theta}^{\pi} U = \theta_1 w + \theta_0 (1 - w)$$

The gradient with respect to w is

$$\nabla \mathbb{E}_{\theta}^{\pi} U = \theta_1 - \theta_0,$$

so we can use the update

$$w^{(n+1)} = w^{(n)} + \delta^{(n)}\theta_1 - \theta_0.$$

However,  $w \in [0, 1]$ , which means our optimisation must be constrained. Then we obtain that w = 1 if  $\theta_1 > \theta_0$  and 0 otherwise.

When  $\theta$  is not known, we can use stochastic gradient descent.

$$\nabla \mathbb{E}_{\beta}^{\pi} U = \int_{\Theta} [\nabla \mathbb{E}_{\theta}^{\pi} U] \, \mathrm{d}\beta(\theta)$$

to obtain:

$$w^{(n+1)} = w^{(n)} + \delta^{(n)}\theta_1^{(n)} - \theta_0^{(n)}.$$

where  $\theta^{(n)} \sim \beta$ .

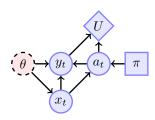


Figure 5.15: Decision problem with covariates.

Example 20. Let  $a_t,x_t=\{0,1\},\,y_t\in\mathbb{R},\,\theta\in\mathbb{R}^4$  and  $y_t\mid a_t=a,x_t=x\sim \text{Bernoulli}(\theta_{a,x})$ 

Then, by estimating  $\theta$ , we can predict the effect of any action.

## 5.4 Individual effects and counterfactuals

Counterfactual analysis is mainly about questions relative to individuals, and specifically about what the effects of alternative actions would have been in specific instances in the past. We will assume a decision-theoretic viewpoint throughout, in order to be as clear as possible and avoiding imprecise language.

## 5.4.1 Disturbances and structural equation models

A structural equation model describes the random variables as deterministic functions of the decisions variables and the random exogenous disturbances. This allows us to separate the unobserved randomness from the known functional relationship between the other variables. Structurally, the model is essentially a variant of decision diagrams, as shown in Figure 5.16.

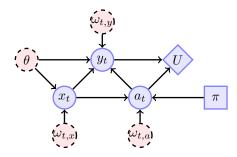


Figure 5.16: Decision diagram with exogenous disturbances  $\omega$ .

We still need to specify particular functional relationships between the variables. Generally speaking, a random variable taking values in  $\mathcal{X}$ , is simply a function  $\Omega \times \Theta \to \mathcal{X}$ . For example, in Figure 5.16  $y_t = f_y(\omega, \theta)$ . Taking into account the dependencies, this can be rewritten in terms of a function of the other random variables, and the local disturbance:  $y_t = f_{y|a,x}(a,x,\omega_{t,y},\theta)$ . The choice of the function, together with the distribution of the parameter  $\theta$  and the disturbances  $\omega$ , fully determines our model.

Example 21 (Structural equation model for Figure 5.16). In structural equation models, the only random variables are the exogenous disturbances. In a fully Bayesian framework,  $\theta$  is also a latent random variable. The remaining variables are deterministic functions.

$$\begin{split} \theta &\sim \mathcal{N}(\mathbf{0}_4, \mathbf{I}_4), \\ x_t &= \theta_0 \omega_{t,x}, \\ y_t &= \theta_1 + \theta_2 x_t + \theta_3 a_t + \omega_{t,y}, \\ a_t &= \pi(x_t) + \omega_{t,a} \mod |\mathcal{A}| \end{split} \qquad \begin{aligned} \omega_{t,x} &\sim \textit{Bernoulli}(0.5) \\ \omega_{t,y} &\sim \mathcal{N}(0,1) \\ \omega_{t,a} &\sim 0.1 \, \textit{D}(0) + 0.9 \, \textit{Unif}(\mathcal{A}), \end{aligned}$$

Structural equation models are particularly interesting in applications such as economics, where there are postulated relations between various economic quantities. If relationships between variables satisfy nice properties, then we can perform counterfactual inferences of the type

: "What if I had not taken an aspirin?" In the example above, if we can infer the noise variables  $\omega$ , we can change the value of some choice variables, i.e.  $a_t$  and see the effect on other variables like  $y_t$  directly.

#### Treatment-unit additivity

An example for a particular assumption for structural equation models is treatment-unit additivity. This states that the outcome depends on the action through a deterministic transformation and some additive noise that is only applied to the specific individual. More specifically, if individual t obtains treatment  $a_t$ , then the outcome only depends on  $a_t$  and individual-specific noise  $\omega_{t,y}$ . The assumption is given more formally below.

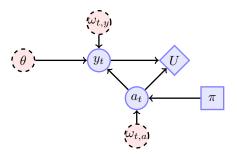


Figure 5.17: Decision diagram for treatment-unit additivity

**Assumption 5.4.1** (TUA). For any given treatment  $a \in \mathcal{A}$ , the response variable satisfies

$$y_t = g(a_t) + \omega_{t,y}$$

This implies that  $\mathbb{E}[y_t \mid a_t = a] = g(a_t) + \mathbb{E}(\omega_{t,y})$ . This implies that we might be able to easily estimate g, for example if the noise is zero-mean. The assumption makes sense for a lot of cases where we do not expect the outcomes of different individuals to be correlated through some confounding variable.

## 5.4.2 Example: Learning instrumental variables

This example is adapted from Hartford et al. <sup>12</sup>, who use a deep learning to infer causal effects in the presence of instrumental variables. They break down the problem in two prediction tasks: the first is treatment prediction, and the second conditional treatment distribution. Unfortunately they do not use a decision-theoretic framework and so the difference between actual prices and policy changes is unclear.

EXAMPLE 22 (Pricing model). In the following pricing model, we wish to understand how sales are affected by different pricing policies. In this example, there is a variable  $z_t$  which is an instrument, as it varies for reasons that are independent of demand and only affects sales through ticket prices.

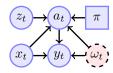


Figure 5.18: Graph of structural equation model for airport pricing policy  $\pi$ :  $a_t$  is the actual price,  $z_t$  are fuel costs,  $x_t$  is the customer type,  $y_t$  is the amount of sales,  $\omega_t$  is whether there is a conference. The dependency on  $\theta$  is omitted for clarity.

There are a number of assumptions we can make on the instrument  $z_t$ .

**Assumption 5.4.2** (Relevance).  $a_t$  depends on  $z_t$ .

In our example, it also depends on  $x_t$ .

**Assumption 5.4.3** (Exclusion).  $z_t \perp \!\!\!\perp y_t \mid x_t, a_t, \omega_t$ .

In other words, the outcome does not depend on the instrument directly. This was also satisfied in our first example of an instrumental variable.

**Assumption 5.4.4** (Unconfounded instrument).  $z_t \perp \!\!\! \perp \omega_t \mid x_t$ .

#### Prediction tasks

We can use the following structural equation model

$$y_t = g_\theta(a_t, x_t) + \omega_t, \qquad \mathbb{E}_\theta \, \omega_t = 0, \qquad \forall \theta \in \Theta$$
 (5.4.1)

There are two slightly different prediction tasks we can think of in this model.

#### Standard prediction

In standard prediction tasks, we just want to estimate the distribution of sales given the characteristics and price. Since the actions are correlated with the outcome through the confounder, this estimate is biased.

$$\mathbb{P}^{\pi}_{\theta}(y_t \mid x_t, a_t), \qquad \mathbb{E}^{\pi}_{\theta}(y_t \mid x_t, a_t) = g_{\theta}(x_t, a_t) + \mathbb{E}^{\pi}_{\theta}(\omega_t \mid x_t, a_t).$$

## Counterfactual prediction

$$\mathbb{E}_{\theta}^{\pi}(y_t \mid x_t, z_t) = \int_{\mathcal{A}} \underbrace{\left[g(a_t \mid x_t, z_t) + \mathbb{E}_{\theta}(\omega \mid x_t)\right]}_{h(a_t, x_t)} d\pi(a_t \mid x_t)$$

#### 5.4.3 Discussion

## Further reading

- Pearl, Causality.
- Dawid 4's decision-theoretic causality
- Halpern's book https://www.cs.cornell.edu/home/halpern/papers/causalitybook-ch1-3. html

• Halpern's video lecture: https://www.cs.cornell.edu/home/halpern/papers/causalitybook-ch1-3 html

#### 5.4.4 Exercises

In the following exercises, we are taking actions  $a_t$  and obtaining outcomes  $y_t$ . Our utility function is simply  $U = y_t$ .

EXERCISE 8. Let us have some data generated from a null treatment policy  $\pi_0$  of the form  $(a_t, y_t)$ . There is a simple model that explains the data of the form

$$y_t \mid a_t = a, \theta \sim \mathcal{N}(a + \theta, 1),$$

where the actions are distribution according to  $\pi(a_t)$ .

- Assume that  $\pi_0 \in [0, 1]$  is given and it is  $a_t \mid \pi = \pi_0 \sim \mathcal{B}emoulli(\pi_0)$ . First, estimate  $\theta$ . Then, calculate the distribution of  $y_t \mid \pi_0, \theta$  for any other policy and plot the resulting mean and variance as  $\pi$  changes. You can do this first in a maximum-likelihood manner. Advanced: estimate the posterior distribution of  $\theta$  for a normal prior on  $\theta$ .
- Now assume that  $\pi_0$  is not given. This means that you must also estimate  $\pi_0$  itself before estimating the effect of any other policy  $\pi$  on the data.
- In this exercise, can you learn about other actions when you are not taking them? Why?

EXERCISE 9. Let us have some data generated from a null treatment policy  $\pi_0$  of the form  $(a_t, y_t)$ . Let us now consider a slightly model where  $\theta \in \mathbb{R}^2$ .

$$y_t \mid a_t = a, \theta \sim \mathcal{N}(\theta_a, 1),$$

where the actions are distribution according to  $\pi(a_t)$ .

- Assume that  $\pi_0 \in [0, 1]$  is given and it is  $a_t \mid \pi = \pi_0 \sim \mathcal{B}emoulli(\pi_0)$ . First, estimate  $\theta$ . Then, calculate the distribution of  $y_t \mid \pi_0, \theta$  for any other policy and plot the resulting mean and variance as  $\pi$  changes. You can do this first in a maximum-likelihood manner. Advanced: estimate the posterior distribution of  $\theta$  for a normal prior on  $\theta$ .
- Now assume that  $\pi_0$  is not given. This means that you must also estimate  $\pi_0$  itself before estimating the effect of any other policy  $\pi$  on the data.
- In this exercise, can you learn about other actions when you are not taking them? Why?

EXERCISE 10. Given your estimates, find the optimal policy for each one of those cases. Measure the quality of this policy on

- The actual data you have already (e.g. using importance sampling)
- On new simulations (using the testing framework).

Advanced: The optimal policy when  $\theta$  is known is to always take the same action. Does that still hold when  $\theta$  is not known and you are estimating it all the time from new data?

EXERCISE 11 (Advanced). Let us have some data generated from a null treatment policy  $\pi_0$  of the form  $(x_t, a_t, y_t)$ , with  $a_t, x_t \in \{0, 1\}$ .

$$y_t \mid a_t = a, x_t = x, \theta \sim \mathcal{N}(\theta_{a,x}, 1),$$

where the actions are distribution according to  $\pi_0(a_t \mid x_t)$ .

- Assume that  $\pi_0$  is given and it is  $a_t \mid x_t = x, \pi = \pi_0 \sim \mathcal{B}ernoulli(w_x)$ . First, estimate  $\theta$ . Repeat your analysis.
- Now assume that  $\pi_0$  is not given. Again, repeat your analysis.
- Is there now globally better action  $a_t$ ? Should it depend on  $x_t$ , like in the observed policy? Can you estimate the optimal policy?

# Appendix A

# Decision and Learning Theory

This chapter deals with simple decision problems, whereby a decision maker (DM) makes a simple choice among many. In some of this problems the DM has to make a decision after first observing some side-information. Then the DM uses a decision rule to assign a probability to each possible decision for each possible side-information. However, designing the decision rule is not trivial, as it relies on previously collected data. A higher-level decision includes choosing the decision rule itself. The problems of classification and regression fall within this framework. While most steps in the process can be automated and formalised, a lot of decisions are actual design choices made by humans. This creates scope for errors and misinterpretation of results.

In this chapter, we shall formalise all these simple decision problems from the point of view of statistical decision theory. The first question is, given a real world application, what type of decision problem does it map to? Then, what kind of machine learning algorithms can we use to solve it? What are the underlying assumptions and how valid are our conclusions?

## A.1 Hierarchies of decision making problems

All machine learning problems are essentially decision problems. This essentially means replacing some human decisions with machine decisions. One of the simplest decision problems is classification, where you want an algorithm to decide the correct class of some data, but even within this simple framework there is a multitude of decisions to be made. The first is how to frame the classification problem the first place. The second is how to collect, process and annotate the data. The third is choosing the type of classification model to use. The fourth is how to use the collected data to find an optimal classifier within the selected type. After all this has been done, there is the problem of classifying new data. In this course, we will take a holistic view of the problem, and consider each problem in turn, starting from the lowest level and working our way up.

## A.1.1 Simple decision problems

#### Preferences

The simplest decision problem involves selecting one item from a set of choices, such as in the following examples

#### Food

Example 23. A McDonald's cheeseburger

- B Surstromming
- C Oatmeal

## Money

- A 10,000,000 CHF
- B 10,000,000 USD
- C 1,000 BTC

#### Entertainment

- A 1 Week in Venice
- B 1 Week Hoch-Tour in the Alps
- C 1 Week Transatlantic Cruise

#### Rewards and utilities

In the decision theoretic framework, the things we receive are called rewards, and we assign a utility value to each one of them, showing which one we prefer.

- Each choice is called a reward  $r \in \mathcal{R}$ .
- There is a utility function  $U: \mathcal{R} \to \mathbb{R}$ , assigning values to reward.
- We (weakly) prefer A to B iff  $U(A) \ge U(B)$ .

In each case, given U the choice between each reward is trivial. We just select the reward:

$$r^* \in \operatorname*{arg\,max}_r U(r)$$

The main difficult is actually selecting the appropriate utility function. In a behavioural context, we simply assume that humans act with respect to a specific utility function. However, figuring out this function from behavioural data is non trivial. ven when this assumption is correct, individuals do not have a common utility function.

EXERCISE 12. From your individual preferences, derive a *common utility function* that reflects everybody's preferences in the class for each of the three examples. Is there a simple algorithm for deciding this? Would you consider the outcome fair?

#### Preferences among random outcomes

Example 24. Would you rather ...

- A Have 100 EUR now?
- B Flip a coin, and get 200 EUR if it comes heads?

#### The expected utility hypothesis

Rational decision makers prefer choice A to B if

$$\mathbb{E}(U|A) \ge \mathbb{E}(U|B),$$

where the expected utility is

$$\mathbb{E}(U|A) = \sum_{r} U(r) \, \mathbb{P}(r|A).$$

In the above example,  $r \in \{0, 100, 200\}$  and U(r) is increasing, and the coin is fair.

## Risk and monetary rewards

When  $r \in \mathbb{R}$ , as in the case of monetary rewards, we can use the shape of the utility function determines the amount of risk aversion. In particular:

- If *U* is convex, we are risk-seeking. In the example above, we would prefer B to A, as the expected utility of B would be higher than A, even though they give the same amount of money on average.
- If *U* is linear, we are risk neutral. In the example above, we would be indifferent between *A* and *B*, as the expected amount of money is the same as the amount of money we get.
- If U is concave, we are risk-averse. Hence, in the example above, we prefer A.

This idea of risk can be used with any other utility function. We can simply replace the original utility function U with a monotonic function f(U) to achieve risk-sensitive behaviour. However, this is not the only risk-sensitive approach possible.

## Uncertain rewards

However, in real life, there are many cases where we can only choose between uncertain outcomes. The simplest example are lottery tickets, where rewards are essentially random. However, in many cases the rewards are not really random, but simply uncertain. In those cases it is useful to represent our uncertainty with probabilities as well, even though there is nothing really random.

- Decisions  $a \in \mathcal{A}$
- Each choice is called a reward  $r \in \mathcal{R}$ .
- There is a utility function  $U: \mathcal{R} \to \mathbb{R}$ , assigning values to reward.
- We (weakly) prefer A to B iff  $U(A) \ge U(B)$ .

EXAMPLE 25. You are going to work, and it might rain. What do you do?

- $a_1$ : Take the umbrella.
- a2: Risk it!
- $\omega_1$ : rain
- $\omega_2$ : dry

$\rho(\omega, a)$	$a_1$	$a_2$
$\omega_1$	dry, carrying umbrella	wet
$\omega_2$	dry, carrying umbrella	dry
$U[\rho(\omega,a)]$	$a_1$	$a_2$
$\omega_1$	0	-10
$\omega_2$	0	1

Table A.1: Rewards and utilities.

$ ho(\omega,a)$	$a_1$	$a_2$
$\omega_1$	dry, carrying umbrella	wet
$\omega_2$	dry, carrying umbrella	dry
$U[\rho(\omega,a)]$	$a_1$	$a_2$
$\omega_1$	0	-10
$\omega_2$	0	1
$\mathbb{E}(U \mid a)$	0	-1.2

Table A.2: Rewards, utilities, expected utility for 20% probability of rain.

- $\max_a \min_{\omega} U = 0$
- $\min_{\omega} \max_{a} U = 0$

#### Expected utility

$$\mathbb{E}(U \mid a) = \sum_{r} U[\rho(\omega, a)] \, \mathbb{P}(\omega \mid a)$$

EXAMPLE 26. 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?

- $a_1$ : Take the umbrella.
- $a_2$ : Risk it!

## A.1.2 Decision rules

We now move from simple decisions to decisions that depend on some observation. We shall start with a simple problem in applied meteorology. Then we will discuss hypothesis testing as a decision making problem. Finally, we will go through an exercise in Bayesian methods for classification.

#### Bayes decision rules

Consider the case where outcomes are independent of decisions:

$$U(\beta,a) \triangleq \sum_{\mu} U(\mu,a)\beta(\mu)$$

This corresponds e.g. to the case where  $\beta(\mu)$  is the belief about an unknown world.

**Definition A.1.1** (Bayes utility). The maximising decision for  $\beta$  has an expected utility equal to:

$$U^*(\beta) \triangleq \max_{a \in \mathcal{A}} U(\beta, a). \tag{A.1.1}$$

#### The *n*-meteorologists problem

Of course, we may not always just be interested in classification performance in terms of predicting the most likely class. It strongly depends on the problem we are actually wanting to solve. In biometric authentication, for example, we want to guard against the unlikely event that an impostor will successfully be authenticated. Even if the decision rule that always says 'OK' has the lowest classification error in practice, the expected cost of impostors means that the optimal decision rule must sometimes say 'Failed' even if this leads to false rejections sometimes.

EXERCISE 13. Assume you have n meteorologists. At each day t, each meteorologist i gives a probability  $p_{t,\mu_i} \triangleq P_{\mu_i}(x_t = \text{rain})$  for rain. Consider the case of there being three meteorologists, and each one making the following prediction for the coming week. Start with a uniform prior  $\beta(\mu) = 1/3$  for each model.

		Т					
CNN	0.5	0.6	0.7	0.9	0.5	0.3	0.1
SMHI	0.3	0.7	0.8	0.9	0.5	0.2	0.1
CNN SMHI YR	0.6	0.9	0.8	0.5	0.4	0.1	0.1
Rain?	Y	Y	Y	N	Y	N	N

Table A.3: Predictions by three different entities for the probability of rain on a particular day, along with whether or not it actually rained.

- 1. What is your belief about the quality of each meteorologist after each day?
- 2. What is your belief about the probability of rain each day?

$$P_{\beta}(x_t = \text{rain} \mid x_1, x_2, \dots x_{t-1}) = \sum_{\mu \in \mathcal{M}} P_{\mu}(x_t = \text{rain} \mid x_1, x_2, \dots x_{t-1}) \beta(\mu \mid x_1, x_2, \dots x_{t-1})$$

3. Assume you can decide whether or not to go running each day. If you go running and it does not rain, your utility is 1. If it rains, it's -10. If you don't go running, your utility is 0. What is the decision maximising utility in expectation (with respect to the posterior) each day?

## A.1.3 Statistical testing\*

A common type of decision problem is a statistical test. This arises when we have a set of possible candidate models  $\mathcal{M}$  and we need to be able to decide which model to select after we see the evidence. Many times, there is only one model under consideration,  $\mu_0$ , the so-called *null hypothesis*. Then, our only decision is whether or not to accept or reject this hypothesis.

## Simple hypothesis testing

Let us start with the simple case of needing to compare two models.

The simple hypothesis test as a decision problem

- $\mathcal{M} = \{\mu_0, \mu_1\}$
- $a_0$ : Accept model  $\mu_0$
- $a_1$ : Accept model  $\mu_1$

$$\begin{array}{c|ccc} U & \mu_0 & \mu_1 \\ \hline a_0 & 1 & 0 \\ a_1 & 0 & 1 \\ \end{array}$$

Table A.4: Example utility function for simple hypothesis tests.

There is no reason for us to be restricted to this utility function. As it is diagonal, it effectively treats both types of errors in the same way.

Example 27 (Continuation of the medium example). •  $\mu_1$ : that John is a medium.

•  $\mu_0$ : that John is not a medium.

Let  $x_t$  be 0 if John makes an incorrect prediction at time t and  $x_t = 1$  if he makes a correct prediction. Let us once more assume a Bernoulli model, so that John's claim that he can predict our tosses perfectly means that for a sequence of tosses  $x = x_1, \ldots, x_n$ ,

$$P_{\mu_1}(\mathbf{x}) = \begin{cases} 1, & x_t = 1 \forall t \in [n] \\ 0, & \exists t \in [n] : x_t = 0. \end{cases}$$

That is, the probability of perfectly correct predictions is 1, and that of one or more incorrect prediction is 0. For the other model, we can assume that all draws are independently and identically distributed from a fair coin. Consequently, no matter what John's predictions are, we have that:

$$P_{\mu_0}(\boldsymbol{x} = 1 \dots 1) = 2^{-n}.$$

So, for the given example, as stated, we have the following facts:

- If John makes one or more mistakes, then  $\mathbb{P}(\boldsymbol{x} \mid \mu_1) = 0$  and  $\mathbb{P}(\boldsymbol{x} \mid \mu_0) = 2^{-n}$ . Thus, we should perhaps say that then John is not a medium
- If John makes no mistakes at all, then

$$\mathbb{P}(\boldsymbol{x} = 1, \dots, 1 \mid \mu_1) = 1,$$
  $\mathbb{P}(\boldsymbol{x} = 1, \dots, 1 \mid \mu_0) = 2^{-n}.$  (A.1.2)

Now we can calculate the posterior distribution, which is

$$\beta(\mu_1 \mid \boldsymbol{x} = 1, \dots, 1) = \frac{1 \times \beta(\mu_1)}{1 \times \beta(\mu_1) + 2^{-n}(1 - \beta(\mu_1))}.$$

Our expected utility for taking action  $a_0$  is actually

$$\mathbb{E}_{\beta}(U \mid a_0) = 1 \times \beta(\mu_0 \mid \boldsymbol{x}) + 0 \times \beta(\mu_1 \mid \boldsymbol{x}), \qquad \mathbb{E}_{\beta}(U \mid a_1) = 0 \times \beta(\mu_0 \mid \boldsymbol{x}) + 1 \times \beta(\mu_1 \mid \boldsymbol{x})$$

#### Null hypothesis test

Many times, there is only one model under consideration,  $\mu_0$ , the so-called *null hypothesis*. This happens when, for example, we have no simple way of defining an appropriate alternative. Consider the example of the medium: How should we expect a medium to predict? Then, our only decision is whether or not to accept or reject this hypothesis.

## The null hypothesis test as a decision problem

- $a_0$ : Accept model  $\mu_0$
- $a_1$ : Reject model  $\mu_0$

Example 28. Construction of the test for the medium

- $\mu_0$  is simply the  $\mathcal{B}emoulli(1/2)$  model: responses are by chance.
- We need to design a policy  $\pi(a \mid x)$  that accepts or rejects depending on the data.
- Since there is no alternative model, we can only construct this policy according to its properties when  $\mu_0$  is true.
- In particular, we can fix a policy that only chooses  $a_1$  when  $\mu_0$  is true a proportion  $\delta$  of the time.
- This can be done by construcing a threshold test from the inverse-CDF.

#### Using p-values to construct statistical tests

**Definition A.1.2** (Null statistical test). A statistical test  $\pi$  is a decision rule for accepting or rejecting a hypothesis on the basis of evidence. A *p*-value test rejects a hypothesis whenever the value of the statistic f(x) is smaller than a threshold. The statistic  $f: \mathcal{X} \to [0, 1]$  is designed to have the property:

$$P_{\mu_0}(\{x \mid f(x) \le \delta\}) = \delta.$$

If our decision rule is:

$$\pi(a \mid x) = \begin{cases} a_0, & f(x) \ge \delta \\ a_1, & f(x) < \delta, \end{cases}$$

the probability of rejecting the null hypothesis when it is true is exactly  $\delta$ .

This is because, by definition, f(x) has a uniform distribution under  $\mu_0$ . Hence the value of f(x) itself is uninformative: high and low values are equally likely. In theory we should simply choose  $\delta$  before seeing the data and just accept or reject based on whether  $f(x) \leq \delta$ . However nobody does that in practice, meaning that p-values are used incorrectly. Better not to use them at all, if uncertain about their meaning.

#### Issues with p-values

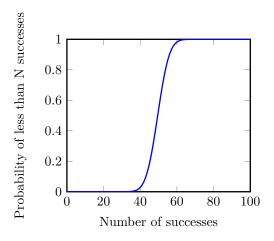
- They only measure quality of fit on the data.
- Not robust to model misspecification. For example, zero-mean testing using the  $\chi^2$ -test has a normality assumption.
- They ignore effect sizes. For example, a linear analysis may determine that there is a significant deviation from zero-mean, but with only a small effect size of 0.01. Thus, reporting only the p-value is misleading
- They do not consider prior information.
- They do not represent the probability of having made an error. In particular, a p-value of  $\delta$  does not mean that the probability that the null hypothesis is false given the data x, is  $\delta$ , i.e.  $\delta \neq \mathbb{P}(\neg \mu_0 \mid x)$ .
- The null-rejection error probability is the same irrespective of the amount of data (by design).

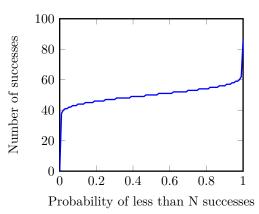
Let us consider the example of the medium.

#### p-values for the medium example

•  $\mu_0$  is simply the Bernoulli(1/2) model: responses are by chance.

- CDF:  $P_{\mu_0}(N \leq n \mid K = 100)$  is the probability of at most N successes if we throw the coin 100 times. This is in fact the cumulative probability function of the binomial distribution. Recall that the binomial represents the distribution for the number of successes of independent experiments, each following a Bernoulli distribution.
- ICDF: the number of successes that will happen with probability at least  $\delta$
- e.g. we'll get at most 50 successes a proportion  $\delta = 1/2$  of the time.
- Using the (inverse) CDF we can construct a policy  $\pi$  that selects  $a_1$  when  $\mu_0$  is true only a  $\delta$  portion of the time, for any choice of  $\delta$ .





## Constructing a Null-Hypothesis test with frequentist properties

#### The test statistic

We want the test to reflect that we don't have a significant number of failures.

$$f(x) = 1 - \text{binocdf}(\sum_{t=1}^{n} x_t, n, 0.5)$$

## 

- It is a **statistic** which is  $\leq \delta$  a  $\delta$  portion of the time when  $\mu_0$  is true.
- It is **not** the probability of observing x under  $\mu_0$ .
- It is **not** the probability of  $\mu_0$  given x.

EXERCISE 14. • Let us throw a coin 8 times, and try and predict the outcome.

- Select a p-value threshold so that  $\delta=0.05$ . For 8 throws, this corresponds to > 6 successes or  $\geq 87.5\%$  success rate.
- Let's calculate the p-value for each one of you
- What is the rejection performance of the test?

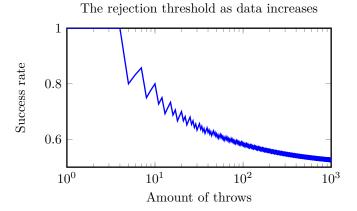
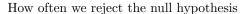


Figure A.1: Here we see how the rejection threshold, in terms of the success rate, changes with the number of throws to achieve an error rate of  $\delta = 0.05$ .

As the amount of throws goes to infinity, the threshold converges to 0.5. This means that a statistically significant difference from the null hypothesis can be obtained, even when the actual model from which the data is drawn is only slightly different from 0.5.



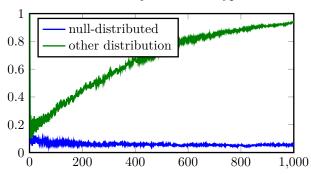


Figure A.2: Here we see the rejection rate of the null hypothesis  $(\mu_0)$  for two cases. Firstly, for the case when  $\mu_0$  is true. Secondly, when the data is generated from  $\mathcal{B}ernoulli(0.55)$ .

As we see, this method keeps its promise: the null is only rejected 0.05 of the time when it's true. We can also examine how often the null is rejected when it is false... but what should we compare against? Here we are generating data from a Bernoulli(0.55) model, and we can see the rejection of the null increases with the amount of data. This is called the power of the test with respect to the Bernoulli(0.55) distribution.

#### ♠ Statistical power and false discovery.

Beyond not rejecting the null when it's true, we also want:

• High power: Rejecting the null when it is false.

• Low false discovery rate: Accepting the null when it is true.

#### Power

The power depends on what hypothesis we use as an alternative. This implies that we cannot simply consider a plain null hypothesis test, but must formulate a specific alternative hypothesis.

#### False discovery rate

False discovery depends on how likely it is a priori that the null is false. This implies that we need to consider a prior probability for the null hypothesis being true.

Both of these problems suggest that a Bayesian approach might be more suitable. Firstly, it allows us to consider an infinite number of possible alternative models as the alternative hypothesis, through Bayesian model averaging. Secondly, it allows us to specify prior probabilities for each alternative. This is especially important when we consider some effects unlikely.

### The Bayesian Null-Hypothesis test

1. Set  $U(a_i, \mu_j) = \mathbb{I}\{i=j\}$ . This choice makes sense if we care equally about either Example 29. type of error.

- 2. Set  $\beta(\mu_i) = 1/2$ . Here we place an equal probability in both models.
- 3.  $\mu_0$ : Bernoulli(1/2). This is the same as the null hypothesis test.
- 4.  $\mu_1$ : Bernoulli( $\theta$ ),  $\theta \sim Unif([0,1])$ . This is an extension of the simple hypothesis test, with an alternative hypothesis that says "the data comes from an arbitrary Bernoulli model".
- 5. Calculate  $\beta(\mu \mid x)$ .
- 6. Choose  $a_i$ , where  $i = \arg \max_i \beta(\mu_i \mid x)$ .

#### Bayesian model averaging for the alternative model $\mu_1$

In this scenario,  $\mu_0$  is a simple point model, e.g. corresponding to a Bernoulli (1/2). However  $\mu_1$  is a marginal distribution integrated over many models, e.g. a *Beta* distribution over Bernoulli parameters.

$$P_{\mu_1}(x) = \int_{\Theta} B_{\theta}(x) \, \mathrm{d}\beta(\theta) \tag{A.1.3}$$

$$P_{\mu_1}(x) = \int_{\Theta} B_{\theta}(x) \, \mathrm{d}\beta(\theta)$$

$$\beta(\mu_0 \mid x) = \frac{P_{\mu_0}(x)\beta(\mu_0)}{P_{\mu_0}(x)\beta(\mu_0) + P_{\mu_1}(x)\beta(\mu_1)}$$
(A.1.4)

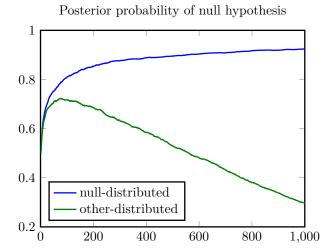


Figure A.3: Here we see the convergence of the posterior probability.

As can be seen in the figure above, in both cases, the posterior converges to the correct value, so it can be used to indicate our confidence that the null is true.

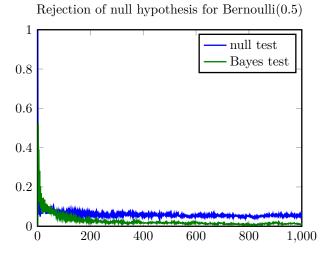


Figure A.4: Comparison of the rejection probability for the null and the Bayesian test when  $\mu_0$  is true.

Now we can use this Bayesian test, with uniform prior, to see how well it performs. While the plain null hypothesis test has a fixed rejection rate of 0.05, the Bayesian test's rejection rate converges to 0 as we collect more data.

Rejection of null hypothesis for Bernoulli(0.55)

Figure A.5: Comparison of the rejection probability for the null and the Bayesian test when  $\mu_1$  is true.

However, both methods are able to reject the null hypothesis more often when it is false, as long as we have more data.

#### Concentration inequalities and confidence intervals

There are a number of inequalities from probability theory that allow us to construct high-probability confidence intervals. The most well-known of those is Hoeffding's inequality , the simplest version of which is the following:

**Lemma A.1.1** (Hoeffding's inequality). Let  $x_1, \ldots, x_n$  be a series of random variables,  $x_i \in [0,1]$ . Then it holds that, for the empirical mean:

$$\mu_n \triangleq \frac{1}{n} sum_{t=1}^n x_t$$

$$\mathbb{P}(\mu_n \geq \mathbb{E} \,\mu_n + \epsilon) \leq e^{-2n\epsilon^2}.$$
(A.1.5)

When  $x_t$  are i.i.d,  $\mathbb{E} \mu_n = \mathbb{E} x_t$ . This allows us to construct an interval of size  $\epsilon$  around the true mean. This can generalise to a two-sided interval:

$$\mathbb{P}(|\mu_n - \mathbb{E}\,\mu_n| \ge \epsilon) \le 2e^{-2n\epsilon^2}.$$

We can also rewrite the equation to say that, with probability at least  $1 - \delta$ 

$$|\mu_n - \mathbb{E}\,\mu_n| \le \sqrt{\frac{\ln 2/\delta}{2n}}$$

## Further reading

#### Points of significance (Nature Methods)

- Importance of being uncertain https://www.nature.com/articles/nmeth.2613
- Error bars https://www.nature.com/articles/nmeth.2659
- P values and the search for significance https://www.nature.com/articles/nmeth. 4120
- Bayes' theorem https://www.nature.com/articles/nmeth.3335
- Sampling distributions and the bootstrap https://www.nature.com/articles/nmeth. 3414

## A.2 Formalising classification problems

One of the simplest decision problems is classification. At the simplest level, this is the problem of observing some data point  $x_t \in \mathcal{X}$  and making a decision about what class  $\mathcal{Y}$  it belongs to. Typically, a fixed classifier is defined as a decision rule  $\pi(a|x)$  making decisions  $a \in \mathcal{A}$ , where the decision space includes the class labels, so that if we observe some point  $x_t$  and choose  $a_t = 1$ , we essentially declare that  $y_t = 1$ .

Typically, we wish to have a classification policy that minimises classification error.

#### Deciding a class given a model

In the simplest classification problem, we observe some features  $x_t$  and want to make a guess  $a_t$  about the true class label  $y_t$ . Assuming we have some probabilistic model  $P_{\mu}(y_t \mid x_t)$ , we want to define a decision rule  $\pi(a_t \mid x_t)$  that is optimal, in the sense that it maximises expected utility for  $P_{\mu}$ .

- Features  $x_t \in \mathcal{X}$ .
- Label  $y_t \in \mathcal{Y}$ .
- Decisions  $a_t \in \mathcal{A}$ .
- Decision rule  $\pi(a_t \mid x_t)$  assigns probabilities to actions.

#### Standard classification problem

In the simplest case, the set of decisions we make are the same as the set of classes

$$A = Y$$
,  $U(a, y) = \mathbb{I}\{a = y\}$ 

EXERCISE 15. If we have a model  $P_{\mu}(y_t \mid x_t)$ , and a suitable U, what is the optimal decision to make?

## Deciding the class given a model family

- Training data  $D_T = \{(x_i, y_i) \mid i = 1, ..., T\}$
- Models  $\{P_{\mu} \mid \mu \in \mathcal{M}\}.$
- Prior  $\beta$  on  $\mathcal{M}$ .

Similarly to our example with the meteorological stations, we can define a posterior distribution over models.

#### Posterior over classification models

$$\beta(\mu \mid D_T) = \frac{P_{\mu}(y_1, \dots, y_T \mid x_1, \dots, x_T)\beta(\mu)}{\sum_{\mu' \in \mathcal{M}} P_{\mu'}(y_1, \dots, y_T \mid x_1, \dots, x_T)\beta(\mu')}$$

This posterior form can be seen as weighing each model according to how well they can predict the class labels. It is a correct form as long as, for every pair of models  $\mu, \mu'$  we have that  $P_{\mu}(x_1, \ldots, x_T) = P_{\mu'}(x_1, \ldots, x_T)$ . This assumption can be easily satisfied without specifying a particular model for the x. If not dealing with time-series data, we assume independence between  $x_t$ :

$$P_{\mu}(y_1, \dots, y_T \mid x_1, \dots, x_T) = \prod_{i=1}^T P_{\mu}(y_i \mid x_i)$$

## The *Bayes rule* for maximising $\mathbb{E}_{\beta}(U \mid a, x_t, D_T)$

The decision rule simply chooses the action:

$$a_t \in \underset{a \in \mathcal{A}}{\arg \max} \sum_{y} \sum_{\mu \in \mathcal{M}} P_{\mu}(y_t = y \mid x_t) \beta(\mu \mid D_T) U(a, y)$$
(A.2.1)

$$= \arg\max_{a \in \mathcal{A}} \sum_{y} \mathbb{P}_{\beta \mid D_T}(y_t = y \mid x_t) U(a, y)$$
(A.2.2)

We can rewrite this by calculating the posterior marginal marginal label probability

$$\mathbb{P}_{\beta \mid D_T}(y_t \mid x_t) \triangleq \mathbb{P}_{\beta}(y_t \mid x_t, D_T) = \sum_{\mu \in \mathcal{M}} P_{\mu}(y_t \mid x_t) \beta(\mu \mid D_T).$$

#### Approximating the model

#### Full Bayesian approach for infinite $\mathcal{M}$

Here  $\beta$  can be a probability density function and

$$\beta(\mu \mid D_T) = P_{\mu}(D_T)\beta(\mu)/\mathbb{P}_{\beta}(D_T), \qquad \mathbb{P}_{\beta}(D_T) = \int_{\mathcal{M}} P_{\mu}(D_T)\beta(\mu) d,$$

can be hard to calculate.

## Maximum a posteriori model

We only choose a single model through the following optimisation:

$$\mu_{\text{MAP}}(\beta, D_T) = \underset{\mu \in \mathcal{M}}{\operatorname{arg \, max}} P_{\mu}(D_T)\beta(\mu) = \underset{\mu \in \mathcal{M}}{\operatorname{arg \, max}} \underbrace{\ln P_{\mu}(D_T)} + \underbrace{\ln \beta(\mu)}_{\text{regulariser}}.$$

You can think of the goodness of fit as how well the model fits the training data, while the regulariser term simply weighs models according to some criterion. Typically, lower weights are used for more complex models.

#### Learning outcomes

## Understanding

- Preferences, utilities and the expected utility principle.
- Hypothesis testing and classification as decision problems.
- How to interpret *p*-values Bayesian tests.
- The MAP approximation to full Bayesian inference.

#### Skills

- Being able to implement an optimal decision rule for a given utility and probability.
- Being able to construct a simple null hypothesis test.

#### Reflection

- When would expected utility maximisation not be a good idea?
- What does a p value represent when you see it in a paper?
- Can we prevent high false discovery rates when using p values?
- When is the MAP approximation good?

## A.3 Beliefs and probabilities

Probability can be used to describe purely chance events, as in for example quantum physics. However, it is mostly used to describe uncertain events, such as the outcome of a dice roll or a coin flip, which only appear random. In fact, one can take it even further than that, and use it to model subjective uncertainty about any arbitrary event. Although probabilities are not the only way in which we can quantify uncertainty, it is a simple enough model, and with a rich enough history in mathematics, statistics, computer science and engineering that it is the most useful.

#### Uncertainty

## Axioms of probability

Let  $\Omega$  be the certain event, and  $\Sigma$  is an appropriate  $\sigma$ -algebra on  $\Omega$ . A probability measure P on  $(\Omega, \Sigma)$  has the following properties:

- 1. The probability of the certain event is  $P(\Omega) = 1$
- 2. The probability of the impossible event is  $P(\emptyset) = 0$

- 3. The probability of any event  $A \in \Sigma$  is  $0 \le P(A) \le 1$ .
- 4. If A,B are disjoint, i.e.  $A\cap B=\emptyset$ , meaning that they cannot happen at the same time, then

$$P(A \cup B) = P(A) + P(B)$$

Sometimes we would like to calculate the probability of some event A happening given that we know that some other event B has happened. For this we need to first define the idea of conditional probability.

**Definition A.3.1** (Conditional probability). The probability of A happening if we know that B has happened is defined to be:

$$P(A \mid B) \triangleq \frac{P(A \cap B)}{P(B)}.$$

Conditional probabilities obey the same rules as probabilities. Here, the probability measure of any event A given B is defined to be the probability of the intersection of of the events divided by the second event. We can rewrite this definition as follows, by using the definition for  $P(B \mid A)$ 

#### Bayes's theorem

For  $P(A_1 \cup A_2) = 1$ ,  $A_1 \cap A_2 = \emptyset$ ,

$$P(A_i \mid B) = \frac{P(B \mid A_i)P(A_i)}{P(B)} = \frac{P(B \mid A_i)P(A_i)}{P(B \mid A_1)P(A_1) + P(B \mid A_2)P(A_2)}$$

EXAMPLE 30 (probability of rain). What is the probability of rain given a forecast  $x_1$  or  $x_2$ ?

$$\begin{array}{c|c} \omega_1 \colon \operatorname{rain} & P(\omega_1) = 80\% \\ \omega_2 \colon \operatorname{dry} & P(\omega_2) = 20\% \end{array}$$

Table A.5: Prior probability of rain tomorrow

$$x_1$$
: rain |  $P(x_1 \mid \omega_1) = 90\%$   
 $x_2$ : dry |  $P(x_2 \mid \omega_2) = 50\%$ 

Table A.6: Probability the forecast is correct

$$P(\omega_1 \mid x_1) = 87.8\%$$
  
 $P(\omega_1 \mid x_2) = 44.4\%$ 

Table A.7: Probability that it will rain given the forecast

#### Classification in terms of conditional probabilities

Conditional probability naturally appears in classification problems. Given a new example vector of data  $x_t \in \mathcal{X}$ , we would like to calculate the probability of different classes  $c \in \mathcal{Y}$  given the data,  $P_{\mu}(y_t = c \mid x_t)$ . If we somehow obtained the distribution of data  $P_{\mu}(x_t \mid y_t)$  for each

possible class, as well as the prior class probability  $P_{\mu}(y_t = c)$ , from Bayes's theorem, we see that we can obtain the probability of the class:

$$P_{\mu}(y_t = c \mid x_t) = \frac{P_{\mu}(x_t \mid y_t = c)P_{\mu}(y_t = c)}{\sum_{c' \in \mathcal{Y}} P_{\mu}(x_t \mid y_t = c')P_{\mu}(y_t = c')}$$

for any class c. This directly gives us a method for classifying new data, as long as we have a way to obtain  $P_{\mu}(x_t \mid y_t)$  and  $P_{\mu}(y_t)$ .

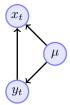


Figure A.6: A generative classification model.  $\mu$  identifies the model (paramter).  $x_t$  are the features and  $y_t$  the class label of the t-th example.

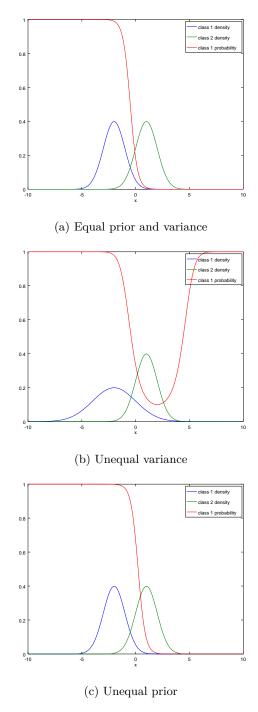


Figure A.7: The effect of changing variance and prior on the classification decision when we assume a normal distribution.

Example 31 (Normal distribution). A simple example is when  $x_t$  is normally distributed in a matter that depends on the class. Figure A.7 shows the distribution of  $x_t$  for two different classes, with means

of -1 and +1 respectively, for three different case. In the first case, both classes have variance of 1, and we assume the same prior probability for both

$$x_t \mid y_t = 0 \sim \mathcal{N}(-1, 1), \qquad x_t \mid y_t = 1 \sim \mathcal{N}(1, 1)$$

$$x_t \mid y_t = 0 \sim \mathcal{N}(-1, 1), \qquad x_t \mid y_t = 1 \sim \mathcal{N}(1, 1)$$

But how can we get a probability model in the first place?

## Subjective probability

While probabilities apply to truly random events, they are also useful for representing subjective uncertainty. In this course, we will use a special symbol for subjective probability,  $\beta$ .

#### Subjective probability measure $\beta$

- If we think event A is more likely than B, then  $\beta(A) > \beta(B)$ .
- Usual rules of probability apply:
  - 1.  $\beta(A) \in [0,1]$ .
  - 2.  $\beta(\emptyset) = 0$ .
  - 3. If  $A \cap B = \emptyset$ , then  $\beta(A \cup B) = \beta(A) + \beta(B)$ .

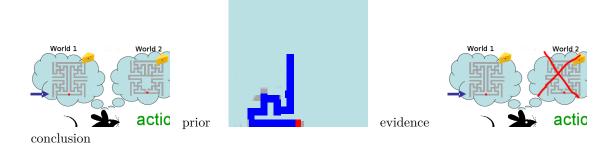
#### Bayesian inference illustration

Here is a simple example to illustrate the idea of Bayesian inference. Imagine you are a rat in the maze and you need to get to the cheese. In the example below, you believe that the maze lay out is only one of two possibilities  $\mu_1, \mu_2$ . Initially you might believe that the first model is more likely than the second, so you decide upon some function  $\beta$  so that  $\beta(\mu_1) > \beta(\mu_2)$ , for example  $\beta(\mu_1) = 2/3$ ,  $\beta(\mu_2) = 1/3$ . Now you start acting in the maze, collecting data h consisting of your history of observation in this maze. The data gives a different amount of evidence for each model of the world, which we can write as  $P_{\mu}(h)$ . In the end, we want to combine this evidence with our initial belief to arrive at a conclusion. In the Bayesian setting, all of these quantities are probabilities, and the conclusion is simply the conditional probability of the model given the data.

## Use a subjective belief $\beta(\mu)$ on $\mathcal{M}$

- Prior belief  $\beta(\mu)$  represents our initial uncertainty.
- We observe history h.
- Each possible  $\mu$  assigns a probability  $P_{\mu}(h)$  to h.
- We can use this to *update* our belief via Bayes' theorem to obtain the *posterior* belief:

$$\beta(\mu \mid h) \propto P_{\mu}(h)\beta(\mu)$$
 (conclusion = evidence × prior)



## A.3.1 Probability and Bayesian inference

One of the most important methods in machine learning and statistics is that of Bayesian inference. This is the most fundamental method of drawing conclusions from data and explicit prior assumptions. In Bayesian inference, prior assumptions are represented as a probabilities on a space of hypotheses. Each hypothesis is seen as a probabilistic model of all possible data that we can see.

Frequently, we want to draw conclusions from data. However, the conclusions are never solely inferred from data, but also depend on prior assumptions about reality.

#### Some examples

EXAMPLE 32. John claims to be a medium. He throws a coin n times and predicts its value always correctly. Should we believe that he is a medium?

- $\mu_1$ : John is a medium.
- $\mu_0$ : John is not a medium.

The answer depends on what we *expect* a medium to be able to do, and how likely we thought he'd be a medium in the first place.

EXAMPLE 33. Traces of DNA are found at a murder scene. We perform a DNA test against a database of  $10^4$  citizens registered to be living in the area. We know that the probability of a false positive (that is, the test finding a match by mistake) is  $10^{-6}$ . If there is a match in the database, does that mean that the citizen was at the scene of the crime?

## Bayesian inference

Now let us apply this idea to our specific problem. We already have the probability of the observation for each model, but we just need to define a *prior probability* for each model. Since this is usually completely subjective, we give it another symbol.

#### Prior probability

The prior probability  $\beta$  on a set of models  $\mathcal{M}$  specifies our subjective belief  $\beta(\mu)$  that each model is true.<sup>a</sup>

<sup>a</sup>More generally  $\beta$  is a probability measure.

This allows us to calculate the probability of John being a medium, given the data:

$$eta(\mu_1 \mid \boldsymbol{x}) = rac{\mathbb{P}(\boldsymbol{x} \mid \mu_1) eta(\mu_1)}{\mathbb{P}_{eta}(\boldsymbol{x})},$$

where

$$\mathbb{P}_{\beta}(\boldsymbol{x}) \triangleq \mathbb{P}(\boldsymbol{x} \mid \mu_1)\beta(\mu_1) + \mathbb{P}(\boldsymbol{x} \mid \mu_0)\beta(\mu_0).$$

The only thing left to specify is  $\beta(\mu_1)$ , the probability that John is a medium before seeing the data. This is our subjective prior belief that mediums exist and that John is one of them. More generally, we can think of Bayesian inference as follows:

- We start with a set of mutually exclusive models  $\mathcal{M} = \{\mu_1, \dots, \mu_k\}$ .
- Each model  $\mu$  is represented by a specific probabilistic model for any possible data x, that is  $P_{\mu}(x) \equiv \mathbb{P}(x \mid \mu)$ .
- For each model, we have a prior probability  $\beta(\mu)$  that it is correct.
- After observing the data, we can calculate a posterior probability that the model is correct:

$$\beta(\mu \mid x) = \frac{\mathbb{P}(x \mid \mu)\beta(\mu)}{\sum_{\mu' \in \mathcal{M}} \mathbb{P}(x \mid \mu')\beta(\mu')} = \frac{P_{\mu}(x)\beta(\mu)}{\sum_{\mu' \in \mathcal{M}} P_{\mu'}(x)\beta(\mu')}.$$

## C Interpretation

- $\mathcal{M}$ : Set of all possible models that could describe the data.
- $P_{\mu}(x)$ : Probability of x under model  $\mu$ .
- Alternative notation  $\mathbb{P}(x \mid \mu)$ : Probability of x given that model  $\mu$  is correct.
- $\beta(\mu)$ : Our belief, before seeing the data, that  $\mu$  is correct.
- $\beta(\mu \mid x)$ : Our belief, aftering seeing the data, that  $\mu$  is correct.

It must be emphasized that  $P_{\mu}(x) = \mathbb{P}(x \mid \mu)$  as they are simply two different notations for the same thing. In words the first can be seen as the probability that model  $\mu$  assigns to data x, while the second as the probability of x if  $\mu$  is the true model. Combining the prior belief with evidence is key in this procedure. Our posterior belief can then be used as a new prior belief when we get more evidence.

EXERCISE 16 (Continued example for medium). Now let us apply this idea to our specific problem. We first make an independence assumption. In particular, we can assume that success and failure comes from a Bernoulli distribution with a parameter depending on the model.

$$P_{\mu}(x) = \prod_{t=1}^{n} P_{\mu}(x_t).$$
 (independence property)

We first need to specify how well a medium could predict. Let's assume that a true medium would be able to predict perfectly, and that a non-medium would only predict randomly. This leads to the following models:

$$P_{\mu_1}(x_t=1)=1,$$
  $P_{\mu_1}(x_t=0)=0.$  (true medium model)  $P_{\mu_0}(x_t=1)=1/2,$   $P_{\mu_0}(x_t=0)=1/2.$  (non-medium model)

The only thing left to specify is  $\beta(\mu_1)$ , the probability that John is a medium before seeing the data. This is our subjective prior belief that mediums exist and that John is one of them.

$$\beta(\mu_0) = 1/2,$$
  $\beta(\mu_1) = 1/2.$  (prior belief)

Combining the prior belief with evidence is key in this procedure. Our posterior belief can then be used as a new prior belief when we get more evidence.

$$\beta(\mu_1 \mid x) = \frac{P_{\mu_1}(x)\beta(\mu_1)}{\mathbb{P}_{\beta}(x)}$$
 (posterior belief)  
$$\mathbb{P}_{\beta}(x) \triangleq P_{\mu_1}(x)\beta(\mu_1) + P_{\mu_0}(x)\beta(\mu_0).$$
 (marginal distribution)

If a classmate correctly predicts 4 coin tosses, what is your belief they are a medium?

#### Sequential update of beliefs

Assume you have n meteorologists. At each day t, each meteorologist i gives a probability  $p_{t,\mu_i} \triangleq P_{\mu_i}(x_t = \text{rain})$  for rain. Consider the case of there being three meteorologists, and each one making the following prediction for the coming week. Start with a uniform prior  $\beta(\mu) = 1/3$  for each model.

	M	T	W	$\Gamma$	F	S	S
CNN	0.5	0.6	0.7	0.9	0.5	0.3	0.1
SMHI	0.3	0.7	0.8	0.9	0.5	0.2	0.1
YR	0.6	0.9	0.8	0.5	0.4	0.1	0.1
Rain?	Y	Y	Y	N	Y	N	N

Table A.8: Predictions by three different entities for the probability of rain on a particular day, along with whether or not it actually rained.

Exercise 17. • n meteorological stations  $\{\mu_i \mid i=1,\ldots,n\}$ 

- The *i*-th station predicts rain  $P_{\mu_i}(x_t \mid x_1, \dots, x_{t-1})$ .
- Let  $\beta_t(\mu)$  be our belief at time t. Derive the next-step belief  $\beta_{t+1}(\mu) \triangleq \beta_t(\mu|x_t)$  in terms of the current belief  $\beta_t$ .
- Write a python function that computes this posterior

$$\beta_{t+1}(\mu) \triangleq \beta_t(\mu|x_t) = \frac{P_{\mu}(x_t \mid x_1, \dots, x_{t-1})\beta_t(\mu)}{\sum_{\mu'} P_{\mu'}(x_t \mid x_1, \dots, x_{t-1})\beta_t(\mu')}$$

## Bayesian inference for Bernoulli distributions

## Estimating a coin's bias

A fair coin comes heads 50% of the time. We want to test an unknown coin, which we think may not be completely fair.

For a sequence of throws  $x_t \in \{0, 1\}$ ,

$$P_{\theta}(x) \propto \prod_{t} \theta^{x_t} (1 - \theta)^{1 - x_t} = \theta^{\text{\#Heads}} (1 - \theta)^{\text{\#Tails}}$$

Say we throw the coin 100 times and obtain 70 heads. Then we plot the *likelihood*  $P_{\theta}(x)$  of different models. From these, we calculate a *posterior* distribution over the correct models.

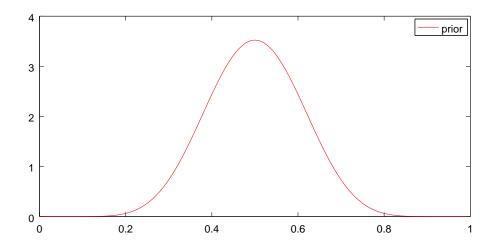


Figure A.8: Prior belief  $\beta$  about the coin bias  $\theta$ .

This represents our conclusion given our prior and the data. If the prior distribution is described by the so-called Beta density

$$f(\theta \mid \alpha, \beta) \propto \theta^{\alpha - 1} (1 - \theta)^{\beta - 1}$$

where  $\alpha, \beta$  describe the shape of the Beta distribution.

#### Riemann and Lebesgue integrals

Since  $\beta$  is a probability measure over models  $\Theta$ , it is always simple to write the posterior probability given some data x in the following terms when  $\Theta$  is discrete (finite or countable):

$$\beta(\theta \mid x) = \frac{P_{\theta}(x)\beta(\theta)}{\sum_{\theta'} P_{\theta'(x)\beta(\theta')}}.$$

However, in many situations  $\Theta$  is uncountable, i.e.  $\Theta \subset \mathbb{R}^k$ . Then, as both the prior  $\beta$  and the posterior  $\beta(\cdot \mid x)$  have to be probability measures on  $\Theta$ , they are defined over subsets of  $\Theta$ . This means that we can write the posterior in terms of Lebesgue integrals:

$$\beta(B \mid x) = \frac{\int_B P_{\theta}(x) \, \mathrm{d}\beta(\theta)}{\int_{\Theta} P_{\theta}(x) \, \mathrm{d}\beta(\theta)}.$$

Alternatively, we can abuse notation and use  $\beta(\theta)$  to describe a density, so that the *posterior density* is written in terms of a Riemann integral:

posterior density

$$\beta(\theta \mid x) \triangleq \frac{P_{\theta}(x)\beta(\theta) d\theta}{\int_{\Theta} P_{\theta'}(x)\beta(\theta')\theta'}.$$

However, the advantage of the Lebesgue integral notation is that it works the same no matter whether  $\Theta$  is discrete or continuous.

## Learning outcomes

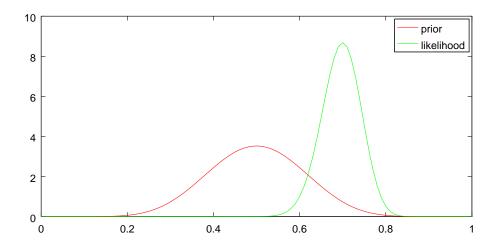


Figure A.9: Prior belief  $\beta$  about the coin bias  $\theta$  and likelihood of  $\theta$  for the data.

#### Understanding

- The axioms of probability, marginals and conditional distributions.
- The philosophical underpinnings of Bayesianism.
- The simple conjugate model for Bernoulli distributions.

## Skills

- Be able to calculate with probabilities using the marginal and conditional definitions and Bayes rule.
- Being able to implement a simple Bayesian inference algorithm in Python.

## Reflection

- How useful is the Bayesian representation of uncertainty?
- How restrictive is the need to select a prior distribution?
- Can you think of another way to explicitly represent uncertainty in a way that can incorporate new evidence?

## A.4 Classification with stochastic gradient descent\*

## Classification as an optimisation problem.

Finding the optimal policy for our belief  $\beta$  is not normally very difficult. However, it requires that we maintain the complete distribution  $\beta$  and the set of models  $P_{\mu}(y \mid x)$ . In simple decision problems, e.g. where the set of actions  $\mathcal{A}$  is finite, it is possible to do this calculation on-the-fly. However, in some cases we might not have a model.

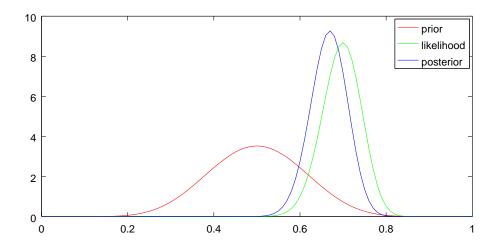


Figure A.10: Prior belief  $\beta(\theta)$  about the coin bias  $\theta$ , likelihood of  $\theta$  for the data, and posterior belief  $\beta(\theta \mid x)$ 

Recall that we wish to maximise expected utility for some policy under some distribution. In general, this has the form

$$\max_{\pi} \mathbb{E}^{\pi}_{\mu}(U).$$

We also know that any expectation can be approximated by sampling. Let  $P_{\mu}(\omega)$  be the distribution on outcomes defined by our model. Then

$$\mathbb{E}_{\mu}^{\pi}(U) = \sum_{\omega} U(a, \omega) P_{\mu}(\omega) \approx T^{-1} \sum_{t=1}^{T} U(a, \omega_{t}), \qquad \omega_{t} \sim P_{\mu}(\omega),$$

i.e. when we can replace the explicit summation over all possible outcomes, weighed by their probability through averaging over T outcomes sampled from the correct distribution. In fact this approximation is *unbiased*, as its expectation is equal to the expected utility.

### The $\mu$ -optimal classifier

Since the performance measure is simply an expectation, it is intuitive to directly optimise the decision rule with respect to an approximation of the expectation

$$\max_{\theta \in \Theta} f(\pi_{\theta}, \mu, U), \qquad f(\pi_{\theta}, \mu, U) \triangleq \mathbb{E}_{\mu}^{\pi_{\theta}}(U) \qquad (A.4.1)$$

$$f(\pi_{\theta}, \mu, U) = \sum_{x,y,a} U(a, y) \pi_{\theta}(a \mid x) P_{\mu}(y \mid x) P_{\mu}(x)$$
(A.4.2)

$$= \sum_{x,y,a} U(a,y)\pi_{\theta}(a \mid x)P_{\mu}(y \mid x)P_{\mu}(x)$$

$$\approx \sum_{t=1}^{T} \sum_{a_{t}} U(a_{t},y_{t})\pi_{\theta}(a_{t} \mid x_{t}),$$

$$(A.4.2)$$

$$(x_{t},y_{t})_{t=1}^{T} \sim P_{\mu}.$$

In practice, this is the empirical expectation on the training set  $\{(x_t, y_t) \mid t = 1, \dots, T\}$ . However, when the amount of data is insufficient, this expectation may be far from reality, and so our classification rule might be far from optimal.

### The Bayes-optimal classifier

An alternative idea is to use our uncertainty to create a distribution over models, and then use this distribution to obtain a single classifier that does take the uncertainty into account.

$$\max_{\theta \in \Theta} f(\pi_{\theta}\beta) \approx \max_{\theta \in \Theta} N^{-1} \sum_{n=1}^{N} \pi_{\theta}(a_t = y_n \mid x_t = x_n), \qquad (x_n, y_n) \sim P_{\mu_n}, \mu_n \sim \beta.$$

In this case, the integrals are replaced by sampling models  $\mu_n$  from the belief, and then sampling  $(x_n, y_n)$  pairs from  $P_{\mu_n}$ .

### Stochastic gradient methods

To find the maximum of a differentiable function q, we can use gradient descent

### Gradient ascent

$$\theta_{i+1} = \theta_i + \alpha \nabla_{\theta} g(\theta_i).$$

When f is an expectation, we don't need to calculate the full gradient. In fact, we only need to take one sample from the related distribution.

### Stochastic gradient ascent

$$g(\theta) = \int_{\mathcal{M}} f(\theta, \mu) \, \mathrm{d}\beta(\mu)$$
  
$$\theta_{i+1} = \theta_i + \alpha \nabla_{\theta} f(\theta_i, \mu_i), \qquad \mu_i \sim \beta.$$

Stochastic gradient methods are commonly employed in neural networks.

### A.4.1 Neural network models

### Two views of neural networks

In the simplest sense a neural network is simply as parametrised functions  $f_{\theta}$ . In classification, neural networks can be used as probabilistic models, so they describes the probability  $P_{\theta}(y|\mathbf{x})$ , or as classification policies so that  $f_{\theta}(x,a)$  describes the probability  $\pi_{\theta}(a \mid x)$  of selecting class label a. Let us begin by describing the simplest type of neural network model, the perceptron.

### Neural network classification model $P_{\theta}(\boldsymbol{y} \mid \boldsymbol{x}_t)$



Objective: Find the best model for  $D_T$ .

In this case, we consider a utility function that minimises the divergence between the network's probability predictions and the labels, such as  $U(\theta, y_t, x_t) = \ln P_{\theta}(y_t|x_t)$ .

### Neural network classification policy $\pi(a_t \mid x_t)$



Objective: Find the best policy for U(a, x).

In this case, we consider a utility function that considers the actual classification decisions, so that if  $u(a_t, y_t)$  is how much we gain by choosing  $a_t$  when the label is  $y_t$ , then the expected utility for a decision a is  $U(a, \mathbf{x}_t) = \sum_y \mathbb{P}(y_t = y | \mathbf{x}_t) u(a, y)$ .

### 

- We can use standard probabilistic methods for P.
- Finding the optimal  $\pi$  is an optimisation problem. However, estimating P can also be formulated as an optimisation.

### Linear networks and the perceptron algorithm



Figure A.11: Abstract graphical model for a neural network

A neural network as used for modelling classification or regression problems, is simply a parametrised mapping  $\mathcal{X} \to \mathcal{Y}$ . If we include the network parameters, then it is instead a mapping  $\mathcal{X} \times \Theta \to \mathcal{Y}$ , as seen in Figure A.13.

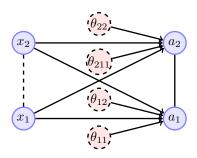


Figure A.12: Graphical model for a linear neural network.

If we see each possible output as a different random variable, this creates a dependence. After all, we are really splitting one variable into many. In particular, if the network's output is the probability of each action, then we must make sure these sum to 1.

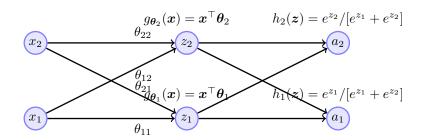


Figure A.13: Architectural view of a linear neural network.

**Definition A.4.1** (Linear classifier). A linear classifier with N inputs and C outputs is parametrised by

$$oldsymbol{\Theta} = egin{bmatrix} oldsymbol{ heta}_1 & \cdots & oldsymbol{ heta}_C \end{bmatrix} = egin{bmatrix} heta_{1,1} & \cdots & heta_{1,C} \ dots & \ddots & dots \ heta_N & \cdots & heta_{N,C} \end{bmatrix}$$

$$\pi_{\Theta}(a \mid \boldsymbol{x}) = \exp\left(\boldsymbol{\theta}_{a}^{\top} \boldsymbol{x}\right) / \sum_{a'} \exp\left(\boldsymbol{\theta}_{a'}^{\top} \boldsymbol{x}\right)$$

Even though the classifier has a linear structure, the final non-linearity at the end is there to ensure that it defines a proper probability distribution over decisions. For classification problems, the observations  $\boldsymbol{x}_t$  are features  $\boldsymbol{x}_t = (x_{t,1} \dots, x_{t,n})$  so that  $\mathcal{X} \subset \mathbb{R}^N$ . It is convenient to consider the network output as a vector on the simplex  $\boldsymbol{y} \in \boldsymbol{\Delta}^A$ , i.e.  $\sum_{i=1}^C y_{t,i} = 1, \ y_{t,i} \geq 0$ . In the neural network model for classification, we typically ignore dependencies between the  $x_{t,i}$  features, as we are not very interested in the distribution of  $\boldsymbol{x}$  itself.

### Gradient ascent for a matrix U

$$\max_{\theta} \sum_{t=1}^{T} \sum_{a_t} U(a_t, y_t) \pi_{\theta}(a_t \mid x_t)$$
 (objective)

$$\nabla_{\theta} \sum_{t=1}^{T} \sum_{a_{t}} U(a_{t}, y_{t}) \pi_{\theta}(a_{t} \mid x_{t})$$
 (gradient)

$$= \sum_{t=1}^{T} \sum_{a_t} U(a_t, y_t) \nabla_{\theta} \pi_{\theta}(a_t \mid x_t)$$
(A.4.4)

We now need to calculate the gradient of the policy.

Chain Rule of Differentiation 
$$f(z), z = g(x), \qquad \qquad \frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx} \qquad \qquad \text{(scalar version)}$$
 
$$\nabla_{\theta} \pi = \nabla_{q} \pi \nabla_{\theta} g \qquad \qquad \text{(vector version)}$$

### Learning outcomes

### Understanding

- Classification as an optimisation problem.
- (Stochastic) gradient methods and the chain rule.
- Neural networks as probability models or classification policies.
- Linear neural netwoks.
- Nonlinear network architectures.

### Skills

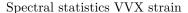
• Using a standard NN class in python.

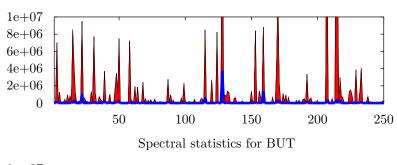
### Reflection

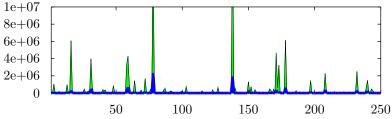
- How useful is the ability to have multiple non-linear layers in a neural network.
- How rich is the representational power of neural networks?
- Is there anything special about neural networks other than their allusions to biology?

### A.5 Nearest neighbours

### Discriminating between diseases



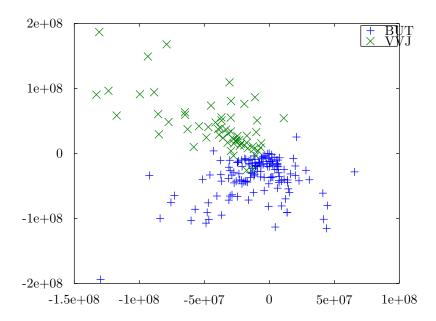




Let's tackle the problem of discriminating between different disease vectors. Ideally, we'd like to have a simple test that tells us what ails us. One kind of test is mass spectrometry. This graph shows spectrometry results for two types of bacteria. There is plenty of variation within

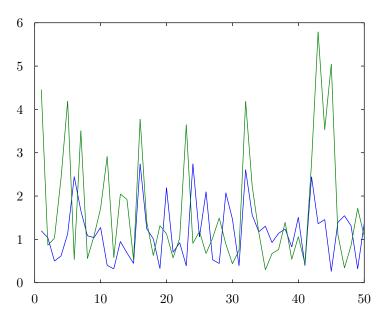
each type, both due to measurement error and due to changes in the bacterial strains. Here, we plot the average and maximum energies measured for about 100 different examples from each strain.

### Nearest neighbour: the hidden secret of machine learning



Now, is it possible to identify an unknown strain based on this data? Actually, this is possible. Sometimes, very simple algorithms work very well. One of the simplest one involves just measuring the distance between the decsription of a new unknown strain and known ones. In this visualisation, I projected the 1300-dimensional data into a 2-dimensional space. Here you can clearly see that it is possible to separate the two strains. We can use the distance to examples VVT and BUT in order to decide the type of an unknown strain.

### Comparing spectral data



The choice of distance in this kind of algorithm is important, particularly for very high dimensions. For something like a spectrogram, one idea is look at the total area of the difference between two spectral lines.

### The nearest neighbour algorithm

The nearest neighbour algorithm for classification (Alg. 2) does not include any complicated learning. Given a training dataset D, it returns a classification decision for any new point x by simply comparing it to its closest k neighbours in the dataset. It then estimates the probability  $p_y$  of each class y by calculating the average number of times the neighbours take the class y.

### Algorithm 2 K-NN Classify

- 1: **Input** Data  $D = \{(x_1, y_1), \dots, (x_T, y_T)\}, k \geq 1, d: \mathcal{X} \times \mathcal{X} \to \mathbb{R}_+, \text{ new point } x \in \mathcal{X}$ 2: D = Sort(D, d) % Sort D so that  $d(x, x_i) \leq d(x, x_{i+1})$ .
  3:  $p_y = \sum_{i=1}^k \mathbb{I}\{y_i = y\}/k \text{ for } y \in \mathcal{Y}$ .
  4: **Return**  $\mathbf{p} \triangleq (p_1, \dots, p_k)$

### **∧** Algorithm parameters

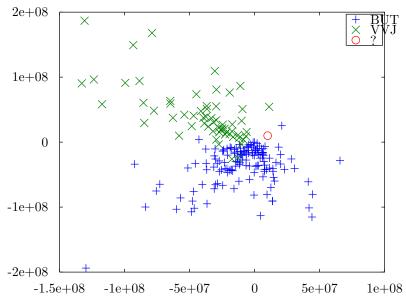
In order to use the algorithm, we must specify some parameters, namely.

- Neighbourhood  $k \geq 1$ . The number of neighbours to consider.
- Distance  $d: \mathcal{X} \times \mathcal{X} \to \mathbb{R}_+$ . The function we use to determine what is a neighbour.



Figure A.14: The nearest neighbours algorithm was introduced by Fix and Hodges  ${\rm Jr}^{\,10},$  who also proved consistency properties.

Nearest neighbour: What type is the new bacterium?



-1.5e+08 -1e+08 -5e+07 0 5e+07 1e+08 Given that the + points represent the BUT type, and the × points the VVJ type, what type of bacterium could the circle point be?

### Separating the model from the classification policy

- The K-NN algorithm returns a model giving class probabilities for new data points.
- It is up to us to decide how to use this model to decide upon a given class. A typical decision making rule can be in the form of a policy  $\pi$  that depends on what the model says. However, the simplest decision rule is to take the most likely class:

$$\pi(a \mid x) = \mathbb{I}\left\{p_a \ge p_y \forall y\right\}, \quad \mathbf{p} = \text{K-NN}(D, k, d, x)$$

### Discussion: Shortcomings of k-nearest neighbour

- Choice of k The larger k is, the more data you take into account when making your decision. This means that you expect your classes to be more spread out.
- Choice of metric d. The metric d encodes prior information you have about the structure of the data.
- Representation of uncertainty. The predictions of kNN models are simply based on distances and counting. This might not be a very good way to represent uncertainty about class label. In particular, label probabilities should be more uncertain when we have less data.
- Scaling with large amounts of data. A naive implementation of kNN requires the algorithm to shift through all the training data to find the k nearest neighbours, suggesting a superlinear computation time. However, advanced data structures such as Cover Trees (or even KD-trees in low dimensional spaces) can be used to find the neighbours in polylog time.

• Meaning of label probabilities. It is best to think of K-NN as a *model* for predicting the class of a new example from a finite set of existing classes. The model itself might be incorrect, but this should nevertheless be OK for our purposes. In particular, we might later use the model in order to derive classification rules.

### Learning outcomes

### Understanding

- · How kNN works
- The effect of hyperameters k, d for nearest neighbour.
- The use of kNN to classify new data.

### Skills

- Use a standard kNN class in python
- Optimise kNN hyperparameters in an unbiased manner.
- Calculate probabilities of class labels using kNN.

### Reflection

- When is kNN a good model?
- How can we deal with large amounts of data?
- How can we best represent uncertainty?

## A.6 Naive Bayes classifiers\*

One special case of this idea is in classification, when each hypothesis corresponds to a specific class. Then, given a new example vector of data  $\boldsymbol{x}$ , we would like to calculate the probability of different classes C given the data,  $\mathbb{P}(C \mid \boldsymbol{x})$ . So here, the class is the hypothesis.

From Bayes's theorem, we see that we can write this as

$$\mathbb{P}(C \mid \boldsymbol{x}) = \frac{\mathbb{P}(\boldsymbol{x} \mid C) \, \mathbb{P}(C)}{\sum_{i} \mathbb{P}(\boldsymbol{x} \mid C_{i}) \, \mathbb{P}(C_{i})}$$

for any class C. This directly gives us a method for classifying new data, as long as we have a way to obtain  $\mathbb{P}(\boldsymbol{x} \mid C)$  and  $\mathbb{P}(C)$ .

But should we use for the probability model  $\mathbb{P}$ ?

### Naive Bayes classifier

Naive Bayes classifiers are one of the simplest classification methods. They can have a full Bayesian interpretation under some assumptions, but otherwise they are too simplistic to be useful.

### Calculating the prior probability of classes

A simple method is to simply count the number of times each class appears in the training data  $D_T = ((x_t, y_t))_{t=1}^T$ . Then we can set

$$\mathbb{P}(C) = 1/T \sum_{t=1}^{T} \mathbb{I}\left\{y_t = C\right\}$$

The Naive Bayes classifier uses the following model for observations, where observations are independent of each other given the class. Thus, for example the result of three different tests for lung cancer (stethoscope, radiography and biopsy) only depend on whether you have cancer, and not on each other.

### Probability model for observations

$$\mathbb{P}(\boldsymbol{x} \mid C) = \mathbb{P}(x(1), \dots, x(n) \mid C) = \prod_{k=1}^{n} \mathbb{P}(x(k) \mid C).$$

There are two different types of models we can have, one of which is mostly useful for continuous attributes and the other for discrete attributes. In the first, we just need to count the number of times each feature takes different values in different classes.

#### Discrete attribute model.

Here we simply count the average number of times that the attribute k had the value i when the label was C. This is in fact analogous to the conditional probability definition.

$$\mathbb{P}(x(k) = i \mid C) = \frac{\sum_{t=1}^{T} \mathbb{I} \{ x_t(k) = i \land y_t = C \}}{\sum_{t=1}^{T} \mathbb{I} \{ y_t = C \}} = \frac{N_k(i, C)}{N(C)},$$

where  $N_k(i, C)$  is the numb 1 l .er of examples in class C whose k-th attribute has the value i, and N(C) is the number of examples in class C.

### 

This estimation is simple maximum likelihood, as it does not maintain a distribution over the parameters.

Sometimes we need to be able to deal with cases where there are no examples at all of one class. In that case, that class would have probability zero. To get around this problem, we add "fake observations" to our data. This is called *Laplace smoothing*.

Remark A.6.1. In Laplace smoothing with constant  $\lambda$ , our probability model is

$$\mathbb{P}(x(k) = i \mid C) = \frac{\sum_{t=1}^{T} \mathbb{I}\{x_t(k) = i \land y_t = C\} + \lambda}{\sum_{t=1}^{T} \mathbb{I}\{y_t = C\} + n_k \lambda} = \frac{N_k(i, C) + \lambda}{N(C) + n_k \lambda}.$$

where  $n_k$  is the number of values that the k-th attribute can take. This is necessary, because we want  $\sum_{i=1}^{n_k} \mathbb{P}(x(k) = i \mid C) = 1$ . (You can check that this is indeed the case as a simple exercise).

Remark A.6.2. In fact, the Laplace smoothing model corresponds to a so-called Dirichelt prior over polynomial parameters with a marginal probability of observation equal to the Laplace smoothing. This is an extension of Beta-Bernoulli example from binary outcomes to multiple outcomes.

### Continuous attribute model.

Here we can use a Gaussian model for each continuous dimension.

$$\mathbb{P}(x(k) = v \mid C) = \frac{1}{\sigma\sqrt{2\pi}}e^{\frac{(v-\mu)^2}{\sigma^2}},$$

where  $\mu$  and  $\sigma$  are the mean and variance of the Gaussian, typically calculated from the training data as:

$$\mu = \frac{\sum_{t=1}^{T} x_t(k) \mathbb{I} \{ y_t = C \}}{\sum_{t=1}^{T} \mathbb{I} \{ y_t = C \}},$$

i.e.  $\mu$  is the mean of the k-th attribute when the label is C and

$$\sigma = \frac{\sum_{t=1}^{T} [x_t(k) - \mu]^2 \mathbb{I} \{ y_t = C \}}{\sum_{t=1}^{T} \mathbb{I} \{ y_t = C \}},$$

i.e.  $\sigma$  is the variance of the k-th attribute when the label is C. Sometimes we can just fix  $\sigma$  to a constant value, i.e.  $\sigma = 1$ .

### 

This estimation is simple maximum likelihood, as it selects a single parameter pair  $\mu = (\mu_1, \dots, \mu_n)$  and  $\sigma = (\sigma_1, \dots, \sigma_n)$  for every class and does not maintain a distribution over the parameters. It also assumes independence between the features. The full Bayesian approach considers an arbitrary covariance matrix  $\Sigma$  and maintains a distribution  $\beta(\mu, \Sigma)$ .

## Appendix B

# Empirical estimates and sampling

Let us fix some probability measure P on  $(\Omega, \Sigma)$ . Let  $P^n$  denote the corresponding product measure.

**Sample.** A single sample from P will be written as  $\omega \sim P$ . A sample of size n from P can be written as

$$(\omega^{(1)},\ldots,\omega^{(n)})\sim P^n,$$

or equivalently as

$$\omega^{(i)} \sim P^n, \quad i = 1, \dots, n.$$

**Expectation.** For any  $\Sigma$ -measurable random variable f, we define the expectation with respect to P as

$$\mathbb{E}_P(f) \triangleq \int_{\Omega} f(\omega) dP(\omega).$$

**Mean.** Let us define the n-sample mean of f as follows:

$$\hat{E}_n(f) \triangleq \frac{1}{n} \sum_{i=1}^n f(\omega^{(n)}), \qquad \omega^{(i)} \sim P, i = 1, \dots, n$$

It is easy to show that

$$\mathbb{E}_{P}[\hat{E}_{n}(f)] = \mathbb{E}_{P}(f).$$

It moreover holds that  $|\mathbb{E}_P(f) - \hat{E}_n(f)| \leq O(1/\sqrt{n})$  with high probability under mild conditions. See for example Hoeffding's inequality (Lemma A.1.1) for bounded random variables.

Importance sampling. When we wish to estimate an expectation under one probability measure, but using data generated from another probability, we can always use importance sampling, also called also called propensity scoring. There are two reasons why we may use importance sampling. The first is that simply the data has been generated from some distribution Q, but we want the expectation under P. The second is that we may want to actively sample from a different distribution than P around more "important" values, hoping to obtain better estimates.

No matter what the reason, we can always calculate the expectation of any random variable f under a measure P by sampling through another measure Q, as follows:

$$\mathbb{E}_{P}(f) = \int_{\Omega} f(\omega) \, dP(\omega) = \int_{\Omega} f(\omega) \frac{dP(\omega)}{dQ(\omega)} \, dQ(\omega) \approx \frac{1}{n} \sum_{i=1}^{n} \int_{\Omega} f(\omega^{(n)}) \frac{dP(\omega^{(n)})}{dQ(\omega^{(n)})} \, dQ(\omega).$$

Here  $\frac{\mathrm{d}P(\omega)}{\mathrm{d}Q(\omega)}$  is the Radon-Nikodym derivative, which for the case of finite  $\Omega$  is simply  $P(\omega)/Q(\omega)$ . Intuitively, it can be thought of as the ratio of densities  $p(\omega)/q(\omega)$  of the two probability measures P,Q.

## Appendix C

# Graphical models

Graphical models are a very useful tool for modelling the relationship between multiple variables. The simplest such models, probabilistic graphical models (otherwise known as Bayesian networks) involve directed acyclic graphs between random variables. There are two other types of probabilistic models, factor graph and undirected graphical models, which are equivalent to each other, though not to directed models.

### Graphical models

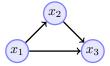


Figure C.1: Graphical model (directed acyclic graph) for three variables.

Consider for example the model in Figure ??. It involves three variables,  $x_1, x_2, x_3$  and there are three arrows, which show how one variable depends on another. Simply put, if you think of each  $x_k$  as a stochastic function, then  $x_k$ 's value only depends on the values of its parents, i.e. the nodes that are point to it. In this example,  $x_1$  does not depend on any other variable, but the value of  $x_2$  depends on the value of  $x_1$ . Such models are useful when we want to describe the joint probability distribution of all the variables in the collection.

The graphical model allows us to factorise the joint probability distribution of these random variables in a simplified manner. First, we define what we mean by a joint probability with respect to some probability measure P.

### Joint probability

Let P is a probability measure on  $(\Omega, \Sigma)$ . Then let the random variable  $\boldsymbol{x} = (x_1, \dots, x_n)$  so that  $\boldsymbol{x} : \Omega \to X$ ,  $X = \prod_i X_i$ . The joint probability of  $\boldsymbol{x}$  can be written in terms of the underlying probability measure P:

$$\mathbb{P}(\boldsymbol{x} \in A) = P(\{\omega \in \Omega \mid \boldsymbol{x}(\omega) \in A\}).$$

When  $X_i$  are finite, we can typically write

$$\mathbb{P}(\boldsymbol{x} = \boldsymbol{a}) = P(\{\omega \in \Omega \mid \boldsymbol{x}(\omega) = \boldsymbol{a}\}),$$

for the probability that  $x_i = a_i$  for all  $i \in [n]$ . Through the definition of conditional probability, we can always factorise the joint distribution of random variables as follows:

#### **Factorisation**

For any subsets  $B \subset [n]$  and its complement C so that  $\mathbf{x}_B = (x_i)_{i \in B}$ ,  $\mathbf{x}_C = (x_i)_{i \notin B}$ 

$$\mathbb{P}(\boldsymbol{x}) = \mathbb{P}(\boldsymbol{x}_B \mid \boldsymbol{x}_C) \, \mathbb{P}(\boldsymbol{x}_C)$$

So we can write any joint distribution as

$$\mathbb{P}(x_1)\,\mathbb{P}(x_2\mid x_1)\,\mathbb{P}(x_3\mid x_1,x_2)\cdots\mathbb{P}(x_n\mid x_1,\ldots,x_{n-1}).$$

Although the above factorisation is always possible to do, sometimes our graphical model has a structure that makes the factors much simpler. In fact, the main reason for introducing graphical models is to represent dependencies between variables. For a given model, we can infer whether some variables are in fact dependent, independent, or conditionally independent.

### Directed graphical models and conditional independence

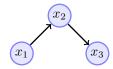


Figure C.2: Graphical model for the factorisation  $\mathbb{P}(x_3 \mid x_2) \mathbb{P}(x_2 \mid x_1) \mathbb{P}(x_1)$ .

### Conditional independence

We say  $x_i$  is conditionally independent of  $x_B$  given  $x_D$  and write  $x_i \perp \!\!\! \perp x_B \mid x_D$  iff

$$\mathbb{P}(x_i, \boldsymbol{x}_B \mid \boldsymbol{x}_D) = \mathbb{P}(x_i \mid \boldsymbol{x}_D) \, \mathbb{P}(\boldsymbol{x}_B \mid \boldsymbol{x}_D).$$

### Directed graphical models

A graphical model is a convenient way to represent conditional independence between variables. There are many variants of graphical models, whose name is context dependent. Other names used in the literature are probabilistic graphical models, Bayesian networks, causal graphs, or decision diagrams. In this set of notes we focus on directed graphical models that depict dependencies between ranom variables.

**Definition C.0.1** (Directed graphical model). A collection of n random variables  $x_i : \Omega \to X_i$ , and let  $X \triangleq \prod_i X_i$ , with underlying probability measure P on  $\Omega$ . Let  $\mathbf{x} = (x_i)_{i=1}^n$  and for any subset  $B \subset [n]$  let

$$\boldsymbol{x}_B \triangleq (x_i)_{i \in B} \tag{C.0.1}$$

$$\boldsymbol{x}_{-i} \triangleq (x_i)_{i \neq i} \tag{C.0.2}$$

In a graphical model, conditional independence is represented through directed edges.

EXAMPLE 34 (Specifying a probability model). A graphical model does not specify the complete distribution, only an allowed factorisation. If we want to specify a complete distribution for the above graphical model of  $x_1, x_2, x_3$ , we can use the following notation:

$$x_1 \sim f$$
 (C.0.3)

$$x_2 \mid x_1 = a \sim g(a) \tag{C.0.4}$$

$$x_3 \mid x_2 = b \sim h(b),$$
 (C.0.5)

where f, g, h are three different distributions, with g and h being specified through a single parameter.

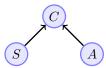


Figure C.3: Smoking and lung cancer graphical model, where S: Smoking, C: cancer, A: asbestos exposure.

EXAMPLE 35 (Smoking and lung cancer). It has been found by Lee <sup>13</sup> that lung incidence not only increases with both asbestos exposure and smoking. This is in agreement with the graphical model shown. The study actually found that there is an amplification effect, whereby smoking and asbestos exposure increases cancer risk by 28 times compared to non-smokers. This implies that the risk is not simply additive. The graphical model only tells us that there is a dependency, and does not describe the nature of this dependency precisely.

### Explaining away

Even though S,A are independent, they become dependent once you know C. For example, let us say we know that you have cancer and that our model says that it's very unlikely to have cancer unless you either smoke or are exposed to asbestos. When we also learn that you do not have asbestos exposure, smoking becomes more likely. In either words, if cancer is caused by either smoking or asbestos, and we rule out asbestos, the only remaining explanation is smoking. This is what is generally called *explaining away*.

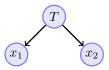


Figure C.4: Time of arrival at work graphical model where T is a traffic jam and  $x_1$  is the time John arrives at the office and  $x_2$  is the time Jane arrives at the office.

EXAMPLE 36 (Time of arrival at work). In this model, the arrival times of John and Jane may seem correlated. However, there is a common cause: The existence of a traffic jam. Whenever there is a traffic jam, both John and Jane are usually late. Whenever there is not a traffic jam, they are both mostly on time.

78%

Treatment B

87%

69%83%

### Conditional independence

Even though  $x_1, x_2$  are correlated, they become independent once you know T.

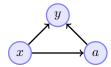


Figure C.5: Kidney treatment model, where x: severity, y: result, a: treatment applied

	Treatment A	Treatment B	Severity	Treatment A
	Heatiment A		Small stones)	93%
Small stones	87	270	_ ′	
Large stones	263	80	Large stones	73%
Large stones	200	00	Average	78%

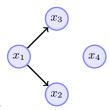
Example 37 (Treatment effects). A curious example is that of applying one of two treatments for kidneys. In the data, it is clear that one treatment is best for both large and small stones. However, when the data is aggregated it appears as though treatment B is best. This is because treatment A is chosen much more frequently when the stones are large, and that's when both treatments perform worse. This apparent discrepancy is called Simpson's paradox



Figure C.6: School admission graphical model, where z: gender, s: school applied to, a: whether you were admitted.

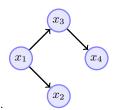
School	Male	Female
A	62%	82%
В	63%	68%
$\mathbf{C}$	37%	34%
D	33%	35%
$\mathbf{E}$	28%	24%
F	6%	7%
$\overline{Average}$	45%	38%

EXAMPLE 38 (School admission). In this example, it appears as though female candidates have a lower acceptance rate than males. However what is missing is the fact that many more males are applying to easier schools. Thus, it is possible that the data is explainable by the fact that admission only reflects the difficulty of each school, and the overall gender imbalance is due to the choices made by the applicants. However, an alternative model is that the admissions process also explicitly takes gender into account. However, both of these models may be inadequate, as we do not have data about each individual applicant, such as their grades. We shall discuss this issue further when we talk about causality, confounding variables and counterfactuals.



EXERCISE 18. Factorise the following graphical model.

$$\mathbb{P}(\boldsymbol{x}) = \mathbb{P}(x_1) \, \mathbb{P}(x_2 \mid x_1) \, \mathbb{P}(x_3 \mid x_1) \, \mathbb{P}(x_4)$$

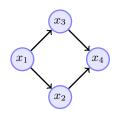


EXERCISE 19. Factorise the following graphical model.

$$\mathbb{P}(\boldsymbol{x}) = \mathbb{P}(x_1) \, \mathbb{P}(x_2 \mid x_1) \, \mathbb{P}(x_3 \mid x_1) \, \mathbb{P}(x_4 \mid x_3)$$

EXERCISE 20. What dependencies does the following factorisation imply?

$$\mathbb{P}(\boldsymbol{x}) = \mathbb{P}(x_1) \, \mathbb{P}(x_2 \mid x_1) \, \mathbb{P}(x_3 \mid x_1) \, \mathbb{P}(x_4 \mid x_2, x_3)$$



### **⚠** Deciding conditional independence

There is an algorithm for deciding conditional independence of any two variables in a graphical model. However, this is beyond the scope of these notes. Here, we shall just use these models as a way to encode dependencies that we assume exist.

## C.1 Decision diagrams

Graphical models can be extended to model causal relations. In particular, we can use decision  $diagrams^1$ , which include not only random variables, but also decision variables, denoted with squares, as well as utility variables, denoted via diamonds. In the following examples, we assume there are some underlying distributions specified by parameters  $\theta$ , which we include in the diagrams for clarity. Even though it may seem intuitively sensible to suppose it, the arrow directions in the diagrams do not indicate direct causes. The only important thing for determining whether some variable influences another is whether or not there is independence between the corresponding decision and random variables.

<sup>&</sup>lt;sup>1</sup>Otherwise called influence diagrams

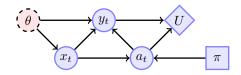


Figure C.7: A typical decision diagram where  $x_t$ : individual information,  $y_t$ : individual result,  $a_t$ : action,  $\pi$ : policy

EXAMPLE 39 (Taking an aspirin). The diagram in Figure C.7 does not completely specify the decision problem. For aspirin taking, we can define the following variables:

- Individual t
- Individual information  $x_t$
- $a_t = 1$  if t takes an aspirin, and 0 otherwise.
- $y_t = 1$  if the headache is cured in 30 minutes, 0 otherwise.
- $\pi$ : intervention policy.

EXAMPLE 40 (A recommendation system). Consider the example of a recommendation system, where we have data of the form  $(x_t, a_t, y_t)$ . The performance of the recommendation system depends not only on the parameter  $\theta$ , but also on the chosen policy  $\pi$ .

- $x_t$ : User information (random variable)
- $a_t$ : System action (random variable)
- $y_t$ : Click (random varaible)
- $\pi$ : recommendation policy (decision variable).

In both cases, there are some questions we can ask using the underlying model. The dependency structure is not enough to know *a priori* whether we can obtain meaningful answers. This depends on the specific assumptions we make about the model.

### Conditional distributions and decision variables.

We begin with a parenthesis on conditional distributions. We normally define the conditional distribution of A given B under a probability measure P as:

$$P(A \mid B) \triangleq \frac{P(A \cap B)}{P(B)}.$$

However, decision variables are outside the scope of this probability measure, and yet we need to define conditional distributions using them.

### The conditional distribution of decisions

If  $\pi \in \Pi$  is a decision variable, we represent the conditional distribution of any random variable a given  $\pi$  simply as a collection of probability measures  $\{\pi(a) \mid \pi \in \Pi\}$ , one for each possible value  $\pi$ . The following notations will be equivalent:

$$\pi(a) \equiv \mathbb{P}^{\pi}(a) \equiv \mathbb{P}(a \mid \pi).$$

The reader should note that the standard definition of a conditional distribution also  $P(A \mid B)$  creates a collection of distributions on A, with elements  $P_B(A)$ . However, it also specifies a rule for doing so from the complete distribution P.

If the random variables a also depends on some probability law  $P_{\theta}$ , then it will be

convenient to use the notation

$$\mathbb{P}_{\theta}^{\pi}(a) \equiv \mathbb{P}(a \mid \theta, \pi).$$

### C.2 Inference and prediction in graphical models

Inference and prediction in graphical models

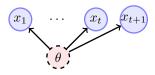


Figure C.8: Inference and prediction in a graphical model

In this example,  $x_1, \ldots, x_t$  are all i.i.d, drawn from the distribution  $P_{\theta}(x_t) = \mathbb{P}(x_t \mid \theta)$ .

### Inference of latent variables

$$\mathbb{P}(\theta \mid x_1, \dots, x_t)$$

Inference in graphical models typically refers to the problem of estimating the distribution of some *hiddeen* or *latent* variable from data. These variables are generally thought of as two types:

- Model parameters. These generally do not change over time. One example is estimating the mean of a Gaussian distribution from data.
- System states. These typically are time-indexed. We will see further examples of such variables when we discuss latent variable models. On example is inferring location from GPS measurements.

### Prediction

$$\mathbb{P}(x_{t+1} \mid x_1, \dots, x_t) = \int_{\Theta} \mathbb{P}(x_{t+1} \mid \theta) d\mathbb{P}(\theta \mid x_1, \dots, x_t)$$

Prediction is a type of inference, but differs in that the variable whose distribution we wish to estimate is only temporarily not observed: we can actually obtain its value in the future. Thus, a prediction is always testable!

### Coin tossing, revisited

Example 41. The Beta-Bernoulli prior



Figure C.9: Graphical model for a Beta-Bernoulli prior

$$\theta \sim \mathcal{B}eta(\beta_1, \beta_2),$$
 i.e.  $\beta$  are Beta distribution parameters (C.2.1)

$$x \mid \theta \sim \mathcal{B}emoulli(\theta),$$
 i.e.  $P_{\theta}(x)$  is a Bernoulli (C.2.2)

In this example, it is obvious why we use the notation above for describing hierarchical models. We simply state what is the distribution on one variable conditioned on the other variables. Here,  $\beta$  is fixed, and it is something we can choose arbitrarily. The data x is observed, while the parameter  $\theta$  remains latent. Using Bayes theorem, we can derive the distribution for  $\beta(\theta \mid x)$ . In this particular case, we elide referring to a sample  $x_1, \ldots, x_t$  as they are all i.i.d.

Example 42. The *n*-meteorologists problem (continuation of Exercise 13)

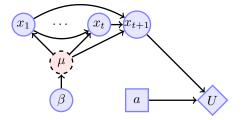


Figure C.10: Inference, prediction and decisions in a graphical model.

In this problem, we allow each model's predictions to have arbitrary dependences with the past i.e.

$$P_{\mu}(x_{t+1} \mid x_t, x_{t-1}, \dots, x_1).$$

For inference, the details of how each model works are not important: just the probability of the next observation given the previous ones.

## C.3 Testing conditional independence

### Measuring independence

The simplest way to measure independence is by looking at whether or not the distribution of the possibly dependent variable changes when we change the value of the other variables.

Theorem C.3.1. If  $x_i \perp \!\!\! \perp x_B \mid x_D$  then

$$\mathbb{P}(x_i \mid \boldsymbol{x}_B, \boldsymbol{x}_D) = \mathbb{P}(x_i \mid \boldsymbol{x}_D)$$

This implies

$$\mathbb{P}(x_i \mid \boldsymbol{x}_B = b, \boldsymbol{x}_D) = \mathbb{P}(x_i \mid \boldsymbol{x}_B = b', \boldsymbol{x}_D)$$

so we can measure independence by seeing how the distribution of  $x_i$  changes when we vary  $x_B$ , keeping  $x_D$  fixed. For any given model, there is either a dependence or there is not. However, sometimes we might be able to tolerate some amount of dependence. Thus, we can simply measure the deviation from independence through a metric or divergence on distributions.

Example 43.

$$\|\mathbb{P}(a \mid y, z) - \mathbb{P}(a \mid y)\|_1$$

which for discrete a, y, z is:

$$\max_{i,j} \| \mathbb{P}(a \mid y = i, z = j) - \mathbb{P}(a \mid y = i)\|_1 = \max_{i,j} \| \sum_k \mathbb{P}(a = k \mid y = i, z = j) - \mathbb{P}(a = k \mid y = i)\|_1.$$

See also src/fairness/ci\_test.py

EXAMPLE 44. An alternative model for coin-tossing This is an elaboration of Example 29 for hypothesis testing.

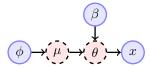


Figure C.11: Graphical model for a hierarchical prior

- $\mu_1$ : A Beta-Bernoulli model with  $\mathcal{B}eta(\beta_1, \beta_2)$
- $\mu_0$ : The coin is fair.

$$\theta \mid \mu = \mu_0 \sim \mathcal{D}(0.5),$$
 i.e.  $\theta$  is always 0.5 (C.3.1)

$$\theta \mid \mu = \mu_1 \sim \mathcal{B}eta(\beta_1, \beta_2),$$
 i.e.  $\theta$  has a Beta distribution (C.3.2)

$$x \mid \theta \sim \mathcal{B}ernoulli(\theta),$$
 i.e.  $P_{\theta}(x)$  is Bernoulli (C.3.3)

Here the posterior over the two models is simply

$$\phi(\mu_0 \mid x) = \frac{P_{0.5}(x)\phi(\mu_0)}{P_{0.5}(x)\phi(\mu_0) + \mathbb{P}_{\mu_1}(x)\phi(\mu_1)}, \qquad \mathbb{P}_{\mu_1(x)} = \int_0^1 P_{\theta}(x) \,\mathrm{d}\beta(\theta).$$

### Bayesian testing of independence

For a given distributional model  $P_{\theta}$ , conditional independence either holds or does not. Consider the set of model parameters  $\Theta_0$  where, for each parameter  $\theta \in \Theta_0$ , we have a conditional independence condition, while  $\Theta_1$  may corresponds to models where there may not be independence. To make this more concrete, let's give an example.

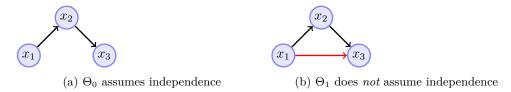


Figure C.12: The two alternative models

EXAMPLE 45. Assume data  $D = \{x_1^t, x_2^t, x_3^t \mid t = 1, \dots, T\}$  with  $x_i^t \in \{0, 1\}$ . First consider model  $\Theta_0$ where the following conditional independence holds

$$P_{\theta}(x_3 \mid x_2, x_1) = P_{\theta}(x_2 \mid x_1), \quad \forall \theta \in \Theta_0.$$

In the alternative model  $\Theta_1$  there is no independence assumption. So the likelihood for either a model in either set is

$$P_{\theta}(D) = \prod_{t} P_{\theta}(x_3^t \mid x_2^t) P_{\theta}(x_2^t \mid x_1^t) P_{\theta}(x_1^t), \qquad \theta \in \Theta_0$$
 (C.3.4)

$$P_{\theta}(D) = \prod_{t} P_{\theta}(x_3^t \mid x_2^t) P_{\theta}(x_2^t \mid x_1^t) P_{\theta}(x_1^t), \qquad \theta \in \Theta_0$$

$$P_{\theta}(D) = \prod_{t} P_{\theta}(x_3^t \mid x_2^t, x_1^t) P_{\theta}(x_2^t \mid x_1^t) P_{\theta}(x_1^t), \qquad \theta \in \Theta_1$$
(C.3.4)

The parameters for this example can be defined as follows

$$\theta_1 \triangleq P_\theta(x_1^t = 1) \tag{\mu_0, \mu_1}$$

$$\theta_{2|1}^i \triangleq P_{\theta}(x_2^t = 1 \mid x_1^t = i)$$
  $(\mu_0, \mu_1)$ 

$$\theta_{3|2}^j \triangleq P_\theta(x_3^t = 1 \mid x_2^t = j) \tag{\mu_0}$$

$$\theta_{3|2,1}^{i,j} \triangleq P_{\theta}(x_3^t = 1 \mid x_2^t = j, x_1^t = i)$$
 (\mu\_1)

We model each one of these parameters with a separate Beta-Bernoulli distribution.

### C.4 Hierarchical Bayesian models

Given some data D, the Bayesian approach would involve specifying a hierarchical prior  $\beta$  so that  $\phi(\mu_1) = 1 - \phi(\mu_0)$  specifies a probability on the two model structures, while for the *i*-th model we define a prior  $\beta_i(\theta)$  over  $\Theta_i$ , so that we obtain the following hierarchical model



Figure C.13: Hierarchical model.

Here the specific model  $\mu$  is unobserved, as well as its parameters  $\theta$ . Only the data D is observed. Our prior distribution is omitted from the graph.

$$\mu_i \sim \phi$$
 (C.4.1)

$$\theta \mid \mu = \mu_i \sim \beta_i \tag{C.4.2}$$

### Marginal likelihood

This gives the the following marginal likelihood for the combined models and each of the models respectively.

$$\mathbb{P}_{\phi}(D) = \phi(\mu_0) \, \mathbb{P}_{\mu_0}(D) + \phi(\mu_1) \, \mathbb{P}_{\mu_1}(D) \tag{C.4.3}$$

$$\mathbb{P}_{\mu_i}(D) = \int_{\Theta_i} P_{\theta}(D) \, \mathrm{d}\beta_i(\theta). \tag{C.4.4}$$

### Model posterior

$$\phi(\mu \mid D) = \frac{\mathbb{P}_{\mu}(D)\phi(\mu)}{\sum_{i} \mathbb{P}_{\mu_{i}}(D)\phi(\mu_{i})}$$
 (C.4.5)

### Calculating the marginal likelihood

Generally speaking, calculating the marginal likelihood for a model with an uncountable parameter set is hard. However, conjugate models admit closed form solutions and efficient calculations. Firstly, let's rewrite the marginal likelihood in terms of an integral Monte-Carlo approximation.

### Monte-Carlo approximation

$$\int_{\Theta} P_{\theta}(D) \, \mathrm{d}\beta(\theta) \approx \frac{1}{N} \sum_{n=1}^{N} P_{\theta_n}(D) + O(1/\sqrt{N}), \qquad \theta_n \sim \beta$$
 (C.4.6)

Even though this approximation is reasonable at first glance, the problem is that the leading constant of the error scales approximately proportionally to the maximum likelihood  $\max_{\theta} P_{\theta}(D)$ . This is a consequence of Hoeffding's inequality (??). Thus, the more data we have the more samples we need to get a good approximation with this simple Monte Carlo approach. For that reason, one typically uses a sample from a proposal distribution  $\psi$  which is different from  $\beta$ . Then it holds

### Importance sampling

For any two measures  $\beta, \psi$  on  $\Theta$ , we can write:

$$\int_{\Theta} P_{\theta}(D) \, \mathrm{d}\beta(\theta) = \int_{\Theta} P_{\theta}(D) \frac{\mathrm{d}\psi(\theta)}{\mathrm{d}\psi(\theta)} \, \mathrm{d}\beta(\theta) \approx \sum_{n=1}^{N} P_{\theta}(D) \frac{\mathrm{d}\beta(\theta_n)}{\mathrm{d}\psi(\theta_n)}, \qquad \theta_n \sim \psi \qquad (C.4.7)$$

This allows us to estimate the marginal likelihood with respect to a belief  $\beta$  by sampling from an alternative belief  $\psi$ .

### Sequential updating of the marginal likelihood

$$\mathbb{P}_{\beta}(D) = \mathbb{P}_{\beta}(x_1, \dots, x_T) \tag{C.4.8}$$

$$= \mathbb{P}_{\beta}(x_2, \dots, x_T \mid x_1) \, \mathbb{P}_{\beta}(x_1) \tag{C.4.9}$$

$$= \prod_{t=1}^{T} \mathbb{P}_{\beta}(x_t \mid x_1, \dots, x_{t-1})$$
 (C.4.10)

$$= \prod_{t=1}^{T} \int_{\Theta} P_{\theta_n}(x_t) \, \mathrm{d} \underbrace{\beta(\theta \mid x_1, \dots, x_{t-1})}_{\text{posterior at time } t}$$
 (C.4.11)

The nice thing about this break down is that for a simple model such as Beta-Bernoulli, the individual datapoint marginal likelihoods are easy to compute

EXAMPLE 46 (Beta-Bernoulli). The marginal predictive distribution for a Beta-Bernoulli prior is

$$\mathbb{P}_{\beta}(x_t = 1 \mid x_1, \dots, x_{t-1}) = \frac{\alpha_t}{\alpha_t + \beta_t},$$

with 
$$\alpha_t = \alpha_0 + \sum_{n=1}^{t-1} x_n$$
,  $\beta_t = \beta_0 + \sum_{n=1}^{t-1} (1 - x_n)$ 

### Further reading

### Python sources

- A simple python measure of conditional independence src/fairness/ci\_test.py
- A simple test for discrete Bayesian network src/fairness/DirichletTest.py
- Using the PyMC package https://docs.pymc.io/notebooks/Bayes\_factor.html

## Index

risk, 89

```
n-meteorologists, 91, 132
                                                 Simpson's paradox, 78, 128
                                                 variable
confounder, 75
                                                     instrumental, 75
gradient ascent, 112
    stochastic, 80, 112
hierarchical Bayesian model, 132, 134
hierarchical Bayesian model, Beta-Bernoulli, 135
Hoeffding inequality, 98
importance sampling, 50, 123, 135
inference, 15, 131
seevariable
    instrumental, 75
latent, 132
MAP inference, 100
Maximum a posteriori, see MP inference100
Null-Hypothesis test, 92
    Bayesian, 96
    frequentist, 94
    classification, 61, 119
    classification accuracy, 62
    default, 77
    evaluation, 78
    exponential mechanism, 40
    for statistical testing, 92
    intervention, 77, 130
    neural network classification, 112
    optimisation, 79
    recommendation, 130
posterior density, 109
prediction, 16, 131
propensity scoring, see importance sampling
```

136 INDEX

# Bibliography

- [1] Craig M. Bennett, George L. Wolford, and Michael B. Miller. The principled control of false positives in neuroimaging. Social cognitive and affective neuroscience, 4 4:417-22, 2009. URL https://pdfs.semanticscholar.org/19c3/d8b67564d0e287a43b1e7e0f496eb1e8a945.pdf.
- [2] Craig M Bennett, Abigail A Baird, Michael B Miller, and George L Wolford. Journal of serendipitous and unexpected results. *Journal of Serendipitous and Unexpected Results (jsur. org)-Vol*, 1(1):1-5, 2012. URL https://teenspecies.github.io/pdfs/NeuralCorrelates.pdf.
- [3] Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123-140, 1996. URL citeseer. nj.nec.com/breiman96bagging.html.
- [4] Philip Dawid. The decision-theoretic approach to causal inference. Causality: Statistical perspectives and applications, pages 25-42, 2012. URL https://arxiv.org/pdf/1405.2292.pdf.
- [5] Christos Dimitrakakis, Blaine Nelson, Zuhe Zhang, Aikaterini Mitrokotsa, and Benjamin I. P. Rubinstein. Differential privacy for bayesian inference through posterior sampling. *Journal of Machine Learning Research*, 18(1):343–381, 2017.
- [6] John C Duchi, Michael I Jordan, and Martin J Wainwright. Local privacy and statistical minimax rates. In 2013 IEEE 54th Annual Symposium on Foundations of Computer Science, pages 429–438. IEEE, 2013.
- [7] Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. Foundations and Trends® in Theoretical Computer Science, 9(3-4):211-407, 2014.
- [8] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [9] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Roth. The reusable holdout: Preserving validity in adaptive data analysis. *Science*, 349 (6248):636–638, 2015.
- [10] Evelyn Fix and Joseph L Hodges Jr. Discriminatory analysis-nonparametric discrimination: consistency properties. Technical report, California Univ Berkeley, 1951.
- [11] Simson L. Garfinkel and Philip Leclerc. Randomness concerns when deploying differential privacy, 2020.

138 BIBLIOGRAPHY

[12] Jason Hartford, Greg Lewis, Kevin Leyton-Brown, and Matt Taddy. Counterfactual prediction with deep instrumental variables networks. Technical Report 1612.09596, arXiv, 2016.

- [13] PN Lee. Relation between exposure to asbestos and smoking jointly and the risk of lung cancer. Occupational and environmental medicine, 58(3):145–153, 2001.
- [14] Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In 48th Annual IEEE Symposium on Foundations of Computer Science (FOCS'07), pages 94–103. IEEE, 2007.
- [15] Pierangela Samarati and Latanya Sweeney. Protecting privacy when disclosing information: k-anonymity and its enforcement through generalization and suppression. Technical report, technical report, SRI International, 1998.
- [16] Stanley L. Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69, 1965. ISSN 01621459. URL http://www.jstor.org/stable/2283137.
- [17] Mengmeng Yang, Lingjuan Lyu, Jun Zhao, Tianqing Zhu, and Kwok-Yan Lam. Local differential privacy and its applications: A comprehensive survey. arXiv preprint arXiv:2008.03686, 2020.