

INTRODUCTION TO STATISTICS
THEORY, METHODS, AND APPLICATIONS

AUTHORED BY
JONAS PETERSEN
December 5, 2024

This page is intentionally left blank

Contents

1	PREFACE	1
1.1	Acknowledgements	2
I	INTRODUCTION	3
2	INTRODUCTION TO SET THEORY	5
3	INTRODUCTION TO PROBABILITY THEORY	13
4	INTRODUCTION TO STATISTICS	37
4.1	Interpretation of a Probability Measure	39
4.2	Relaxation of Notation	41
5	ASSIGNING PROBABILITY FUNCTIONS	43
5.1	The Principle of Maximum Entropy	44
6	FRAMING OF STATISTICS	49
6.1	Assigning a Cost Function	52
6.1.1	Continuous Action Space	52
6.1.2	Discrete Action Space	56
6.2	Statistical Paradigms	56
II	FREQUENTIST STATISTICS	57
7	FREQUENTIST STATISTICS INTRODUCTION	59
8	PARAMETER ESTIMATION	61
III	BAYESIAN STATISTICS	69
9	BAYESIAN STATISTICS INTRODUCTION	71
10	REGRESSION	73
11	CLASSIFICATION	81
12	MAKING INFERENCE ABOUT THE MODEL OF NATURE	103
12.1	Selecting the Robot's Model	103
12.2	Parameter Estimation	104
IV	REFLECTION	107
13	REFLECTIONS ON STATISTICAL PARADIGM	109

V	APPENDIX	111
A	HAMILTONIAN MONTE CARLO	113
B	NESTED SAMPLING	121
	BIBLIOGRAPHY	127

CHAPTER 1

Preface

Statistics is a mathematical discipline that use probability theory (which in turn require set theory) to extract insights from information (data). Probability theory is a branch of pure mathematics – probabilistic questions can be posed and solved using axiomatic reasoning, and therefore there is one correct answer to any probability question. Statistical questions can be converted to probability questions by the use of probability models. Given certain assumptions about the mechanism generating the data, statistical questions can be answered using probability theory. This highlights the dual nature of statistics, comprised of two integral parts.

1. The first part involves the formulation and evaluation of probabilistic models, a process situated within the realm of the philosophy of science. This phase grapples with the foundational aspects of constructing models that accurately represent the problem at hand.
2. The second part concerns itself with extracting answers after assuming a specific model. Here, statistics becomes a practical application of probability theory, involving not only theoretical considerations but also numerical analysis in real-world scenarios.

This duality underscores the interdisciplinary nature of statistics, bridging the gap between the conceptual and the applied aspects of probability theory. Although probabilities are well defined (see Chapter 3), their interpretation is not defined beyond their definition. This ambiguity has given birth to two

competing interpretations of probability, leading to two competing branches of statistics; Frequentist and Bayesian Statistics. This book aims to explain how these competing branches of statistics fit together as well as providing a non-exhaustive presentation of some of the methods within both branches. The philosophy of the book is rather straight to the point, but with a lot of examples both big and small. Some of these are anonymized versions of projects from industry. The book is split into three parts; introduction (Part i), Frequentist statistics (Part ii) and Bayesian statistics (Part iii).

1.1 ACKNOWLEDGEMENTS

The philosophy of the book is similar to [1], a few exercises from [2] used as examples, the idea of phrasing decision theory as "Robot vs Nature" from [3] and the review of probability theory is inspired by [4].

Part I

INTRODUCTION

CHAPTER 2

Introduction to Set Theory

Set theory is a fundamental branch of mathematical logic that provides a foundation for much of mathematics, including probability theory. At its core, set theory deals with the concept of a set, which is a collection of distinct objects or elements. In this introduction, the essential properties and operations of sets are explored in order to lay the groundwork for the axiomatic formation of probability theory and statistics.

Definition 1 (Membership). *In set theory, the membership relation between an object o and a set A is fundamental. $o \in A$ denotes that o is an element or member of A .*

Definition 2 (Set). *A set is a collection of distinct objects, considered as an object in its own right. Sets are typically denoted using curly braces $\{\}$ and can be described in two primary ways:*

1. *By listing its elements separated by commas, e.g., $A = \{a_1, a_2, a_3\}$.*
2. *By specifying a characterizing property of its elements, e.g., $A = \{x \mid x \text{ is a natural number}\}$.*

Sets can also be illustrated graphically, as shown in Figure 1.

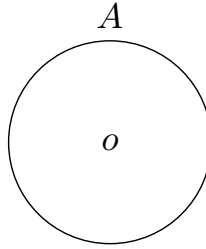


Figure 1: The graphical representation of a generic set A with generic elements o .

Definition 3 (Subset). *A set A is called a subset of a set B , denoted $A \subseteq B$, if every element of A is also an element of B . Formally, $A \subseteq B$ if $\forall x \in A, x \in B$. By this definition, a set is always a subset of itself.*

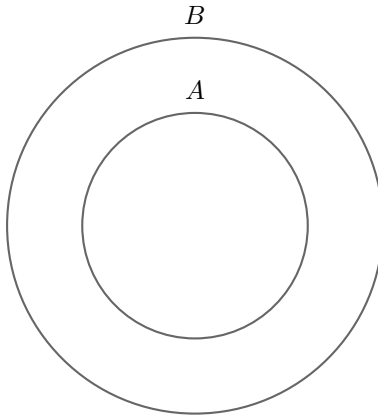


Figure 2: The graphical representation of $A \subseteq B$.

Definition 4 (Proper Subset). *A set A is called a proper subset of a set B , denoted $A \subset B$, if $A \subseteq B$ and $A \neq B$. This means that A is a subset of B but A is not equal to B ; there is at least one element in B that is not in A .*

Example 2.1.

Suppose $A = \{\text{🍌}, \text{🍏}, \text{🍇}\}$, then $\{\text{🍌}, \text{🍏}\}$ and $\{\text{🍏}\}$ are proper subsets of

A , meaning $\{\text{🍌}, \text{🍎}\}, \{\text{🍎}\} \subset A$. $\{\text{🍌}, \text{🍌}\}$, on the other hand, is not a subset of A , meaning $\{\text{🍌}, \text{🍌}\} \not\subset A$.

Example 2.2.

🍌, 🍎, and 🍇 are members (elements) of the set $\{\text{🍌}, \text{🍎}, \text{🍇}\}$, but are not subsets of it; and in turn, the subsets, such as $\{\text{🍌}\}$, are not members of the set $\{\text{🍌}, \text{🍎}, \text{🍇}\}$.

Definition 5 (Empty Set). The empty set, denoted by \emptyset or $\{\}$, is the set that contains no elements.

Definition 6 (Universal Set). The universal set, denoted by Ω , is the set that contains all the objects or elements under consideration in a particular discussion or problem. It is the largest set in the context of a given study.

Definition 7 (Closure). A set A is said to be closed under a certain operation if, for every pair of elements x and y in A , the result of applying the operation to x and y is also in A .

Definition 8 (Union). The union of sets A and B , denoted by $A \cup B$, is defined as the set containing all elements that are in A or B (or both). Figure 3 provide a graphical representation of $A \cup B$.

Definition 9 (Intersection). The intersection of sets A and B , denoted by $A \cap B$, is defined as the set containing all elements that are common to both A and B . Figure 4 provide a graphical representation of $A \cap B$.

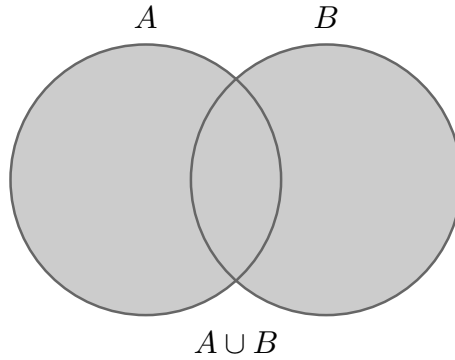


Figure 3: The figure show the union of sets A and B . Each circle represent the sets and the colored region represent the result of the result of the binary operation.

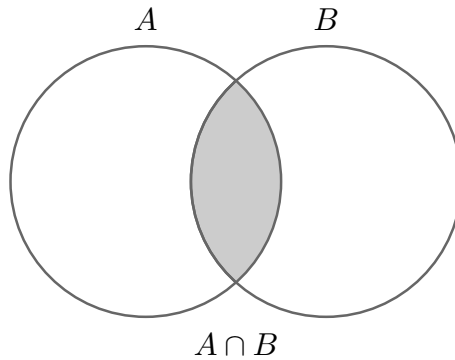


Figure 4: The figure show the intersection of sets A and B . Each circle represent the sets and the colored region represent the result of the result of the binary operation.

Definition 10 (Disjoint). *Two sets A and B are said to be disjoint if their intersection is the empty set, i.e., $A \cap B = \emptyset$. Figure 5 provide a graphical representation of $A \cap B = \emptyset$.*

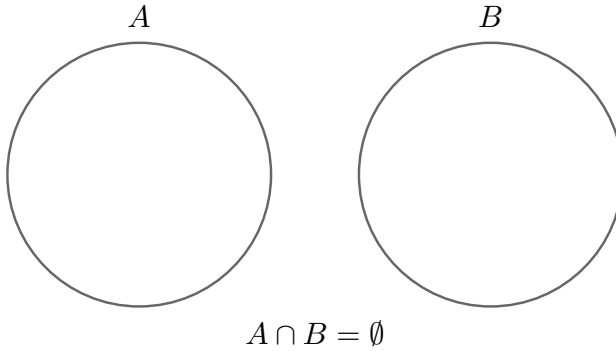


Figure 5: The figure show the case where the intersection of sets A and B is the empty set. Each circle represent the sets and the colored region represent the result of the result of the binary operation.

Definition 11 (Complementation). *The complement of set A , denoted by A^c , is defined as the set containing all elements in the universal set Ω that are not in A . Figure 6 provide a graphical representation of $(A \cap B)^c$.*

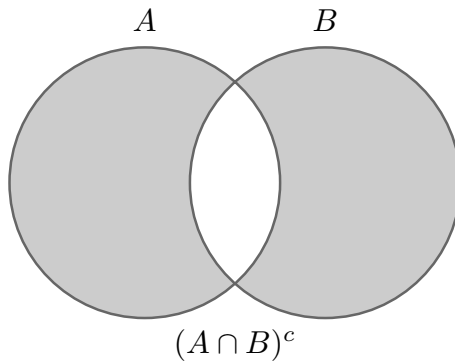


Figure 6: The figure show the complementary of the intersection of sets A and B . Each circle represent the sets and the colored region represent the result of the result of the binary operation.

Definition 12 (Difference). *The difference between set A and B , denoted by $A \setminus B = A \cap B^c$, is defined as the set containing all elements in A that are not in B . Figure 7 provide a graphical representation of $A \setminus B$ and $B \setminus A$.*

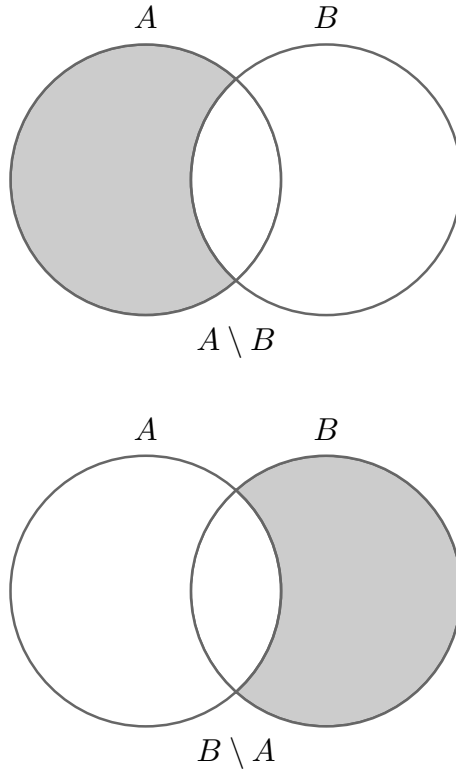


Figure 7: (left) show A minus B and (right) show B minus A . Each circle represent the sets and the colored region represent the result of the result of the binary operation.

Definition 13 (Power Set). *The power set of a set A , denoted by 2^A , is defined as the set containing all possible subsets of A , including A itself and the empty set.*

Example 2.3.

Suppose $A = \{a_1, a_2, a_3\}$, then

$$2^A = \{\emptyset, \{a_1\}, \{a_2\}, \{a_3\}, \{a_1, a_2\}, \{a_1, a_3\}, \{a_2, a_3\}, \{a_1, a_2, a_3\}\}. \quad (1)$$

Definition 14 (Symmetric Difference). *The symmetric difference of sets A and B , denoted by $A \Delta B$, is defined as the set containing all elements that are in either A or B but not in both, meaning $A \Delta B = (A \cap B)^c$. Figure 6 show the symmetric difference between sets A and B .*

Definition 15 (Finite and Infinite Unions). *For a collection $\{A_i\}$, the union is denoted by $\bigcup_i A_i$ and is defined as the set containing all elements that are in at least one of the sets A_i .*

Definition 16 (Partition). *A collection of non-empty subsets $\{A_i\}$ of a set A is called a partition of A if the following conditions are satisfied:*

1. *The subsets A are pairwise disjoint, i.e., $A_i \cap A_j = \emptyset$ for all $i \neq j$.*
2. *The union of all subsets A_i is equal to the set A , i.e., $\bigcup_{i \in I} A_i = A$.*

A graphical representation of the set $A = \{A_1, A_2, A_3\}$, where A_j are partitions, is shown in Figure 8.

Definition 17 (Finite and Infinite Intersections). *For a collection $\{A_i\}$, the intersection is denoted by $\bigcap_i A_i$ and is defined as the set containing all elements that are common to all sets A_i .*

Definition 18 (Cartesian Product). *The Cartesian product of sets A and B , denoted by $A \times B$, is defined as the set containing all ordered pairs (a, b) , where a is in A and b is in B .*

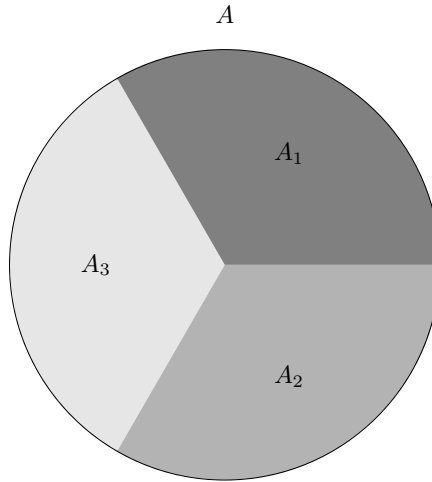


Figure 8: The figure show $A = \{A_1, A_2, A_3\}$ where A_j are partitions.

Example 2.4.

Suppose $A = \{a_1, a_2\}$ and $B = \{b_1, b_2, b_3\}$, then

$$A \times B = \{(a_1, b_1), (a_1, b_2), (a_1, b_3), (a_2, b_1), (a_2, b_2), (a_2, b_3)\} \quad (2)$$

CHAPTER 3

Introduction to Probability Theory

Probability theory aims to provide a mathematical framework for analyzing random experiments, where outcomes cannot be predicted with certainty beforehand. Its objective is to systematically study and understand the potential outcomes of these experiments.

Definition 19 (Sample Space). *The sample space, denoted by Ω , represents the set of all possible outcomes of a random experiment. It encompasses every conceivable result that could occur, serving as the foundation for analyzing probabilities associated with different outcomes.*

Definition 20 (Event). *An event, E , is a subset of the sample space, denoted by $E \subseteq \Omega$, that corresponds to a specific collection of possible outcomes in a random experiment. Events may consist of single or multiple outcomes and are defined by the occurrence or non-occurrence of particular conditions.*

Example 3.1.

Consider the roll of a fair six-sided die. The sample space for this experiment is given by $\Omega = \{\square, \blacksquare, \blacklozenge, \boxtimes, \boxplus, \boxminus\}$. $E = \{\blacksquare, \boxtimes, \boxplus\}$, is the event of rolling an even number.

Definition 21 (Event Space). *The set containing all valid possible events for a random experiment is referred to as the event space, \mathcal{F} . The notion of "all valid possible events for a random experiment" is formally defined by requiring \mathcal{F} to be a σ -algebra satisfying the following properties:*

1. \mathcal{F} is the set of all subsets of the sample space Ω , including the empty set \emptyset and Ω itself, along with various combinations of outcomes.
2. Closure under complementation: If E is in the σ -algebra ($E \in \mathcal{F}$), then its complement E^c is also in the σ -algebra.
3. Closure under countable union and intersection: If the events E_1, E_2, E_3, \dots are in the σ -algebra ($E_i \in \mathcal{F}$ for all i), then their countable union $\bigcup_{i=1}^{\infty} E_i$ and intersection $\bigcap_{i=1}^{\infty} E_i$ are also in the σ -algebra.

In the case where the outcomes of the random experiment can take discrete values, these properties are sufficient. However, in the case where the outcomes are continuous, \mathcal{F} is required to be a Borel σ -algebra, meaning it must further fulfill the closure property under countable intersection with open sets. This ensures that \mathcal{F} contains all sets that can be formed by taking unions, intersections, and complements of open sets, which are essential for defining probabilities in continuous spaces.

Example 3.2.

For the roll with the fair die considered in Example 3.1, the sample space is $\Omega = \{\square, \blacksquare, \boxtimes, \boxplus, \boxminus, \boxdot\}$ and the event space (the set of all possible events) is given by

$$\begin{aligned} \mathcal{F} &= \{\emptyset, \{\square\}, \{\square, \blacksquare\}, \{\square, \boxtimes\}, \{\square\}, \{\square, \blacksquare, \boxtimes, \boxplus\}, \{\boxplus\}, \dots\} \\ &= 2^{\Omega}. \end{aligned} \tag{3}$$

Definition 22 (Measurable Space). The pair (Ω, \mathcal{F}) is called a measurable space.

Probability can loosely be defined [4] as a measure of the size of an event (a set) relative to the sample space (another set), meaning it is a function that operates on an event (a set). In particular the probability measure maps any valid event, i.e. any $E \in \mathcal{F}$, to a number between 0 and 1, representing the relative size of the event to the sample space.

Definition 23 (Probability Measure). *A Probability measure, \mathbb{P} , is a set function defined on a measurable space (Definition 22) (Ω, \mathcal{F})*

$$\mathbb{P} : \mathcal{F} \mapsto [0, 1] \quad (4)$$

that obey [5] Axiom 1-Axiom 3.

Axiom 1 (Non-negativity). *For any event $E \in \mathcal{F}$, the probability measure $\mathbb{P}(E)$ is non-negative, satisfying*

$$\mathbb{P}(E) \geq 0 \quad \forall E \in \mathcal{F}. \quad (5)$$

Axiom 2 (Normalization). *The probability of the universal set Ω is 1, satisfying*

$$\mathbb{P}(\Omega) = 1. \quad (6)$$

Axiom 3 (Additivity). *For any countable sequence of mutually exclusive events $E_1, E_2, \dots \in \mathcal{F}$, the probability of their union is the sum of their individual probabilities, such that*

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} E_i\right) = \sum_{i=1}^{\infty} \mathbb{P}(E_i) \quad \forall E_i \in \mathcal{F} \text{ where } \bigcap_{i=1}^{\infty} E_i = \emptyset. \quad (7)$$

Together, the probability measure, the sample space and the algebra form the tuple $(\Omega, \mathcal{F}, \mathbb{P})$ which define what a probability space. The non-negativity and normalization axioms are largely matters of convention, although it is non-trivial that probability measures take at least the two values 0 and 1, and that they have a maximal value (unlike various other measures, such as length, volume, and so on, which are unbounded). The axioms are supplemented by two definitions.

Definition 24 (Conditional Probability). *For events E_1 and E_2 in a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ with $\mathbb{P}(E_2) > 0$, the conditional probability of E_1 given E_2 is defined viz*

$$\mathbb{P}(E_1|E_2) \equiv \frac{\mathbb{P}(E_1, E_2)}{\mathbb{P}(E_2)}, \quad (8)$$

where $\mathbb{P}(E_1, E_2) = \mathbb{P}(E_1 \cap E_2)$ to ease the notation.

Definition 25 (Independence). *Events E_1 and E_2 in a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ are said to be conditionally independent if*

$$\mathbb{P}(E_1, E_2) = \mathbb{P}(E_1)\mathbb{P}(E_2). \quad (9)$$

From Axiom 1-Axiom 3, Definition 8 and Definition 9, the chain rule, the concept of marginalization, conditional independence and the law of total probability can be derived.

Theorem 1 (Chain Rule). *Given $\{E_1, E_2, \dots, E_n\} \subseteq \mathcal{F}$ denotes a set of events in a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, the chain rule for this set of events can be written*

$$\mathbb{P}(E_1, \dots, E_n) = \mathbb{P}(E_1) \prod_{j=2}^n \mathbb{P}(E_j | E_1, \dots, E_{j-1}). \quad (10)$$

Proof. From the definition of conditional probability in Equation 8

$$\mathbb{P}(E_1, E_2, \dots, E_n) = \mathbb{P}(E_1 | E_2, \dots, E_n) \mathbb{P}(E_2, \dots, E_n). \quad (11)$$

Using the definition of conditional probability again

$$\mathbb{P}(E_2, \dots, E_n) = \mathbb{P}(E_2 | \dots, E_n) \mathbb{P}(\dots, E_n). \quad (12)$$

Continuing in this way, Equation 10 follows. \square

Equation 10 illustrates how to decompose the joint probability of multiple events into a product of conditional probabilities. The idea is to calculate the probability of each event in the sequence conditioned on the occurrence of the previous events in the chain. The chain rule is particularly powerful when dealing with complex systems where events may be interdependent. It allows breaking down joint probabilities into more manageable conditional probabilities, making it easier to analyze and model intricate relationships between events. Whether in the context of statistical modeling or machine learning, the chain rule plays a key role in calculating the joint probability of multiple events and provides a foundation for more advanced probabilistic reasoning.

Theorem 2 (Bayes theorem). *For events $E_1, E_2, E_3 \in \mathcal{F}$ in a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, Bayes theorem can be formulated viz*

$$\mathbb{P}(E_1|E_2, E_3) = \frac{\mathbb{P}(E_2|E_1, E_3)\mathbb{P}(E_1, E_2)}{\mathbb{P}(E_2, E_3)}. \quad (13)$$

Proof. Bayes theorem follows directly from applying the chain rule and applying the concept of symmetry viz

$$\begin{aligned} \mathbb{P}(E_1, E_2, E_3) &= \mathbb{P}(E_1|E_2, E_3)\mathbb{P}(E_2, E_3) \\ &= \mathbb{P}(E_2|E_1, E_3)\mathbb{P}(E_1, E_3) \end{aligned} \quad (14)$$

from which

$$\mathbb{P}(E_1|E_2, E_3) = \frac{\mathbb{P}(E_2|E_1, E_3)\mathbb{P}(E_1, E_2)}{\mathbb{P}(E_2, E_3)} \quad (15)$$

which is Bayes theorem. \square

Theorem 3 (Law of Total Probability). *Let $\{E_1, E_2, \dots, E_n\}$ be a partition of the sample space Ω of the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, then for any $A \subseteq \Omega$,*

$$\mathbb{P}(A) = \sum_i \mathbb{P}(A, E_i). \quad (16)$$

In continuous cases, the summation is replaced by integration.

Proof. Consider an event $A \subseteq \Omega$ and a partition $\{E_1, E_2, \dots, E_n\}$ of Ω such that $\cup_i E_i = \Omega$. For mutually exclusive events (which a partition by definition is), finite additivity can be used such that

$$\sum_i \mathbb{P}(A, E_i) = \mathbb{P}(\bigcup_i (A, E_i)). \quad (17)$$

$\bigcup_i (A, E_i)$ is the union of all intersections between A and the E 's. However, since the E 's form a partition of Ω , they together form Ω and the intersection between Ω and A is A , meaning

$$\begin{aligned} \bigcup_i (A, E_i) &= (A, \bigcup_i E_i) \\ &= (A, \Omega) \\ &= A. \end{aligned} \quad (18)$$

Combining Equation 17–Equation 18 then yields

$$\mathbb{P}(A) = \sum_i \mathbb{P}(A, E_i). \quad (19)$$

□

Example 3.3.

Consider the roll of a fair six-sided die. The sample space for this experiment is given by $\Omega = \{\square, \blacksquare, \blacklozenge, \blacksquare, \blacksquare, \blacksquare\}$. Let $E_1 = \{\blacksquare, \blacksquare, \blacksquare\}$ and $E_2 = \{\blacksquare\}$ be two events, then from Equation 8

$$\begin{aligned} \mathbb{P}(E_1|E_2) &= \frac{\mathbb{P}(E_1, E_2)}{\mathbb{P}(E_2)} \\ &= 1 \end{aligned} \quad (20)$$

where $\mathbb{P}(E_1, E_2) = \frac{1}{6}$ since $E_1, E_2 = E_1 \cap E_2 = E_2 = \{\blacksquare\}$ is one of 6 possible values and $\mathbb{P}(E_2) = \frac{1}{6}$. Intuitively this makes sense because E_2 is a set with one member and since E_2 is known, the outcome of the experiment is known with certainty in this case.

Definition 26 (Random Variable). A random variable X is a function

$$X : \Omega \mapsto \Omega_X \quad (21)$$

that maps outcomes from a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ to a measurable space $(\Omega_X, \mathcal{F}_X)$, where Ω_X is the codomain of X and \mathcal{F}_X is a σ -algebra on Ω_X . The σ -algebra \mathcal{F}_X ensures that X is measurable, meaning that for any set $x \in \mathcal{F}_X$, the preimage $X^{-1}(x)$ must belong to \mathcal{F} . Formally, this can be written as

$$X^{-1}(x) = \{\omega \in \Omega | X(\omega) = x\} \in \mathcal{F} \quad \forall x \in \mathcal{F}_X. \quad (22)$$

Random variables are classified as either discrete or continuous, based on the discrete or continuous nature of their sample space. Discrete random variables have countable sample spaces, while continuous random variables have uncountable sample spaces, often modeled as intervals on the real line. The role of random variables is to provide

a numerical representation of the outcomes of a random experiment, allowing quantification and analysis of the likelihood of different numerical outcomes.

Definition 27 (Expected value). Let X be a real-valued random variable defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, then the expected value of X , denoted by $\mathbb{E}[X]$, is defined by the Lebesgue integral [6]

$$\mathbb{E}_X[X] \equiv \int_{\Omega} X(\omega) d\mathbb{P}(\{\omega\}). \quad (23)$$

Theorem 4 (Non-negativity of expected value). If $X \geq 0$ for a random variable X , then $\mathbb{E}_X[X] \geq 0$.

Theorem 5 (Linearity of expected value). The expectation is a linear operator meaning $\mathbb{E}_X[a + X] = a + \mathbb{E}_X[X]$ and $\mathbb{E}_X[aX] = a\mathbb{E}_X[X]$ for any constant a .

Theorem 6 (The law of the unconscious statistician). The law of the unconscious statistician generalize the expectation of a random variable to the expectation of a function $g : \Omega \mapsto \mathbb{R}$ of a random variable $X(\omega) \in \Omega_X \quad \forall \omega \in \Omega$ defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ such that

$$\mathbb{E}[g(X)] \equiv \int_{\Omega} g(X(\omega)) d\mathbb{P}(\{\omega\}). \quad (24)$$

Definition 28 (Variance). Let X be a real-valued random variable defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, then the variance of X , denoted by $\text{Var}[X]$, is defined viz

$$\begin{aligned} \text{Var}[X] &\equiv \mathbb{E}_X[(X - \mathbb{E}_X[X])^2] \\ &= \mathbb{E}_X[X^2] - \mathbb{E}_X[X]^2. \end{aligned} \quad (25)$$

Theorem 7 (Non-linearity of variance). The variance is a non-linear operator, where $\text{Var}[a + X] = \text{Var}[X]$ and $\text{Var}[aX] = a^2 \text{Var}[X]$ for any constant a .

Definition 29 (Image Measure). Let $X : \Omega \mapsto \Omega_X$ be a random variable that maps from the probability space $(\Omega, \mathcal{F}, \mathbb{P})$ to a measurable space $(\Omega_X, \mathcal{F}_X)$. Then [7]

$$\mathbb{P} \circ X^{-1} : \mathcal{F}_X \mapsto [0, 1] \quad (26)$$

defines a probability measure on $(\Omega_X, \mathcal{F}_X)$. $\mathbb{P} \circ X^{-1}$ is called the image measure or the push forward measure of \mathbb{P} .

Definition 30 (Probability Mass Function). In case of a discrete random variable $X : \Omega \mapsto \Omega_X$ that maps from the probability space $(\Omega, \mathcal{F}, \mathbb{P})$ to a measurable space $(\Omega_X, \mathcal{F}_X)$, the image measure is defined as the probability mass function

$$\begin{aligned} p(X = x) &\equiv \mathbb{P} \circ X^{-1}(x) \\ &= \mathbb{P}(X^{-1}(x)). \end{aligned} \quad (27)$$

According to Axiom 1-Axiom 3 $\sum_{all\ x} p(X = x) = 1$ and $p(X = x) \geq 0 \quad \forall x \in \Omega_X$.

Theorem 8 (Expected value of discrete random variable). The expected value of a discrete random variable X with probability mass function p can be written

$$\mathbb{E}_X[X] = \sum_i x_i p(X = x_i). \quad (28)$$

Definition 31 (Probability Density Function). Let $X : \Omega \mapsto \Omega_X$ be a continuous random variable that maps from the probability space $(\Omega, \mathcal{F}, \mathbb{P})$ to a measurable space $(\Omega_X, \mathcal{F}_X)$. Probabilities are described using a probability density function (PDF)

$$f : \Omega_X \mapsto \mathbb{R}_{\geq 0}, \quad (29)$$

which is related to the probability measure viz

$$\mathbb{P}(\{\omega \in \Omega | X(\omega) \leq x\}) = \int_{-\infty}^x f(X = t) dt, \quad (30)$$

where $\int_{-\infty}^{\infty} f(X = t) dt = 1$ and $f(X = x) = 0$ for any individual point $x \in \mathbb{R}$.

Theorem 9 (Expected value of continuous random variable). *The expected value of a continuous random variable X with probability density function f can be written*

$$\mathbb{E}_X[X] = \int_{\Omega_X} xf(X=x)dx. \quad (31)$$

Theorem 10 (Total expectation). *The expectation of a random variable X can be expressed in terms of another random variable Y viz*

$$\mathbb{E}_X[X] = \mathbb{E}_Y[\mathbb{E}_{X|Y}[X|Y=y]], \quad (32)$$

where the subscript specify the probability distribution the expectation is with respect to.

Proof.

$$\begin{aligned} \mathbb{E}_X[X] &= \int_{\Omega_X} dx xf(X=x) \\ &= \int_{\Omega_Y} dy \int_{\Omega_X} dx xf(X=x, Y=y) \\ &= \int_{\Omega_Y} dy f(Y=y) \int_{\Omega_X} dx xf(X=x|Y=y) \\ &= \int_{\Omega_Y} dy f(Y=y) \underbrace{\int_{\Omega_X} dx xf(X=x|Y=y)}_{=\mathbb{E}_{X|Y}[X|Y=y]} \\ &= \mathbb{E}_Y[\mathbb{E}_{X|Y}[X|Y=y]]. \end{aligned} \quad (33)$$

□

Theorem 11 (Expectation of product of independent random variables). *Let $X(\omega) \in \Omega_X$ and $Y(\omega) \in \Omega_Y$ be independent continuous random variables defined on the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, such that $f(X=x, Y=y) = f(X=x)f(Y=y)$, then $\mathbb{E}_{XY}[XY] = \mathbb{E}_X[X]\mathbb{E}_Y[Y]$.*

Proof.

$$\begin{aligned}
 \mathbb{E}_{XY}[XY] &= \int_{\Omega_X} \int_{\Omega_Y} xyf(X=x, Y=y)dxdy \\
 &= \int_{\Omega_X} xf(X=x)dx \int_{\Omega_Y} yf(Y=y)dy \\
 &= \mathbb{E}_X[X]\mathbb{E}_Y[Y]
 \end{aligned} \tag{34}$$

□

Definition 32 (Covariance). *Let X and Y be a real-valued random variables defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, then the covariance of X and Y , denoted by $\text{Cov}[X, Y]$, is defined viz*

$$\begin{aligned}
 \text{Cov}[X, Y] &= \mathbb{E}_{XY}[(X - \mathbb{E}_X[X])(Y - \mathbb{E}_Y[Y])] \\
 &= \mathbb{E}_{XY}[XY] - \mathbb{E}_X[X]\mathbb{E}_Y[Y],
 \end{aligned} \tag{35}$$

Theorem 12 (Covariance of independent random variables). *For independent random variables $X \in \Omega_X$ and $Y \in \Omega_Y$ the covariance is given by $\text{Cov}[X, Y] = 0$.*

Proof. Using $\mathbb{E}_{XY}[XY] = \mathbb{E}_X[X]\mathbb{E}_Y[Y]$ (Theorem 11) in Definition 32 yield $\text{Cov}[X, Y] = 0$. □

Definition 33 (Correlation). *Let X and Y be real-valued random variables defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. The correlation between X and Y , denoted by $\text{Corr}[X, Y]$, is defined as*

$$\begin{aligned}
 \text{Corr}[X, Y] &= \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X]\text{Var}[Y]}} \\
 &= \frac{\mathbb{E}_{XY}[XY] - \mathbb{E}_X[X]\mathbb{E}_Y[Y]}{\sqrt{(\mathbb{E}_X[X^2] - \mathbb{E}_X[X]^2)(\mathbb{E}_Y[Y^2] - \mathbb{E}_Y[Y]^2)}}.
 \end{aligned} \tag{36}$$

Correlation and covariance are both measures of the relationship between two random variables. While covariance indicates the extent to which two variables change together, correlation provides a standardized measure of this relationship, taking into account the scales of the variables. In particular, the correlation between two variables, denoted by $\text{Corr}[X, Y]$, is the

covariance of X and Y divided by the product of their standard deviations. This normalization makes correlation a unitless quantity that ranges between -1 and 1 , where -1 indicates a perfect negative linear relationship, 1 indicates a perfect positive linear relationship, and 0 indicates no linear relationship. In essence, correlation provides a more interpretable measure of the strength and direction of the linear association between two variables compared to covariance.

Definition 34 (Change of Variables for PDFs). *Let X be a continuous random variable with probability density function (PDF) $f(X = x)$, defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. Suppose $Y = g(X)$ is a continuous and differentiable function of X , and let g^{-1} denote the inverse function of g . If $Y = g(X)$ and the inverse function g^{-1} exists and is differentiable, the PDF of the random variable Y , denoted $f(Y = y)$, can be obtained by the change of variables formula viz [1]*

$$f(Y = y) = f(X = g^{-1}(y)) \left| \frac{d}{dY} (g^{-1}(Y)) \right|_{Y=y}. \quad (37)$$

Example 3.4.

Let X be a continuous random variable with PDF $f(X = x)$, and let $Y = g(X) = aX + b$, where $a \neq 0$ and b are constants. The inverse function is given by

$$g^{-1}(y) = \frac{y - b}{a} \quad (38)$$

Using Definition 34

$$\begin{aligned} f(Y = y) &= f(X = g^{-1}(y)) \left| \frac{d}{dY} (g^{-1}(Y)) \right|_{Y=y} \\ &= f\left(X = \frac{y - b}{a}\right) \left| \frac{d}{dY} \left(\frac{Y - b}{a}\right) \right|_{Y=y} \\ &= f\left(X = \frac{y - b}{a}\right) \left| \frac{1}{a} \right|. \end{aligned} \quad (39)$$

Thus, the PDF of Y is

$$f(Y = y) = \frac{1}{|a|} f\left(X = \frac{y - b}{a}\right). \quad (40)$$

Example 3.5.

Let $X = \ln\left(\frac{Y}{1-Y}\right)$ be a continuous random variable with PDF $f(X = x) \propto \text{const}$. The inverse function is given by

$$g^{-1}(y) = \ln\left(\frac{y}{1-y}\right). \quad (41)$$

Using Definition 34

$$\begin{aligned} f(Y = y) &= f\left(X = g^{-1}(y)\right) \left| \frac{d}{dY} \left(g^{-1}(Y)\right) \right|_{Y=y} \\ &= \text{const} \cdot \frac{1-Y}{Y} \left(\frac{1}{1-Y} + \frac{Y}{(1-Y)^2} \right) \Big|_{Y=y} \\ &= \text{const} \cdot Y^{-1}(1-Y)^{-1} \Big|_{Y=y} \\ &= \text{Beta}(Y = y | a = 0, b = 0). \end{aligned} \quad (42)$$

Definition 35 (Error-Propagation). Let X_1, \dots, X_n be continuous random variables with means $\mathbb{E}[X_1], \dots, \mathbb{E}[X_n]$ and variances denoted $\text{Var}[X_1], \dots, \text{Var}[X_n]$, defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. Consider a function $g(X_1, \dots, X_n)$ of these random variables. The variance of g , which quantifies the uncertainty in g due to the uncertainties in the X_i , can be written

$$\begin{aligned} \text{Var}[g(X_1, \dots, X_n)] &\equiv \mathbb{E}[(g(X_1, \dots, X_n) - \mathbb{E}[g(X_1, \dots, X_n)])^2] \\ &= \mathbb{E}[g(X_1, \dots, X_n)^2] - (\mathbb{E}[g(X_1, \dots, X_n)])^2. \end{aligned} \quad (43)$$

In practice, the variance is often analytically intractable, in which case, a linear approximation of the variance can be used. This involves expanding g about the means using a first-order Taylor series expansion around the point $= \{X_1 = \mathbb{E}[X_1], \dots, X_n = \mathbb{E}[X_n]\}$

$$g(X_1, \dots, X_n) = g(\text{point}) + \sum_{i=1}^n (X_i - \mu_i) \frac{\partial g}{\partial X_i} \Big|_{\text{point}} + \mathcal{O}(\partial^2 g) \quad (44)$$

with this approximation, the variance of g can be approximated viz

$$\begin{aligned}
 \text{Var}[g] &= \mathbb{E}[(g - \mathbb{E}[g])^2] \\
 &= \mathbb{E}\left[\left(\sum_{i=1}^n (X_i - \mu_i) \frac{\partial g}{\partial X_i} \Big|_{\text{point}} + \mathcal{O}(\partial^2 f)\right)^2\right] \\
 &= \sum_{i=1}^n \left(\frac{\partial g}{\partial X_i} \Big|_{\text{point}}\right)^2 \text{Var}[X_i] + \sum_{i \neq j} \frac{\partial g}{\partial X_i} \frac{\partial g}{\partial X_j} \Big|_{\text{point}} \text{Cov}[X_i, X_j] \\
 &\quad + \mathcal{O}(\partial^2 g).
 \end{aligned} \tag{45}$$

where it has been used that $\mathbb{E}[g] = g(\text{point}) + \mathcal{O}(\partial^2 g)$ and it is understood all derivatives are evaluated at the means of the random variables. When the random variables are independent, $\text{Cov}[X_i, X_j] = 0$ for all $i \neq j$ (see Theorem 12), and the formula simplifies to

$$\text{Var}[g] \approx \sum_{i=1}^n \left(\frac{\partial g}{\partial X_i} \Big|_{\text{point}}\right)^2 \text{Var}[X_i]. \tag{46}$$

Example 3.6.

A company produce square plates. Let the plate dimensions be characterized by two independent random variables $X \sim \mathcal{N}(2m, (0.01m)^2)$ and $Y \sim \mathcal{N}(3m, (0.02m)^2)$ and the area given by XY . Determine the variance of XY . From Definition 35, the exact variance is

$$\begin{aligned}
 \text{Var}[XY] &= \mathbb{E}[(XY)^2] - (\mathbb{E}[XY])^2 \\
 &= \left(\text{Var}[X] + \mathbb{E}[X]^2\right) \left(\text{Var}[Y] + \mathbb{E}[Y]^2\right) - \mathbb{E}[X]^2 \mathbb{E}[Y]^2 \\
 &= \mathbb{E}[Y]^2 \text{Var}[X] + \mathbb{E}[X]^2 \text{Var}[Y] + \text{Var}[X] \text{Var}[Y]
 \end{aligned} \tag{47}$$

where it has been used that X and Y are independent, such that $\mathbb{E}[(XY)^2] = \mathbb{E}[X^2] \mathbb{E}[Y^2]$. Via the linear approximation

$$\begin{aligned}
 \text{Var}[XY] &\approx \sum_{i=X,Y} \left(\frac{\partial(XY)}{\partial i} \Big|_{X=\mu_X, Y=\mu_Y}\right)^2 \text{Var}[i] \\
 &= \mathbb{E}[Y]^2 \text{Var}[X] + \mathbb{E}[X]^2 \text{Var}[Y]
 \end{aligned} \tag{48}$$

Comparing Equation 47 and Equation 48 the relative difference can be written

$$\frac{\text{Var}[XY] - \text{Var}[XY]|_{\text{linear approximation}}}{\text{Var}[XY]} = \frac{\text{Var}[X]\text{Var}[Y]}{\text{Var}[XY]} \quad (49)$$

$$\simeq 1.6 \cdot 10^{-5}.$$

Example 3.7.

Consider a thought experiment in which a father with amnesia is told he has two children, but does not know the sex of them. The sample space can be constructed from the sample space for each child

$$\begin{aligned} \Omega_{\text{child } 1} &= \{(\text{♂}, \text{♂}), (\text{♂}, \text{♀}), (\text{♀}, \text{♂}), (\text{♀}, \text{♀})\}, \\ \Omega_{\text{child } 2} &= \{(\text{♂}, \text{♂}), (\text{♂}, \text{♀}), (\text{♀}, \text{♂}), (\text{♀}, \text{♀})\} \end{aligned} \quad (50)$$

such that

$$\begin{aligned} \Omega &= \Omega_{\text{child } 1} \times \Omega_{\text{child } 2} \\ &= \{(\text{♂}, \text{♂}), (\text{♂}, \text{♀}), (\text{♀}, \text{♂}), (\text{♀}, \text{♀})\}. \end{aligned} \quad (51)$$

Assuming the sex of a child is like a coin flip, it is most likely, *a priori*, that the father has one boy and one girl with probability $\frac{1}{2}$, i.e. $\mathbb{P}(\{(\text{♂}, \text{♀})\}) = \frac{1}{2}$. The other possibilities (two boys or two girls) have probability $\frac{1}{4}$, meaning $\mathbb{P}(\{(\text{♂}, \text{♂})\}) = \frac{1}{4}$ and $\mathbb{P}(\{(\text{♀}, \text{♀})\}) = \frac{1}{4}$. In order to simplify the formalism, define the random variables $B : \Omega \mapsto \{0, 1, 2\}$ and $G : \Omega \mapsto \{0, 1, 2\}$ that maps the events in \mathcal{F} to a number of boys $B(E) \forall E \in \mathcal{F}$ and girls $G(E) \forall E \in \mathcal{F}$. The probability mass function associated to B and G is given by Equation 27, such that e.g.

$$p(B = 1, G = 1) = \mathbb{P}(\{(\text{♂}, \text{♀})\}). \quad (52)$$

1. Suppose the father ask his wife whether he has any boys, and she says yes. What is the probability that one child is a girl?

The exact framing of the question is important here; "any boys" means "at least one boy"

$$p(G = 1, B \geq 1) = \frac{p(B \geq 1|G = 1)p(G = 1)}{p(B \geq 1)}. \quad (53)$$

Given the father has two children, if he has exactly one girl, then the other must be a boy, so $p(B \geq 1|G = 1) = 1$. $p(G = 1) = \frac{1}{2}$ since it is a priori assumed to be equally likely to be a boy or girl. $p(B \geq 1) = 1 - p(G = 2, B = 0) = \frac{3}{4}$, so

$$p(G = 1|B \geq 1) = \frac{2}{3}. \quad (54)$$

2. Suppose instead the father meets one of his children and it is a boy. What is the probability that the other is a girl?

Since one child is known to be a boy, what is asked about is $p(G = 1|B = 1) = \frac{1}{2}$.

Example 3.8.

Suppose a crime has been committed. Blood is found at the crime scene for which there is no innocent explanation. It is of the type which is present in 1% of the population.

1. The prosecutor claims: "There is a 1% chance that the defendant would have the crime blood type if he were innocent. Thus there is a 99% chance that he is guilty". This is known as the prosecutors fallacy. What is wrong with this argument?

Let E denote the event of having the blood type found at the crime scene, then "there is a 1% chance that the defendant would have the crime blood type if he were innocent" means

$$\mathbb{P}(E|\text{innocent}) = 0.01. \quad (55)$$

This is not the relevant quantity, rather

$$\mathbb{P}(\text{innocent}|E) = \frac{\mathbb{P}(E|\text{innocent})\mathbb{P}(\text{innocent})}{p(E)}. \quad (56)$$

Since

$$\mathbb{P}(\text{innocent}|E) + \mathbb{P}(\text{guilty}|E) = 1 \quad (57)$$

and so $\mathbb{P}(\text{innocent}|E) = 0.01$ means $\mathbb{P}(\text{guilty}|E) = 0.99$, which is what is stated in the exercise, however, in general $\mathbb{P}(E|\text{innocent}) \neq \mathbb{P}(\text{innocent}|E)$.

2. The defender claims: "The crime occurred in a city of 800 000 people. Hence, the blood type found at the crime scene would be found in $800\,000 \cdot 0.01 = 8\,000$ people". The evidence has thus provided a probability of $\frac{1}{8\,000}$ that the defendant is guilty, and therefore has no relevance". This is known as the defendants fallacy. What is wrong with this argument?

$$\mathbb{P}(\text{guilty}|E) = \frac{\mathbb{P}(E|\text{guilty})\mathbb{P}(\text{guilty})}{\mathbb{P}(E)}, \quad (58)$$

with $\mathbb{P}(E|\text{guilty}) = 1$, $\mathbb{P}(\text{guilty}) = \frac{1}{8\,000}$ and

$$\mathbb{P}(E) = \mathbb{P}(E|\text{guilty})\mathbb{P}(\text{guilty}) + \mathbb{P}(E|\text{innocent})p(\text{innocent}) \quad (59)$$

where $\mathbb{P}(E|\text{innocent}) = 0.01$ and $\mathbb{P}(\text{innocent}) = 1 - \mathbb{P}(\text{guilty})$, meaning

$$\mathbb{P}(\text{guilty}|E) = \frac{100}{800\,099}. \quad (60)$$

$\frac{100}{800\,099}$ is very close to $\frac{1}{8\,000}$, however, this assumes the only evidence against the defendant is the blood type found at the crime scene. If this changes, the calculation can change significantly, depending on the evidence.

Example 3.9.

Show that the variance of a sum is $\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y]$.

$$\begin{aligned} \text{Var}[X + Y] &= \mathbb{E}_{XY}[(X + Y - \mathbb{E}_{XY}[X + Y])^2] \\ &= \mathbb{E}_X[(X - \mathbb{E}_X[X])^2] + \mathbb{E}_Y[(Y - \mathbb{E}_Y[Y])^2] \\ &\quad + 2\mathbb{E}_{XY}[(X - \mathbb{E}_X[X])(Y - \mathbb{E}_Y[Y])] \\ &= \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y]. \end{aligned} \quad (61)$$

Example 3.10.

After your yearly checkup, the doctor has bad news and good news. The bad news is that you tested positive for a serious disease, and that the test is 99% accurate (i.e. the probability of testing positive given that you have the disease is 99%, as is the probability of testing negative given that you don't have the disease). The good news is that this is a rare disease, striking only one in 10 000 people. What are the chances that you actually have the disease?

Let "s" denote the event of being sick, "h" the event of being healthy, "p" the event of a positive test and "n" the event of a negative test, then

$$\begin{aligned}\mathbb{P}(s|p) &= \frac{\mathbb{P}(p|\text{sick})\mathbb{P}(s)}{\mathbb{P}(p)} \\ &= \frac{\mathbb{P}(p|s)\mathbb{P}(s)}{\mathbb{P}(p|s)\mathbb{P}(s) + \mathbb{P}(p|h)\mathbb{P}(h)}\end{aligned}\tag{62}$$

where $\mathbb{P}(p|s) = 0.99$, $\mathbb{P}(s) = \frac{1}{10\,000}$, $\mathbb{P}(p|h) = 1 - \mathbb{P}(n|h)$, $\mathbb{P}(n|h) = 0.99$ and $\mathbb{P}(h) = 1 - \mathbb{P}(s)$. This means

$$\mathbb{P}(s|p) \simeq 0.0098.\tag{63}$$

Example 3.11.

On a game show, a contestant is told the rules as follows: There are 3 doors labeled 1, 2, 3. A single prize has been hidden behind one of them. You get to select one door. Initially your chosen door will not be opened, instead, the gameshow host will open one of the other two doors in such a way as not to reveal the prize. For example, if you first choose door 1, the gameshow host will open one of doors 2 and 3, and it is guaranteed that he will choose which one to open so that the prize will not be revealed. At this point you will be given a fresh choice of door: You can either stick with your first choice, or you can switch to the other closed door. All the doors will then be opened and you will receive whatever is behind your final choice of door.

Imagine that the contestant chooses first door 1; then the gameshow host opens door 3, revealing nothing. Should the contestant a) stick with door 1, b) switch to door 2 or c) it does not matter? You may assume that initially, the prize is equally likely to be behind any of the 3 doors.

Let z_i denote the prize being behind the i 'th door, o_i the action of opening the i 'th door and c_i the action of choosing the i 'th door. The door with the largest probability of containing the prize should be picked, meaning

$$z^* = \underset{z}{\operatorname{argmax}}(\mathbb{P}(z|o_3, c_1)). \quad (64)$$

Since the host cannot open the door containing the prize, $\mathbb{P}(z_3|o_3, c_1) = 0$ and only $\mathbb{P}(z_1|o_3, c_1)$ and $\mathbb{P}(z_2|o_3, c_1)$ will have to be considered. For z_1

$$\mathbb{P}(z_1|o_3, c_1) = \frac{\mathbb{P}(o_3|c_1, z_1)\mathbb{P}(c_1, z_1)}{\mathbb{P}(o_3, c_1)} \quad (65)$$

with

$$\begin{aligned} \mathbb{P}(o_3, c_1) &= \sum_i \mathbb{P}(o_3, c_1, z_i) \\ &= \mathbb{P}(o_3, c_1, z_1) + \mathbb{P}(o_3, c_1, z_2) + \mathbb{P}(o_3, c_1, z_3) \\ &= \mathbb{P}(o_3|c_1, z_1)\mathbb{P}(c_1, z_1) + \mathbb{P}(o_3|c_1, z_2)\mathbb{P}(c_1, z_2) \\ &\quad + \mathbb{P}(o_3|c_1, z_3)\mathbb{P}(c_1, z_3). \end{aligned} \quad (66)$$

$\mathbb{P}(o_3|c_1, z_3) = 0$ since the host will not open the door with the prize. $p(o_3|c_1, z_2) = 1$ since the host has no other option in this case. $\mathbb{P}(o_3|c_1, z_1) = \frac{1}{2}$ since the host has two options in this case. There is no connection between the choice of door and position of the prize, so $\mathbb{P}(c_1, z_j) = \mathbb{P}(c_1)\mathbb{P}(z_j)$ and initially $\mathbb{P}(z_j) = \mathbb{P}(z_k) \forall j, k \in \{1, 2, 3\}$. Hence

$$\begin{aligned} \mathbb{P}(z_1|o_3, c_1) &= \frac{\mathbb{P}(o_3|c_1, z_1)}{\sum_i \mathbb{P}(o_3|c_1, z_i)} \\ &= \frac{1}{3}. \end{aligned} \quad (67)$$

Similarly

$$\begin{aligned}\mathbb{P}(z_2|o_3, c_1) &= \frac{\mathbb{P}(o_3|c_1, z_2)}{\sum_i \mathbb{P}(o_3|c_1, z_i)} \\ &= \frac{2}{3}.\end{aligned}\tag{68}$$

Since $\mathbb{P}(z_2|o_3, c_1) > \mathbb{P}(z_1|o_3, c_1) > \mathbb{P}(z_3|o_3, c_1)$, door number 2 is the optimal choice. Hence, answer "b" is correct. The intuition behind the answer is the information the contestant has at the time of making the decision; initially, there is no a priori information and so $\mathbb{P}(z_1|o_3, c_1) = \frac{1}{3}$. At this time, there is $\frac{2}{3}$ probability that the prize is behind doors 2, 3. When the gameshow host open door 3, this probability converge on door 2.

Example 3.12.

Let $X \sim \text{Unif}(a = -1, b = 1)$ and $Y = X^2$. Clearly Y is dependent on X (in fact Y is uniquely determined by X). However, show that $\text{Corr}[X, Y] = 0$.

$$\begin{aligned}\text{Corr}[X, Y] &= \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X]\text{Var}[Y]}} \\ &= \frac{\mathbb{E}_{XY}[XY] - \mathbb{E}_X[X]\mathbb{E}_Y[Y]}{\sqrt{\text{Var}[X]\text{Var}[Y]}}\end{aligned}\tag{69}$$

In this case for the nonimator

$$\begin{aligned}\text{Cov}[X, Y] &= \int dx x^3 p(x) - \int dx' x' p(x') \int dx'' x''^2 p(x'') \\ &= \frac{1}{b-a} \int_a^b x^3 dx - \frac{1}{(b-a)^2} \int_a^b dx' x' \int_a^b dx'' x''^2 \\ &= \frac{1}{12} (a-b)^2 (a+b) \\ &= 0\end{aligned}\tag{70}$$

where the last equality comes from the fact that $a + b = 0$ in this case. However, we need to make sure the denominator does not diverge

$$\begin{aligned} \text{Var}[X]\text{Var}[X^2] &= (\mathbb{E}_X[X^2] - \mathbb{E}_X[X]^2)(\mathbb{E}_X[X^4] - \mathbb{E}_X[X^2]^2) \\ &= \frac{1}{540}(b-a)^4(4a^2 + 7ab + 4b^2) \\ &\neq 0. \end{aligned} \quad (71)$$

It denominator does not diverge, so the factorized $a + b$ from the nominator makes $\text{Corr}[X, X^2] = 0$.

Example 3.13.

Let $X \sim N(\mu = 0, \sigma^2 = 1)$ and $Y = WX$, where W is a discrete random variable defined by $p(W = -1) = p(W = 1) = \frac{1}{2}$. It is clear that X and Y are not independent, since Y is a function of X .

1. Show $Y \sim N(\mu = 0, \sigma^2 = 1)$.

To show that $Y \sim N(\mu = 0, \sigma^2 = 1)$, show that Y has zero mean and unity variance.

$$\begin{aligned} \mathbb{E}_Y[Y] &= \mathbb{E}_{WX}[WX] \\ &= \mathbb{E}_W[W]\mathbb{E}_X[X] \rightarrow 0 \\ &= 0. \end{aligned} \quad (72)$$

The variance

$$\begin{aligned} \text{Var}[Y] &= \mathbb{E}_Y[Y^2] - \mathbb{E}_Y[Y]^2 \rightarrow 0 \\ &= \mathbb{E}_{WX}[W^2X^2] \\ &= \mathbb{E}_W[W^2]\mathbb{E}_X[X^2] \\ &= \mathbb{E}_W[W^2]\text{Var}[X] \end{aligned} \quad (73)$$

since $\text{Var}[X] = \mathbb{E}_X[X^2] - \mathbb{E}_X[X]^2 = 1$. Now

$$\begin{aligned}\mathbb{E}_W[W^2] &= \frac{1}{n} \sum_{i=1}^n w_i^2 p(W = w_i) \\ &= \frac{1}{2} [(-1)^2 \frac{1}{2} + 1^2 \frac{1}{2}] \\ &= 1\end{aligned}\tag{74}$$

so $\text{Var}[Y] = 1$.

2. Show $\text{Cov}[X, Y] = 0$. Thus X and Y are uncorrelated but dependent, even though they are Gaussian.

$$\begin{aligned}\text{Cov}[X, Y] &= \text{Cov}[X, WX] \\ &= \mathbb{E}_{WX}[WX^2] - \mathbb{E}_X[X] \mathbb{E}_{WX}[WX] \\ &= \mathbb{E}_W[W] \mathbb{E}_X[X^2] - \mathbb{E}_W[W] \mathbb{E}_X[X]^2 \\ &= \mathbb{E}_W[W] \text{Var}[X] \\ &= 0\end{aligned}\tag{75}$$

where for the last equality it has been used that

$$\begin{aligned}\mathbb{E}_W[W] &= \frac{1}{n} \sum_{i=1}^n w_i p(W = w_i) \\ &= \frac{1}{2} [(-1) \frac{1}{2} + 1 \frac{1}{2}] \\ &= 0\end{aligned}\tag{76}$$

Example 3.14.

Prove that $-1 \leq \text{Corr}[X, Y] \leq 1$.

Since the variance is defined as positive definite

$$\begin{aligned}
 0 &\leq \text{Var} \left[\frac{X}{\sigma_X} \pm \frac{Y}{\sigma_Y} \right] \\
 &= \frac{\text{Var}[X]}{\sigma_X^2} + \frac{\text{Var}[Y]}{\sigma_Y^2} \pm \frac{2}{\sigma_X \sigma_Y} \text{Cov}[X, Y] \\
 &= \frac{\text{Var}[X]}{\sigma_X^2} + \frac{\text{Var}[Y]}{\sigma_Y^2} \pm 2\text{Corr}[X, Y] \\
 &= 2 \pm 2\text{Corr}[X, Y]
 \end{aligned} \tag{77}$$

where for the last equality it has been used that $\sigma_i^2 = \text{Var}[i]$. $0 \leq 2 \pm 2\text{Corr}[X, Y] \Leftrightarrow -1 \leq \text{Corr}[X, Y] \leq 1$.

Example 3.15.

Show that if $Y = aX + b$ for some parameters $a > 0$ and b , then $\text{Corr}[X, Y] = 1$. Similarly show that if $a < 0$, then $\text{Corr}[X, Y] = -1$.

$$\text{Corr}[X, Y] = \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X]\text{Var}[Y]}} \tag{78}$$

$$\begin{aligned}
 \text{Cov}[X, Y] &= \mathbb{E}_{XY}[XY] - \mathbb{E}_X[X]\mathbb{E}_Y[Y] \\
 &= \mathbb{E}_X[X(aX + b)] - \mathbb{E}_X[X]\mathbb{E}_X[aX + b] \\
 &= a\mathbb{E}_X[X^2] + b\mathbb{E}_X[X] - a\mathbb{E}_X[X]^2 - b\mathbb{E}_X[X] \\
 &= a\text{Var}[X]
 \end{aligned} \tag{79}$$

$$\begin{aligned}
 \text{Var}[Y] &= \text{Var}[aX + b] \\
 &= a^2\text{Var}[X] + \cancel{\text{Var}[b]}^0 + \cancel{2\text{Cov}[aX, b]}^0 \\
 &= a^2\text{Var}[X]
 \end{aligned} \tag{80}$$

$$\begin{aligned}
 \text{Corr}[X, Y] &= \frac{a\text{Var}[X]}{\sqrt{a^2\text{Var}[X]\text{Var}[X]}} \\
 &= \frac{a}{|a|}
 \end{aligned} \tag{81}$$

Hence, the sign of "a" determine if $\text{Corr}[X, Y] = \pm 1$ for the particular Y of this example.

CHAPTER 4

Introduction to Statistics

Let the observed outcome of a statistical experiment be described by the probability space $(\Omega, \mathcal{F}, \mathbb{P})$ (see Chapter 3), where as opposed to the case in probability theory, \mathbb{P} is now unknown. A generic number of random variables are defined on the sample space viz [4, 7–9]

$$X_i : \Omega \mapsto \Omega_{X_i}, \quad (82)$$

where Ω_{X_i} is part of the probability space $(\Omega_{X_i}, \mathcal{F}_{X_i}, \mathbb{P}_{X_i})$, where

$$\mathbb{P}_{X_i} = \mathbb{P} \circ X_i^{-1} \quad (83)$$

is the push forward measure (see Definition 29) of \mathbb{P} with respect to X_i . The joint probability measure can be defined viz

$$\mathbb{P}_{X_1, \dots, X_n} = \mathbb{P} \circ (X_1, \dots, X_n)^{-1}. \quad (84)$$

on the measurable space

$$(\Omega_{X_1} \cdots \times \Omega_{X_n}, \mathcal{F}_{X_1} \cdots \otimes \mathcal{F}_{X_n}) \quad (85)$$

which for brevity will be written $(\Omega_{X_{1:n}}, \mathcal{F}_{X_{1:n}})$. Depending on the discrete or continuous nature of the different random variables, there are discrete (PMF, see Definition 30) or continuous probability distributions (PDF, see Definition 31) associated to the joint probability measure. All probability distributions related to the random variables can be derived from the joint probability distribution via marginalization (see Theorem 3).

Definition 36 (Set of Probability Measures). Let \mathcal{P} be the set of all probability measures on $(\Omega_{X_{1:n}}, \mathcal{F}_{X_{1:n}})$. It is assumed, often based on prior information, that $\mathbb{P}_{X_1, \dots, X_n} \in \mathcal{P}' \subseteq \mathcal{P}$, which is described in parametric form viz

$$\mathcal{P}' = \{\mathbb{P}_{X_1, \dots, X_n}(w) | w \in \Omega_W\}, \quad (86)$$

where Ω_W is called the parameter space.

Definition 37 (Parameter Space). $\mathbb{P}_{X_1, \dots, X_n}(w) \in \mathcal{P}'$ is specified by parameters $w \in \Omega_W$, where Ω_W is the parameter space.

Definition 38 (Identifiable statistical model). A statistical model is identifiable if $w \in \Omega_W \mapsto \mathbb{P}_{X_1, \dots, X_n}(w) \in \mathcal{P}'$ is injective (one-to-one).

The parameters $w \in \Omega_W$ can either be viewed as fixed constants or the realization of a random variable.

Axiom 4 (Parameter Fixedness). The parameter $w \in \Omega_W$ is treated as a fixed but unknown constant in the statistical model.

Axiom 5 (Parameter as a Random Variable). The parameter $w \in \Omega_W$ is treated as a realization of a random variable. In this case, the parameter space must be endowed with a σ -algebra (\mathcal{F}_W) and a probability measure (\mathbb{P}_W) that must be the result of another measure pushed forward (see Definition 29) with respect to the random variable W . This means

$$W : \Omega \mapsto \Omega_W \quad (87)$$

is defined as a random variable that maps from the probability space $(\Omega, \mathcal{F}, \mathbb{P})$ to the probability space $(\Omega_W, \mathcal{F}_W, \mathbb{P}_W)$, and where

$$\mathbb{P}_W : \mathcal{F}_W \mapsto [0, 1], \quad (88)$$

is called the prior measure, which is the push forward measure of \mathbb{P} with respect to W , i.e.

$$\mathbb{P}_W = \mathbb{P} \circ W^{-1}. \quad (89)$$

For both Axiom 4 and Axiom 5, the value of a parameter is considered fixed. Axiom 5 introduces a random variable W not to add randomness to the parameter w but to model uncertainty or variability about the fixed but unknown parameter value. Observations of the random variables X_1, \dots, X_n are used to a) estimate the parameters if they are fixed and b) estimate the joint probability distribution of the parameters if they are random variables. Hence, given a set of observations of the random variables X_1, \dots, X_n and defining an appropriate subset \mathcal{P}' for the joint probability measure, probability theory can be used to answer statistical questions. This highlights the dual nature of statistics, comprised of two integral parts.

1. The first part involves the formulation and evaluation of probabilistic models, a process situated within the realm of the philosophy of science. This phase grapples with the foundational aspects of constructing models that accurately represent the problem at hand.
2. The second part concerns itself with extracting answers after assuming a specific model. Here, statistics becomes a practical application of probability theory, involving not only theoretical considerations but also numerical analysis in real-world scenarios.

This duality underscores the interdisciplinary nature of statistics, bridging the gap between the conceptual and the applied aspects of probability theory.

4.1 INTERPRETATION OF A PROBABILITY MEASURE

Although probability measures are well defined (see Chapter 3), their interpretation is not defined beyond their definition. For this reason there are two broadly accepted interpretations of probability; objective and subjective.

Definition 39 (Objective Probability Measure). Let \mathbb{P} denote a generic probability measure defined on the generic probability space $(\Omega, \mathcal{F}, \mathbb{P})$. The "objective probability measure"-interpretation define \mathbb{P} as the long-run or limiting frequency of an event, E . That is, let m be the number of occurrences of E , and let n be the number of experiments, then [10]

$$\mathbb{P}(E) \equiv \lim_{n \rightarrow \infty} \left(\frac{m}{n} \right) \quad (90)$$

define the probability measure as the limit of a relative frequency.

Definition 40 (Sugeno Measure). Let (Ω, \mathcal{F}) be a measurable space (Definition 22) and $\text{Bel} : \mathcal{F} \rightarrow [0, 1]$ a Sugeno measure iff [11]

1. **Non-negativity:** $\text{Bel}(\emptyset) = 0$,
2. **Normalization:** $\text{Bel}(\Omega) = 1$,
3. **Monotonicity:** For all $A, B \in \mathcal{F}$, if $A \subseteq B$, then $\text{Bel}(A) \leq \text{Bel}(B)$.

Definition 41 (Subjective Probability Measure). A subjective probability measure is a numerical representation of rational beliefs. Formally, it is a probability measure (Definition 23) \mathbb{P} on a measurable space (Ω, \mathcal{F}) that fulfills the definition of a Sugeno measure (Definition 40) [11, 12].

Theorem 13. Any probability measure \mathbb{P} on (Ω, \mathcal{F}) is a Sugeno measure.

Proof. Let \mathbb{P} be a probability measure on (Ω, \mathcal{F}) . By definition, \mathbb{P} satisfies:

1. $\mathbb{P}(\emptyset) = 0$ and $\mathbb{P}(\Omega) = 1$ (Boundary Conditions).
2. If $A, B \in \mathcal{F}$ and $A \subseteq B$, then $\mathbb{P}(A) \leq \mathbb{P}(B)$ (Monotonicity).

Thus, \mathbb{P} is a Sugeno measure. □

Corollary 1. *Since a probability measure \mathbb{P} satisfies the axioms of a Sugeno measure, it can be interpreted as a belief function.*

Definition 42 (Frequentist Statistics). *Frequentist statistics is a paradigm that adopts Axiom 4 and Definition 39 of probability.*

Definition 43 (Bayesian Statistics). *Bayesian statistics is a paradigm that adopts Axiom 5 and definition Definition 41 of probability.*

Example 4.1.

In the Frequentist approach one can say; if an experiment is repeated many times, in (e.g.) 95% of these cases the computed confidence interval will contain the true parameter value.

In the Bayesian approach one can say; given the observed data, there is a 95% probability that the value of the true parameter lies within the Bayesian interval.

Note how in the Frequentist approach the true parameter is fixed and the confidence interval is varying. In the Bayesian approach the interval is fixed and the true parameter is varying.

Example 4.2.

Consider a Bayesian statistical model involving both a normal distribution with parameters μ, σ and a beta distribution with parameters a, b , then

$$W = \left(W_\mu \quad W_\sigma \quad W_a \quad W_b \right)^T, \quad (91)$$

such that each individual parameter has an associated probability distribution.

4.2 RELAXATION OF NOTATION

Fortunately, a lot of the details around probability spaces and measures can be abstracted in the practical application of statistics. For this reason, in the remainder of the book, where the practical application of statistics is considered, the notation and

formalization especially around probability spaces, algebras, probability measures ect. is relaxed considerably – which is the norm, by the way. Specifically, in the rest of this book, p will be used to denote anything related to probability distributions or measures and the probability for a random variable to take on a specific value, e.g. $p(X = x)$, will usually be denoted $p(x)$ for shorthand. This relaxation of notation facilitates advanced manipulation of probabilities, which would otherwise be incredibly cumbersome. It is, however, beneficial to have some background knowledge about the formal definitions, hence this introduction.

CHAPTER 5

Assigning Probability Functions

The axioms and definitions (Axiom 1-Axiom 3, Definition 8 and Definition 9) of probability theory can be used to define and relate probability measures, however, they are not sufficient to conduct inference because, ultimately, the probability measure or relevant probability functions (density or mass) needs to be specified. Thus, the rules for manipulating probability functions must be supplemented by rules for assigning probability functions. To assign any probability function, there is ultimately only one way, logical analysis, i.e., non-self-contradictory analysis of the available information. The difficulty is to incorporate only the information one actually possesses without making gratuitous assumptions about things one does not know. A number of procedures have been developed that accomplish this task: Logical analysis may be applied directly to the sum and product rules to yield probability functions [13]. Logical analysis may be used to exploit the group invariances of a problem [14]. Logical analysis may be used to ensure consistency when uninteresting or nuisance parameter are marginalized from probability functions [15]. And last, logical analysis may be applied in the form of the principle of maximum entropy to yield probability functions [14, 16–19]. Of these techniques the principle of maximum entropy is probably the most powerful.

5.1 THE PRINCIPLE OF MAXIMUM ENTROPY

The principle of maximum entropy, first proposed by Jaynes [20], considers the issue of assigning a probability distribution to a random variable. Let Z be a generic random variable that describes an abstract experiment. Z follow a distribution $p(z|\lambda, I)$ with associated parameters $\lambda = \{\lambda_0, \dots, \lambda_n\}$. The principle of maximum entropy propose that the probability distribution, $p(z|\lambda, I)$, which best represents the current state of knowledge about a system is the one with largest constrained entropy [1], defined by the Lagrangian

$$\mathcal{L} = \int F dz, \quad (92)$$

with

$$F = -p(z|\lambda, I) \ln \frac{p(z|\lambda, I)}{m(z)} - \lambda_0 p(z|\lambda, I) - \sum_{j=1}^n \lambda_j C_j(z). \quad (93)$$

m – called the Lebesgue measure – ensures the entropy, given by $-\int p(z|\lambda, I) \ln \frac{p(z|\lambda, I)}{m(z)} dz$, is invariant under a change of variables and $C_j(z)$ represent the constraints beyond normalization. The constraint beyond normality depend on the background information related to the random variable, X . In variational calculus the Lagrangian is optimized via solving the Euler-Lagrange equation

$$\frac{\partial F}{\partial p(z|\lambda, I)} - \frac{d}{dx} \frac{\partial F}{\partial p(z|\lambda, I)'} = 0, \quad (94)$$

where $\frac{\partial p(z|\lambda, I)}{\partial x} = p(z|\lambda, I)'$ for shorthand. Since $p(z|\lambda, I)' \notin F$, the Euler-Lagrange equation simplify to simply

$$\frac{\partial F}{\partial p(z|\lambda, I)} = 0. \quad (95)$$

Combining Equation 92 and Equation 95

$$\begin{aligned}\frac{\partial F}{\partial p(z|\lambda, I)} &= -\ln \left(\frac{p(z|\lambda, I)}{m(z)} \right) - 1 - \sum_j \lambda_j C_j(z) \\ &= 0\end{aligned}\quad (96)$$

and so

$$\begin{aligned}p(z|\lambda, I) &= m(z)e^{-1-\sum_j \lambda_j C_j(z)} \\ &= \tilde{m}(z)e^{-\sum_j \lambda_j C_j(z)},\end{aligned}\quad (97)$$

where $\tilde{m}(z) \equiv m(z)e^{-1}$. Using that $\int p(z|\lambda, I)dx = 1$

$$p(z|\lambda, I) = \frac{\tilde{m}(z)e^{-\sum_j \lambda_j C_j(z)}}{\int \tilde{m}(z')e^{-\sum_j \lambda_j C_j(z')}dz'},\quad (98)$$

where m is a reference distribution that is invariant under parameter transformations. λ_j are determined from the additional constraints, e.g. on the mean or variance.

Example 5.1.

Consider a random variable, Z , with unlimited support, $z \in [-\infty, \infty]$, assumed to be symmetric around a single peak defined by the mean μ , standard deviation σ . In this case $\lambda = \{\lambda_0, \lambda_1, \lambda_2\}$, where it will be shown that λ_1, λ_2 are related to μ, σ . In this case F can be written

$$\begin{aligned}F &= -\int p(z|\lambda, I) \ln \left(\frac{p(z|\lambda, I)}{m(z)} \right) dz - \lambda_0 \int p(z|\lambda, I) dz \\ &\quad - \lambda_1 \int p(z|\lambda, I) z dz - \lambda_2 \int p(z|\lambda, I) z^2 dz\end{aligned}\quad (99)$$

with the derivative

$$\begin{aligned}\frac{\partial F}{\partial p(z|\lambda, I)} &= -1 - \ln \left(\frac{p(z|\lambda, I)}{m(z)} \right) - \lambda_1 z - \lambda_2 z^2 \\ &= 0,\end{aligned}\quad (100)$$

meaning

$$p(z|\lambda, I) = m(z)e^{-1-\lambda_0-\lambda_1 z-\lambda_2 z^2}.\quad (101)$$

Taking a uniform measure ($m = \text{const}$) and imposing the normalization constraint

$$\begin{aligned} \int p(z|\lambda, I) dz &= me^{-1-\lambda_0} \int e^{-\lambda_1 z - \lambda_2 z^2} dz \\ &= me^{-1-\lambda_0} \sqrt{\frac{\pi}{\lambda_2}} e^{\frac{\lambda_1^2}{4\lambda_2}} \\ &= 1. \end{aligned} \quad (102)$$

Defining $K^{-1} = me^{-1-\lambda_0}$ yields

$$\begin{aligned} p(z|\lambda, I) &= \frac{e^{-\lambda_1 z - \lambda_2 z^2}}{K} \\ &= \sqrt{\frac{\lambda_2}{\pi}} e^{-\frac{\lambda_1^2}{4\lambda_2} - \lambda_1 z - \lambda_2 z^2}. \end{aligned} \quad (103)$$

Now, imposing the mean constraint

$$\begin{aligned} \int zp(z|\lambda, I) dz &= \frac{\int ze^{-\lambda_1 z - \lambda_2 z^2} dz}{K} \\ &= -\frac{\lambda_1}{2\lambda_2} \\ &= \mu. \end{aligned} \quad (104)$$

Hereby

$$\begin{aligned} p(z|\lambda, I) &= \sqrt{\frac{\lambda_2}{\pi}} e^{-\mu^2 \lambda_2 + 2\mu \lambda_2 z - \lambda_2 z^2} \\ &= \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2} \left(\frac{\mu - z}{\sigma} \right)^2}, \end{aligned} \quad (105)$$

where $\sigma \equiv \frac{1}{2\lambda_2}$ has been defined. Hence, it is clear that the normal distribution can be derived from general constraints via the principle of maximum entropy.

Example 5.2.

Consider a random variable, Z , with limited support, $z \in [0, 1]$. In order to impose the limited support, require that $\ln(z)$ and $\ln(1 - z)$ be well defined. In this case F can be written

$$F = -p(z|\lambda, I) \ln \left(\frac{p(z|\lambda, I)}{m(z)} \right) - \lambda_0 p(z|\lambda, I) - \lambda_1 p(z|\lambda, I) \ln(z) - \lambda_2 p(z|\lambda, I) \ln(1 - z) \quad (106)$$

with the derivative

$$\begin{aligned} \frac{\partial F}{\partial p(z|\lambda, I)} &= -1 - \ln \left(\frac{p(z|\lambda, I)}{m(z)} \right) - \lambda_1 \ln(z) - \lambda_2 \ln(1 - z) \\ &= 0, \end{aligned} \quad (107)$$

meaning

$$p(z|\lambda, I) = m(z) e^{-1 - \lambda_0 - \lambda_1 \ln(z) - \lambda_2 \ln(1 - z)}. \quad (108)$$

Taking a uniform measure ($m = \text{const}$) and imposing the normalization constraint

$$\begin{aligned} \int p(z|\lambda, I) dz &= m e^{-1 - \lambda_0} \int z^{-\lambda_1} (1 - z)^{-\lambda_2} dz \\ &= m e^{-1 - \lambda_0} \frac{\Gamma(1 - \lambda_1) \Gamma(1 - \lambda_2)}{\Gamma(2 - \lambda_1 - \lambda_2)} \\ &= 1. \end{aligned} \quad (109)$$

Now define $\alpha \equiv 1 - \lambda_1$ and $\beta \equiv 1 - \lambda_2$. Hereby

$$p(z|\alpha, \beta, I) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \Gamma(\beta)} x^{\alpha-1} (1 - x)^{\beta-1}, \quad (110)$$

which is the beta distribution.

CHAPTER 6

Framing of Statistics

In this book, the field of statistics will be framed as a game against Nature, as is conventionally done in decision theory. In this game there are two players or decision makers

1. **Robot:** This is the name given to the primary decision maker.
2. **Nature:** This decision maker is a mysterious entity that is unpredictable to the Robot. It has its own set of actions, and it can choose them in a way that interferes with the achievements of the Robot. Nature can be considered as a synthetic decision maker that is constructed for the purposes of modeling uncertainty in the decision-making or planning process.

The game is described by the interaction between the Robot and Nature, characterized by the probability space, $(\Omega, \mathcal{F}, \mathbb{P})$, the parameter space Ω_W , and the set of probability distributions \mathcal{P} parameterized by the parameters $w \in \Omega_W$. Imagine that the Robot and Nature each make a decision by choosing an action from a set, $u \in \Omega_U$ and $s \in \Omega_S$, respectively. Ω_U is referred to as the action space, and Ω_S as the Nature action space. The Robot receives a numerical penalty, assigned by a cost function, depending on the two decisions made.

Definition 44 (Cost Function). *A cost function associates a numerical penalty depending on decision $u \in \Omega_U$ and $s \in \Omega_S$,*

$$C : \Omega_U \times \Omega_S \mapsto \mathbb{R}. \quad (111)$$

Given the observation $X = x$ as well as a set of past observations and matching actions of Nature $D = \{(x_i, s_i) | i = 1 : n\}$, the Robot's objective is to formulate a decision rule that minimize the expected cost associated with its decisions.

Definition 45 (Decision Rule). *A decision rule is a function U that maps from the observation space Ω_X and past observations and decisions $\Omega_X^n \times \Omega_S^n$ to a set of possible actions Ω_U , meaning*

$$U : \Omega_X \times \Omega_S \mapsto \Omega_U. \quad (112)$$

Example 6.1.

Suppose the Robot has an umbrella and considers if it should bring it on a trip outside, i.e.

$$\mathbb{U} = \{\text{"bring umbrella"}, \text{"don't bring umbrella"}\}. \quad (113)$$

Nature have already picked whether or not it will rain later, i.e.

$$\Omega_S = \{\text{"rain"}, \text{"no rain"}\}, \quad (114)$$

so the Robot's task is to estimate Nature's decision regarding rain later and either bring the umbrella or not. The Robot's decision rule, denoted as U , maps the available information $X = x$ (possibly $X =$ weather forecasts, current weather conditions, etc.) to one of its possible actions. For instance, $U(\text{weather forecast})$ might map to the action "bring umbrella" if rain is predicted and "don't bring umbrella" otherwise.

The random variable $X : \Omega \mapsto \Omega_X$ represent the information available (the information may be missing or null) to the Robot regarding the decision Nature will make, while $S : \Omega \mapsto \Omega_S$ represent the different possible decisions of Nature. Ω_X and Ω_S have associated σ -algebras and probability measures, however, such details are assumed *to be understood* in the practical application of statistics. Given the observation $X = x$ as well as a set of past observations

$$D = \{(X = x_1, S = s_1), \dots (X = x_n, S = s_n)\}, \quad (115)$$

the objective of the Robot is to minimize the expected cost associated with its decisions [2]

$$\begin{aligned}\mathbb{E}[C(U, S)|I] &= \int dD dx ds C(U(x, D), s) p(X = x, S = s, D|I) \\ &= \int d\tilde{D} ds C(U(\tilde{D}), s) p(S = s, \tilde{D}|I)\end{aligned}\quad (116)$$

where $\tilde{D} = \{D, X = x\}$ and the Robot aims to find the decision rule which minimizes Equation 116, meaning

$$U^* = \arg \min_U \mathbb{E}[C(U, S)|I]. \quad (117)$$

From Theorem 10

$$\mathbb{E}[C(U, S)|I] = \mathbb{E}_{\tilde{D}}[\mathbb{E}_{S|\tilde{D}}[C(U, S)|\tilde{D}, I]]. \quad (118)$$

Using Equation 118 in Equation 117

$$\begin{aligned}U^* &= \arg \min_U \mathbb{E}_{\tilde{D}}[\mathbb{E}_{S|\tilde{D}}[C(U, S)|\tilde{D}, I]] \\ &= \arg \min_U \int dx p(\tilde{D}|I) \mathbb{E}_{S|\tilde{D}}[C(U, S)|\tilde{D}, I].\end{aligned}\quad (119)$$

Since $p(\tilde{D}|I)$ is a non-negative function, the minimizer of the integral is the same as the minimizer of the conditional expectation, meaning

$$\begin{aligned}U^*(\tilde{D}) &= \arg \min_{U(\tilde{D})} \mathbb{E}_{S|\tilde{D}}[C(U(\tilde{D}), S)|\tilde{D}, I] \\ &= \arg \min_{U(\tilde{D})} \int ds C(U(\tilde{D}), s) p(S = s|X = x, D, I).\end{aligned}\quad (120)$$

Example 6.2.

In general the random variable X represent the observations the Robot has available that are related to the decision Nature is going to make. However, this information may not be given, in which case $\{x, D_x\} = \emptyset$ and consequently

$$\begin{aligned}\tilde{D} &= \{S_1 = s_1, \dots, S_n = s_n\} \\ &\equiv D_s.\end{aligned}\quad (121)$$

In this case, the Robot is forced to model the decisions of Nature with a probability distribution with associated parameters without observations. From Equation 234 the optimal action for the Robot can be written

$$U^*(D_s) = \arg \min_{U(D_s)} \mathbb{E}_{S|\tilde{D}}[C(U(\tilde{D}), S)|\tilde{D}, I] \quad (122)$$

6.1 ASSIGNING A COST FUNCTION

The cost function (see definition 44) associates a numerical penalty to the Robot's action and thus the details of it determine the decisions made by the Robot. Under certain conditions, a cost function can be shown to exist [3], however, there is no systematic way of producing or deriving the cost function beyond applied logic. In general, the topic can be split into considering a continuous and discrete action space, Ω_U .

6.1.1 Continuous Action Space

In case of a continuous action space, the cost function is typically picked from a set of standard choices.

Definition 46 (Linear Cost Function). *The linear cost function is defined viz*

$$C(U(\tilde{D}), s) \equiv |U(\tilde{D}) - s|. \quad (123)$$

Theorem 14 (Median Decision Rule). *Assuming the cost function of Definition 46*

$$\begin{aligned} \mathbb{E}_{S|\tilde{D}}[C(U(\tilde{D}), S)|\tilde{D}, I] &= \int_{-\infty}^{\infty} ds |U(\tilde{D}) - s| p(s|\tilde{D}, I) \\ &= \int_{-\infty}^{U(\tilde{D})} (s - U(\tilde{D})) p(s|\tilde{D}, I) ds \\ &\quad + \int_{U(\tilde{D})}^{\infty} (U(\tilde{D}) - s) p(s|\tilde{D}, I) ds \end{aligned} \quad (124)$$

$$\begin{aligned}
0 &= \frac{d\mathbb{E}_{s|\tilde{D}}[C(U(\tilde{D}), s)|\tilde{D}, I]}{dU(\tilde{D})} \Big|_{U(\tilde{D})=U^*(\tilde{D})} \\
&= (U^*(\tilde{D}) - U^*(\tilde{D}))p(U^*(\tilde{D})|\tilde{D}, I) + \int_{-\infty}^{U^*(\tilde{D})} p(s|\tilde{D}, I)ds \\
&\quad + (U^*(\tilde{D}) - U^*(\tilde{D}))p(U^*(\tilde{D})|\tilde{D}, I) - \int_{U^*(\tilde{D})}^{\infty} p(s|\tilde{D}, I)ds
\end{aligned} \tag{125}$$

$$\begin{aligned}
\int_{-\infty}^{U^*(\tilde{D})} p(s|\tilde{D}, I)ds &= \int_{U^*(\tilde{D})}^{\infty} p(s|\tilde{D}, I)ds \\
&= 1 - \int_{-\infty}^{U^*(\tilde{D})} p(s|\tilde{D}, I)ds
\end{aligned} \tag{126}$$

$$\int_{-\infty}^{U^*(\tilde{D})} p(s|\tilde{D}, I)ds = \frac{1}{2} \tag{127}$$

which is the definition of the median.

Definition 47 (Quadratic Cost Function). *The quadratic cost function is defined as*

$$C(U(\tilde{D}), s) \equiv (U(\tilde{D}) - s)^2. \tag{128}$$

Theorem 15 (Expectation Decision Rule). *Assuming the cost function of Definition 47*

$$\begin{aligned}
 \mathbb{E}_{S|\tilde{D}}[C(U(\tilde{D}), S)|\tilde{D}, I] &= \int ds (U(\tilde{D}) - s)^2 p(s|\tilde{D}, I) \\
 &\Downarrow \\
 \frac{d\mathbb{E}_{S|\tilde{D}}[C(U(\tilde{D}), S)|\tilde{D}, I]}{dU(\tilde{D})} \Big|_{U(\tilde{D})=U^*(x)} &= 2U^*(\tilde{D}) - 2 \int dss p(s|\tilde{D}, I) \\
 &= 0 \\
 &\Downarrow \\
 U^*(\tilde{D}) &= \int dss p(s|\tilde{D}, I) \\
 &= \mathbb{E}[S|\tilde{D}, I]
 \end{aligned} \tag{129}$$

which is the definition of the expectation value.

Definition 48 (0-1 Cost Function). *The 0-1 cost function is defined viz*

$$C(U(\tilde{D}), s) \equiv 1 - \delta(U(\tilde{D}) - s). \tag{130}$$

Theorem 16 (MAP Decision Rule). *The maximum a posteriori (MAP) follows from assuming 0-1 loss viz*

$$\mathbb{E}_{S|\tilde{D}}[C((\tilde{D}), S)|\tilde{D}, I] = 1 - \int ds \delta(U(\tilde{D}) - s) p(s|\tilde{D}, I) \tag{131}$$

meaning

$$\begin{aligned}
 \frac{d\mathbb{E}_{S|\tilde{D}}[C(U(\tilde{D}), S)|\tilde{D}, I]}{dU(\tilde{D})} \Big|_{U(\tilde{D})=U^*(\tilde{D})} &= - \frac{dp(s|\tilde{D}, I)}{ds} \Big|_{s=U^*(\tilde{D})} \\
 &= 0
 \end{aligned} \tag{132}$$

which is the definition of the MAP.

Example 6.3.

Take

$$C(U(x), s) = \alpha \cdot \text{swish}(U(x) - s, \beta) + (1 - \alpha) \cdot \text{swish}(s - U(x), \beta) \quad (133)$$

where

$$\text{swish}(z, \beta) = \frac{z}{1 + e^{-z\beta}}. \quad (134)$$

and $z \equiv U(x) - s$. Taking $\alpha \ll 1$, then $z < 0$ will be penalized relatively more than $z > 0$. $z < 0$ corresponds to underestimation, so this is penalized greater relative to overestimation. Now

$$\begin{aligned} \mathbb{E}[C|\dots] &= \int dsp(s|\dots) \left(\alpha \cdot \text{swish}(U(x) - s, \beta) \right. \\ &\quad \left. + (1 - \alpha) \cdot \text{swish}(s - U(x), \beta) \right) \end{aligned} \quad (135)$$

Let $z \equiv U(x) - s$, then

$$\begin{aligned} \frac{dC}{dU(x)} &= \frac{dC}{dz} \frac{dz}{dU(x)} \\ &= \left(\frac{\alpha}{1 + e^{-\beta z}} - \frac{1 - \alpha}{1 + e^{\beta z}} \right. \\ &\quad \left. + \frac{\alpha \beta e^{-\beta z} z}{(1 + e^{-\beta z})^2} + \frac{(1 - \alpha) \beta e^{\beta z} z}{(1 + e^{\beta z})^2} \right) \frac{dz}{dU(x)} \quad (136) \\ &= \frac{\beta z e^{\beta z} - e^{\beta z} - 1}{(1 + e^{\beta z})^2} + \alpha + \mathcal{O}(\alpha^2) \\ &\approx \alpha - \frac{1}{(1 + e^{\beta z})^2} \end{aligned}$$

$$\begin{aligned} \frac{d\mathbb{E}[C|\dots]}{dU(x)} &\approx \int dsp(s|\dots) \left(\alpha - \frac{1}{(1 + e^{\beta z})^2} \right) \\ &= \alpha - \int dsp(s|\dots) \frac{1}{(1 + e^{\beta z})^2} \quad (137) \\ &= 0 \end{aligned}$$

$\frac{1}{(1+e^{\beta z})^2}$ approximate a unit step which is 1 for $z < 0$ and 0 otherwise.
 $z < 0 \Rightarrow s > U(x)$. This means

$$\int_{-\infty}^{\infty} dsp(s|\dots) \frac{1}{(1+e^{\beta z})^2} \approx \int_{U(x)}^{\infty} dsp(s|\dots) \quad (138)$$

This means

$$\alpha \approx \int_{U(x)}^{\infty} dsp(s|\dots). \quad (139)$$

6.1.2 Discrete Action Space

In case of a continuous action space, the conditional expected loss can be written

$$\mathbb{E}_{S|\tilde{D}}[C(U(\tilde{D}), S)|\tilde{D}, I] = \sum_{s \in S}^n C(U(\tilde{D}), s)p(s|\tilde{D}, I), \quad (140)$$

where the cost function is typically represented in matrix form viz

		S		
		s_1	\dots	$s_{\dim(\Omega_S)}$
$U(x)$	u_1	$C(u_1, s_1)$	\dots	$C(u_1, s_{\dim(\Omega_S)})$
	\vdots	\vdots	\vdots	\vdots
	$u_{\dim(\Omega_U)}$	$C(u_{\dim(\Omega_U)}, s_1)$	\dots	$C(u_{\dim(\Omega_U)}, s_{\dim(\Omega_S)})$

6.2 STATISTICAL PARADIGMS

So far in this chapter, there has been no reference to the statistical paradigms (Bayesian and Frequentist). This is because all so far is valid for both the Bayesian (Definition 4.3) and Frequentist (Definition 4.1) paradigms. The difference between the two comes to light when considering the parameters of Nature's model.

Part II

FREQUENTIST STATISTICS

CHAPTER 7

Frequentist Statistics Introduction

Frequentist statistics is based on Definition 4.1, which follows the definition of objective probability (Definition 39) and the principle of fixed, unknown parameters (Axiom 4). The foundations of Frequentist statistics trace back to seminal works such as those of Neyman and Pearson [21] and Fisher [22], who laid the groundwork for much of its methodology. Subsequent developments by Wald [23], Neyman [24], and Lehmann [25] further refined its theories and techniques.

In the Frequentist paradigm, it is assumed that Nature's decisions can be captured by a model with unknown, fixed parameters w . This means everything in Chapter 6 becomes conditioned on w , and the focus becomes estimating w via an estimator $\hat{w}(\tilde{D})$ and subsequently deciding the optimal decision rule

$$U^*(X = x, \hat{w}(\tilde{D})) \tag{141}$$

according to a cost function, as specified in Chapter 6.

Example 7.1.

Let X and S be continuous random variables with the relationship [26]

$$S = f(X, w) + \epsilon \tag{142}$$

where ϵ is a random variable representing noise, with $\mathbb{E}[\epsilon] = 0$ and f is some model with parameters w . Using the quadratic cost function of Definition 47, Theorem ?? yields

$$\begin{aligned} U^*(X = x, w) &= \mathbb{E}[S|X = x, w, I] \\ &= f(x, w), \end{aligned} \tag{143}$$

where it has been used that f does not depend on past data. Hence, if w was known, the Robot could map any observation $X = x$ using the optimal decision rule Equation 143. In reality w is not known and must be estimated.

CHAPTER 8

Parameter Estimation

Given a decision rule, $U(x, w)$, the unknown parameter w need to be estimated. This is unique to frequentist statistics as the parameters w are considered realizations of random variables in Bayesian statistics. w is estimated by re-applying decision theory. To distinguish this scenario from previous ones, denote the decision rule in this case \hat{w} .

Definition 49 (Fisher Information). *The Fisher information is a way of measuring the amount of information about an unknown parameter a random variable contains. Let w be an unknown parameter, $p(X|w)$ a probability distribution for the generic random variable $X : \Omega \mapsto \Omega_X$ and define $l(X|w, I) = \frac{\partial}{\partial w} \ln p(X|w, I)$. The Fisher information can then be written*

$$\begin{aligned} \mathcal{I}(w) &\equiv \mathbb{E}[l(X|w)^2|w, I] \\ &= \text{Var}[l(X|w)|w, I]. \end{aligned} \tag{144}$$

Proof. In general

$$\mathbb{E}[l(X|w)^2|w, I] = \text{Var}[l(X|w)|w, I] + \mathbb{E}[l(X|w)|w, I]^2 \tag{145}$$

however

$$\begin{aligned} \mathbb{E}[l(X|w, I)|w] &= \int \frac{\partial}{\partial w} \ln p(X = x|w, I) p(X = x|w, I) dx \\ &= \frac{\partial}{\partial w} \int p(X = x|w, I) dx \\ &= 0. \end{aligned} \tag{146}$$

□

Theorem 17 (Fisher information for sample). *Let X_1, X_2, \dots, X_n be a set of independent and identically distributed random variables from the measurable space (Ω, \mathcal{F}) . The Fisher information in a sample is*

$$\mathcal{I}(w) = n\mathcal{I}_1(w), \quad (147)$$

where $\mathcal{I}_1(w)$ is the Fisher information of any one of the random variables.

Definition 50 (Maximum Likelihood Estimator (MLE) Decision Rule). *The Maximum Likelihood Estimator (MLE) decision rule \hat{w}_{MLE} is defined as the decision rule that maximizes the likelihood $p(D_s|D_x, w)$ given the data D_x and past Nature decisions D_s*

$$\hat{w}_{MLE}(\tilde{D}) \equiv \arg \max_w p(D_s|D_x, w, I). \quad (148)$$

Theorem 18 (Unbiasedness of the MLE Decision Rule). *Under certain regularity conditions, the MLE decision rule \hat{w}_{MLE} is asymptotically unbiased, meaning*

$$\sqrt{n}(\hat{w}_{MLE} - w) \xrightarrow{d} N(0, I(w)^{-1}), \quad (149)$$

where $I(w)^{-1}$ is the Fisher information matrix at w and \xrightarrow{d} represents convergence in distribution.

Definition 51 (Minimax Decision Rule). *A decision rule \hat{w}' is said to be minimax if it minimize the maximum expected cost, meaning*

$$\hat{w}' \equiv \inf_{\hat{w}} \sup_{w \in \Omega_W} \mathbb{E}[C(\hat{w}, w)|w, D, I]. \quad (150)$$

Theorem 19 (Mean Squared Error (MSE)). *The expectation of the quadratic cost function (Definition 47) can be written*

$$\begin{aligned} \mathbb{E}[C(\hat{w}, w)|w, D, I] &= \mathbb{E}[(\hat{w} - w)^2|w, D, I] \\ &= \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^2] + (w - \mathbb{E}[\hat{w}])^2 \\ &= \text{Var}[\hat{w}] + \text{Bias}[\hat{w}]^2 \end{aligned} \quad (151)$$

where conditions have been suppressed in the second line (to fit to the page) and the bias of the estimator of \hat{w} is defined viz

$$\text{Bias}[\hat{w}] \equiv w - \mathbb{E}[\hat{w}|w, D, I]. \quad (152)$$

If $\mathbb{E}[C(\hat{w}, w)|w, D, I] \xrightarrow{\text{data} \rightarrow \infty} 0$ then $C(\hat{w}, w)$ is a weakly consistent estimate of w . There can be different consistent estimates that converge towards w at different speeds. It is desirable for an estimate to be consistent and with small (quadratic) cost, meaning that both the bias and variance of the estimator should be small. In many cases, however, there is bias-variance which means that both cannot be minimized at the same time.

Corollary 2 (MLE is Approximately Minimax for quadratic Loss). *Under certain regularity conditions, the Maximum Likelihood decision rule (MLE) \hat{w}_{MLE} is approximately minimax for the quadratic cost function (Definition 47), meaning it approximately minimizes the maximum expected cost.*

Proof. From theorem Theorem 19

$$\mathbb{E}[(\hat{w} - w)^2] = \text{Var}[\hat{w}] + \text{Bias}[\hat{w}]^2. \quad (153)$$

Under the regularity conditions where the MLE is unbiased and has asymptotically minimal variance, the bias term vanishes, meaning $\text{Bias}[\hat{w}_{\text{MLE}}] = 0$ and the variance term $\text{Var}[\hat{w}_{\text{MLE}}]$ is minimized among a class of estimators. Thus, the expected quadratic cost for the MLE can be approximated by

$$\begin{aligned} \mathbb{E}[(\hat{w}_{\text{MLE}} - w)^2] &\approx \text{Var}[\hat{w}_{\text{MLE}}] \\ &\approx \frac{\text{tr}[I(w)^{-1}]}{n}, \end{aligned} \quad (154)$$

where Theorem 18 was used for the second line. The Cramer-Rao lower bound [27] for variance states that

$$\text{Var}[\hat{w}] \geq \frac{\text{tr}[I(w)^{-1}]}{n}, \quad (155)$$

implying that the MLE decision rule achieves the smallest possible variance asymptotically and therefore that

$$\sup_{w \in \Omega_W} \mathbb{E}[(\hat{w}_{\text{MLE}} - w)^2] \approx \inf_{\hat{w}} \sup_{w \in \Omega_W} \mathbb{E}[(\hat{w} - w)^2], \quad (156)$$

meaning the MLE decision rule is approximately the minimax decision rule under quadratic cost. \square

Example 8.1.

The bias-variance decomposition (Theorem 19) is only relevant for frequentist statistics and Bayesian statistics does not struggle with this tradeoff. It relates to overfitting and underfitting. Bayesians do not fit in the same way as frequentists do. They do not determine a single set of parameters, rather they use a set of parameters due to integration. In that sense, they are protected against overfitting and underfitting and they do not struggle with hyperparameter finetuning as well.

Example 8.2.

Consider $X_1, \dots, X_n \sim \text{Ber}(w)$. Determine the (quadratic) cost of three different decision rules for the mean; the arithmetic sample mean, the number 0.5 and the first data entry and X_1 .

- For the arithmetic mean

$$\hat{w} = \frac{1}{n} \sum_{i=1}^n X_i \quad (157)$$

meaning

$$\begin{aligned} \mathbb{E}[\hat{w}] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] \\ &= w, \\ \text{Var}[\hat{w}] &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] \\ &= \frac{w(1-w)}{n}, \\ \mathbb{E}[(\hat{w} - w)^2] &= \frac{w(1-w)}{n}. \end{aligned} \quad (158)$$

- For the number 0.5

$$\hat{w} = 0.5 \quad (159)$$

meaning

$$\begin{aligned} \mathbb{E}[\hat{w}] &= 0.5, \\ \text{Var}[\hat{w}] &= 0, \\ \mathbb{E}[(\hat{w} - w)^2] &= (0.5 - w)^2. \end{aligned} \quad (160)$$

- For the first entry, X_1 ,

$$\hat{w} = \frac{1}{n} \sum_{i=1}^n X_i \quad (161)$$

meaning

$$\begin{aligned} \mathbb{E}[\hat{w}] &= \mathbb{E}[X_1] \\ &= w, \\ \text{Var}[\hat{w}] &= \text{Var}[X_1] \\ &= w(1 - w), \\ \mathbb{E}[(\hat{w} - w)^2] &= w(1 - w). \end{aligned} \quad (162)$$

The arithmetic mean minimizes the quadratic cost over the entire range of w , while the constant value 0.5 performs better for a specific range of w . The cost for X_1 is independent of n , making it less favorable as n increases.

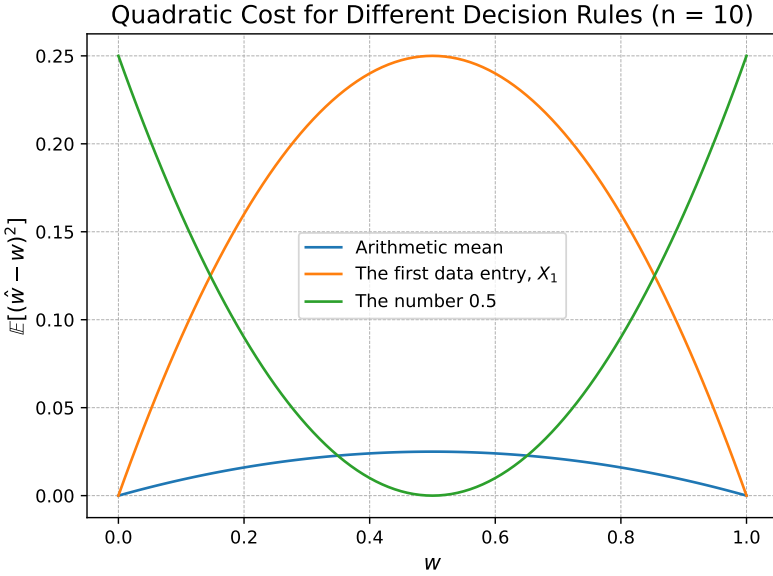


Figure 9: The quadratic cost, $\mathbb{E}[(\hat{w} - w)^2]$, for three different decision rules: the arithmetic mean (blue), the first data entry X_1 (orange), and the constant value 0.5 (green).

Example 8.3.

Determine the maximum likelihood estimate of w for the model $(\{0, 1\}, \{Ber(w)\}_{w \in (0,1)})$.

In this case

$$p(D_s | D_x, w, I) = \prod_{i=1}^n w^{x_i} (1 - w)^{1-x_i}. \quad (163)$$

Let $l(w) \equiv \ln p(D_s | D_x, w, I)$, then

$$\begin{aligned}
 \operatorname{argmax}_w l(w) &= \operatorname{argmax}_w p(D_s | D_x, w, I) \\
 &= \operatorname{argmax}_w \ln \left(\prod_{i=1}^n w^{x_i} (1-w)^{1-x_i} \right) \\
 &= \operatorname{argmax}_w \left[\ln w \sum_{i=1}^n x_i + \ln(1-w) \sum_{i=1}^n (1-x_i) \right]
 \end{aligned} \tag{164}$$

Now

$$\frac{d}{dw} l(w) = \frac{\sum_{i=1}^n x_i}{w} - \frac{n - \sum_{i=1}^n x_i}{1-w} \tag{165}$$

Requiring the derivative to vanish means the maximum likelihood estimate of w is given by

$$\hat{w}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i. \tag{166}$$

Example 8.4.

Determine the maximum likelihood estimate of w for the model $([0, \infty), \{\operatorname{Exp}(w)\}_{w>0})$.

In this case

$$p(D_s | D_x, w, I) = \prod_{i=1}^n w e^{-w x_i}. \tag{167}$$

Let $l(w) \equiv \ln p(D_s | D_x, w, I)$, then

$$\frac{d}{dw} l(w) = \frac{n}{w} - \sum_{i=1}^n x_i \tag{168}$$

Requiring the derivative to vanish means the maximum likelihood estimate of w is given by

$$\hat{w}_{MLE} = \frac{1}{\frac{1}{n} \sum_{i=1}^n x_i}. \tag{169}$$

Part III

BAYESIAN STATISTICS

CHAPTER 9

Bayesian Statistics Introduction

Bayesian statistics is based on Definition 43, which follows the definition of subjective probability (Definition 41) and the treating the parameters as realizations of a random variable (Axiom 5). The Bayesian framework originally come from the work of Bayes [28] and Laplace [29] with much of the modern discussions and formalism created later by Finetti [30] and Jeffreys [31] and Savage [32].

In the Bayesian paradigm, it is assumed that Natures decisions can be captured by a statistical model with parameters that are modeled as realizations of random variables. This means that the probability $p(S = s|X = x, D, I)$ in equation Equation 120 depend on the parameters w_1, \dots, w_n of the statistical model. Introducing the shorthand notation $W = w_1 \dots W = w_n \rightarrow w$, $dw_1 \dots dw_n \rightarrow dw$ and $X = x \rightarrow x$, then

$$\begin{aligned} p(s|x, D, I) &= \int dw p(w, s|x, D, I) \\ &= \int dw p(s|w, x, D, I) p(w|x, D, I) \end{aligned} \tag{170}$$

Example 9.1.

Writing out the shorthand notation

$$\begin{aligned} p(W = w_1, \dots, W = w_n, S = s|X = x, D, I) &\rightarrow p(w, s|x, D, I), \\ dw_1 \dots dw_n &\rightarrow dw. \end{aligned} \tag{171}$$

To evaluate $p(w|D, I)$ a combination of the chain rule (Theorem 1), Bayes' theorem (Theorem 2) and marginalization (Theorem 3) can be employed viz

$$\begin{aligned} p(w|x, D, I) &= p(w|D, I) \\ &= \frac{p(D_s|w, D_x, I)p(w|I)}{p(D_s|D_x, I)}, \end{aligned} \quad (172)$$

where $D_s = \{S = s_1 \dots S = s_n\}$, $D_x = \{X = x_1, \dots X = x_n\}$ and $p(D_s|D_x, I)$ can be expanded via marginalization and Axiom 6 has been used for the first and second equality.

Axiom 6 (Relevance of Observations). *The Robot's observations are relevant for estimating Nature's model only when they map to known actions of Nature.*

$p(w|I)$ is the Robot's prior belief about w . $p(D_s|w, D_x, I)$ is the likelihood of the past observations of Nature's actions, and $p(w|D, I)$ called the posterior distribution represent the belief of the Robot after seeing data. The prior distribution depends on parameters that must be specified and cannot be learned from data since it reflects the Robot's belief before observing data. These parameters are included in the background information, I . From Equation 172, it is evident that, given the relevant probability distributions are specified, the probability of a parameter taking a specific value follows deductively from probability theory. The subjectivity arises from the assignment and specification of probability distributions which depend on the background information.

CHAPTER 10

Regression

Regression involves the Robot building a model, $f : \Omega_W \times \Omega_X \mapsto \mathbb{R}$, with associated parameters $w \in \Omega_W$, that estimates Nature's actions S based on observed data X . Note that the output of f is \mathbb{R} implying that S is assumed continuous. The model f acts as a proxy for the Robot in that it on behalf of the Robot estimates the action of Nature given an input. Hence, in providing an estimate, the model must make a choice, similar to the Robot and thus the Robot must pick a cost function for the model. In this study, the quadratic cost function from Definition 47 will be considered to review the subject. From Theorem 15 the best action for the Robot can be written

$$U^*(x) = \int ds p(s|x, D, I) \quad (173)$$

Assuming the actions of Nature follow a normal distribution with the function f as mean and an unknown variance, $\xi \in \Omega_W$

$$p(s|x, w, \xi, I) = \sqrt{\frac{\xi}{2\pi}} e^{-\frac{\xi}{2}(f(w,x)-s)^2}. \quad (174)$$

Using Equation 174 and marginalizing over ξ, w

$$\begin{aligned} p(s|x, D, I) &= \int p(s, w, \xi|x, D, I) dw d\xi \\ &= \int p(s|x, w, \xi, D, I) p(w, \xi|x, D, I) dw d\xi \quad (175) \\ &= \int p(s|x, w, \xi, I) p(w, \xi|D, I) dw d\xi, \end{aligned}$$

where it has been used that $p(s|w, \xi, x, D, I) = p(s|w, \xi, x, I)$ since by definition f produce a $1 - 1$ map of the input x (Equation 174) and $p(w, \xi|x, D, I) = p(w, \xi|D, I)$ from Axiom 6. Using Equation 175 in Equation 173¹

$$\begin{aligned} U^*(x) &= \int f(w, x) p(w, \xi|D, I) dw d\xi, \\ &= \mathbb{E}[f|x, D, I] \end{aligned} \quad (176)$$

where it has been used that

$$\begin{aligned} \mathbb{E}[S|x, w, \xi, I] &= \int s p(s|x, w, \xi, I) dy \\ &= f(w, x) \end{aligned} \quad (177)$$

according to Equation 174. Using Bayes theorem (Theorem 2)

$$p(w, \xi|D, I) = \frac{p(D_s|D_x, w, \xi, I) p(w, \xi|D_x, I)}{p(D_s|D_x, I)} \quad (178)$$

where from marginalization (Theorem 3)

$$p(D_s|D_x, I) = \int p(D_s|D_x, w, \xi, I) p(w, \xi|D_x, I) dw d\xi. \quad (179)$$

Assuming the past actions of Nature are independent and identically distributed, the likelihood can be written (using equation Equation 174)

$$p(D_s|D_x, w, \xi, I) = \left(\frac{\xi}{2\pi} \right)^{\frac{n}{2}} \prod_{i=1}^n e^{-\frac{\xi}{2} (f(w, x_i) - s_i)^2} \quad (180)$$

From the chain rule (see Theorem 1) and Theorem 6

$$p(w, \xi|D_x, I) = p(w|\xi, I) p(\xi|I). \quad (181)$$

¹ Note that a function of a random variable is itself a random variable, so f is a random variable.

Assuming the distributions of the w 's are i) independent of ξ and ii) normally distributed² with zero mean and a precision described by a hyperparameter, λ .

$$\begin{aligned} p(w|\xi, I) &= p(w|I) \\ &= \int p(w|\lambda, I) p(\lambda|I) d\lambda \end{aligned} \quad (182)$$

The precision is constructed as a wide gamma distribution so as to approximate an objective prior

$$p(w|\lambda, I) p(\lambda|I) = \prod_{q=1}^{\tilde{n}} \frac{\lambda_q^{\frac{n_q}{2}}}{(2\pi)^{\frac{n_q}{2}}} e^{-\frac{\lambda_q}{2} \sum_{l=1}^{n_q} w_l^2} \frac{\beta_q^{\alpha_q}}{\Gamma(\alpha_q)} \lambda_q^{\alpha_q-1} e^{-\beta_q \lambda_q} \quad (183)$$

where α_q, β_q are prior parameters (a part of the background information) and \tilde{n} is the number of hyper parameters. In the completely general case \tilde{n} would equal the number of parameters w , such that each parameter has an independent precision. In practice, the Robot may consider assigning some parameters the same precision, e.g. for parameters in the same layer in a neural network. Since $p(\xi|I)$ is analogous to $p(\lambda|I)$ – in that both are prior distributions for precision parameters – $p(\xi|I)$ is assumed to be a wide gamma distribution, then

$$\begin{aligned} p(\xi|I) &= \text{Ga}(\xi|\tilde{\alpha}, \tilde{\beta}) \\ &= \frac{\tilde{\beta}^{\tilde{\alpha}}}{\Gamma(\tilde{\alpha})} \xi^{\tilde{\alpha}-1} e^{-\tilde{\beta}\xi}. \end{aligned} \quad (184)$$

At this point equation Equation 173 is fully specified (the parameters $\alpha, \beta, \tilde{\alpha}, \tilde{\beta}$ and the functional form of $f(w, x)$ are assumed specified as part of the background information) and can be approximated by obtaining samples from $p(w, \xi, \lambda|D, I)$ via HMC [34–37] (see Appendix A for a review of HMC). The

² The normally distributed prior is closely related to weight decay [33], a principle conventionally used in frequentist statistics to avoid the issue of overfitting.

centerpiece in the HMC algorithm is the Hamiltonian defined viz [36, 37]

$$H \equiv \sum_{q=1}^{\tilde{n}} \sum_{l=1}^{n_q} \frac{p_l^2}{2m_l} - \ln[p(w, \xi, \lambda | D, I)] + \text{const}, \quad (185)$$

where

$$p(w, \xi | D, I) = \int d\lambda p(w, \xi, \lambda | D, I). \quad (186)$$

Besides its function in the HMC algorithm, the Hamiltonian represent the details of the Bayesian model well and should be a familiar sight for people used to the more commonly applied frequentist formalism (since, in this case, it is in form similar to a cost function comprised of a sum of squared errors, weight decay on the coefficients and further penalty terms [38–40]). Using Equation 178–Equation 186 yields

$$\begin{aligned} H = & \sum_{q=1}^{\tilde{n}} \sum_{l=1}^{n_q} \frac{p_l^2}{2m_l} + \frac{n}{2} [\ln(2\pi) - \ln(\xi)] + \frac{\xi}{2} \sum_{i=1}^n (f(w, x_i) - s_i)^2 \\ & + \sum_{q=1}^{\tilde{n}} \left(\ln(\Gamma(\alpha_q)) - \alpha_q \ln(\beta_q) + (1 - \alpha_q) \ln(\lambda_q) + \beta_q \lambda_q \right. \\ & \quad \left. + \frac{n_q}{2} (\ln(2\pi) - \ln(\lambda_q)) + \frac{\lambda_q}{2} \sum_{l=1}^{n_q} w_l^2 \right) \\ & + \ln(\Gamma(\tilde{\alpha})) - \tilde{\alpha} \ln(\tilde{\beta}) + (1 - \tilde{\alpha}) \ln(\tilde{\xi}) + \tilde{\beta} \tilde{\xi} + \text{const}. \end{aligned} \quad (187)$$

Example 10.1.

Let $\xi \equiv e^{\zeta}$, such that $\zeta \in [-\infty, \infty]$ maps to $\xi \in [0, \infty]$ and ξ is ensured to be positive definite regardless of the value of ζ . Using the differential $d\xi = \xi d\zeta$ in Equation 176 means $p(\theta, \xi, \lambda|D, I)$ is multiplied with ξ . Hence, when taking $-\ln(p(\theta, \xi, \lambda|D, I))$ according to Equation 185, $a - \ln(\xi)$ is added to the Hamiltonian. In practice this means

$$(1 - \tilde{\alpha}) \ln(\xi) \in H \Rightarrow -\tilde{\alpha} \ln(\xi). \quad (188)$$

Example 10.2.

Suppose there is a game between a Robot and Nature in which the Robot objective is to guess the position of the Robot. The Robot is given measurements of its velocity and acceleration, at any time step and must formulate a belief about its position. Nature decides the true position of the Robot and will penalize the Robot according to the deviation between the Robots estimate of its position and the true position viz (Definition 47)

$$C(U(x), s) = (U(x) - s)^2, \quad (189)$$

where s is the true position $U(x)$ is the Robots estimate based on data x containing velocity and acceleration measurements.

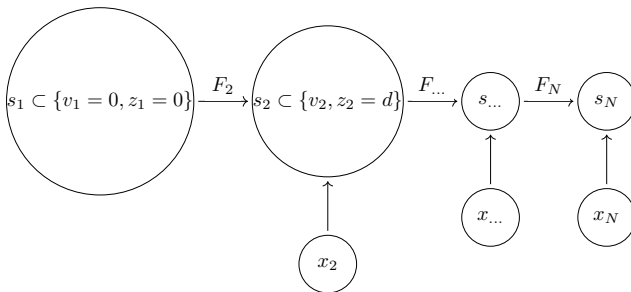


Figure 10

Suppose the Robot is given all historical data, then the expected cost can be written

$$\mathbb{E}_{S|X}[C(U(X), S|s_1, \{x\}_{2:N}, I)] = \int ds (U(x_N) - s)^2 p(s|s_1, \{x\}_{2:N}, I). \quad (190)$$

The optimal decision is defined by

$$\left. \frac{d}{dU(X)} \left(\mathbb{E}_{S|X}[C(U(X), S|s_1, \{x\}_{2:N}, I)] \right) \right|_{U(x_N)=U^*(x_N)} = 0 \quad (191)$$

leading to (Theorem 15)

$$U^*(x_N) = \mathbb{E}[S_N|s_1, \{x\}_{2:N}, I]. \quad (192)$$

The inuitive interpretation of Equation 192 is that the optimal decision of the Robot is to estimate the position that it expects Nature to have chosen. Now

$$\mathbb{E}[S_N|s_1, \{x\}_{2:N}, I] = \int ds_N s_N p(s_N|s_1, \{x\}_{2:N}, I). \quad (193)$$

Assume that

$$\begin{aligned} p(s_i|s_{i-1}, I) &= N(s_i|\mu_i = F_i s_{i-1} + b_i, \Sigma_i = Q_i), \\ p(x_i|s_i, I) &= N(x_i|\mu_i = H_i s_i + d_i, \Sigma_i = R_i) \end{aligned} \quad (194)$$

Hence, $p(s_N|s_1, \{x\}_{2:N}, I)$ must be reformulated such that the above can be utilized. Now

$$\begin{aligned} p(s_N|s_1, \{x\}_{2:N}, I) &= \frac{p(x_N|s_N, s_1, \{x\}_{2:N-1}, I) p(s_N|s_1, \{x\}_{2:N-1}, I)}{p(x_N|s_1, \{x\}_{2:N-1}, I)} \\ &= \frac{p(x_N|s_N, I) p(s_N|s_1, \{x\}_{2:N-1}, I)}{p(x_N|s_1, \{x\}_{2:N-1}, I)} \end{aligned} \quad (195)$$

where

$$\begin{aligned} p(s_N|s_1, \{x\}_{2:N-1}, I) &= \int ds_{N-1} p(s_N, s_{N-1}|s_1, \{x\}_{2:N-1}, I) \\ &= \int ds_{N-1} p(s_N|s_{N-1}) p(s_{N-1}|s_1, \{x\}_{2:N-1}, I) \\ &= N(s_N|\mu_{N|N-1}, \Sigma_{N|N-1}) \end{aligned}$$

$$(196)$$

where

$$\begin{aligned}\mu_{N|N-1} &= F_N \mu_{N-1} + b_N \\ \Sigma_{N|N-1} &= F_N \Sigma_{N-1} F_N^T + Q_N\end{aligned}\quad (197)$$

Using Equation 196 and Equation 194 in Equation 195 then yields [2]

$$p(s_N | s_1, \{x\}_{2:N}, I) = N(s_N | \mu_N, \Sigma_N) \quad (198)$$

where

$$\begin{aligned}\mu_N &= \mu_{N|N-1} + K_N(x_N - \hat{x}_N), \\ \Sigma_N &= \Sigma_{N|N-1} - K_N S_N K_N^T, \\ K_N &= \Sigma_{N|N-1} H_N^T S_N^{-1}, \\ \hat{x}_N &= H_N \mu_{N|N-1} + d_N, \\ S_N &= H_N \Sigma_{N|N-1} H_N^T + R_N\end{aligned}\quad (199)$$

Combining Equation 198 with Equation 193 then yields the optimal decision for the Robot

$$\mathbb{E}[S_N | s_1, \{x\}_{2:N}, I] = \mu_N. \quad (200)$$

Example 10.3.

Consider a Robot moving in one dimension. Every time interval Δt , the Robot will sample a wheel counter and an accelerometer. The wheel counter will be incremented every distance d . The Robot is interested in knowing its position and velocity. Expand the position viz

$$z(t + \Delta t) = z(t) + \Delta t \frac{dz(t)}{dt} + \frac{1}{2} (\Delta t)^2 \frac{d^2 z(t)}{dt^2} + \mathcal{O}(\Delta t^3) \quad (201)$$

which discretize to

$$z_k \simeq z_{k-1} + \Delta t v_{k-1} + \frac{1}{2} (\Delta t)^2 a_{k-1}, \quad (202)$$

where $\Delta t = \text{const}$ and

$$\begin{aligned} v_k &\simeq v_{k-1} + \Delta t a_{k-1}, \\ a_k &\simeq a_{k-1}, \end{aligned} \tag{203}$$

This means

$$\begin{aligned} s_k &= \begin{pmatrix} z_k \\ v_k \\ a_k \end{pmatrix} \\ &\simeq \underbrace{\begin{pmatrix} 1 & \Delta t & \frac{1}{2}\Delta t^2 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{pmatrix}}_{=F_k} \begin{pmatrix} z_{k-1} \\ v_{k-1} \\ a_{k-1} \end{pmatrix} \end{aligned} \tag{204}$$

where $b_k = \emptyset$ for simplicity. Now take

$$\begin{aligned} x_k &= \begin{pmatrix} c_k \\ a_k \end{pmatrix} \\ &= \underbrace{\begin{pmatrix} d^{-1} & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}}_{=H_k} \begin{pmatrix} z_{k-1} \\ v_{k-1} \\ a_{k-1} \end{pmatrix} + r_k \end{aligned} \tag{205}$$

where $r_k \sim N(0, R_k)$ with

$$R_k = \begin{pmatrix} \sigma_z^2 & \sigma_z \sigma_a \\ \sigma_z \sigma_a & \sigma_a^2 \end{pmatrix}, \tag{206}$$

where σ_z and σ_a are estimated from observations. The process noise can be determined viz Kalman link

$$Q_k = \begin{pmatrix} \frac{\Delta t^4}{4} & \frac{\Delta t^3}{2} & \frac{\Delta t^2}{2} \\ \frac{\Delta t^3}{2} & \Delta t^2 & \Delta t \\ \frac{\Delta t^2}{2} & \Delta t & 1 \end{pmatrix} \sigma^2, \tag{207}$$

where $\sigma \sim \Delta a$.

CHAPTER 11

Classification

Classification is the discrete version of regression, meaning it involves the Robot building a model, $f : \Omega_W \times \Omega_X \mapsto [0, 1]$, with associated parameters $w \in \Omega_W$, that estimates Nature's actions S based on observed data X . As opposed to regression, the random variable S is now discrete and the function is identified with the probability of each action

$$p(S = s|x, w, I) = f_{S=s}(w, x), \quad (208)$$

with

$$\sum_{s \in S} p(S = s|x, w, I) = 1. \quad (209)$$

In this case, the Robot's action space is equal to Nature's action space, with the possible addition of a reject option, $\Omega_U = \Omega_S \cup \text{"Reject"}$. To review this subject the Robot will be considered to be penalized equally in case of a classification error, which corresponds to the 0 – 1 cost function, with the addition of a reject option at cost λ . This means

$$C(U(x), s) = 1 - \delta_{U(x), s} + (\lambda - 1)\delta_{U(x), \text{"Reject"}}. \quad (210)$$

The optimal decision rule for the robot can be written

$$\begin{aligned}
 U^*(x) &= \arg \min_{U(x)} \mathbb{E}[C(U(X), S)|x, D, I] \\
 &= \arg \min_{U(x)} \left(\sum_s C(U(x), s) p(S = s|x, D, I) \right. \\
 &\quad \left. + (\lambda - 1) \delta_{U(x), \text{"Reject"}} \right) \\
 &= \arg \min_{U(x)} \left(1 - p(S = U(x)|x, D, I) \right. \\
 &\quad \left. + (\lambda - 1) \delta_{U(x), \text{"Reject"}} \right).
 \end{aligned} \tag{211}$$

In absence of the reject option, the optimal decision rule is to pick the MAP, similar to Theorem 16. Using Equation 236 and marginalizing over w

$$\begin{aligned}
 p(S = U(x)|x, D, I) &= \int p(S = U(x), w|x, D, I) dw \\
 &= \int p(S = U(x)|x, w, D, I) p(w|x, D, I) dw \\
 &= \int p(S = U(x)|x, w, I) p(w|D, I) dw \\
 &= \int f_{S=U(x)}(w, x) p(w|D, I) dw \\
 &= \mathbb{E}[f_{S=U(x)}(w, x)|D, I],
 \end{aligned} \tag{212}$$

where for the second to last equality it has been assumed that $p(S = U(x)|w, x, D, I) = p(S = U(x)|w, x, I)$ since by definition f (see Equation 236) produce a 1 – 1 map of the input x and $p(w|x, D, I) = p(w|D, I)$ from Axiom 6. From Bayes theorem

$$p(w|D, I) = \frac{p(D_s|D_x, w, I) p(w|D_x, I)}{p(D_s|D_x, I)}, \tag{213}$$

where from Axiom 6 $p(w|D_x, I) = p(w|I)$. Assuming the distribution over w is normally distributed with zero mean and a precision described by a hyperparameter, λ ,

$$p(w|I) = \int p(w|\lambda, I)p(\lambda|I)d\lambda. \quad (214)$$

where $p(w|\lambda, I)p(\lambda|I)$ is given by Equation 241. Assuming the past actions of Nature are independent and identically distributed, the likelihood can be written [41]

$$\begin{aligned} p(D_s|D_x, w, I) &= \prod_{i=1}^n p(S = s_i|X = x_i, w, I) \\ &= \prod_{i=1}^n f_{s_i}(w, x_i) \end{aligned} \quad (215)$$

At this point Equation 211 is fully specified and can be approximated by HMC similarly to the regression case. In this case, the model can be represented by the Hamiltonian

$$H \equiv \sum_q \sum_l \frac{p_l^2}{2m_l} - \ln(p(w, \lambda|D, I)) + \text{const} \quad (216)$$

where

$$p(w|D, I) = \int d\lambda p(w, \lambda|D, I). \quad (217)$$

Using Equation 235-Equation 242 in equation (243) yields the Hamiltonian

$$\begin{aligned} H &= \sum_{q=1}^{\tilde{n}} \sum_{l=1}^{n_q} \frac{p_l^2}{2m_l} - \sum_{i=1}^n \ln(f_{s_i}(w, x_i)) + \text{const} \\ &+ \sum_{q=1}^{\tilde{n}} \left(\ln(\Gamma(\alpha_q)) - \alpha_q \ln(\beta_q) + (1 - \alpha_q) \ln(\lambda_q) + \beta_q \lambda_q \right. \\ &\quad \left. + \frac{n_q}{2} (\ln(2\pi) - \ln(\lambda_q)) + \frac{\lambda_q}{2} \sum_{l=1}^{n_q} w_l^2 \right) \end{aligned}$$

(218)

Sampling Equation 245 yields a set of coefficients which can be used to compute $\mathbb{E}[f_s(w, x)|D, I]$ which in turn (see Equation 235) can be used to compute $U^*(x)$.

Example 11.1.

Consider a discrete action space with an observation $X = x$ and available data D . Picking a class corresponds to an action, so classification can be viewed as a game against nature, where nature has picked the true class and the robot has to pick a class as well. Suppose there are only two classes and the cost function is defined by the matrix

$$\begin{array}{cc}
 & S \\
 & \begin{array}{cc} s_1 & s_2 \end{array} \\
 U(x) \quad u_1 & \begin{array}{cc} 0 & \lambda_{01} \end{array} \\
 & u_2 \quad \begin{array}{cc} \lambda_{10} & o \end{array}
 \end{array}$$

1. Show that the decision u that minimizes the expected loss is equivalent to setting a probability threshold w and predicting $U(x) = u_1$ if $p(S = s_1|x, D, I) < w$ and $U(x) = u_2$ if $p(S = s_2|x, D, I) \geq w$. What is w as a function of λ_{01} and λ_{10} ?

The conditional expected cost

$$\begin{aligned}
 \mathbb{E}_{S|X}[C(u, S)|x, D, I] &= \sum_s C(u, S = s)p(S = s|x, D, I) \\
 &= C(u, S = s_1)p(S = s_1|x, D, I) \\
 &\quad + C(u, S = s_2)p(S = s_2|x, D, I)
 \end{aligned} \tag{219}$$

For the different possible actions

$$\begin{aligned}
 \mathbb{E}_{S|X}[C(u_1, S)|x, D, I] &= \lambda_{01}p(S = s_2|x, D, I), \\
 \mathbb{E}_{S|X}[C(u_2, S)|x, D, I] &= \lambda_{10}p(S = s_1|x, D, I),
 \end{aligned} \tag{220}$$

$U(x) = u_1$ iff

$$\mathbb{E}_{S|X}[C(u_1, S)|x, D, I] < \mathbb{E}_{S|X}[C(u_2, S)|x, D, I] \tag{221}$$

meaning

$$\begin{aligned}\lambda_{01}p(S = s_2|x, D, I) &< \lambda_{10}p(S = s_1|x, D, I) \\ &= \lambda_{10}(1 - p(S = s_2|x, D, I))\end{aligned}\quad (222)$$

meaning $U(x) = u_0$ iff

$$p(S = s_2|x, D, I) < \frac{\lambda_{10}}{\lambda_{01} + \lambda_{10}} = w \quad (223)$$

2. Show a loss matrix where the threshold is 0.1.

$$w = \frac{1}{10} = \frac{\lambda_{10}}{\lambda_{01} + \lambda_{10}} \Rightarrow \lambda_{01} = 9\lambda_{10} \text{ yielding the loss matrix}$$

		S	
		s ₁	s ₂
U(x)	u ₁	0	9λ ₁₀
	u ₂	λ ₁₀	0

You may set $\lambda_{10} = 1$ since only the relative magnitude is important in relation to making a decision.

Example 11.2.

In many classification problems one has the option of assigning x to class $k \in K$ or, if the robot is too uncertain, choosing a reject option. If the cost for rejection is less than the cost of falsely classifying the object, it may be the optimal action. Define the cost function as follows

$$C(u, s) = \begin{cases} 0 & \text{if correct classification } (u = s) \\ \lambda_r & \text{if reject option } u = \text{reject} \\ \lambda_s & \text{if wrong classification } (u \neq s) \end{cases} \quad (224)$$

1. Show that the minimum cost is obtained if the robot decides on class u if $p(S = u|x, D, I) \geq p(S \neq u|x, D, I)$ and if $p(S = u|x, D, I) \geq 1 - \frac{\lambda_r}{\lambda_s}$.

The conditional expected cost if the robot does not pick the reject option, meaning $u \in \mathbb{U} \setminus \text{reject}$

$$\begin{aligned}
 \mathbb{E}_{S|X}[C(u, S)|x, D, I] &= \sum_s C(u, S = s)p(S = s|x, D, I) \\
 &= \sum_{s \neq u} \lambda_s p(S = s|x, D, I) \\
 &= \lambda_s(1 - p(S = u|x, D, I))
 \end{aligned} \tag{225}$$

where for the second equality it has been used that the cost of a correct classification is 0, so the case of $S = u$ does not enter the sum. For the third equality it has been used that summing over all but $S = u$ is equal to $1 - p(S = u|x, D, I)$. The larger $p(S = u|x, D, I)$, the smaller loss (for $\lambda_s > 0$), meaning the loss is minimized for the largest probability. The conditional expected loss if the robot picks the reject option

$$\begin{aligned}
 \mathbb{E}_{S|X}[C(\text{reject}, S)|x, D, I] &= \lambda_r \sum_s p(S = s|x, D, I) \\
 &= \lambda_r.
 \end{aligned} \tag{226}$$

Equation (225) show picking $\arg \max_{u \in \mathbb{U} \setminus \text{reject}} p(S = u|x, D, I)$ is the best option among classes $u \neq \text{reject}$. To be the best option overall, it also needs to have lower cost than the reject option. Using equations (225) and (226) yields

$$(1 - p(S = u|x, D, I))\lambda_s < \lambda_r \tag{227}$$

meaning

$$p(S = u|x, D, I) \geq 1 - \frac{\lambda_r}{\lambda_s}. \tag{228}$$

2. Describe qualitatively what happens as $\frac{\lambda_r}{\lambda_s}$ is increased from 0 to 1.

$\frac{\lambda_r}{\lambda_s} = 0$ means rejection is rated as a successful classification – i.e. no cost associated – and this become the best option (rejection that is) unless $p(y = j|x) = 1$, corresponding to knowing the correct class with absolute certainty. In other words; in this limit rejection is best unless the robot is certain of the correct class. $\frac{\lambda_r}{\lambda_s} = 1$ means rejection is rated a misclassification – i.e. $\lambda_r = \lambda_s$ – and thus and "automatic cost". Hence, in this case rejection is never chosen. In between the limits, an interpolation of interpretations apply.

Example 11.3.

SETUP: Consider a farmer who wishes to retire and therefore would like to sell their set of live animals. For simplicity, assume that animals make up a simple group and a given person can either be interested in purchasing an animal or not – the simplification consist of not differentiating between different animal types. In order to sell their animals, the farmer needs to contact people with a sale in mind. For simplicity, the contact will be assumed to be via a telephone call only. The farmer can call (or not) with the intent to sell an animal to the receipient of the call. The receipient of the call can (or not) be interested in purchasing an animal (Natures decision). Let Ω_U denote the set of the farmers actions and Ω_S the set of Natures actions, then

$$\begin{aligned}\Omega_U &= \{u_1 = \text{call}, u_2 = \text{don't call}\}, \\ \Omega_S &= \{s_1 = \text{interested}, s_2 = \text{not interested}\}.\end{aligned}\tag{229}$$

A nuisance for the contacted people is associated to the call which is represented by the abstract monetary loss, $\lambda \in \mathbb{R}^+$. The degree to which people are annoyed by a sales call is independent in general and the monetary loss represents the average animosity generated and the

associated monetary loss connected to a worsened reputation. Aside from the nuisance associated with a sales call, there is also a monetary reward for a successful sale, ψ . If the farmer cannot sell his animals, he will have them terminated with no associated cost in order not to spend additional time or money on them. Given these assumptions, the cost function can be represented by the matrix

	$s_1 = \text{Interested}$	$s_2 = \text{Not interested}$
$u_1 = \text{Call}$	$\lambda - \psi$	λ
$u_2 = \text{Don't call}$	0	0

The farmer has available to them observations $X = x$ that contain information regarding the decision Nature is going to make, $S = s$. The farmer also have a collection of past observations and resulting decisions of Nature, i.e. $D = \{(X = x_1, S = s_1), (X = x_2, S = s_2), \dots (X = x_n, S = s_n)\} = D_s \times D_s$.

OPTIMAL DECISION RULE: The optimal decision for the farmer to call the i 'th person can then be written viz

$$U^*(x) = \arg \min_{U(x)} \mathbb{E}_{S|X}[C(U(x), S) | X = x, D, I], \quad (230)$$

where

$$\mathbb{E}_{S|X}[C(U(x), S) | X = x, D, I] = \sum_{s \in S} C(U(x), s) p(S = s | X = x, D, I). \quad (231)$$

Writing out the conditional expectation

$$\begin{aligned} \mathbb{E}[C(u_1, S)] &= \sum_s C(u_1, s) p(s) \\ &= C(u_1, s_1) p(s_1 | x, D, I) + C(u_1, s_2) p(s_2 | x, D, I) \\ &= (\lambda - \psi) p(s_1 | x, D, I) + \lambda_i p(s_2 | x, D, I), \\ \mathbb{E}[C(u_2, S)] &= \sum_s C(u_2, s) p(s | x, D, I) \\ &= C(u_2, s_1) p(s_1 | x, D, I) + C(u_2, s_2) p(s_2 | x, D, I), \\ &= 0 \end{aligned}$$

(232)

where the notation has been compressed to fit the equations to the page. The optimal decision rule $U^*(x)$ can be implicitly specified as picking u_1 (call) iff $\mathbb{E}_{S|X}[C(u_1, S)|x, D, I] < \mathbb{E}_{S|X}[C(u_2, S)|x, D, I]$, corresponding to picking u_1 (call) iff

$$(\lambda - \psi)p(S = s_1|x, D, I) + \lambda p(S = s_2|x, D, I) < 0 \quad (233)$$

Since $p(S = s_1|x, D, I) + p(S = s_2|x, D, I) = 1$

$$\frac{\lambda}{\psi} < p(S = s_1|x, D, I) \quad (234)$$

meaning the farmer should call (action u_1) iff the probability for the recipient of the call to be interested in at least one animal is larger than the penalty of calling divided by the gain of calling.

THE PROBABILITY: Equation 234 implicitly specify the decision rule for the farmer. λ, ψ is assumed specified, so only the probability $p(S = s_1|x, D, I)$ remain to be specified. Suppose now a model, $f : \Omega_W \times \Omega_X \mapsto [0, 1]$, with associated parameters $w \in \Omega_W$, that estimates Nature's actions S based on observed data X is introduced. Using marginalization and assuming independence

$$\begin{aligned} p(S = s|x, D, I) &= \int p(S = s, w|x, D, I)dw \\ &= \int p(S = s|x, w, D, I)p(w|x, D, I)dw \quad (235) \\ &= \int p(S = s|x, w, I)p(w|D, I)dw. \end{aligned}$$

The random variable S is discrete and the function is identified with the probability of each action

$$p(S = s|x, w, I) = f_{S=s}(w, x), \quad (236)$$

with

$$\sum_{s \in S} p(S = s|x, w, I) = 1. \quad (237)$$

Combining Equation 235 and Equation 236

$$\begin{aligned} p(S = s|x, D, I) &= \int f_{S=s}(w, x) p(w|D, I) dw \\ &= \mathbb{E}[f_{S=s}(w, x)|D, I]. \end{aligned} \quad (238)$$

From Bayes theorem

$$p(w|D, I) = \frac{p(D_s|D_x, w, I)p(w|D_x, I)}{p(D_s|D_x, I)}, \quad (239)$$

where $p(w|D_x, I) = p(w|I)$. Assuming the distribution over w is normally distributed with zero mean and a precision described by a hyperparameter, λ ,

$$p(w|I) = \int p(w|\lambda, I)p(\lambda|I)d\lambda. \quad (240)$$

The precision is constructed as a wide gamma distribution so as to approximate an objective prior

$$p(w|\lambda, I)p(\lambda|I) = \prod_{q=1}^{\tilde{n}} \frac{\lambda_q^{\frac{n_q}{2}}}{(2\pi)^{\frac{n_q}{2}}} e^{-\frac{\lambda_q}{2} \sum_{l=1}^{n_q} w_l^2} \frac{\beta_q^{\alpha_q}}{\Gamma(\alpha_q)} \lambda_q^{\alpha_q-1} e^{-\beta_q \lambda_q} \quad (241)$$

Assuming the past actions of Nature are independent and identically distributed, the likelihood can be written

$$\begin{aligned} p(D_s|D_x, w, I) &= \prod_{i=1}^n p(S = s_i|X = x_i, w, I) \\ &= \prod_{i=1}^n f_{s_i}(w, x_i) \end{aligned} \quad (242)$$

Aside from the specification of the model f , $p(S = s|x, D, I)$ is at this point fully specified and can be approximated by HMC similarly to the regression case. In this case, the model can be represented by the Hamiltonian

$$H \equiv \sum_q \sum_l \frac{p_l^2}{2m_l} - \ln(p(w, \lambda|D, I)) + \text{const} \quad (243)$$

where

$$p(w|D, I) = \int d\lambda p(w, \lambda|D, I). \quad (244)$$

Using Equation 235-Equation 242 in Equation 243 yields the Hamiltonian

$$\begin{aligned} H = & \sum_{q=1}^{\tilde{n}} \sum_{l=1}^{n_q} \frac{p_l^2}{2m_l} - \sum_{i=1}^n \ln(f_{s_i}(w, x_i)) + \text{const} \\ & + \sum_{q=1}^{\tilde{n}} \left(\ln(\Gamma(\alpha_q)) - \alpha_q \ln(\beta_q) + (1 - \alpha_q) \ln(\lambda_q) + \beta_q \lambda_q \right. \\ & \left. + \frac{n_q}{2} (\ln(2\pi) - \ln(\lambda_q)) + \frac{\lambda_q}{2} \sum_{l=1}^{n_q} w_l^2 \right) \end{aligned} \quad (245)$$

SIMPLE MODEL: Let

$$f_{S=s}(w, x_i) = \frac{e^{b_s + \sum_q a_{sq} x_{iq}}}{\sum_{k \in S} e^{b_k + \sum_q a_{kq} x_{iq}}}, \quad (246)$$

where $w = \{b, a\}$.

MANUAL HMC: The Hamiltonian is given by

$$\begin{aligned} H = & \sum_{q=1}^2 \sum_{l=1}^2 \frac{p_{ql}^2}{2m_{ql}} - \sum_{i=1}^n \ln(f_{s_i}(w, x_i)) \\ & + \ln(\Gamma(\alpha_a)) - \alpha_a \ln(\beta_a) + (1 - \alpha_a) \ln(\lambda_a) + \beta_a \lambda_a \\ & + \frac{1}{2} (\ln(2\pi) - \ln(\lambda_a)) + \frac{\lambda_a}{2} \sum_{j,q} a_{jq}^2 \\ & + \ln(\Gamma(\alpha_b)) - \alpha_b \ln(\beta_b) + (1 - \alpha_b) \ln(\lambda_b) + \beta_b \lambda_b \\ & + \frac{1}{2} (\ln(2\pi) - \ln(\lambda_b)) + \frac{\lambda_b}{2} \sum_j b_j^2 \end{aligned} \quad (247)$$

λ_j is positive definite. In order to uphold this numerically, let $\lambda_j = e^{\tau_j}$. When making this transformation, the integration measure of Equation 240 has to be transformed as well. This proceeds viz

$$d\lambda_j = \lambda_j d\tau_j, \quad (248)$$

meaning effectively λ_j is multiplied on $p(w, \lambda | D, I)$ such that $H \rightarrow H - \ln(\lambda_j)$. This means

$$(1 - \alpha_j) \ln(\lambda_j) \in H \Rightarrow -\alpha_j \ln(\lambda_j). \quad (249)$$

Additionally, it is convenient to pick out the s_i via a one-hot target vector such that

$$\begin{aligned} H = & \sum_{q=1}^2 \sum_{l=1}^2 \frac{p_{ql}^2}{2m_{ql}} - \sum_{j \in S} \sum_{i=1}^n s_{ij} \ln(f_j(w, x_i)) \\ & + \ln(\Gamma(\alpha_a)) - \alpha_a \ln(\beta_a) - \alpha_a \tau_a + \beta_a e^{\tau_a} \\ & + \frac{1}{2} (\ln(2\pi) - \tau_a) + \frac{e^{\tau_a}}{2} \sum_{j,q} a_{jq}^2 \\ & + \ln(\Gamma(\alpha_b)) - \alpha_b \ln(\beta_b) - \alpha_b \tau_b + \beta_b e^{\tau_b} \\ & + \frac{1}{2} (\ln(2\pi) - \tau_b) + \frac{e^{\tau_b}}{2} \sum_j b_j^2 \end{aligned} \quad (250)$$

The derivatives are needed for the HMC algorithm

$$\frac{\partial H}{\partial a_{ml}} = - \sum_{i,j} \frac{s_{ij}}{f_{ij}} \frac{\partial f_{ij}}{\partial a_{ml}} + e^{\tau_a} a_{ml}, \quad (251)$$

$$\begin{aligned} \frac{\partial f_{ij}}{\partial a_{ml}} &= \frac{e^{b_j + \sum_{q_1} a_{jq_1} x_{iq_1}}}{\sum_{k \in S} e^{b_k + \sum_{q_2} a_{kq_2} x_{iq_2}}} \sum_{q_3} \delta_{jm} \delta_{q_3 l} x_{iq_3} \\ &\quad - \frac{e^{b_j + \sum_{q_4} a_{jq_4} x_{iq_4}}}{(\sum_{k \in S} e^{b_k + \sum_{q_5} a_{kq_5} x_{iq_5}})^2} \sum_{k' \in S} e^{b_{k'} + \sum_{q_6} a_{k'q_6} x_{iq_6}} \sum_{q_7} \delta_{k'm} \delta_{q_7 l} x_{iq_7} \\ &= f_{ij} \delta_{jm} x_{il} - f_{ij} f_{im} x_{il} \end{aligned} \quad (252)$$

where it has been used that

$$\frac{\partial a_{jq_3}}{\partial a_{ml}} = \delta_{jm} \delta_{q_3l} \quad (253)$$

$$\begin{aligned} \frac{\partial H}{\partial a_{ml}} &= - \sum_{i,j} \frac{s_{ij}}{f_{ij}} (f_{ij} \delta_{jm} x_{il} - f_{ij} f_{im} x_{il}) + e^{\tau_a} a_{ml} \\ &= \sum_i x_{il} (f_{im} - s_{im}) + e^{\tau_a} a_m \end{aligned} \quad (254)$$

$$\frac{\partial H}{\partial b_m} = \sum_i (f_{im} - s_{im}) + e^{\tau_b} b_m. \quad (255)$$

$$\frac{\partial H}{\partial \tau_m} = -\alpha_m + \beta_m e^{\tau_m} - \frac{1}{2} + \frac{e^{\tau_m}}{2} \sum_j m_j^2. \quad (256)$$

The masses for the HMC algorithm can be set by approximating the second order derivatives as fixed. Let

$$\begin{aligned} \frac{\partial^2 H}{\partial a_{ml}^2} &= \sum_i x_{il} \frac{\partial f_{im}}{\partial a_{ml}} + e^{\tau_a} \\ &= \sum_i x_{il}^2 f_{im} (1 - f_{im}) + e^{\tau_a} \end{aligned} \quad (257)$$

then taking $x_{il}^2 \sim 1$, $f_{im} \sim \frac{1}{2}$ and any parameter ~ 0 , meaning $e^{\tau_a} \sim 1$ yield the mass approximation

$$\begin{aligned} m_{ml}^{(a)} &\sim \left. \frac{\partial^2 H}{\partial a_{ml}^2} \right|_{\text{fixed approximation}} \\ &\sim N \cdot 1^2 \cdot \frac{1}{2} (1 - \frac{1}{2}) + 1 \\ &= \frac{N}{4} + 1, \end{aligned} \quad (258)$$

where N is the number of data samples in D_x . Similarly

$$\begin{aligned} \frac{\partial^2 H}{\partial b_m^2} &= \sum_i \frac{\partial f_{im}}{\partial a_{ml}} + e^{\tau_b} \\ &= \sum_i f_{im} (1 - f_{im}) + e^{\tau_b}, \end{aligned} \quad (259)$$

meaning (since $x_{il}^2 \sim 1$)

$$m_{ml}^{(a)} \sim \left. \frac{\partial^2 H}{\partial b_m^2} \right|_{\text{fixed approximation}} m_{ml}^{(b)}. \quad (260)$$

The precision parameter

$$\frac{\partial^2 H}{\partial \tau_q^2} = \beta_q e^{\tau_q} + \frac{e^{\tau_q}}{2} \sum_j q_j^2. \quad (261)$$

Take $\beta_q = 3$, then

$$m_q^{(\tau)} \sim \left. \frac{\partial^2 H}{\partial \tau_q^2} \right|_{\text{fixed approximation}} \sim 3. \quad (262)$$

DATA: Take $x = (\text{area}, \text{number of animals})^T$ and $s = 2d$ one-hot vector, with $\dim(D) = 1000$. D is split into two sets $D^{(\text{training})}$ and $D^{(\text{test})}$, with $\dim(D^{(\text{training})}) = \gamma \dim(D)$ and $\dim(D^{(\text{test})}) = (1 - \gamma) \dim(D)$ and $\gamma = 0.6$. $D^{(\text{training})}$ will be used to train the model and $D^{(\text{test})}$ to evaluate the quality of the trained model. The underlying truth of Nature (unbeknownst to the model) is that an animal will be purchased iff

$$\text{total area} - 3.2 \cdot \text{number of animals} \geq 3.2. \quad (263)$$

TRAINING: Using $D^{(\text{training})}$ as input, the algorithms the algorithms are trained for 2000 iterations. The first 500 iterations are taken as burn in to be conservative. The coefficients, w , for iterations [500, 2000] are used to make a model prediction viz

$$p(S = s | x, D, I) = \frac{1}{1500} \sum_{i=500}^{2000} f(w_i, x) \quad (264)$$

The accuracy of the modeled probabilities can be gauged by considering the case where $\psi = 2\lambda$ such that the decision rule (Equation 234) becomes

$$\frac{1}{2} < p(S = s_1|x, D, I), \quad (265)$$

and the classification is driven by the probabilities alone.

PYMC HMC ALGORITHM: PyMC is a probabilistic programming library for Python that allows users to build Bayesian models with a Python API and fit them using Markov Chain Monte Carlo methods PyMC link. Using this API it is possible to create, train and test a PyMC-equivalent of the model described in the previous sections. The general approach to building a model using PyMC consists of stating the data generating process, specifying a likelihood, and related prior distributions for any parameters involved. While modeling the data generating process and framing the statistical problem correctly are never completely trivial and require some effort from the user it is rather straight forward to perform the Markov Chain Monte Carlo sampling with PyMC. The user do not need to do any calculations related to the Hamiltonian Monte Carlo method nor do they need to specifically handle any integrals. In addition, there is a range of pre-defined probability distributions both discrete and continuous readily available in the library such as the Gamma (figure 11) and Normal distribution used for modeling the priors.

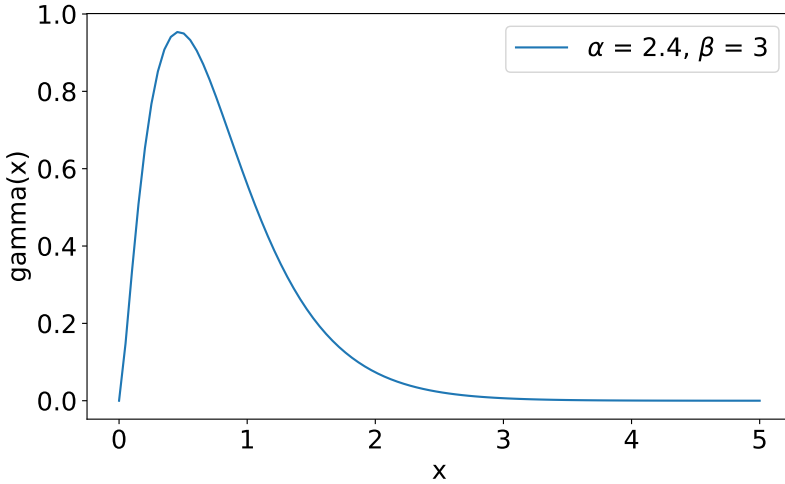


Figure 11: Plot of Gamma probability distribution used as priors for the precision parameters in Normal distributions.

The python code in Algorithm 1 shows how the simple model can be declared using PyMC. All models in PyMC are declared using a "with pm.Model() as modelname"-statement, where pm is the abbreviated form of the PyMC package. Within this statement at range of random variables, data and their relations can be declared. The syntax for declaring a random variable by its probability distribution (such as pm.Normal) is to pass a name for the random variable as a string for the first argument. The rest of the arguments are typically values for the parameters specific to the distribution. PyMC random variables can often be used for stating parameters in other random variables forming a hierarchy of distributions as is the case for the precision variables used to model the variance for the normally distributed parameters a and b (see Equation 246) with zero mean. In the code example, these parameters are then combined deterministically with D_x to constitute the data generating process of the simple model. D_x is declared as a PyMC mutable object so that it will be possible later to change the values to generate predictions. As a last step the likelihood is declared. In PyMC this is the step where the relationship between

the observed conditional (D_s) data and parameters is stated. As the output is the probabilities for the two classes a Multinomial distribution with the number of independent trials (n) equal to 1 is used as the likelihood.

Algorithm 1 PyMC Python Code

```

import pymc as pm
import numpy as np
with pm.Model() as classifier_model:
    covars = pm.MutableData('covars', data_x_training)
    # Priors for precision
    precision_a = pm.Gamma('precision_a', alpha=2.4, beta=3)
    precision_b = pm.Gamma('precision_b', alpha=2.4, beta=3)
    # Priors for parameters
    param_a = pm.Normal('param_a',
        0,
        sigma=1/np.sqrt(precision_a),
        shape=(2,2))
    param_b = pm.Normal('param_b',
        0,
        sigma=1/np.sqrt(precision_b),
        shape=(2,))
    # Data generating process
    T = pm.Deterministic('T', pm.Math.exp(param_b + pm.math.dot(
        covars, param_a.T)))
    class_conditional_probability = pm.Deterministic('
        class_conditional_probability', T
        /T.sum(axis=1, keepdims=True))
    # Likelihood
    obs = pm.Multinomial('obs', n=1, p=
        class_conditional_probability,
        observed=data_s_training, shape=
        class_conditional_probability.
        shape)
  
```

After declaring the model it is possible to sample the posterior by calling the `pm.sample()` method. Burn-in can be controlled by setting the `tune` argument. The `draw` argument determines how many sam-

ples are being drawn while a number of chains can be run in parallel by setting the *chains* and *cores* (computational) arguments.

Algorithm 2 PyMC Posterior Python Code

```
with classifier_model:
    posterior = pm.sample(tune=512, draws=1024, chains=4, cores=
                        4)
```

Using the result of drawing samples from the distribution of model parameters (posterior) it is possible to draw from the posterior predictive distribution using the `pm.sample_posterior_predictive()` method. Without changing the input data this is equivalent to obtaining the result of applying the trained model on $D_x^{(training)}$. By changing the input data to $D_x^{(test)}$ using the `pm.set_data()` method it is possible to obtain the posterior predictive distribution.

Algorithm 3 PyMC Posterior Python Code

```
with classifier_model:
    posterior_predictive = pm.sample_posterior_predictive(
        posterior)
    pm.set_data({'covars': data_x_test})
    posterior_predictive_test = pm.sample_posterior_predictive(
        posterior)
```

Using Equation 265 for $D^{(training)}$, the PyMC model correctly classify 598 of 600 data points. For $D^{(test)}$, the model correctly classify 398.

RESULTS: The HMC algorithm have parameters "step_scale" and "number_of_steps_scale", which adjust the overall scale of the step lengths and number of steps in phase space. Ideally, the distance between points should be large, so that step scale should be small (what the step length is divided by should be small) and the number of

steps should be large. Numerical stability only exist for $\text{step_scale} \gtrsim 5$ (given accurate mass estimation) and thus only the number of steps remain as a variable to tune. In this study $\text{step_scale} = 10$ and $\text{number_of_steps_scale} = 1500$, where the latter is limited by reasonable computation time (to match approximately the pymc computation time). Given these parameters, the manual HMC algorithm misclassify a single training data point

$$x_{\text{misclassified } 1}^{(\text{training})} = \begin{pmatrix} 3.21798365 \\ 0 \end{pmatrix} \quad (266)$$

and two test data points

$$x_{\text{misclassified } 1}^{(\text{test})} = \begin{pmatrix} 6.40501793 \\ 1 \end{pmatrix}, \quad x_{\text{misclassified } 2}^{(\text{test})} = \begin{pmatrix} 9.60627827 \\ 2 \end{pmatrix}. \quad (267)$$

With Equation 263 in mind, it is clear that the misclassifications of Equation 266 and Equation 267 are close to the limit with respect to purchasing an animal.

The PyMC HMC Algorithm obtain misclassify two training data points; Equation 266 and

$$x_{\text{misclassified } 2}^{(\text{training})} = \begin{pmatrix} 12.80826256 \\ 3 \end{pmatrix} \quad (268)$$

and two test data points (Equation 267).

SUMMARY AND DISCUSSION: It has been shown how decision theory can be used in conjunction with statistics to make theoretically optimal decisions based on a user specified set of preferences (cost function). Using mock data, a "manual HMC algorithm" written by hand and a standard Python "PyMC HMC algorithm" have been compared. The two yield identical results on test data with the manual model yielding marginally better results on training data. Overall the performance is deemed equivalent both in terms of accuracy and computational speed. The manual HMC algorithm require the user to derive

the gradients, write the sampling algorithm in Python and tune the algorithm, whereas the latter only require a specification of the model via a standardized PyMC interface. Hence, from a user complexity perspective, the PyMC algorithm has a significant advantage.

CHAPTER 12

Making Inference About the Model of Nature

In some instances, the robot is interested in inference related to the model of Nature. The observation $X = x$ by definition does not have an associated known action of Nature and thus by Axiom 6 is disregarded in this context. From Equation 120

$$U^*(D) = \arg \min_{U(D)} \mathbb{E}_{S|D}[C(U(D), S)|D, I] \quad (269)$$

where $S = s$ is interpreted as an action related to the model of Nature, e.g. Nature picking a given systematic that generates data.

12.1 SELECTING THE ROBOT'S MODEL

Suppose the Robot must choose between two competing models, aiming to select the one that best represents Nature's true model. The two competing models could e.g. be two different functions f in regression or two different probability distribution assignments. In this case the Robot has actions u_1 and u_2 representing picking either model and Nature has two actions s_1 and s_2 which represent which model that in truth fit Nature's true model best. From Equation 269

$$\begin{aligned} \mathbb{E}[C(u_1, S)|D, I] &= \sum_{S=s_1, s_2} C(u_1, s)p(S = s|D, I), \\ \mathbb{E}[C(u_2, S)|D, I] &= \sum_{S=s_1, s_2} C(u_2, s)p(S = s|D, I), \end{aligned} \quad (270)$$

where in this case $u_i = s_i \quad \forall (u_i, s_i) \in \mathbb{U} \times \mathbb{S}$ but the notational distinction is kept to avoid confusion. Since there is no input

$X = x$ in this case, the decision rule U is fixed (i.e. it does not depend on x). $U = u_1$ is picked iff $\mathbb{E}[C(U = u_1, S)|D, I] < \mathbb{E}[C(U = u_2, S)|D, I]$, meaning

$$\frac{p(s_1|D, I)}{p(s_2|D, I)} > \frac{C(u_1, s_2) - C(u_2, s_2)}{C(u_2, s_1) - C(u_1, s_1)}. \quad (271)$$

The ratio $\frac{p(s_1|D, I)}{p(s_2|D, I)}$ is referred to as the posterior ratio. Using Bayes theorem it can be re-written viz

$$\begin{aligned} \text{posterior ratio} &= \frac{p(s_1|D, I)}{p(s_2|D, I)} \\ &= \frac{p(D_s|s_1, D_x, I)p(s_1|I)}{p(D_s|s_2, D_x, I)p(s_2|I)}, \end{aligned} \quad (272)$$

where for the second equality it has been used that the normalization $p(D|I)$ cancels out between the denominator and nominator and Axiom 6 has been employed. Given there is no a priori bias towards any model, $p(s_1|I) = p(s_2|I)$

$$\text{posterior ratio} = \frac{p(D_s|s_1, D_x, I)}{p(D_s|s_2, D_x, I)}. \quad (273)$$

$p(D_s|s_1, D_x, I)$ and $p(D_s|s_2, D_x, I)$ can then be expanded via marginalization, the chain rule and Bayes theorem until they can be evaluated either analytically or numerically. Equation 273 is referred to as Bayes factor and as a rule of thumb

Definition 52 (Bayes Factor Interpretation Rule of Thumb). *If the probability of either of two models being the model of Nature is more than 3 times likely than the other, the likelier model is accepted. Otherwise the result does not significantly favor either model.*

12.2 PARAMETER ESTIMATION

Let $w_j \in \Omega_W$ represent the j 'th parameter with the associated random variable W_j . In case of parameter estimation, the action of Nature is identified with the parameter of interest from the

model of Nature's and the Robot's action with the act of estimating the parameters value, meaning

$$U^* = \arg \min_U \mathbb{E}[C(U, W_j)|D, I], \quad (274)$$

with

$$\mathbb{E}[C(U, W_j)|D, I] = \int dw_j C(U, w_j) p(w_j|D, I). \quad (275)$$

At this point, the Robot can select a cost function like in Section 6.1 and proceed by expanding $p(w_j|D, I)$ similarly to Equation 172. Picking the quadratic cost (Definition 47) yields

$$U^* = \mathbb{E}[w_j|D, I] \quad (276)$$

$p(w_j|D, I)$ in Equation 276 can be expanded as shown in Equation 172.

Example 12.1.

Consider the scenario where two sets of costumers are subjected to two different products, A and B. After exposure to the product, the customer will be asked whether or not they are satisfied and they will be able to answer "yes" or "no" to this. Denote the probability of a customer liking product A/B by w_A/w_B , respectively. In this context, the probabilities w_A/w_B are parameters of Nature's model (similar to how the probability is a parameter for a binomial distribution). What will be of interest is the integral of the joint probability distribution where $w_B > w_A$, meaning

$$p(w_B > w_A | D, I) = \int_0^1 \int_{w_A}^1 p(w_A, w_B | D, I) dw_A dw_B. \quad (277)$$

Assuming the customer sets are independent

$$\begin{aligned} p(w_A, w_B | D, I) &= p(w_B | w_A, D, I) p(w_A | D, I) \\ &= p(w_B | D_A, I) p(w_A | D_A, I), \end{aligned} \quad (278)$$

with

$$p(w_i | D_i, I) = \frac{p(D_i | w_i, I) p(w_i | I)}{p(D_i | I)}. \quad (279)$$

Assuming a beta prior and a binomial likelihood yields (since the binomial and beta distributions are conjugate)

$$p(w_i | D_i, I) = \frac{w_i^{\alpha_i-1} (1 - w_i)^{\beta_i-1}}{B(\alpha_i, \beta_i)}, \quad (280)$$

where $\alpha_i \equiv \alpha + s_i$, $\beta_i \equiv \beta + f_i$ and s_i/f_i denotes the successes/failure, respectively, registered in the two sets of costumers. Evaluating Equation 277 yields

$$p(w_B > w_A | D, I) = \sum_{j=0}^{\alpha_B-1} \frac{B(\alpha_A + j, \beta_A + \beta_B)}{(\beta_B + j) B(1 + j, \beta_B) B(\alpha_A, \beta_A)}. \quad (281)$$

Part IV

REFLECTION

CHAPTER 13

Reflections on Statistical Paradigm

Bayesian statistics offers a mathematically rich framework for modeling uncertainty, akin to the depth of general relativity in physics. Just as general relativity provides a more comprehensive understanding of gravity—encompassing and extending the concepts of Newtonian physics—Bayesian methods offer a flexible, coherent way to model uncertainty and incorporate prior knowledge. However, like general relativity, Bayesian statistics is computationally intensive and, in practice, often requires more sophisticated methods and resources.

Frequentist statistics, on the other hand, serves as a simpler, more practical tool—similar to Newtonian physics. It provides intuitive and computationally efficient techniques that are easier to implement, particularly in large datasets or when computational power is limited. This is why Frequentist methods remain dominant in many areas, despite the richer theoretical framework offered by Bayesian methods.

As technology continues to advance, computational techniques such as Markov Chain Monte Carlo (MCMC) and variational inference are making Bayesian statistics more accessible, potentially transforming it from a complex theoretical framework into a practical tool. If these trends continue, we may see a greater shift toward Bayesian approaches in applied statistics. However, rather than one paradigm replacing the other, a more likely outcome is that both approaches will continue to coexist, with each method being applied where it is most suited. Bayesian methods will become more prevalent as computational

resources improve, but Frequentist methods will likely remain a key tool in many domains due to their simplicity and efficiency.

Ultimately, the future of statistical analysis may not be about replacing one paradigm with the other, but rather about combining the strengths of both. As such, the increasing integration of Bayesian and Frequentist methods in modern statistical practice will continue to shape the field in meaningful ways.

Part V

APPENDIX

APPENDIX A

Hamiltonian Monte Carlo

This appendix is taken from Petersen [53]. The Hamiltonian Monte Carlo Algorithm (HMC algorithm) is a Markov Chain Monte Carlo (MCMC) algorithm used to evaluate integrals on the form

$$\begin{aligned}\mathbb{E}[f] &= \int f(\theta)g(\theta)d\theta \\ &\approx \frac{1}{N} \sum_{j \in g} f(\theta_j),\end{aligned}\tag{282}$$

with f being a generic function and N denoting the number of samples from the posterior distribution, g . The sample $\{j\}$ from g can be generated via a MCMC algorithm that has g as a stationary distribution. The Markov chain is defined by an initial distribution for the initial state of the chain, θ , and a set of transition probabilities, $p(\theta'|\theta)$, determining the sequential evolution of the chain. A distribution of points in the Markov Chain are said to comprise a stationary distribution if they are drawn from the same distribution and that this distribution persist once established. Hence, if g is the a stationary distribution of the Markov Chain defined by the initial point θ and the transition probability $p(\theta'|\theta)$, then [36]

$$g(\theta') = \int p(\theta'|\theta)g(\theta)d\theta.\tag{283}$$

Equation 283 is implied by the stronger condition of detailed balance, defined viz

$$p(\theta'|\theta)g(\theta) = p(\theta|\theta')g(\theta').\tag{284}$$

A Markov chain is ergodic if it has a unique stationary distribution, called the equilibrium distribution, to which it converges from any initial state. $\{i\}$ can be taken as a sequential subset (discarding the part of the chain before the equilibrium distribution) of a Markov chain that has $g(\theta)$ as its equilibrium distribution.

The simplest MCMC algorithm is perhaps the Metropolis-Hastings (MH) algorithm [54, 55]. The MH algorithm works by randomly initiating all coefficients for the distribution wanting to be sampled. Then, a loop runs a subjective number of times in which one coefficient at a time is perturbed by a symmetric proposal distribution. A common choice of proposal distribution is the normal distribution with the coefficient value as the mean and a subjectively chosen variance. If $g(\theta') \geq g(\theta)$ the perturbation of the coefficient is accepted, otherwise the perturbation is accepted with probability $\frac{g(\theta')}{g(\theta)}$.

The greatest weaknesses of the MH algorithm is i) a slow approach to the equilibrium distribution, ii) relatively high correlation between samples from the equilibrium distribution and iii) a relatively high rejection rate of states. ii) can be rectified by only accepting every n 'th accepted state, with n being some subjective number. For $n \rightarrow \infty$ the correlation naturally disappears, so there is a trade off between efficiency and correlation. Hence, in the end the weaknesses of the MH algorithm can be boiled down to inefficiency. This weakness is remedied by the HCM algorithm [35] in which Hamiltonian dynamics are used to generate proposed states in the Markov chain and thus guide the journey in parameter space. Hamiltonian dynamics are useful for proposing states because [37] 1) the dynamics are reversible, implying that detailed balance is fulfilled and so there exist a stationary distribution, 2) the Hamiltonian (H) is conserved during the dynamics if there is no explicit time dependence in the Hamiltonian ($\frac{dH}{dt} = \frac{\partial H}{\partial t}$), resulting in all proposed states being accepted in the case the dynamics are exact and 3) Hamiltonian dynamics preserve the volume in phase space (q_i, p_i -space),

which means that the Jacobian is unity (relevant for Metropolis updates that succeeds the Hamiltonian dynamics in the algorithm). By making sure the algorithm travel (in parameter space) a longer distance between proposed states, the proposed states can be ensured to have very low correlation, hence alleviating issues 1) and 2) of the MH algorithm. The price to pay for using the HMC algorithm relative to the MH algorithm is a) the HMC algorithm is gradient based meaning it requires the Hamiltonian to be continuous and b) the computation time can be long depending on the distribution being sampled (e.g. some recurrent ANNs are computationally heavy due to extensive gradient calculations).

As previously stated, the HMC algorithm works by drawing a physical analogy and using Hamiltonian dynamics to generate proposed states and thus guide the journey in parameter space. The analogy consists in viewing g as the canonical probability distribution describing the probability of a given configuration of parameters. In doing so, g is related to the Hamiltonian, H , viz

$$g = e^{\frac{F-H}{k_B T}} \Rightarrow H = F - k_B T \ln[g], \quad (285)$$

where $F = -k_B T \ln[Z]$ denotes Helmholtz free energy of the (fictitious in this case) physical system and Z is the partition function. $\ln[g(\theta)]$ contain the position (by analogy) variables of the Hamiltonian and so Z must contain the momentum variables. Almost exclusively [56] $Z \sim \mathcal{N}(0, \sqrt{m_i})$ is taken yielding the Hamiltonian

$$H = -k_B T \left[\ln[g] - \sum_i \frac{p_i^2}{2m_i} \right] + \text{const}, \quad (286)$$

where i run over the number of variables and "const" is an additive constant (up to which the Hamiltonian is always defined). $T = k_b^{-1}$ is most often taken [37], however, the temperature can be used to manipulate the range of states which can be accepted

e.g. via simulated annealing [57]. Here $T = k_b^{-1}$ will be adopted in accordance with [36, 37] and as such

$$H = \sum_i \frac{p_i^2}{2m_i} - \ln[g]. \quad (287)$$

The dynamics in parameter space are determined by Hamiltons equations

$$\dot{\theta}_i = \frac{\partial H}{\partial p_i}, \quad \dot{p}_i = -\frac{\partial H}{\partial \theta_i}, \quad (288)$$

with θ_i denoting the different variables (coefficients). In order to implement Hamiltons equations, they are discretized via the leap frog method [36, 37] viz

$$\begin{aligned} p_i \left(t + \frac{\epsilon}{2} \right) &= p_i(t) - \frac{\epsilon}{2} \frac{\partial H(\theta_i(t), p_i(t))}{\partial \theta_i}, \\ \theta_i(t + \epsilon) &= \theta_i(t) + \frac{\epsilon}{m_i} p_i \left(t + \frac{\epsilon}{2} \right), \\ p_i(t + \epsilon) &= p_i \left(t + \frac{\epsilon}{2} \right) - \frac{\epsilon}{2} \frac{\partial H(\theta_i(t + \frac{\epsilon}{2}), p_i(t + \frac{\epsilon}{2}))}{\partial \theta_i}, \end{aligned} \quad (289)$$

with ϵ being an infinitesimal parameter. In the algorithm the initial state is defined by a random initialization of coordinates and momenta, yielding $H_{initial}$. Subsequently Hamiltonian dynamics are simulated a subjective (L loops) amount of time resulting in a final state, H_{final} , the coordinates of which take the role of proposal state. The loop that performs L steps of ϵ in time is here referred to as the dive. During the dive, the Hamiltonian remains constant, so ideally $H_{initial} = H_{final}$, however, imperfections in the discretization procedure of the dynamics can result in deviations from this equality (for larger values of ϵ , as will be discussed further later on). For this reason, the proposed state is accepted as the next state in the Markov chain with probability

$$\mathbb{P}(\text{transition}) = \min [1, e^{H_{initial} - H_{final}}]. \quad (290)$$

Whether or not the proposed state is accepted, a new proposed state is next generated via Hamiltonian dynamics and so the loop goes on for a subjective amount of time.

Most often, the HMC algorithm will be ergodic, meaning it will converge to its unique stationary distribution from any given initialization (i.e. the algorithm will not be trapped in some subspace of parameter space), however, this may not be so for a periodic Hamiltonian if $L\epsilon$ equal the periodicity. This potential problem can however be avoided by randomly choosing L and ϵ from small intervals for each iteration. The intervals are in the end subjective, however, with some constraints and rules of thumb; the leap frog method has an error of $\mathcal{O}(\epsilon^2)$ [36] and so the error can be controlled by ensuring that $\epsilon \ll 1$. A too small value of ϵ will waste computation time as a correspondingly larger number of iterations in the dive (L) must be used to obtain a large enough trajectory length $L\epsilon$. If the trajectory length is too short the parameter space will be slowly explored by a random walk instead of the otherwise approximately independent sampling (the advantage of non-random walks in HMC is a more uncorrelated Markov chain and better sampling of the parameter space). A rule of thumb for the choice of ϵ can be derived from a one dimensional Gaussian Hamiltonian

$$H = \frac{q^2}{2\sigma^2} + \frac{p^2}{2}. \quad (291)$$

The leap frog step for this system is a linear map from $t \rightarrow t + \epsilon$. The mapping can be written

$$\begin{bmatrix} q(t + \epsilon) \\ p(t + \epsilon) \end{bmatrix} = \begin{bmatrix} 1 - \frac{\epsilon^2}{2\sigma^2} & \epsilon \\ \epsilon(\frac{1}{4}\epsilon^2\sigma^{-4} - \sigma^{-2}) & 1 - \frac{1}{2}\epsilon^2\sigma^{-2} \end{bmatrix} \begin{bmatrix} q(t) \\ p(t) \end{bmatrix} \quad (292)$$

The eigenvalues of the coefficient matrix represent the powers of the exponentials that are the solutions to the differential equation. They are given by

$$\text{Eigenvalues} = 1 - \frac{1}{2}\epsilon^2\sigma^{-2} \pm \epsilon\sigma^{-1}\sqrt{\frac{1}{4}\epsilon^2\sigma^{-2} - 1}. \quad (293)$$

In order for the solutions to be bounded, the eigenvalues must be imaginary, meaning that

$$\epsilon < 2\sigma. \quad (294)$$

In higher dimensions a rule of thumb is to take $\epsilon \lesssim 2\sigma_x$, where σ_x is the standard deviation in the most constrained direction, i.e. the square root of the smallest eigenvalue of the covariance matrix. In general [56] a stable solution with $\frac{1}{2}p^T\Sigma^{-1}p$ as the kinetic term in the Hamiltonian require

$$\epsilon_i < 2\lambda_i^{-\frac{1}{2}}, \quad (295)$$

for each eigenvalue λ_i of the matrix

$$M_{ij} = (\Sigma^{-1})_{ij} \frac{\partial^2 H}{\partial q_i \partial q_j}, \quad (296)$$

which means that in the case of $\Sigma^{-1} = \text{diag}(m_i^{-1})$;

$$\epsilon_i < 2\sqrt{\frac{m_i}{\frac{\partial^2 H}{\partial q_i^2}}}. \quad (297)$$

Setting ϵ according to Equation 295 can however introduce issues for hierarchical models (models including hyper parameters) since the reversibility property of Hamiltonian dynamics is broken if ϵ depend on any parameters. This issue can be alleviated by using the MH algorithm on a subgroup of parameters [36, 37] (which are then allowed in the expression for ϵ) that is to be included in ϵ . However, unless the MH algorithm is used for all parameters, some degree of approximation is required.

Algorithm 4 Hamiltonian Monte Carlo Algorithm in pseudo code

```

1: Save:  $q$  and  $V(q)$ , with  $q$  randomly initialized
2: for  $i \leftarrow 1$  to  $N$  do
3:    $p \leftarrow$  Sample from standard normal distribution
4:    $H_{\text{old}} \leftarrow H(q, p)$ 
5:    $p \leftarrow p - \frac{\epsilon}{2} \frac{\partial H(q, p)}{\partial q}$ 
6:    $L \leftarrow$  Random integer between  $L_{\text{lower}}$  and  $L_{\text{upper}}$ 
7:   for  $j \leftarrow 1$  to  $L$  do
8:      $q \leftarrow q + \epsilon \frac{p}{\text{mass}}$ 
9:     if  $j \neq L$  then
10:       $p \leftarrow p - \epsilon \frac{\partial H(q, p)}{\partial q}$ 
11:    end if
12:  end for
13:   $p \leftarrow p - \frac{\epsilon}{2} \frac{\partial H(q, p)}{\partial q}$ 
14:   $H_{\text{new}} \leftarrow H(q, p)$ 
15:   $u \leftarrow$  Sample from uniform distribution
16:  if  $u < \min(1, e^{-(H_{\text{new}} - H_{\text{old}})})$  then
17:     $H_{\text{old}} \leftarrow H_{\text{new}}$ 
18:    Save:  $q$  and  $V(q)$ 
19:  end if
20: end for

```

APPENDIX B

Nested Sampling

A major challenge in estimating the evidence via conventional Monte Carlo Methods is that generally the prior is a very broad and regular distribution whereas the likelihood is a very narrow and irregular distribution. This poses a challenge when the evidence is estimated conventionally, i.e. as the mean of the likelihood evaluated at points in parameter space corresponding to samples from the prior distribution. For a reasonable number of samples, the conventional procedure has a relatively high likelihood of relatively poor sampling in regions near the peaks in the likelihood distribution. This means a conventional estimate of the evidence via Monte Carlo Methods has a high variance. Nested Sampling [58] (NS) address this challenge by accounting for the likelihood distribution when sampling the prior distribution. Consider the integral

$$Z = \int L(\theta)\pi(\theta)d\theta, \quad (298)$$

with L being the likelihood distribution and π the prior distribution. Conventional Monte Carlo methods approximate this integral via importance sampling, meaning

$$\begin{aligned} Z &= \mathbb{E}_{\pi}[L] \\ &\approx \frac{1}{N} \sum_{i \in \pi} L(\theta_i)' \end{aligned} \quad (299)$$

where the second equality become exact for $N \rightarrow \infty$. NS project the integral down into one dimension viz¹

$$\begin{aligned} Z &= \int_0^1 L(\xi) d\xi \\ &\approx \sum_i L(\xi_i) \Delta \xi_i' \end{aligned} \quad (300)$$

where

$$\xi(\lambda) = \int_{L>\lambda} \pi(\theta) d\theta, \quad (301)$$

is the proportion of the prior with likelihood greater than λ and $\Delta \xi_i \equiv \xi_{i-1} - \xi_i$. Due to the constraint $L > \lambda$ on the integral bound of ξ , $L(\xi)$ is a decreasing function of ξ , meaning $L(\xi_1) > L(\xi_2)$ if $\xi_1 < \xi_2$. The sum in Equation 300 can then be evaluated by generating a sequence

$$\{\{L(\xi_m), \xi_m\}, \{L(\xi_{m-1}), \xi_{m-1}\}, \dots, \{L(\xi_1), \xi_1\}\}, \quad (302)$$

with $\xi_1 < \xi_2 < \dots < \xi_m$. The sorting operation eliminate coordinate dependent complications of geometry, topology and dimensionality [59]. A sequence upholding Equation 302 can be generated as follows; consider n random draws from g with corresponding values of L and ξ . Let $L(\xi^*)$ denote the minimum value of L in the sample with ξ^* the corresponding value of ξ in the sample. $\{L(\xi^*), \xi^*\}$ is replaced by another set which is sampled from g with the constraint that $\xi_{new} < \xi^*$ and stored in a list of discarded states. Continuing this sequence again and again will fill the list of discarded states that uphold Equation 302. In practice $L(\xi)$ is not readily available, so instead L can be generated from values of θ . The value of ξ_k can be determined by using that [58]

$$\xi_k = \xi_0 \prod_{i=1}^k t_i, \quad (303)$$

¹ Attempting a higher accuracy via better numerical approximations of the integral is mute since the uncertainty in ξ dominate the approximation [58].

with $t_i = \frac{\xi_k}{\xi_{k-1}}$, called the shrinkage ratio. The shrinkage ratio follow a beta distribution

$$p(t) = nt^{n-1}, \quad (304)$$

with n being the number of initially samples from g (the number of live points), such that

$$\begin{aligned} \langle \ln(t) \rangle &= \mathbb{E}[\ln(t)] \pm \sqrt{V[\ln(t)]} \\ &= \int_0^1 nt^{n-1} \ln(t) dt \pm I_2, \\ &= \frac{1}{n}(-1 \pm 1) \end{aligned} \quad (305)$$

with

$$I_2 = \sqrt{\int_0^1 nt^{n-1} \ln(t)^2 dt - \left(\int_0^1 nt^{n-1} \ln(t) dt \right)^2}. \quad (306)$$

Using $\ln(\xi_k) = \sum_{i=1}^k \ln(t_i)$ and taking t_i to be i.i.d. yield

$$\begin{aligned} \langle \ln(\xi_k) \rangle &= k\mathbb{E}[\ln(t)] \pm \sqrt{kV[\ln(t)]} \\ &= \frac{1}{n}(-k \pm \sqrt{k}). \end{aligned} \quad (307)$$

Ignoring uncertainty ξ_k can be approximated by the mean viz

$$\xi_k \approx e^{-\frac{k}{n}}, \quad (308)$$

meaning

$$\Delta \xi_i \approx e^{-\frac{i}{n}} \left(e^{\frac{1}{n}} - 1 \right). \quad (309)$$

A heuristic measure for terminating the collection of samples is to require that the maximum likelihood collected make up only a small fraction, B , of the evidence, meaning

$$\max(\{L\})\xi_j < BZ, \quad (310)$$

for iteration j . Another approach to terminating the collection of samples is to use that most of the area in the $L\xi$ -plane is usually found in the region [58, 59] $\xi \sim e^{-\mathcal{H}} \sim e^{-\frac{i}{n}}$, meaning the collection of samples can be terminated when

$$i \gg n\mathcal{H}, \quad (311)$$

with \mathcal{H} being the information [58]

$$\begin{aligned} \mathcal{H} &= \int \frac{L(\xi)}{Z} \ln \left(\frac{L(\xi)}{Z} \right) d\xi \\ &\approx \sum_i \frac{L(\xi_i)}{Z} \ln \left(\frac{L(\xi_i)}{Z} \right) \Delta\xi_i. \end{aligned} \quad (312)$$

Terminating at $i \sim n\mathcal{H}$ yield (Equation 307) an uncertainty $\delta(\langle \ln(\xi_i) \rangle) = \sqrt{\frac{\mathcal{H}}{n}}$ meaning

$$\ln(Z) \approx \ln \left(\sum_i L(\xi_i) \Delta\xi_i \right) \pm \sqrt{\frac{\mathcal{H}}{n}}. \quad (313)$$

The NS algorithm with Equation 311 as termination criterion is shown in Algorithm 5. A and B are parameters of the algorithm. The "Remainder" in the second to last line in Algorithm 5 fills in the missing band $0 < \xi < e^{-\frac{k+1}{n}}$ with the average value of the remaining values of L . Due to the chosen stopping criterion, the "Remainder" will be construction be small.

Algorithm 5 Nested Sampling Algorithm in pseudo code

```

1: Import:  $S = n$  samples  $\theta_1, \theta_2, \dots, \theta_n$  from the prior distribu-
   tion with  $L$  being the corresponding likelihoods
2: Initialize:  $k \leftarrow 0, a \leftarrow 0, B \leftarrow 1, Z \leftarrow$  Empty list
3: while  $f > B$  do
4:   Let  $L^* \equiv \min(L)$  and  $S^* \triangleq L^*$ 
5:    $S2 \leftarrow S \setminus S^*$  and  $L2 \leftarrow L \setminus L^*$ 
6:   Define  $\Delta\zeta_k = e^{\frac{k+1}{n}} (e^{\frac{1}{n}} - 1)$ 
7:   Store  $L^* \Delta\zeta_k$  in  $Z$ 
8:    $S_{new}, L_{new} \leftarrow \text{proposer}(\text{random}(S2), L^*)$ 
9:    $S \leftarrow S2 \cup S_{new}$  and  $L \leftarrow L2 \cup L_{new}$ 
10:   $f \leftarrow \frac{\max(L)e^{-\frac{k+1}{n}}}{\sum_{s=0}^k Z_s}$ 
11:  if  $a == A$  then
12:    Display status, e.g.  $f, n\mathcal{H} - k, k, \sum_{s=0}^k Z_s, \dots$ 
13:     $a \leftarrow 0$ 
14:  end if
15:   $k \leftarrow k + 1$ 
16:   $a \leftarrow a + 1$ 
17: end while
18:  $\text{Remainder} \leftarrow \frac{1}{n} \sum_i L_i e^{-\frac{k+1}{n}}$ 
19:  $Z \approx \sum_{s=0}^k Z_s + \text{Remainder}$ 

```

Bibliography

- [1] D. S. Sivia and J. Skilling. *Data Analysis - A Bayesian Tutorial*. 2nd. Oxford Science Publications. Oxford University Press, 2006.
- [2] Kevin P. Murphy. *Probabilistic Machine Learning: Advanced Topics*. MIT Press, 2023. URL: <http://probml.github.io/book2>.
- [3] Steven M. Lavalle. *Planning Algorithms*. Cambridge University Press, 2006. ISBN: 0521862051.
- [4] S.H. Chan. *Introduction to Probability for Data Science*. Michigan Publishing, 2021. ISBN: 9781607857464. URL: <https://books.google.dk/books?id=GR2jzgEACAAJ>.
- [5] Andrey Kolmogorov. *Foundations of the Theory of Probability*. Providence, RI, USA: Chelsea Publishing Company, 1950.
- [6] Marco Taboga. *Expected value and the Lebesgue integral*. Online appendix. 2021. URL: <https://www.statlect.com/fundamentals-of-probability/expected-value-and-Lebesgue-integral>.
- [7] Alexander Drewitz. *Introduction to Probability and Statistics*. Preliminary version, February 1. University of Cologne, 2019.
- [8] Peter Orbanz. *Functional Conjugacy in Parametric Bayesian Models*. Technical Report. University of Cambridge, 2009.
- [9] Daniel V. Tausk. *A Basic Introduction to Probability and Statistics for Mathematicians*. Date: January 24th, 2023. 2023.
- [10] Edward E. Leamer. *Specification Searches: Ad Hoc Inference with Nonexperimental Data*. Wiley, 1978, p. 25.

- [11] Glenn Shafer. "BELIEF FUNCTIONS AND POSSIBILITY MEASURES." English (US). In: *Anal of Fuzzy Inf.* CRC Press Inc, 1987, pp. 51–84. ISBN: 0849362962.
- [12] Peter D. Hoff. *A First Course in Bayesian Statistical Methods.* Springer Texts in Statistics. Springer, 2009. DOI: 10.1007/978-0-387-92407-6.
- [13] E. T. Jaynes. "Probability Theory - The Logic of Science."
- [14] E. T. Jaynes. "Prior Probabilities." In: *IEEE Transactions on Systems Science and Cybernetics* SSC-4 (1968), pp. 227–241.
- [15] E. T. Jaynes. "Marginalization and Prior Probabilities." In: *Bayesian Analysis in Econometrics and Statistics.* Ed. by A. Zellner. Reprinted in [**jaynes_maximum_entropy_formalism**]. Amsterdam: North-Holland Publishing Company, 1980.
- [16] A. Zellner. *An Introduction to Bayesian Inference in Econometrics.* New York: John Wiley and Sons, 1971.
- [17] E. T. Jaynes. "Where Do We Stand On Maximum Entropy?" In: *The Maximum Entropy Formalism.* Ed. by R. D. Levine and M. Tribus. MIT Press, 1978, pp. 15–118.
- [18] J. E. Shore and R. W. Johnson. "Axiomatic Derivation of the Principle of Maximum Entropy and the Principle of Minimum Cross-Entropy." In: *IEEE Transactions on Information Theory* IT-26.1 (1980), pp. 26–37.
- [19] J. E. Shore and R. W. Johnson. "Properties of Cross-Entropy Minimization." In: *IEEE Transactions on Information Theory* IT-27.4 (1981), pp. 472–482.
- [20] E. T. Jaynes. "Information Theory and Statistical Mechanics." In: *Phys. Rev.* 106.4 (May 1957), pp. 620–630. DOI: 10.1103/PhysRev.106.620. URL: http://prola.aps.org/abstract/PR/v106/i4/p620_1.

- [21] J. NEYMAN and E. S. PEARSON. "ON THE USE AND INTERPRETATION OF CERTAIN TEST CRITERIA FOR PURPOSES OF STATISTICAL INFERENCE." In: *Biometrika* 20A.3-4 (Dec. 1928), pp. 263–294. ISSN: 0006-3444. DOI: 10.1093/biomet/20A.3-4.263. eprint: <https://academic.oup.com/biomet/article-pdf/20A/3-4/263/1037410/20A-3-4-263.pdf>. URL: <https://doi.org/10.1093/biomet/20A.3-4.263>.
- [22] R.A. Fisher. *Statistical methods for research workers*. Edinburgh Oliver & Boyd, 1925.
- [23] A. Wald. "Sequential Tests of Statistical Hypotheses." In: *The Annals of Mathematical Statistics* 16.2 (1945), pp. 117 – 186. DOI: 10.1214/aoms/1177731118. URL: <https://doi.org/10.1214/aoms/1177731118>.
- [24] Jerzy Neyman and Elizabeth Letitia Scott. "Consistent Estimates Based on Partially Consistent Observations." In: *Econometrica* 16 (1948), p. 1. URL: <https://api.semanticscholar.org/CorpusID:155631889>.
- [25] E.L. Lehmann. *Testing Statistical Hypotheses*. Probability and Statistics Series. Wiley, 1986. ISBN: 9780471840831. URL: <https://books.google.dk/books?id=jexQAAAAAAAJ>.
- [26] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. New York, NY, USA: Springer New York Inc., 2001.
- [27] C. Radhakrishna Rao. *Linear Statistical Inference and Its Applications*. 2nd. See Chapter 3 for the Cram r-Rao inequality and its applications. New York: John Wiley & Sons, 1973. ISBN: 978-0-471-34969-5.
- [28] T. Bayes. "An essay towards solving a problem in the doctrine of chances." In: *Phil. Trans. of the Royal Soc. of London* 53 (1763), pp. 370–418.

- [29] Pierre-Simon Laplace. *Théorie analytique des probabilités*. Paris: Courcier, 1812. URL: <http://gallica.bnf.fr/ark:/12148/bpt6k88764q>.
- [30] Bruno de Finetti. "La prévision : ses lois logiques, ses sources subjectives." fr. In: *Annales de l'institut Henri Poincaré* 7.1 (1937), pp. 1–68. URL: http://www.numdam.org/item/AIHP_1937__7_1_1_0.
- [31] Harold Jeffreys. *The Theory of Probability*. Oxford Classic Texts in the Physical Sciences. 1939. ISBN: 978-0-19-850368-2, 978-0-19-853193-7.
- [32] L. Savage. *The Foundations of Statistics*. New York: Wiley, 1954.
- [33] D. C. Plaut, S. J. Nowlan, and G. E. Hinton. *Experiments on learning back propagation*. Tech. rep. CMU-CS-86-126. Pittsburgh, PA: Carnegie-Mellon University, 1986.
- [34] J. M Hammersley and D. C. Handscomb. *Monte Carlo Methods*. London, Methuen., 1964.
- [35] S. Duane, A.D. Kennedy, B.J. Pendleton, and D. Roweth. "Hybrid Monte Carlo." In: *Phys. Lett. B* 195 (1987), pp. 216–222. DOI: 10.1016/0370-2693(87)91197-X.
- [36] Radford M. Neal. Berlin, Heidelberg: Springer-Verlag, 1996. ISBN: 0387947248.
- [37] Radford M. Neal. "MCMC using Hamiltonian dynamics." In: (2012). cite arxiv:1206.1901. URL: <http://arxiv.org/abs/1206.1901>.
- [38] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining, inference and prediction*. 2nd ed. Springer, 2009. URL: <http://www-stat.stanford.edu/~tibs/ElemStatLearn/>.

- [39] Kevin P. Murphy. *Machine learning : a probabilistic perspective*. Cambridge, Mass. [u.a.]: MIT Press, 2013. ISBN: 9780262018029 0262018020. URL: https://www.amazon.com/Machine-Learning-Probabilistic-Perspective-Computation/dp/0262018020/ref=sr_1_2?ie=UTF8&qid=1336857747&sr=8-2.
- [40] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. The MIT Press, 2016. ISBN: 0262035618.
- [41] Manfred Fischer and Petra Staufer-Steinnocher. "Optimization in an Error Backpropagation Neural Network Environment with a Performance Test on a Spectral Pattern Classification Problem." In: *Geographical Analysis* 31 (Jan. 1999), pp. 89–108. DOI: 10.1111/gean.1999.31.1.89.
- [42] Renzo Sancisi. "The visible matter – dark matter coupling." In: *Symposium - International Astronomical Union* 220 (June 2004), p. 233. DOI: 10.1017/S0074180900183299.
- [43] R.A. Swaters, R.H. Sanders, and S.S. McGaugh. "Testing Modified Newtonian Dynamics with Rotation Curves of Dwarf and Low Surface Brightness Galaxies." In: *Astrophys. J.* 718 (2010), pp. 380–391. DOI: 10.1088/0004-637X/718/1/380. arXiv: 1005.5456 [astro-ph.CO].
- [44] R.A. Swaters, R. Sancisi, J.M. van der Hulst, and T.S. van Albada. "The Link between the Baryonic Mass Distribution and the Rotation Curve Shape." In: *Mon. Not. Roy. Astron. Soc.* 425 (2012), p. 2299. DOI: 10.1111/j.1365-2966.2012.21599.x. arXiv: 1207.2729 [astro-ph.CO].
- [45] Benoît Famaey and Stacy S. McGaugh. "Modified Newtonian Dynamics (MOND): Observational Phenomenology and Relativistic Extensions." In: *Living Reviews in Relativity* 15.1, 10 (Sept. 2012), p. 10. DOI: 10.12942/lrr-2012-10. arXiv: 1112.3960 [astro-ph.CO].

- [46] Jonas Petersen and Mads T Frandsen. "A method for discriminating between dark matter models and MOND modified inertia via galactic rotation curves." In: *Monthly Notices of the Royal Astronomical Society* 496.2 (June 2020), pp. 1077–1091. ISSN: 0035-8711. DOI: 10.1093/mnras/staa1541. eprint: <https://academic.oup.com/mnras/article-pdf/496/2/1077/33419277/staa1541.pdf>. URL: <https://doi.org/10.1093/mnras/staa1541>.
- [47] James Binney and Scott Tremaine. *Galactic Dynamics: Second Edition*. 2. 2008.
- [48] Federico Lelli, Stacy S. McGaugh, James M. Schombert, and Marcel S. Pawlowski. In: *Astrophys. J.* 836.2 (2017), p. 152. DOI: 10.3847/1538-4357/836/2/152. arXiv: 1610.08981 [astro-ph.GA].
- [49] Stacy McGaugh, Federico Lelli, and Jim Schombert. In: *Phys. Rev. Lett.* 117.20 (2016), p. 201101. DOI: 10.1103/PhysRevLett.117.201101. arXiv: 1609.05917 [astro-ph.GA].
- [50] Y. Bengio, P. Simard, and P. Frasconi. "Learning long-term dependencies with gradient descent is difficult." In: *IEEE Transactions on Neural Networks* 5.2 (1994), pp. 157–166.
- [51] Y. Bengio and P. Frasconi. "Diffusion of Context and Credit Information in Markovian Models." In: *arXiv e-prints*, cs/9510101 (Sept. 1995), cs/9510101. arXiv: cs/9510101 [cs.AI].
- [52] Oriol Abril-Pla et al. "PyMC: a modern, and comprehensive probabilistic programming framework in Python." In: *PeerJ Computer Science* 9 (Sept. 2023), e1516. ISSN: 2376-5992. DOI: 10.7717/peerj-cs.1516. URL: <https://doi.org/10.7717/peerj-cs.1516>.
- [53] J. Petersen. "The Missing MAss Problem on Galactic Scales." PhD thesis.

- [54] Nicholas Metropolis, Arianna W. Rosenbluth, Marshall N. Rosenbluth, Augusta H. Teller, and Edward Teller. "Equation of State Calculations by Fast Computing Machines." In: *The Journal of Chemical Physics* 21.6 (1953), pp. 1087–1092. DOI: 10.1063/1.1699114. URL: <http://link.aip.org/link/?JCP/21/1087/1>.
- [55] W. K. Hastings. "Monte Carlo sampling methods using Markov chains and their applications." In: *Biometrika* 57.1 (1970), pp. 97–109. DOI: 10.1093/biomet/57.1.97. eprint: <http://biomet.oxfordjournals.org/cgi/reprint/57/1/97.pdf>.
- [56] M. Betancourt and Mark Girolami. "Hamiltonian Monte Carlo for Hierarchical Models." In: (Dec. 2013). DOI: 10.1201/b18502-5.
- [57] David J. C. MacKay. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2002.
- [58] John Skilling. "Nested Sampling." In: *Bayesian Inference and Maximum Entropy Methods in Science and Engineering: 24th International Workshop on Bayesian Inference and Maximum Entropy Methods in Science and Engineering*. Ed. by Rainer Fischer, Roland Preuss, and Udo Von Toussaint. Vol. 735. American Institute of Physics Conference Series. Nov. 2004, pp. 395–405. DOI: 10.1063/1.1835238.
- [59] John Skilling. "Skilling, J.: Nested sampling for general Bayesian computation. *Bayesian Anal.* 1(4), 833–860." In: *Bayesian Analysis* 1 (Dec. 2006), pp. 833–860. DOI: 10.1214/06-BA127.

Index

- Axioms of probability theory, 9
- Bayes factor, 74
- Bayes theorem, 11
- Chain rule, 10
- Coin experiment, 16
- Conditional probability, 10
- Example: Bad news from the doctor, 17
- Example: Bayes naive classifier, 51
- Example: Bayesian decision theory, 49, 50
- Example: Correlation coefficient, 19–21
- Example: Crime, 16
- Example: Future event counts, 64
- Example: Gameshow, 18
- Example: HMC Hamiltonian variable change, 44
- Example: Maximum entropy beta distribution, 36
- Example: Maximum entropy normal distribution, 35
- Example: Maximum likelihood, 76
- Example: Normal distribution, 76
- Example: The neighbours children, 16
- Example: Variable transformation, 75
- Example: Variance of a sum, 17
- Expectation value, 64
- Frequentist statistics, 41
- Gamma distribution, 40, 64, 66, 67
- Marginalization, 11
- Marginalized, 64
- Maximum entropy, 33, 36, 64, 67
- Negative binomial distribution, 66
- Normal distribution, 36, 67
- Parameter space, 23
- Poisson distribution, 64, 66
- Posterior ratio, 73
- Probability space, 10
- Random variable, 12
- Set function, 9

Uniform distribution, 76

Variance, 64