

Applied Philosophy of Science and Data Ethics



Dr. Carlos Vega

Applied Philosophy of Science and Data Ethics

Dr. Carlos Vega

2021 - 2025: Last compiled: 2025-11-10

Contents

Preface	7
Course presentation	7
Course at UL	8
Target audience	9
Course objectives	9
Course structure	11
About this course	11
Roadmap for instructors	13
Disclaimer	13
About the author	13
About this class book	14
License	15
Contributing to the book	15
1 Scientific Goals, Methods and Knowledge	17
1.1 What is Science?	19
1.1.1 Scientific Goals and Knowledge	19
1.1.1.1 Data, information and knowledge	23
1.1.2 What is Philosophy of Science?	25
1.2 The scientific method	25
1.3 Methodology	26
1.4 Examples	28
1.4.1 Neptune and Vulcan	28
1.4.2 The most famous “failed” experiment	28
1.4.3 Eddington expeditions	29
1.4.4 The smoke debate	32
1.4.5 Kekulé’s dream	33
1.4.6 Why Most Published Research Findings Are False	34
2 Scientific Inference	39
2.1 Overview	39
2.2 Types of inferences	40

2.2.1	Deduction and Induction	40
2.2.2	Modus ponens and Modus tollens	42
2.3	The problem(s) of induction	44
2.3.1	David Hume’s Problem of Induction	44
2.4	The Hypothetico-deductive Method	48
2.4.1	A good hypothesis	49
2.4.2	Falsification	50
2.4.3	Confirmation	53
2.5	Beyond the Hypothetico-Deductive Method: Bayesian Inference and the Logic of Uncertainty	57
2.5.1	Limitations of the HD method in probabilistic contexts	57
2.5.2	Bayesian Epistemology and Conditionalization	58
2.5.3	Bayesian vs Frequentist Views	58
2.5.4	Bayesian Inference and the Problem of Induction	59
2.5.5	Underdetermination and Bayesian Model Selection	60
2.6	Other types of inference	60
2.6.1	Abduction and Inference to the Best Explanation (IBE)	60
2.6.2	Monotonic and Non-monotonic Logic	61
2.6.3	Defeasible Reasoning	62
2.7	Summary of the different inference methods	63
2.8	Explanation	64
2.8.1	Limitations of Hempel’s model	65
2.8.2	Explanation and causality	67
2.8.3	Summary and modern views	68
2.9	Examples	71
2.9.1	The problem is in your hands!	71
2.9.1.1	How a hypothesis is tested	74
2.9.2	Wason selection task	75
2.9.3	U.S.A. Presidents	78
2.9.4	Yersinia pestis	79
2.9.5	Risks of induction and non-epistemic values in ML	81
3	Empirical Practices and Models	83
3.1	Overview	83
3.2	Experiments	84
3.2.1	Observational studies	85
3.2.1.1	Natural experiments	85
3.2.1.2	Observability	87
3.2.1.3	Indicators	88
3.2.1.4	Data and Evidence	88
3.2.2	Field, laboratory and simulation experiments	89
3.2.2.1	Field experiments	89
3.2.2.2	Laboratory experiments	90

3.2.2.3	Simulation experiments	91
3.2.3	Summary of the different experiments	91
3.2.4	How to evaluate experiment success	94
3.3	Scientific models	95
3.3.1	The models of the atom	98
3.3.2	The models of benzene	99
3.3.3	Models as analogies	99
3.3.4	Differences between Models and Experiments	101
3.3.5	What makes a good model?	102
3.3.6	Models as mirrors	106
3.3.7	Models as isolations	107
3.3.8	Summary of the different model views	108
3.4	A tool for scientific reasoning	110
3.4.1	A Case of Negative Evidence. “Gene Analysis Upsets Turtle Theory”	113
3.4.2	A Case of Positive Evidence. “New View of the Mind Gives Unconscious an Expanded Role”	115
3.4.3	Inconclusive Data. “Was That a Greenhouse Effect? It Depends on Your Theory”	116
3.5	Examples	118
3.5.1	Willow tree experiment	118
3.5.2	1854 Broad Street cholera outbreak	120
3.5.3	Causal models: Estimating treatment effect in the presence of a confounder	123
3.5.3.1	Backdoor criterion	127
3.5.4	Causal models: Detecting an overestimation	129
4	Experimental Control and Statistical Abuse	135
4.1	Overview	135
4.2	The smoke debate - Part II	136
4.3	Experimental Control	139
4.3.1	Other experimental control techniques	139
4.4	Randomised Control Trials	141
4.4.1	Origins of RCTs	143
4.4.2	Validity	144
4.5	Cross-validation in Machine Learning	145
4.6	Surrogates, Proxies, Confounders and Colliders	149
4.6.1	Surrogates and Proxies	149
4.6.2	Confounding factors	150
4.6.3	Collider bias and M-bias	152
4.7	Data alone is not enough	153
4.8	Examples	156

4.8.1	Covid-19: How can efficacy versus severe disease be strong when 60% of hospitalized are vaccinated?	156
4.8.2	Misinterpretations of hurricane forecast maps	157
5	Ethics and Responsibility	161
5.1	Overview	161
5.2	Morality and Ethics	162
5.3	Ethical Frameworks	164
5.3.1	Consequentialism	164
5.3.2	Deontology	165
5.3.3	Virtues	166
5.4	Values in Science	167
5.5	Ethics in action	169
5.6	Data Ethics	172
5.6.1	Origins	172
5.6.2	What are data ethics?	173
5.7	General Data Protection Regulation (GDPR)	176
5.7.1	Personal data	176
5.7.1.1	Special Category Data	177
5.7.1.2	Children's data	177
5.7.2	Principles GDPR	177
5.7.2.1	Lawfulness, Fairness and Transparency	177
5.7.2.2	Purpose Limitation	178
5.7.2.3	Data Minimization	178
5.7.2.4	Storage Limitation	178
5.7.2.5	Integrity and Confidentiality	178
5.7.2.6	Accountability	179
5.7.3	Controllers, Processors and Subjects.	179
5.7.4	Rights of the data subject	179
5.7.5	International transfers	180
5.7.6	Data breaches	180
5.8	Examples	181
5.8.1	A Genocide Incited on Facebook	181
5.8.2	The Theranos scandal - A drop of blood for hundreds of different assays	182
6	Extra Material	185
6.1	History of science	186
6.2	Theory-relatedness of observations	186
6.3	Thomas Kuhn and the idea of scientific revolutions	186
6.4	Gettier problems	187
6.5	Realism and anti-realism	187
6.6	Pessimistic meta-induction	187

Preface

This class book gathers the contents addressed in the course **Applied Philosophy of Science and Data Ethics** from the Master of Data Science¹ at the University of Luxembourg (UL). This course will introduce basic philosophical and scientific concepts supported by examples and discussion. The course expects pro-active participation from the students in the form of presentations and essays as well as open debates.

Course presentation

As much as data science involves automating tasks, we should avoid falling for the automatization of the mental processes we undergo to solve new problems. We data scientists are used to dealing with data as much as we do with the methods used to process such data. Still, rarely, if ever, we stop to question the methodologies that justify employing one method or another. Surprisingly for many, data science is as much about science as about data. We have refined the data acquisition tools and technology stacks used for accomplishing all tasks covered under the data science umbrella. However, below the layers of tools and frameworks that facilitate our daily work, there lies a foundational ground of assumptions, principles, and constraints that shape the way we do data science.

During the last few years, I have taught the students of the Master of Data Science to question such foundations and the reasons behind many tools' weaknesses and strengths. In this course, we delve into philosophical issues, such as Hume's problem of induction, and link them with the recent challenges of artificial intelligence (AI) in the domains of my working experience, like healthcare, systems biology, and computer networks.

We tackle scientific explanation, causation, and association, learning the strengths of Hempel's covering-law model and how it fails to generalize once causality enters

¹https://wwwfr.uni.lu/formations/fstm/master_of_data_science

into play. We discuss the issues of empiricism with causality and showcase the power of causal inference in epidemiology. This leads us to the role of biases, confounders, and surrogate attributes in AI performance. The core of the course is the scientific method, we jump from Popper’s attempt to make science purely deductive to the risks of confirmation, and how diversity plays a crucial role in shaping the quality of datasets. The students get acquainted with the power of randomized controlled trials and learn to relate the strengths of randomization and stratification to current AI techniques such as cross-validation. This is accompanied by historical cases from scientists and doctors like Semmelweis, Fisher, Yersin, et al. They showed us how to properly set scientific hypotheses and how experimental control can play both in favor and against us. This leads to the final part of the course, data ethics, where we learn how cases like the Tuskegee tragedy shaped current regulations for informed consent, and how we can turn ethical values into actions to prevent our solutions from discriminating or wrongly modeling the phenomena at issue.

I believe this course should be taught in every computer science faculty. It requires constant updates to follow the trending topics. One year is COVID-19 and the misleading statistics of vaccine effectiveness put by anti-vax supporters; the next year is the issues of social media applications among teenagers; and today the consequences of generative AI. However changing the present may be, past theories provide a strong base material for this course that students learn to appreciate. The topics learned in this course will remain useful independently of technological advancements. They are essential knowledge for future data scientists to face the upcoming challenges and build the solutions to overcome them.

Course at UL

Most of the contents of this book are taught in the course “Applied Philosophy of Science and Data Ethics” in the Master of Data Science at the University of Luxembourg (UL). The course goal is to provide the students with guidelines and methodologies to identify epistemic and ethical issues present in data science. We expect the students to develop a critical eye that helps them mitigate such problems in their daily work as data scientists.

During this course, students will learn by example different layers of the scientific method and how they relate to data science and data ethics. In particular, they will learn how the mechanisms behind the data affect the data analysis, and how the different types of scientific inference condition method choice and affect the conclusions drawn from the analysis. In this sense, examples of statistical abuse, misconduct and bad visualization will be shown together with their, sometimes catastrophic, collat-

eral consequences.

Target audience

Before we describe the course topics and objectives, we must acknowledge the characteristics of the students who register for the Master of Data Science at the University of Luxembourg as they will shape the course contents and structure. The study program receives students from the European Union (EU) and abroad, adding several dimensions of diversity, such as native languages, study frameworks, and cultural differences, which increase the richness of the class debates but also showcase differences in the student's background knowledge. Of course, we have students from neighboring countries like Germany, France, or Belgium who are also part of the Bologna zone, ensuring comparability in the standards and quality of higher education qualifications. Students from other continents (e.g., Asia, Africa, and South America) can sometimes constitute most enrolled students depending on the year. On top of these differences, students come from different bachelor backgrounds (e.g., computer science, statistics, biology, physics), and some have had professional experience in their fields. The Master in Data Science is a program from the Department of Mathematics, with half of the courses in the first two semesters consisting of mathematics courses. By way of example, one key difference that can be appreciated during the course in students both from the EU and abroad is their knowledge regarding propositional logic, which is usually taught (if at all) in high school curriculums.

Course objectives

In line with the European Quality Framework, Bachelor degrees require a critical understanding of theories and principles, while Master degrees involve higher specialised knowledge and critical awareness of knowledge issues in a field. In this case, the field at issue is data science and the contents will tackle philosophical and ethical issues concerning data science.

Yet, the aim of this course must not be regarded as offering a complete course on the philosophy of science and data ethics but rather as a leveling course to acquaint the students with the concepts of the philosophy of science and how they relate to the challenges of data science and the current ethical debates of the information era. Importantly, as Dr Grüne-Yanoff highlights, teaching the evolution of scientific methods is pedagogical, especially if done through strategies of hyperbole, skepticism, and discussing historical scientific and ethical errors (Grüne-Yanoff, 2014).

The course objectives include:

1. Critically appraise the goals of science, so students can distinguish prediction, explanation, understanding, and design, and grasp how different philosophical views shape what counts as scientific success.
2. Justify methodological choices, because data science involves implicit assumptions and conventions, and students must be able to explain their decisions and recognize how methods evolve over time.
3. Understand scientific inference, to differentiate between deductive, inductive, and abductive reasoning and relate these to reasoning strategies used in both science and AI systems.
4. Grasp the epistemology of uncertainty, to understand how science moves from data to belief and why traditional models like the Hypothetico-Deductive method fall short in data science, where probabilistic reasoning, model comparison, and frameworks like Bayesian and frequentist inference better capture uncertainty.
5. Acknowledge the limits of data, so students appreciate the role of domain knowledge, assumptions, and the iterative nature of inquiry when moving from data to generalizations.
6. Discriminate between correlation and causation, as confusing the two remains a widespread and harmful error in data-driven fields.
7. Recognize and mitigate biases, because biases, whether from measurement, sampling, confounding, or algorithms, can distort every stage of the data pipeline.
8. Understand the role and limits of models and experiments, since modeling, experimental control, and design decisions shape what can be inferred, especially regarding causality and generalizability.
9. Integrate non-epistemic values, because ethical, legal, and social considerations directly impact how data science is practiced and evaluated in real-world contexts.
10. Connect concepts, as a meta-goal students should relate philosophical, methodological, and ethical dimensions into a coherent understanding of data science practice.

The justification of these objectives derives from the work these future data scientists will do. Data science practice is driven by legal, ethical, and resource pressures, so students must see how these forces shape their methods (Obj 9-10) and where each tool reaches its limits (Obj 2; Obj 8). Spotting bias, spurious links, and causal traps

in messy data (Obj 6-7) depends on a strong understanding of science's aims, the logic of inference, and the role of domain knowledge (Obj 1; Obj 3-5). The key takeaway message from this course is to question the data science methods, acknowledge that data without domain knowledge is insufficient, and recognize the current ethical debates and how they guide and constrain the data science developments of our time.

Course structure

The course is part of the first semester of a two-year master's program. It consists of five chapters covered over 14 weeks, with a 90-minute class each week. These sessions are supplemented by weekly online quizzes, readings, and essay assignments. Depending on the pacing and in-class participation, topics may shift between weeks. To address potential scheduling issues, some video lectures cover topics in detail and can replace in-class lectures if needed to allow more time for discussions, exams and debates. These videos also include foundational content, such as an introduction to basic concepts of propositional logic in Chapter 2.

The current chapters and the order they are taught in class are:

- Chapter 1: Scientific Goals, Methods, and Knowledge (Week 1)
- Chapter 2: Scientific Inference (Weeks 2-6)
- Chapter 3: Empirical Practices and Models (Weeks 7-9)
- Chapter 4: Experimental Control and Statistical Abuse (Weeks 10-12)
- Chapter 5: Ethics and Responsibility (Weeks 13-14)

About this course

Despite deep historical ties between philosophy and computer science, especially in areas such as logic, the two disciplines have been largely taken apart largely in recent decades. Computer science curricula have increasingly focused on technical and practical applications, such as programming, algorithms, and systems design, while philosophy courses tackling theoretical inquiry and critical thinking have been relegated to other departments. This separation has diminished opportunities for students to explore foundational questions about the nature of computation, the ethics of technological innovation, and the philosophical implications of artificial intelligence. Below, a selection of courses and resources is available as guide for future instructors in the philosophy of science for data science. Moreover, the rest of the paper is accompanied with examples and references related to the course content which further extend the materials listed below.

While optional and online courses in topics like philosophy of science and data ethics are available, these subjects are not typically integrated into the core courses of computer science departments. Philosophy of science, in particular, is rarely positioned as a fundamental component of computer science education. This gap calls for better interdisciplinary in computer science programs to bridge technical skills with philosophical inquiry.

Some of these online courses which align with the philosophy of science and data ethics and which have influenced this course, are described below.

The first chapters of this course are inspired by the book from Prof. Dr Lars-Göran Johansson ([Johansson et al., 2016](#)) “Philosophy of science for scientists”. For several years, Prof. Dr. Till Grüne-Yanoff taught an online course (now unavailable) titled “Philosophy of Science for Engineers and Scientists” on the edX platform, which combined theoretical concepts with practical examples. Currently, he teaches courses on the methodology of science at KTH Royal Institute of Technology in Stockholm, Sweden. He advocates for the revision of the traditional curriculum of Philosophy of Science to better align with modern scientific disciplines ([Grüne-Yanoff, 2014](#)). Regarding the second part of the course, which covers data ethics, I would like to thank the University of Michigan, for their online courses (from which I was already a fan during my PhD). Especially, Prof. Dr. H. V. Jagadish’s Data Science Ethics course, offered on Coursera by the University of Michigan, which has influenced this course’s data ethics content. The University of Oxford’s Department for Continuing Education provides short courses on Philosophy of Science and Data Ethics, which offer valuable opportunities for expanding knowledge in these areas.

However, the integration of Philosophy of Science into computer science curricula remains limited and the author has struggled to find dedicated resources and courses. Some exceptions exist, such as the University of Bayreuth in Germany, which offers a Master’s program in Philosophy & Computer Science, and the University of Oxford, with an undergraduate program in Computer Science and Philosophy. More recently, the University of Cordoba in Spain will debut a new bachelor in Mathematics and Philosophy in 2025. These programs highlight the potential for interdisciplinary approaches, though they remain rare in broader computer science education.

Additionally, a great variety of books are available and many have shaped the course materials, including: The Book of Why ([Pearl and Mackenzie, 2018](#)), Ethics and Data Science ([Loukides et al., 2018](#)), Understanding philosophy of science ([Ladyman, 2012](#)), Philosophy of Natural Science ([Hempel, 1966](#)). Again, one key work is Philosophy of Science for Scientists ([Johansson et al., 2016](#)), which emphasizes that lessons from natural sciences are relevant to students and researchers in social and

human sciences, to which I must add data scientists. Other books aimed toward the general public such as Philosophy of science: A very short introduction (Okasha, 2016) or Philosophy of science : a new introduction (Barker and Kitcher, 2014) offer a variety of examples to illustrate the topics of Philosophy of Science. Some more have influenced examples and ethics lessons such as How charts lie (Cairo, 2019) and Automating inequality (Eubanks, 2018).

Since 2021, the author regularly updates this online class book that suits both students and instructors with topics taught in class and additional material.

I hope any resemblance or imitation is seen as an act of flattery.

Roadmap for instructors

As a brief roadmap for instructors without formal training in philosophy, this author believes (Okasha, 2016) and (Pearl and Mackenzie, 2018) (chapters 1 to 5) offer an accessible way into core issues of Philosophy of Science and causal inference. This can be complemented with (Johansson et al., 2016) (chapters 2, 3, 4 and 13) and (Ladyman, 2012) (chapter 2) as reference books to understand and link relevant concepts. Though, the most important task is to connect such issues with the instructor expertise in data science to enhance the applied aspect of the course.

Disclaimer

Although the impact and extension of the topics addressed in this class book (and the course) are broad and diverse, its length is limited. Hence the scope and depth of the contents are restricted. Consequently, several topics on Philosophy of Science are tackled superficially while some others are completely ignored. Such philosophical questions are handled from a practical data science point of view. Similarly, Data Ethics is a relatively new matter in continuous evolution. Therefore we will try to cope with the main issues in the most practical way.

About the author

Since October 2023 Carlos Vega² works as a Research Engineer in Digital Health at the **Luxembourg Institute of Health** (LIH) in the Department of Medical Informatics led by the CMIO (Chief Medical Information Officer) Dr. Maximilian Fünfgeld. Since September 2023, he is a fellow of MIT Catalyst Europe innovation

²<https://researchportal.lih.lu/en/persons/carlos-vega>

program, co-leading two teams: (1) on improving emergency response for cardiac device patients and (2) tackling school absenteeism linked to period poverty in Ghana.

From 2018 until September 2023, he worked as a postdoctoral researcher in the Bioinformatics Core Group led by Prof. Dr Reinhard Schneider at the **Luxembourg Centre for Systems Biomedicine** (LCSB) at the University of Luxembourg.

Previously, he worked at the Autonomous University of Madrid (UAM) researching high-performance solutions for big data analysis as well as anomaly detection methodologies. He received his B.Sc (2013), M.Sc (2014) and PhD (2018 with *cum laude* and industrial mentions) degrees in Computer Science Engineering from UAM. His research career started in 2012 when he joined the High-Performance Computing and Networking Research Group (HPCN) led by Prof. Dr Javier Aracil, first as a student and later as a researcher as part of the Network of Excellence InterNet Science European project. During his PhD (2014 - 2017), he continued his work at the HPCN group as a technical researcher for the project TRÁFICA and the European projects Fed4Fire and dReDBox, among others. At the same time, he worked at Naudit HPCN³ (2015 - 2018) applying his research in computer network auditing projects with different enterprises.

In 2022, he became a senior member of the Institute of Electrical and Electronics Engineers (IEEE⁴). In 2023, he became a fellow of the MIT Catalyst Europe program⁵.

Additionally, his past teaching experience includes several courses on computer networks, such as Multimedia Networking, Network Planning and Network Management, taught during his time at UAM. Since 2021, he is the main teacher of the course “Applied Philosophy of Science and Data Ethics” for the Master of Data Science at the Faculty of Science, Technology and Medicine (FSTM) of the University of Luxembourg.

About this class book

This class book was made thanks to the great tutorial available on the book “Open tools for writing open interactive textbooks (and more)”⁶.

³<https://www.naudit.es/en>

⁴<https://www.ieee.org>

⁵<https://www.catalysteurope.eu>

⁶https://www.crumplab.com/OER_bookdown

License

Licensed under CC BY-NC-ND 4.0

The book is released under CC BY-NC-ND 4.0⁷ license. This means that you are free to:

- **Share:** copy and redistribute the material in any medium or format.

Under the following terms:

- **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
- **Non-Commercial:** You may not use the material for commercial purposes.
- **No Derivatives:** If you remix, transform, or build upon the material, you may not distribute the modified material.
- **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributing to the book

Contributions are welcomed, feel free to open a pull-request in github for small changes and it will be reviewed as soon as possible. However, for larger contributions (e.g. sections, chapters) please contact the main author at carlos.vega [at] lih.lu .

⁷<https://creativecommons.org/licenses/by-nc-nd/4.0/>

Chapter 1

Scientific Goals, Methods and Knowledge



Figure 1.1: Allegory of the cave¹ from the Greek philosopher Plato. Source: Wikipedia user 4edges².

To gain knowledge about the world, the universe, and the rules behind the phenomena that shape them, we often require taking measurements about the entities we seek to understand. We can indeed play with many of them: we can take a rock and throw it, we can dissect a mouse or mix multiple liquids. However, to compare rocks, mice or liquids, we must take measurements that allow us to tell their differences and commonalities. This information, these data, are not the entities themselves, and when we work with them, we stop playing with the *real* entities. We begin to play with shadows, and two different objects can certainly project the same shadow. The relationships showcased by these entities are not contained in the data

but in the reality they live on, of which the data is merely a blurry shadow of reality. Therefore, data science works with the shadows of real objects, and this practice involves assessing risks, identifying mistakes and avoiding situations of appearance of knowledge.

In this course, we will introduce central topics of the scientific method and translate practices employed in natural sciences to the domain of data science. Such practices are meant to control the plethora of considerations that must be bore in mind when playing with shadows. From the confirmation of hypothesis, asymmetries of explanation and causality, to the experimental control techniques and their relationship with surrogate variables, confounders and colliders.

Illustrated with examples, this course will tackle the scientific goals, the definition of knowledge, the criteria for a relevant hypothesis, the conditions for experiment success and requirements for good scientific models. Once the empirical cycle is clear, the last chapters are dedicated to experimental control, statistical abuse and data ethics. These chapters provide equivalents of techniques employed in natural sciences for their use in data science, relating experimental designs like Randomised Control Trials to causal inference or cross-validation in Machine Learning. Finally, the chapter on data ethics does not just introduce ethical frameworks but delves into how to bring such values to action.

Allegory of the cave from the Greek philosopher Plato

“Plato describes a group of people who have lived chained to the wall of a cave all their lives, facing a blank wall. The people watch shadows projected on the wall from objects passing in front of a fire behind them and give names to these shadows. The shadows are the prisoners’ reality, but are not accurate representations of the real world. The shadows represent the fragment of reality that we can normally perceive through our senses, while the objects under the sun represent the true forms of objects that we can only perceive through reason. Three higher levels exist: the natural sciences; mathematics, geometry, and deductive logic; and the theory of forms.” — Wikipedia^a

^ahttps://en.wikipedia.org/wiki/Allegory_of_the_cave

Course objectives

This chapter supports course objectives 1–3, by introducing key concepts through simple messages and analogies that build understanding without overwhelming students. Relevant examples are resumed in later chapters to connect contents. Higher-level debates (e.g., Thomas Kuhn, realism vs. antirealism) are omitted to focus on the core ideas of the first course weeks.

1.1 What is Science?

[...] questions about scientific methodology and knowledge in philosophy of science are really continuous with questions in cognitive science about how human beings reason and form beliefs. However, one need not imagine an absolute distinction between philosophy and empirical forms of inquiry to appreciate the broad differences between the latter and the study of philosophical questions that arise when we reflect on science. Of course, this characterisation is of little use unless we know what science is, so perhaps the most fundamental task for the philosophy of science is to answer the question, ‘what is science?’. — ([Ladyman, 2012](#))

This question attempts to answer what common features share subjects such as physics or biology to be called sciences, i.e. what it is that which *makes* something a science. The problem of saying what is scientific and what is not is called **the demarcation problem**. Among other things, science aims to understand, explain and predict the world we live in. But also religions, astrology or alchemy attempt to understand, explain or predict our world. What makes them different from science?

Four historical elements are essential for the development of a scientific approach. Namely: to seek explanations of natural phenomena; to argue; to investigate the rules of argumentation and logical validity; to build them into a logically consistent system. ([Johansson et al., 2016](#))

Rather than finding a proper definition of science, which many have struggled with, we will focus on what makes science different and why its methods are called scientific.

1.1.1 Scientific Goals and Knowledge

In this course we consider prediction, explanation, understanding and design as the most accessible goals of science for the students. Typically, scientific inquiry begins with observation which allow us to make predictions, develop explanations and achieve understanding of nature phenomena. Though, as we will see, these do not necessarily follow a fixed order. Thanks to our understanding we can as well design experiments, instruments and solutions that help us better gather data and further explain, understand, and predict our world. However, the pursuit of one goal does not necessarily contribute to others, with trade-offs often at play ([Potochnik, 2015](#)). The goals of data science mirror this pluralism but all these goals share a common ingredient: knowledge.

Predicting X means knowing that at time t , X will happen. Explaining X means to know the cause(s) that produced X . Designing X requires knowing that artifact X

will satisfy certain functions F . Again, all these goals share a common ingredient, scientific knowledge. Scientists arrive to such knowledge by applying the scientific method (see § 1.2). The goals of science are achieved through a series of activities that constitute the scientific method which include systematic observation and experimentation, inductive and deductive reasoning, and the formation and testing of hypotheses and theories.

Knowledge is justified true belief — Plato (428 - 348 BC)

The most popular definition of knowledge was given by philosopher Plato in the above's quote. This definition specifies that a statement must meet three criteria to be considered knowledge. This definition of knowledge is sufficiently good for this course. However, the definition of knowledge is an ongoing debate among epistemologists. Although these criteria are necessary conditions, they are not sufficient as there are situations that satisfy all these conditions and yet don't constitute knowledge (see Gettier cases³) but such cases are rather philosophical and will not be discussed during this course.

- **True** because statements must refer to an actual state of the world.
 - A wet sidewalk does not necessarily imply it rained even if you believe so.
 - Even if we are justified to believe that something is true, it might not be true.
- **Justified** because you need proper proof, evidence or reasons to defend our statement.
 - Even if it actually rained, a wet sidewalk caused by a sprinkler is not good justification for you to believe it rained.
- **Belief** because even under justified reasons about true facts, people can choose not to believe such knowledge. We define belief as to the state of mind of a person that thinks something is the case. This state of mind is of course *tied to the individual* and *comes in degrees*. We act based on our beliefs and values, and new knowledge can affect these.

³https://en.wikipedia.org/wiki/Gettier_case

Content note

This definition of knowledge serves as a simple but useful entry point into epistemological thinking for those unfamiliar with these concepts. Although it does not fully capture more current epistemological views like Bayesian epistemology, it is a pedagogical and mnemonic tool to introduce several concepts at once. First, values like honesty and transparency from true. Then, justification, which is central in method choice and model selection. Last, belief, which reminds us that subjective priors and cognitive biases influence reasoning. Through the course we will also see their shortcomings.

Certainty of belief and truth are different. Is possible to have certain beliefs about false claims. Similarly, we can have uncertain beliefs about true claims. From tossing a coin, we can expect a fair probability in which head and tails have the same probability. But we cannot know for sure if the coin is biased or not until it lands. Similarly, is possible that even our best theories are wrong or partially wrong. Even after many successful experiments, they might be proved wrong (see 1.4.1). In fact, scientific hypotheses can rarely if ever be proved right, they can, however, be proven wrong.

“We never are definitely right, we can only be sure we are wrong” —
Richard Feynman

Below you can find a clip from the last lecture⁴ of a series of 7 special Messenger Lectures given by the renowned American theoretical physicist Richard Phillips Feynman. The transcription can also be found online⁵.

While is relatively easy to determine cases of failed justification, is much harder to identify what suffices to justify a belief. Few claims can be conclusively proven so that no doubt remains. An ideal justification of a belief would consider all relevant reasons for and against believing a statement. This is why science is a human enterprise where justifications, hypotheses and experiments are made public for review, replication or rebuttal.

Definitions should not be accepted without reason, and instead, we should attend to the arguments that support such definitions. Certain definitions may have widespread popularity but that doesn't make them any more true. For example, a dolphin is a mammal even if many people consider it a fish. In the same way, tomatoes and cucumbers are fruits for botanists even if we daily sort them as

⁴<https://youtu.be/ECY-4Ng9Nkc?t=1190>

⁵<https://sites.google.com/site/barrykort/feynman-on-the-scientific-method>

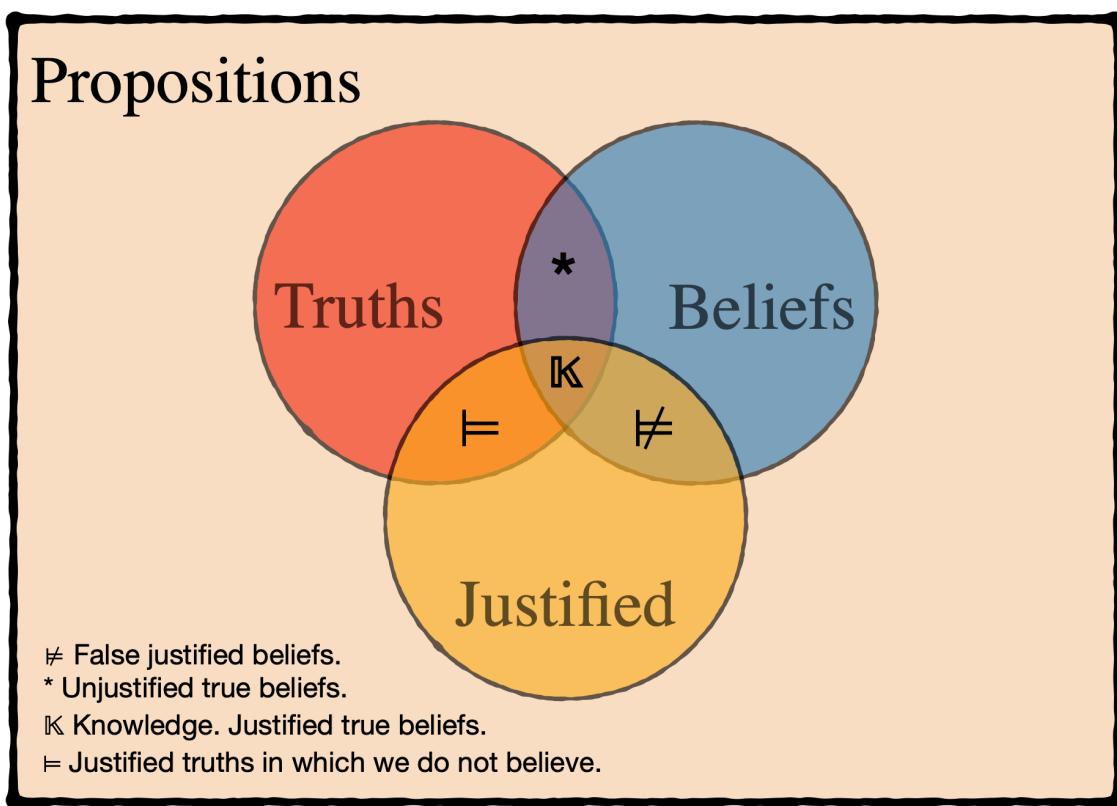


Figure 1.2: A Venn diagram illustrating the classical theory of knowledge.

vegetables. And sometimes, even the EU⁶ and the Supreme Court of the United States of America⁷ need to act to set certain market debates.

As another example, the first astronomers who lacked the telescope believed on the geocentric model because their observations did not suffice to reject it. These first astronomers had **false justified beliefs**. After the advent of the telescope in 1609, the geocentric model was rejected. But how do we know we are not in a similar situation that of the pre-telescope astronomers? Events as recent as the Michelson-Morley (see 1.4.2) experiment, the expeditions of Sir Arthur S. Eddington or the predictions of Urbain Le Verrier (see 1.4.1) have changed our conceptions of the universe and physics forcing scientists to re-formulate models and theories.

In the next chapter we will see how knowledge is obtained.

One [fundamental question of philosophy of science] is ‘how can we have knowledge as opposed to mere belief or opinion?’, and one very general answer to it is ‘**follow the scientific method**’. [...] The branch of philosophy that inquires into knowledge and justification is called epistemology. The central questions of epistemology include: what is knowledge as opposed to mere belief?; can we be sure that we have any knowledge?; what things do we in fact know?. — (Ladyman, 2012)

1.1.1.1 Data, information and knowledge

Nowadays, technology allows us to collect data into datasets, transform datasets into information and arrive at new knowledge. Such processes have always been crucial in science but computer science comes to question concepts such as data, information and knowledge. (Johansson et al., 2016)

By **knowledge**, we can understand three different things. First, knowledge of truths, e.g. we know that the sun rises on the east. Such knowledge can be obtained by reading a book or listening to the radio. The second category of knowledge consists of skills, such as riding a bike or speaking a foreign language. However, this knowledge requires more than language to be communicated. It requires practice. Finally, the third category is the knowledge of objects, what Bertrand Russell called knowledge by acquaintance (Russell, 1912). This knowledge is obtained through experience.

In common English, we can't distinguish between knowledge of truths and objects. However, languages such as German, French or Spanish make clear this distinction by using different verbs, Rusell proposed to use the word “acquaintance”.

⁶<https://eur-lex.europa.eu/LexUriServ/LexUriServ.do?uri=CELEX:32001L0113:En:HTML#d1e32-72-1>

⁷<https://www.nationalgeographic.com/culture/article/fruit-or-vegetable>

- German: wissen vs kennen.
- French: savoir vs connaître.
- Spanish: saber vs conocer.

Example from William James (1890) — “I am acquainted with many people and things, which I know very little about, except their presence in the places where I have met them. I know the colour blue when I see it, and the flavour of a pear when I taste it; [...]”

What is the difference between data and information?

The following excerpt from ([Johansson et al., 2016](#)) may clarify this question:

But why call the input ‘information’? The reason seems to be that we can describe the input as being about something, often the state of the environment. It has content. Or rather, when we humans describe the input and the workings of the system we find it natural to talk as if the information-containing system consciously sent messages to us humans; we say that the systems obtain information, transmit information or store information about something, as if it were like a human mind. The core feature of this use of the word ‘information’ is thus its aboutness, its intentionality.

Finally data. It is common in computer science to say that information is data with meaning. This is ok as far as it goes, but what is ‘meaning’? And how do data acquire meaning? It seems that minimally it means that meaningful data becomes information when we have been able to formulate declarative sentences expressing the information that is obtained from a data set. Almost anything can be data. In order to obtain data from e.g. a story, from light from distant stars, or from the result of an experiment, we need to divide the stream of sounds, lights, or states of detectors into distinct items. When using written text as data source one must divide the string of linguistic signs into distinct items, such as words or longer or shorter expressions. [...] In short, in order to obtain a data set, we need to define a principle for dividing up something into distinct pieces. Hence from a conceptual point of view, discerning data and collecting a data set presupposes that we have a prior principle of making distinctions within a phenomenon. [...] Sometimes we have lot of background knowledge from start.

In short, a piece of knowledge is a piece of information for which the knower can provide good reasons.

1.1.2 What is Philosophy of Science?

One of the tasks of philosophy of science is to question assumptions that scientists take for granted. For example, suppose a scientist conducts an experiment that yields a particular result. The scientist then repeats the experiment a couple of times more obtaining the same result. The scientist then stops repeating the experiment, convinced that repeating it under the same conditions will produce the same result. But *why* does the scientist assume that future repetitions will provide the same outcome? How are we sure this is true?

Therefore, one of the main objectives of philosophy of science is to study the methods and methodologies of enquiry used in the sciences, understanding how techniques like experimentation, observation and theory building enable scientists to reveal new knowledge. Philosophy of science asks questions such as: What is knowledge? What is a scientifically acceptable observation? What makes an explanation scientific? What is a scientific theory?

Finally, the philosophy of science tackles a wide range of topics that would require its own master. Moreover, not all topics are directly related to the aims of this course and the scope of the master. For this reason, a brief summary of the topics left behind is included in § 6. Of course, the curriculum is subject to change in the future and the list might change too.

It is usually thought that if there is anything of which science consists it is a method or set of methods, so the study of scientific method (known as methodology of science) is at the centre of the philosophy of science.
— ([Ladyman, 2012](#))

1.2 The scientific method

The scientific method is the main pillar of science. All science begins with *observation*, as this is the **first step** of the scientific method. Moreover, such observation must be *repeatable*, either actually or potentially. Once an observation has been made, the **second step** involves the definition of a *problem*, or in other words, asking a question about the observation. However, such a question needs to be valuable scientifically, it must be *relevant* and must be *testable*. Questions need to be reformulated until they become testable. The **third step** may seem a rather unscientific procedure as it involves guessing what the answer to the question might be by postulating a *hypothesis*. The **fourth step** will tell the scientist if the *hypothesis* is correct through *experimentation*, which tests the validity of a scientific guess. Notwithstanding, experiments do not guarantee a scientific conclusion. Ex-

periment results represent *evidence*, i.e. the hypothesis in answer to the question is confirmed as correct or invalidated. Given the latter, a new hypothesis with new experiments might be needed. Finally, experimental evidence is key for the **fifth step** of the scientific method, the formulation of a *theory*. A good theory has a *predictive* value, usually predicting that something is *likely* to happen with a certain degree of probability ([Nidditch, 1968](#)).

1.3 Methodology

A method is a particular tool to reach a particular goal (e.g. statistical test). Methodology is the systematic assessment and justification of method choice. Scientists often need to choose between alternative methods in order to reach a particular scientific goal. But specifying a goal does not directly determine what method to choose. We need to consider the reasons why some method is better than another for a particular goal. This process could require a better definition of the initial goal or learning more about the context and domain where the methods will be applied. Methodology must be distinguished from describing methods, which usually concerns the design and implementation of particular research approaches and focus on the technical aspects (e.g. how to program simulations or set up instruments).

For example, a laboratory experiment can be advantageous because the test conditions can be controlled but laboratory experiments might not be realistic enough for certain tests. On the other hand, a field experiment provides more realistic test conditions but is difficult to control all variables.

Similar considerations may be necessary for other seemingly trivial questions such as model choice or data visualization. Should we use a significance test or a Bayesian approach? Should we present our results using a bar chart or a violin plot? Should we use a structural model or a quantum model? Methodology asks questions such as: What methods are available to reach a particular goal? What reasons speak for or against the alternatives? How should be weight the reasons to form a final decision?

How do we decide between alternative methods? Is there a way to determine what is rational to choose? Traditionally there are three ways to choose between alternative methods.

By convention, The methods are chosen because you have been taught to, or because is an established convention between your peers. Conventionalisms create long-term issues when methods become dominant in a field. A good example is the use of p-value in hypothesis significance testing. Similarly, accuracy and precision metrics in Machine Learning can be considered conventionalism. More problems arise when

different disciplines have different conventions, hindering inter-disciplinary work.

Outcome-oriented. While choosing the method that yields the best results may seem well-intended and appropriate, this certainly sounds very vague too. The intention is to find a method that serves some purpose best, but this purpose is sometimes not sufficiently clear. Science frequently involves long-term projects where the final material outcome is uncertain or unknown. For example, the International Space Station or the Large Hadron Collider. This methodology raises the question of how to measure the outcome. For instance, is speed the best way to assess which car is best? Should we focus on fuel autonomy or pollution instead? What about combining all of them?

Reason-based. Choosing the method based on the overall best reasons seems the best option, particularly when the reasons include considerations that justify choosing a method over others for a given scientific goal (e.g. prediction). But sometimes there are methods that despite providing more valuable results could be unethical and/or illegal. For example, randomized control trials (RCT) are often employed to test the effectiveness of a new drug. Participants are divided *at random* into two groups (treatment and control), eliminating the effect of confounding factors on the outcome of interest. However, RCTs are not always feasible, for either practical or ethical reasons. For instance, it won't be ethical to assign people to smoke for decades in order to study if cigarette smoking causes cancer. These other aspects need to be weighted together with the scientific reasons during method choice.

See 1.4.4 for an example of how reason-based methods are not always easy to implement while at the same time, outcome-oriented methods led the mainstream of an important debate.

1.4 Examples

1.4.1 Neptune and Vulcan

Newton's gravitational theory predicted the paths the planets should follow as they orbit the sun. Most of these were confirmed by observation, but the orbit of Uranus differed from Newton's predictions. In 1846 John Adams in England and Urbain Le Verrier in France solved the mystery. Both of them suggested that another planet, yet undiscovered, was the cause of an additional gravitational force exerted on Uranus. These scientists calculated the mass and position that this planet would need to have to explain Uranus' orbit. The planet Neptune was indeed found close to the location predicted by Adams and Le Verrier.

So, instead of rejecting Newton's theory right away (see 2.4.2), these scientists stuck to it and tried to find another missing factor that could explain the difference. When the motion of Uranus was found not to match the predictions of Newton's laws, the theory "There are seven planets in the solar system" was rejected, and not Newton's laws themselves.

However, Le Verrier also found irregularities in the motion of the planet Mercury and tried to explain them as resulting from the gravitational pull of an, again, yet undetected planet Vulcan. This hypothetical planet would have to be a very dense and small object between the sun and Mercury. In this case, no planet was found between Mercury and the sun. A satisfactory explanation was provided much later by the general theory of relativity, which justified irregularities through a new system of laws. In this case, the hypothesis or theory had to be reformulated or replaced by new one.

Below you can find a clip from lecture "The Law of Gravitation"⁸, from the Messenger Lectures given by the renowned American theoretical physicist Richard Phillips Feynman.

1.4.2 The most famous "failed" experiment

The Michelson-Morley experiment (1887) was designed to detect the motion of the Earth through the luminiferous aether. XIX century physicists used aether to explain how light could be transmitted through empty space between the Sun and the Earth. The result of this experiment is considered to be the first strong evidence against the then-prevalent aether theory, and the beginning of a new line of research that eventually led to special relativity, which rules out a stationary aether.

⁸<https://www.youtube.com/watch?v=j3mhkYbznBk&t=1830s>

To the ancients, the concept of a void universe was impossible. Aristotle arrived at the hypothesis of the aether to explain the cosmos and several natural phenomena such as the movement of the planets. By the XIX century, the aether became more than a philosophical need. Whenever there is a wave, something must be waving. But what waves when light waves travel from the Sun? For XIX physicists, the aether was the medium through which light waves from the Sun would propagate.

Michelson and Morley attempted to detect the absolute motion of Earth through space. For that, they set an experiment in which a beam of light was sent through a half-silvered mirror used to split the light beam into two beams travelling at right angles to one another. The beams were then reflected back to the half-silvered mirror by two respective mirrors and recombined into a single beam. The experiment can be seen as a race between two light beams. If the beams arrive in a tie, the result is a bright spot at the centre of the interference pattern, otherwise, a destructive interference would make the centre of the image dark. The hypothesis foretold that a tie was not possible since the two beams were racing on a moving track. It was assumed that the Earth was moving through the aether and therefore the beam should trace different paths with respect to the aether.

The extent to which the negative result of the Michelson–Morley experiment influenced Einstein is disputed. However, the null result helped the notion of the constancy of the speed of light gain acceptance in the physics community. This example shows the impact a well-designed experiment can have.

For a longer and deeper explanation of the experiment, its historical context and consequences, watch episode 41 from *The Mechanical Universe*⁹. The timeline of luminiferous aether can be found at the Wikipedia¹⁰.

1.4.3 Eddington expeditions

The following example also relates to falsification (see 2.4.2) which is taught in the next Chapter. However, this is also a good example of how a good theory should make definite predictions such as those from Einstein's theory of general relativity.

Figure 1.4 shows the positions of different stars during the eclipse. Such stars are not normally visible in the daytime due to the brightness of the Sun but become visible during the moment when the Moon fully covers the solar disc. A difference in the observed position of the stars during the eclipse, compared to their normal position at night, indicates that the light from these stars had bent as it passed close to the Sun.

⁹https://www.youtube.com/watch?v=Ip_jdcA8fcw

¹⁰https://en.wikipedia.org/wiki/Timeline_of_luminiferous_aether

Einstein's theory made a clear prediction: light rays from distant stars would be deflected by the gravitational field of the sun. Normally this effect would be impossible to observe — except during a solar eclipse. In 1919 the English astrophysicist Sir Arthur Eddington organized two expeditions to observe the solar eclipse of that year, one to Brazil and one to the island of Principe off the Atlantic coast of Africa, with the aim of testing Einstein's prediction. The expeditions found that starlight was indeed deflected by the sun, by almost exactly the amount Einstein had predicted. Popper was very impressed by this. Einstein's theory had made a definite, precise prediction, which was confirmed by observations. Had it turned out that starlight was not deflected by the sun, this would have shown that Einstein was wrong. So Einstein's theory satisfies the criterion of falsifiability. — ([Okasha, 2016](#))

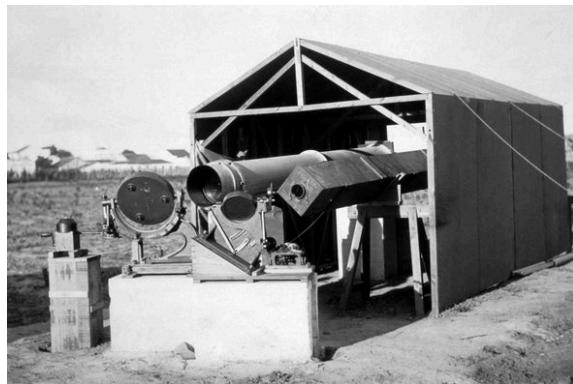


Figure 1.3: Instruments used in the 1919 observations to test Einstein's predictions about warped spacetime. Credit: Getty Images

Einstein published his general theory of relativity in 1915. The total solar eclipse of 1919 offered the perfect opportunity to test it experimentally, by exploring whether — and how — the immense gravity of the Sun bends and distorts incoming light from more distant stars, as predicted by Einstein's theory. For a brief moment during the eclipse, the Moon would block the Sun's light in the sky and make visible some of the stars that lie close to the line of sight of the Sun, not normally visible during the daytime. By measuring the positions of these stars during the eclipse and comparing them to their positions at night, when the sun is not in the field of view, it would be possible to determine whether their light rays bends while passing close to the Sun. — European Southern Observatory

One of the interesting facts from Stanley's account is that Einstein had

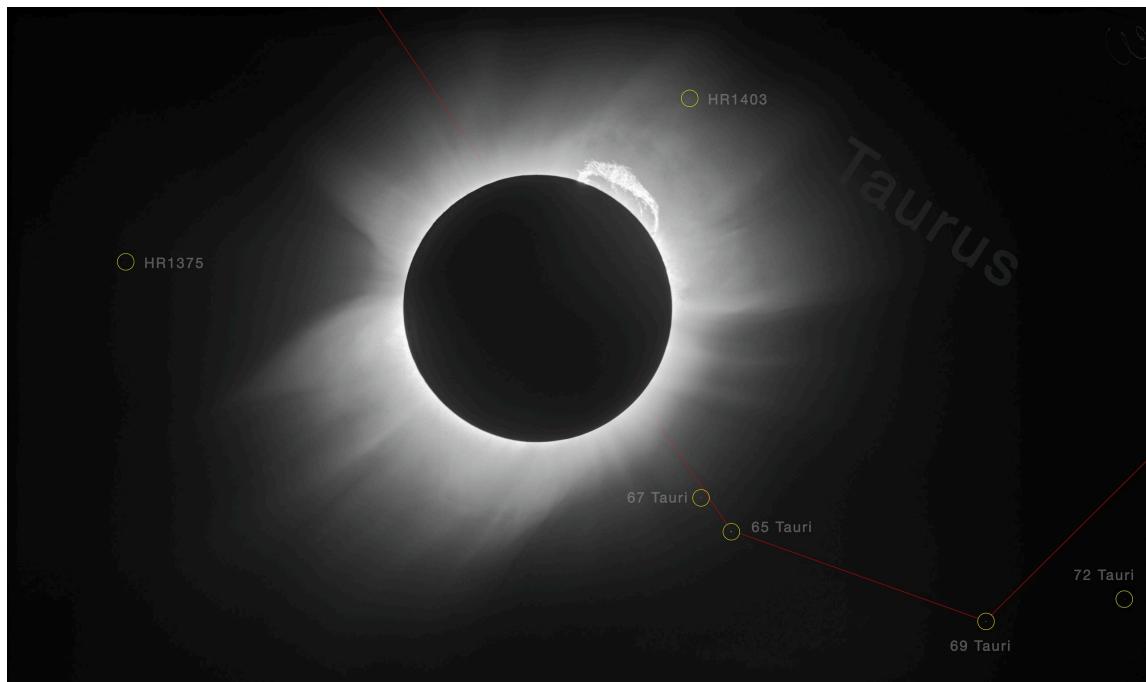


Figure 1.4: Eddington and Crommelin imaged the eclipse using the technology of the time: photographic plates made of glass. Sadly, the original plates from the 1919 expedition have been lost — but, luckily, copies of one of the plates were made and sent to observatories around the world to allow scientists everywhere to see the evidence in support of relativity with their own eyes. Source: European Southern Observatory¹¹.

made a stab at calculating the bending of light back in 1911, before he had formulated the full general theory of relativity. His result was precisely the same as the Newtonian value. I was left wondering what would have happened to his reputation if measurements had been taken then. Would they have been a setback? Or would they just have driven him harder to produce the full theory, with its crucial factor of two? — (Coles, 2019)

1.4.4 The smoke debate

In the mid-1700s, James Lind discovered that citrus fruits prevent scurvy, while in the mid-1800s, John Snow figured out that water contaminated with faecal matter caused cholera. These two examples share a common fortunate one-to-one relation between cause and effect. Deficiency of vitamin C is necessary to produce scurvy. Similarly, cholera bacillus is the only cause of cholera.

However, during the late 1950s and early 1960s, whether or not smoking caused lung cancer was not clear. The subject of the debate wasn't tobacco or cancer but rather the word *caused* as one of the most important arguments against the smoking-cancer hypothesis was the possible existence of confounding factors that may cause lung cancer and nicotine dependency. Many smokers live long lives without getting lung cancer while others develop cancer without ever smoking. Plotting the rates of lung cancer and tobacco consumption makes the connection impossible to miss (See Figure 1.5). However, time-series data are poor evidence for causality. Researchers already knew about RCT though its use was unethical in this case.



Sir Arthur Stanley Eddington (1882–1944).

Austin B. Hill proposed to compare patients already diagnosed with cancer to a control group of healthy volunteers. The results showed that all but two of the 649 lung cancer patients had been smokers. This type of study is today called a case-control study because it compares cases to controls. However, **this method has some drawbacks too**. First, the study is retrospective, meaning that participants known to have cancer are considered and researchers look back to understand why. Second, the probability logic is backwards, as the data tell us the probability that a cancer patient is a smoker instead of the probability that a smoker will get cancer. Moreover, case-control studies admit several possible sources of bias such as recall bias or selection bias. Hospitalised cancer patients were not a representative sample of the population, not even from the smoke population. Researchers were careful to

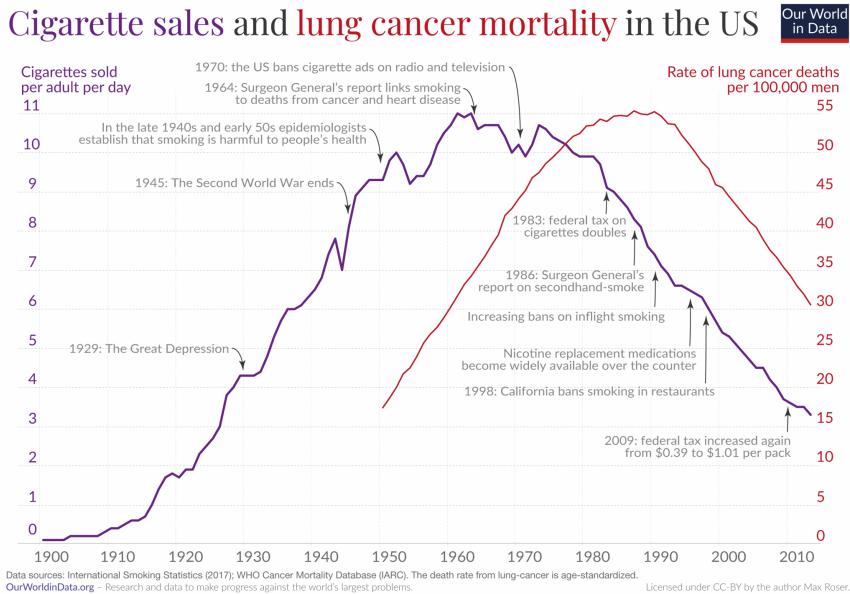


Figure 1.5: Source: Our World in Data¹².

call their results an “association”. Later on, the study was replicated with similar results. Deniers such as R. A. Fischer were right to point out that repeating a biased study doesn’t make it any better as is still biased.

We won’t focus on how this story ends here but **is important to notice how methods chosen based on scientific reasons are sometimes tough to implement and often need to fight against outcome-oriented studies** such as those sponsored by leading tobacco companies.

In the end, many subsequent studies settled the smoking-cancer debate. We will come back to this example in upcoming sections of the course. If you can’t wait, read Chapter 5 from the Book of Why, by Judea Pearl and Dana Mackenzie ([Pearl and Mackenzie, 2018](#)).

1.4.5 Kekulé’s dream

The third step of the scientific method (see § 1.2) requires guessing an answer (or hypothesis) to a previously determined question. There is no clear method to arrive at a hypothesis. Experience, historical context, and previously failed hypothesis condition how a hypothesis is conceived. But sometimes hypotheses can be reached in the most unlikely and unconventional of ways. It makes no difference as long as the hypothesis is then scientifically tested before its acceptance. One of the most famous examples is the structural model of the benzene molecule. In 1865 the chemist August Kekulé hit on the hypothesis of the structure after dreaming of a snake trying

to bite its tail (See Figure 1.6).

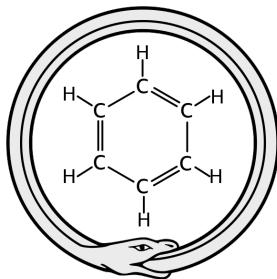


Figure 1.6: Source and credits to: Haltopub, from Wikimedia¹³.

1.4.6 Why Most Published Research Findings Are False

In his Why Most Published Research Findings Are False¹⁴, John P. A. Ioannidis (Prof. at the Stanford School of Medicine) argues that “the probability that a research claim is true may depend on study power (see note) and bias, the number of other studies on the same question, and, importantly, the ratio of true to no relationships among the relationships probed in each scientific field” (Ioannidis, 2005).

In the form of corollaries, the author enumerates several factors making a paper with a positive result more likely to enter the literature and suppress negative-result papers. For instance, the popularity of a field or a trendy topic has more chances to have more publications about. Financial or political interests can act as prejudices affecting research and publication. Other factors such as sample size and effect sizes (e.g. the strength of the relationship between two variables in a population), experimenter bias, white hat bias (e.g. cherry picking the evidence and publication bias); flexibility in designs, definitions, outcomes, and analytical modes.

These corollaries consider each factor separately, but these factors often influence each other. For example, investigators working in fields where true effect sizes are perceived to be small may be more likely to perform large studies than investigators working in fields where true effect sizes are perceived to be large. — (Ioannidis, 2005)

In summary, the issue is that academia tends to value and promote “positive results” and discard negative results as not interesting to be published. This biases scientific literature towards certain topics, designs and knowledge.

Thus, each team may prioritize on pursuing and disseminating its most impressive “positive” results. “Negative” results may become attractive

¹⁴<https://doi.org/10.1371/journal.pmed.0020124>

for dissemination only if some other team has found a “positive” association on the same question. In that case, it may be attractive to refute a claim made in some prestigious journal. — ([Ioannidis, 2005](#))

Consider the following example. Let R be the ratio of the number of “true relationships” to “no relationships” among those tested in the field. Let u be the proportion of probed analyses that would not have been “research findings”, but nevertheless end up presented and reported as such, because of bias. Importantly, “Bias should not be confused with chance variability that causes some findings to be false by chance even though the study design, data, analysis, and presentation are perfect. Bias can entail manipulation in the analysis or reporting of findings. Selective or distorted reporting is a typical form of such bias” ([Ioannidis, 2005](#)).

Let us assume that a team of investigators performs a whole genome association study to test whether any of 100,000 gene polymorphisms are associated with susceptibility to schizophrenia. Based on what we know about the extent of heritability of the disease, it is reasonable to expect that probably around ten gene polymorphisms among those tested would be truly associated with schizophrenia, with relatively similar odds ratios around 1.3 for the ten or so polymorphisms and with a fairly similar power to identify any of them. Then $R = 10/100,000 = 10^{-4}$, and the pre-study probability for any polymorphism to be associated with schizophrenia is also $R/(R + 1) = 10^{-4}$. Let us also suppose that the study has 60% power to find an association with an odds ratio of 1.3 at $\alpha = 0.05$. Then it can be estimated that if a statistically significant association is found with the p -value barely crossing the 0.05 threshold, the post-study probability that this is true increases about 12-fold compared with the pre-study probability, but it is still only 12×10^{-4} .

Now let us suppose that the investigators manipulate their design, analyses, and reporting so as to make more relationships cross the $p = 0.05$ threshold even though this would not have been crossed with a perfectly adhered to design and analysis and with perfect comprehensive reporting of the results, strictly according to the original study plan. Such manipulation could be done, for example, with serendipitous inclusion or exclusion of certain patients or controls, post hoc subgroup analyses, investigation of genetic contrasts that were not originally specified, changes in the disease or control definitions, and various combinations of selective or distorted reporting of the results. Commercially available “data mining” packages actually are proud of their ability to yield statistically significant results through data dredging. In the presence of bias

with $u = 0.10$, the post-study probability that a research finding is true is only 4.4×10^{-4} . Furthermore, even in the absence of any bias, when ten independent research teams perform similar experiments around the world, if one of them finds a formally statistically significant association, the probability that the research finding is true is only 1.5×10^{-4} , hardly any higher than the probability we had before any of this extensive research was undertaken! — ([Ioannidis, 2005](#))

Although most of the arguments and warnings of the paper are shared by most researchers, the work of Ioannidis is not exempt from critics. Researchers Steven Goodman and Sander Greenland published an analysis of Ioannidis' approach ([Goodman and Greenland, 2007](#)). Ironically, this latter work is much less cited than the original work. Goodman and Greenland agree on the conclusions and recommendations but reject the exaggerated language of the paper regarding the falsity of most published research. Researchers Jager and Leek criticized the model as being based on justifiable but arbitrary assumptions rather than empirical data. They calculated that the false positive rate in biomedical studies was estimated to be around 14% instead of over 50% as Ioannidis asserted ([Jager and Leek, 2014](#)).

Whether the model is correct or not, Ioannidis' claims are reasonable meta-scientific research continues to increase, providing stronger knowledge and more credible scientific literature.

Note for data scientists!

The **study power** is defined as “the ability of a study to detect an effect or association if one really exists in a wider population”. Clinical studies are conducted on a subset of the patient population because it is not possible to measure a characteristic in the entire population. Whenever a statistical inference is made from a sample, it is subject to some error. Researchers attempt to reduce systematic errors through proper design so that only random errors remain. There are two types of random errors to be considered before making inferences about the studied population: **type I and type II errors**. To make a statistical inference, 2 hypotheses must be set: the **null hypothesis** (there is no difference) and **alternate hypothesis** (there is a difference). The probability of reaching a statistically significant result if in truth there is no difference or of rejecting the null hypothesis when it should have been accepted is denoted as α , or the **probability of type I error**. It is similar to the false positive result of a clinical test. The probability of not detecting a minimum clinically important difference if in truth there is a difference or of accepting the null hypothesis when it should have been rejected is denoted as β , or the **probability of type II error**. It is similar to the false negative result of a clinical test. Properly, researchers choose the size of α and β before gathering data so that their choices cannot be influenced by study results. The typical value of α is set at 0.05, and the significance level (p value) determined from the data is compared with α to decide on statistical significance. The typical value of β is set at 0.2. The power of the study, its complement, is $1 - \beta$ and is commonly reported as a percentage. — Adapted from ([Cadeddu et al., 2008](#))

Chapter 2

Scientific Inference

“Knowledge can be communicated, but not wisdom. One can find it, live it, be fortified by it, do wonders through it, but one cannot communicate and teach it.” — Hermann Hesse

2.1 Overview

One of the fundamental questions of the philosophy of science is “How can we obtain knowledge as opposed to mere belief or opinion?” ([Ladyman, 2012](#)). Humans have a natural ability to conjecture and spot relationships about the world, *jumping* from hypotheses to conclusions. The scientific attitude is to keep such *jumps* under control and use a well defined procedure to arrive to a conclusion from an hypothesis. This chapter taughts how conclusions can be reached from known facts in different ways. For this, a sound basis of logic is needed to understand more complex concepts from this course. The reader will tackle deductive and inductive arguments and link them to data science applications. With this the reader will be equipped to tackle other topics of scientific inference such as the problem of induction, the Hypothetico-deductive method and Falsificationism.

But our sun is only one of a billion-trillion stars within the observable universe. And those countless suns all obey natural laws some of which are already known to us. How did we discover that there are such laws? If we lived on a planet where nothing ever changed, there wouldn’t be much to do, there’d be nothing to figure out. There’d be no impetus for science. And if we lived in an unpredictable world where things changed in random or very complex ways, we wouldn’t be able to figure things out. And again, there’d be no such thing as science. But we live in an in

between universe where things change alright, but according to patterns, rules, or as we call them, laws of nature. If I throw a stick up in the air, it always falls down. If the sun sets in the west, it always rises again the next morning in the east. And so it's possible to figure things out. We can do science. And with it we can improve our lives. — Carl Sagan

Course objectives

A foundational grasp of logic is essential to understand scientific inference. This chapter content supports objectives 3, 4 and 6 through examples that motivate in-class debate (e.g., Hume, Popper, Hempel) with a narrative that moves from initial naïve support of theories, to a deeper reflection of their weaknesses and later emphasis on the relevant pieces that affect the practice of data science.

2.2 Types of inferences

Most scientific conclusions are uncanny at first glance and difficult to believe without more information and proper explanations about them (e.g. expansion of the universe, electromagnetism, etc.). How do scientists reach such unlikely conclusions? An inference is the act of reaching a conclusion from known facts but there are multiple types as we will see below.

2.2.1 Deduction and Induction

A good argument is one whose conclusions follow from its premises. But how do we tell if the conclusion is a consequence of its premises? Is often assumed that as long as the premises are valid, the conclusions will be valid too. This does not imply that the conclusion is also true. The premises might not be true, but if they are true, then the conclusion will also be true. However, is the truth of the premises always *necessarily sufficient* for the truth of the conclusions? Logicians distinguish between deductive and inductive inference. ([Douven, 2021](#))

Below there is an example of a deductive inference with two premises followed by a conclusion.

All Frenchmen like cheese
Loubin is a Frenchman

Therefore, Loubin likes cheese

All As are Bs

a is an A

Therefore, a is a B

We call an inference *deductive* whenever the conclusion *necessarily* follows from the premises. **The truth of the premises guarantees the truth of the conclusion.** Or in other words, what is inferred is *necessarily* true if the premises from which it is inferred are true. We call this type of inferences *explicative*.

Not all inferences are deductive. For example:

The first five eggs in the box were good.

All the eggs have the same best-before date stamped on them.

Therefore, the next egg will be good too.

In this case, the premises do not entail the conclusion. Even if the previous eggs were good, it is possible that the next egg will be rotten. In this case, is logically possible for the premises to be true and yet the conclusion false. We call this type of inferences *inductive*. Contrary to deduction, where the truth of the premises guarantees the truth of the conclusion, **inductive inferences are *ampliative* — since whose conclusions go beyond what is contained in their premises** — and their conclusions could be totally wrong even if infinitely many examples confirm them. ([Bergadano, 1991](#))

In these regards, deduction seems safer than induction. Whenever we reason deductively we can be sure that given true premises we will reach true conclusions. On the other hand, **inductive reasoning can take us from true premises to a false conclusion.** Notwithstanding, we rely on inductive reasoning every day. For instance, every day we turn on our computers and we are confident they will not explode in our faces. ([Okasha, 2016](#)) But why? Simply because we do it every morning and it has never exploded up to now.

We are sure that the sun will rise tomorrow, and if we are asked why we believe so, we will naturally answer “Because it always does”. We believe that it will rise in the future because it has risen in the past. Of course, when we are challenged to answer what *justifies* our belief we can refer to the laws of motion and nature. But will the laws of motion remain the same tomorrow? ([Russell, 1912](#))

2.2.2 Modus ponens and Modus tollens

Course Note:

The following content relates to deduction and is usually taught in high school philosophy courses as part of propositional calculus. It will help getting a better understanding of the deductive inference rules. If this is already clear to you, feel free to jump to the problem(s) of induction [2.3](#).

There are two rules of inference in deductive reasoning. Deduction constitutes top-down logic because particular conclusions are drawn from general premises. Whereas in bottom-up logic the conclusion is reached by generalizing from specific cases.

- Modus ponens: P implies Q. P is true. Therefore Q must also be true.
- Modus tollens: If P, then Q. Not Q. Therefore, not P.

The form of a **modus ponens** argument looks like a syllogism consisting of two premises and a conclusion. The first premise is a conditional if-then claim (e.g. P implies Q). The second premise is an assertion that *P* (the antecedent of the first premise) is indeed true. From these two premises, it can be concluded that *Q*, (the consequent of the first premise) must be true as well.

If P, then Q.

P.

Therefore, Q.

The next example fits the form of *modus ponens*.

If today rains, John will take the umbrella.

Today is raining.

Therefore, John will take the umbrella.

The argument is valid but it doesn't matter if the statements in the argument are actually true. An argument can be valid but nonetheless unsound if their premises are false. *Modus ponens* rule can be written as $P \rightarrow Q, P \vdash Q$. In logic, an argument is sound if it is both valid in form and its premises are true.

On the other hand, the form of a **modus tollens** argument also consists of two premises and a conclusion. The first premise is a conditional if-then claim (e.g. P implies Q). The second premise is an assertion that *Q* (the consequent of the conditional claim) is not the case. From these two premises, it can be concluded that *P*

is also not the case. *Modus tollens* rule can be written as $P \rightarrow Q, \neg Q \vdash \neg P$.

If P , then Q .

Not Q .

Therefore, not P .

Modus tollens is specially important in falsification (see 2.4.2). For instance, we take our hypothesis H to test and assume that is true. If H is true, then consequent C is true. We make an observation and see that C is false. Therefore, we conclude that H is false.

If H , then C .

C is false.

Therefore, H is false.

Other forms of arguments are apparently **similar but invalid forms**.

Affirming the consequent. This formal fallacy consists of taking a true conditional statement $P \rightarrow Q$ and invalidly inferring its converse $Q \rightarrow P$. For example, the statement “if the light is broken, the room would be dark” does not justify inferring the converse “the room is dark, therefore the lamp is broken”. This situations may arise when a consequent has more than one possible antecedent.

Denying the antecedent. This fallacy is committed by reasoning in the form: If P , then Q . Therefore, if not P , then not Q . This kind of arguments can seem valid at first glance. Consider this famous example from Alan Turing:

If each man had a definite set of rules of conduct by which he regulated his life he would be no better than a machine. But there are no such rules, so men cannot be machines. — Alan Turing

Men could still be machines that do not follow a definite set of rules.

Another trivial example of this second fallacy.

If you are a bus driver, then you have a job.

You are not a bus driver.

Therefore, you have no job.

2.3 The problem(s) of induction

Do scientists use induction? Pretty much all the time. Whenever scientists move from limited data to general conclusions scientists reason inductively. **In inductive arguments, the truth of the premises is never sufficient for the truth of the conclusion.** For instance, a newspaper may run the headline “scientists find experimental proof that transgenic maize is safe to eat”. This means scientists tested transgenic maize on a large number of people without finding any issues. Does this *strictly prove* that transgenic maize is safe? Is this prove as strong as the proof of the Pythagoras’ theorem? Going from “the transgenic maize didn’t harm any of the people on whom it was tested” to “the transgenic maize will not harm anyone” is an inductive inference, not deductive.

Writing Note:

Suppose the following inductive inference I : If the probability of observing R , given that H is true, is smaller than a significance level of 0.05, then reject H . Is important to distinguish between the two following things:

- Justification *with* an inference rule: Justifying the conclusion by pointing to the premise and the inference rule. Inference rules justify conclusions.
- Justification *of* an inference rule: What makes I a good inductive inference? Why not choosing other parameters? The choice of a particular inference rule must be justified.

2.3.1 David Hume’s Problem of Induction

In his work, David Hume begins with a critique of causal inference, arguing that we cannot make a causal inference by purely a priori means, as we can always conceive a different effect. Later, he exposes the broader epistemological problem, challenging us to justify any inference that moves from the observed to the unobserved.

We use induction to justify our statements but **how do we justify induction itself?** How would you convince someone else that induction is a good inference method? The Scottish philosopher David Hume (1711-76) argued that the use of induction cannot be rationally justified at all. In 1739, still under the shadow of the bubonic plague in Europe, David Hume publishes *A Treatise of Human Nature*, presumably without knowing that his work would not only continue to be debated more than 200 years later, but also still remarkably relevant in the technological advances of our time. In *the problem of induction* Hume argues that we cannot make a causal inference just by *a priori* means, and poses the question of how we can

conclude from the observed to the unobserved.

Hume admitted that we use induction all the time in everyday life and science but insisted that this is just a matter of brute animal habit. What does he mean by that? Bertrand Russell (1872-1970) gives us a good example on this. He argues that the inductive association is also present in animals.

“And this kind of association is not confined to men; in animals also it is very strong. A horse which has been often driven along a certain road resists the attempt to drive him in a different direction. Domestic animals expect food when they see the person who usually feeds them. We know that all these rather crude expectations of uniformity are liable to be misleading. The man who has fed the chicken every day throughout its life at last wrings its neck instead, showing that more refined views as to the uniformity of nature would have been useful to the chicken. [...] The mere fact that something has happened a certain number of times causes animals and men to expect that it will happen again. Thus our instincts certainly cause us to believe that the sun will rise to-morrow, but we may be in no better a position than the chicken which unexpectedly has its neck wrung.” — (Russell, 1912)

Hume arrived to this conclusion by noting that whenever we make inductive inferences we presuppose the *uniformity of nature*. Remember the eggs box example in § 2.2.1 ? Our reasoning depends on the assumption that objects that we have not examined yet will resemble those objects that we have already examined. Then, Hume argues that we cannot prove the truth of the uniformity assumption. Basically, from the mere act of being able to imagine a world where nature is not uniform but changes at random it follows that we cannot prove that the uniformity assumption is true. Also, if we try to argue for the uniformity assumption on empirical grounds, we end up reasoning in a circle.

Hume’s argument concerns specific inductive inferences such as All observed instances of A have been B and The next instance of A will be B. His argument proceeds as follows:

- Suppose an inductive inference I as: All observed instances of A have been B;



Portrait of David Hume by Allan Ramsay.

The next instance of A will be B.

- To justify I, it must be inferred from some premises in an argument.
- Every argument is either inductive or deductive.
- Inference I presupposes the uniformity of nature (UP), i.e., the resemblance between observed and unobserved regularities; but this cannot be proven by deduction because we can conceive situations where nature will not be uniform any more without causing a logical contradiction. Hence, a deductive argument is not possible ([Henderson, 2024](#)).
- Similarly, any inductive argument for UP would presuppose UP itself, producing a vicious circle. However, an argument cannot rely on the very principle it aims to justify. Thus, there is no inductive argument for UP.

Therefore, there is no non-circular argument for UP and thus, inference I is not justified. The reader is left to debate this unsettling conclusion and its consequences. For Hume, the conclusion is that our tendency to project past regularities into the future is not underpinned by reason.

The problem now becomes how to find a way to avoid this conclusion, despite Hume's argument ([Henderson, 2024](#)). **Hume's problem of induction is still an active area of research for philosophers.** We can frame this problem with different lenses.

- **Infinitism** embraces an epistemic view that challenges the premise that inference I presupposes UP, since, after all, the future only resembles the past in some respects, but not others. Maybe inductive inferences do not rely in a common UP, but rather, each inductive inference depends on separate empirical presuppositions ([Henderson, 2024](#)), **replacing the justification circularity problem with an infinite regress.**
- **Foundationalism** argues that the justification chain should end with foundational beliefs that do not require inferential justification.
- **Coherentism** proposes that beliefs are justified if they cohere within a web of beliefs, with a chain of reasons that is never-ending but which does not involve infinite beliefs ([Olsson, 2025](#)).

There are many different ways to respond to Hume's argument, yet none is fully convincing. Peter Strawson (1950s) used the following analogy: **justifying induction is like asking whether the law is itself legal.** This is rather odd, since the law is the standard against which the legality of other things is judged. Others, like Karl Popper (1902-1994) argued that science is not in fact based on inductive inferences at all and presented a deductivist view of science. We will study this in detail in § **2.4.2. Bayesian approaches frame induction probabilistically.** Schurz's

meta-induction justifies using induction not by logic, but by tracking which predictive strategy has been most successful over time. This is linked to on-line machine learning methods, where the learner adapts based on past success ([Ortner, 2023](#)).

Note on uniformity of nature:

Notice how ML models can be regarded as inductive machines performing inductive inferences based on previous observations. For the ML model to perform well on novel data, it is often assumed that novel data will resemble past data. But it is easy to conceive that data would stop resembling the past (see data drift and concept drift).

Hume refers to this assumption as the Principle of Uniformity of Nature: “*If reason determined us, it would proceed upon that principle, that instances, of which we have had no experience, must resemble those, of which we have had experience, and that the course of nature continues always uniformly the same.*”

And it continues: “*Our foregoing method of reasoning will easily convince us, that there can be no demonstrative arguments to prove, that those instances, of which we have had no experience, resemble those, of which we have had experience. We can at least conceive a change in the course of nature; which sufficiently proves, that such a change is not absolutely impossible. To form a clear idea of any thing, is an undeniable argument for its possibility, and is alone a refutation of any pretended demonstration against it.*”

([Hume, 1739](#)) T. 1.3.6.4

As scientists, Hume’s problem of induction may leave a huge void in our heart. An empty feeling that science is indeed fallible and the sudden realisation of the impossibility of establishing the truth or falsity of scientific laws ([Rosenberg and McIntyre, 2019](#)). But perhaps there is a way to fill such gap, and perhaps big part of the effort of science is put on filling this void with as much certainty as possible. Just because an inference rule has yield true conclusions in the past does not necessarily imply that it will do so in the future. Consequently, Hume concludes that no inductive inference rule can be justified. But, does this mean all scientific inductive inferences are not justified?

Note for data scientists!

If we visualise the data as points in a plane; every set of finite points belongs to infinite functions or curves. The problem of induction, in this case, consists in establishing criteria that allow us to say that the finite series of data confirms only one of the functions, or less dramatically but just as problematic, that one is more confirmed than the others (Díez and Moulines, 1997). (See the problem of underdetermination in §2.4.3).

Notice how the problem of induction relates to model choice given different data-generation processes. Statistical learning theory focuses on how model complexity and sample size affect the reliability of generalization from finite data, motivating practices like regularization and validation. Formal learning theory tackles which learners can converge to the truth given enough evidence. These learning-theoretic perspectives show how the problem of induction affects the practice of data science.

Ancient views on the regress argument

Pyrrhonist philosopher Sextus Empiricus (mid-late 2nd century CE) raised concerns which applied to all types of knowledge and doubted the validity of induction long before David Hume, raising the regress argument against all forms of reasoning (Wikipedia^a). This view is known as Pyrronian skepticism.

Those who claim for themselves to judge the truth are bound to possess a criterion of truth. This criterion, then, either is without a judge's approval or has been approved. But if it is without approval, whence comes it that it is trustworthy? For no matter of dispute is to be trusted without judging. And, if it has been approved, that which approves it, in turn, either has been approved or has not been approved, and so on ad infinitum. – Sextus Empiricus

^ahttps://en.wikipedia.org/wiki/Sextus_Empiricus#Philosophy

2.4 The Hypothetico-deductive Method

In the section about the **scientific method**, we learnt how scientists begin proposing (or guessing) unproven hypotheses. After an initial consideration of the problem and collection of data a conjecture or hypothesis to explain a particular phenomena is formulated. Afterwards, deduction is used to derive consequences or observable implications $\{C_i\}$ from such hypotheses H . These consequences should be relevant for H and observable directly or with the help of instruments (e.g. microscope, MRI, etc.). Next, hypotheses are put to test and either based on the results scientists

decrease or increase the confidence over the hypotheses.

1. Propose a hypothesis H .
2. Deduce observable consequences $\{C_i\}$ from H .
3. Test. Look for evidence that conflicts with the predicted consequences $\{C_i\}$ in order to disprove H .
4. If $\{C_i\}$ is false, infer that H is false, reformulate H . (See § 2.4.2)
5. If $\{C_i\}$ is true, increase confidence in H . (See § 2.4.3)

For relevant examples, check 2.9.1 and 2.9.2. Importantly, read the limitations of the HD Method at the end of this section 2.5. Notably, the HD method assumes a strict logical entailment between theory and data. However, in data science, models yield probabilistic predictions rather than deductive consequences (Mayo, 1996) (chapter VI).

2.4.1 A good hypothesis

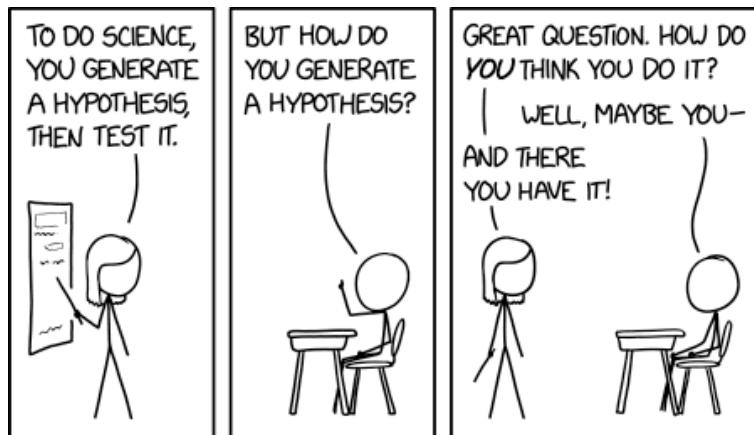


Figure 2.1: The key aspect being conveyed in this simple exchange is that one of the many good practices in science (no matter the aspect, though the specifics may change according to the precise field of study) is that one should first have an idea of what you can test and then perform the test to confirm (or rule out) your idea. Title text: “Frazzled scientists are requesting that everyone please stop generating hypotheses for a little bit while they work through the backlog”. Source: xkcd.com

Still, there are some criteria for a good hypothesis¹. Apart from criteria such as parsimony, scope, fruitfulness and conservatism, these are other criteria to recall.

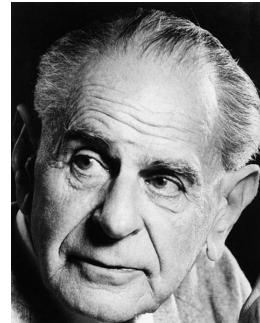
¹<https://opentext.wsu.edu/carriecuttler/chapter/developing-a-hypothesis/>

- It should be a statement that can be either true or false (e.g. “Boiling point of a liquid increases with increase in pressure”). In other words, it should be **testable and falsifiable**. We must be able to test the hypothesis using the methods of science and according to Popper’s falsifiability criterion, it must be possible to gather evidence that will reject the hypothesis if it is indeed false.
- A hypothesis must not be a tautology (i.e. claims that are necessarily true or false; e.g. “Either it will rain tomorrow or it will not rain.” or “all bachelors are unmarried”).
- Hypotheses should be informed by previous theories or observations and logical reasoning.
- Finally, the hypothesis should be positive. That is, the hypothesis should make a positive statement about the existence of a relationship or effect, rather than a statement that a relationship or effect does not exist.
- Finally, it should have some generality (e.g. “things of certain type...”) or be about some non-directly observable property of a particular.

2.4.2 Falsification

According to the Hypothetico-deductive method (H-D), a hypothesis is formulated, then relevant consequences are deduced, and finally we observe whether these consequences are false or true. Depending on these observations the hypothesis will be either falsified or confirmed.

It is important to note a key difference between confirmation and falsification. In step 4 of the **H-D method** we can infer the falsity of the hypothesis from the falsity of even a single one of the expected consequences. In contrast, in step 5 confirmation of the hypothesis is not inferred from the truth of even a large set of the consequences. Instead, we only increase our confidence on the hypothesis after finding that many consequences of the hypothesis are true. This difference is referred as the **asymmetry between confirmation and falsification**. Although a scientific theory can never be proved true by a finite amount of data, it can be proved false, or refuted by a single experiment.



Karl Popper in the 1980's.

Note for data scientists!

Importantly, we must note that in statistical hypothesis testing, data do not logically falsify a hypothesis; rather, they can provide a severe test of it ([Mayo-Wilson, 2021](#)) (see Section ??[bayesian-inference](#))). Falsification, in the strict Popperian sense, refers to universal (deterministic) hypotheses, i.e., those that can be refuted by a single counterexample. The kind of hypothesis that can be refuted by a single failed experiment is a universal hypothesis, or deterministic law-like hypothesis. Probabilistic statements like “The mean of population A equals that of population B” do not entail any single observation with logical necessity. Hence, no single data point can logically falsify such a hypothesis. We can only evaluate the hypothesis’ consistency with data.

“No amount of experimentation can ever prove me right; a single experiment can prove me wrong.” — Albert Einstein

This asymmetry forms the basics of Karl Popper’s (1902-1994) falsificationism.

- Propose falsifiable hypotheses.
- Try to falsify these hypotheses with observable evidence.
- Reject any falsified hypothesis as false.
- Never accept any hypothesis as true - consider non-falsified hypotheses as “not-rejected yet”.

One objection to this [the asymmetry between confirmation and falsification] holds that the asymmetry is an illusion, because whenever we refute a universal statement we thereby verify its negation. A universal statement “All x are y” is equivalent to “There is no non-y x.” Therefore, when we refute “All apples are green” we automatically verify “There is a non-green apple.” — ([Percival, 2015](#))

Popper is quite radical in this last step. For him, confirmation places no role at all. One can never infer the truth of hypotheses - Popper argues - from the observations regarding their implications. Not even increase the confidence in the truth of the hypothesis. Popper hoped to avoid Hume’s problem of induction by not employing induction in science. Popper thought that science was and should be deductive, and therefore that the lack of justification for inductive inferences was not as damaging for science. Below example is illustrative.

Suppose a scientist is testing the hypothesis that all pieces of metal conduct electricity. Even if every piece of metal they examine conducts electricity, this doesn’t prove that the hypothesis is true, for reasons that we’ve seen. But if the scientist finds even one piece of metal that fails to

conduct electricity, this conclusively refutes the theory. For the inference from ‘this piece of metal does not conduct electricity’ to ‘it is false that all pieces of metal conduct electricity’ is a deductive inference—the premise entails the conclusion. So if a scientist were trying to refute their theory, rather than establish its truth, their goal could be accomplished without the use of induction. — ([Okasha, 2016](#))

However, this view of the scientific process could be rather limiting with respect to the actual scientific practice. First, it does not allow to distinguish between non-falsified hypotheses. Popper argues that obtaining evidence in favour of a given theory is generally easy, and holds that such *corroboration* should count scientifically only if it is the positive result of a genuinely *risky* prediction, which might conceivably have been false.

It is logically impossible to verify a universal proposition by reference to experience (as Hume saw clearly), but a single genuine counter-instance falsifies the corresponding universal law. In a word, an exception, far from “proving” a rule, conclusively refutes it. — ([Thornton, 2021](#))

Second, in scientific practice hypotheses rarely have immediate observable consequences, they often require measurements or experiments to do so. For instance, the hypothesis “this liquid contains 3 substances” does not entail any direct observable consequence. We might use distillation or chromatography to test such hypothesis but this requires relying on **auxiliary hypothesis** (e.g. the distillation machine works properly). This consideration quite changes the **H-D method** steps. Moreover, we never test a single hypothesis alone, but only in conjunction with various auxiliary hypotheses (Duhem-Quine Thesis). One relevant example is the work of Galileo Galilei and his reports of mountains on the moon and Jupiter satellites. Philosophers such as Cesare Cremonini refused to look through the telescope, arguing that the instrument itself might have introduced artefacts, producing a visual illusion. Therefore, Duhem-Quine thesis states that in order to falsify a hypothesis we must be confident that the responsible for falsity of the consequence are not the auxiliary hypotheses but the main hypothesis.

- 1. Propose a hypothesis H .
- 2. Deduce observable consequences $\{C_i\}$ from H in conjunction with auxiliary hypotheses AH_j
- 3. Test. Look for evidence that conflicts with the predicted consequences $\{C_i\}$ in order to disprove H .
- 4. If $\{C_i\}$ is false, infer that $H \& \{AH_j\}$ is false, reformulate H .
- 5. If $\{C_i\}$ is true, increase confidence in $H \& \{AH_j\}$.

Semantic Note:

Note the difference between *falsifiable* and *falsified*.

Falsifiability is a quality of a hypothesis or a theory. Is the quality of a conjecture or hypothesis to be proven wrong. Some theories have no empirical implications. Popper claimed that astrology and Freud's psychoanalysis were not falsifiable. He argued that *falsifiability* demarcates whether a theory is scientific or not (see the demarcation problem^a (Hansson, 2021)). Similarly, some hypotheses might be more falsifiable than others because they have more empirically testable implications. For example, Newton's law of gravitation is falsifiable (e.g. it is falsified by "The brick fell upwards when released").

Falsification is the observation that an implication of a hypothesis is not true which implies (by *modus tollens*) the falsity of the hypothesis. Hypothesis can only be falsified if they are falsifiable.

Falsification uses the valid inference modus tollens: if from a statement P we logically deduce Q , but what is observed is $\neg Q$, we infer that P is false. For example, given the statement "all swans are white" and the initial condition "there is a swan here", we can deduce "the swan here is white", but if what is observed is "the swan here is not white" (say black), then "all swans are white" is false, or it was not a swan.

^a<https://plato.stanford.edu/entries/pseudo-science/>

The take-away message from falsification is that despite proposing an unrealistically restrictive practice of science, it might be a useful inference method for scientists. However, they should be aware of its limitations and for instance, bear in mind the pitfalls of *ad-hoc* modifications. (See negative weight in phlogiston theory (Grünbaum, 1976)). An *ad-hoc* hypothesis is added to a theory to save it from being falsified. A modification is considered *ad-hoc* if it reduces the falsifiability of the hypothesis in question. Again, remember that in statistical hypothesis testing, data do not logically falsify a hypothesis; rather, they can provide a severe test of it (Mayo-Wilson, 2021) (see Section ??bayesian-inference)).

2.4.3 Confirmation

Confirmation is the act of using evidence to justify increasing the confidence in the hypothesis. Confirmation is not based on deductively valid inferences. For instance, in the **H-D method** we identify some C that is an implication of H . H implies C , then if H is true we conclude (by *modus ponens*) that C is also true. Moreover, if we observe that C is false, then we conclude (by *modus tollens*) that H is false as well.

Modus ponens	Modus tollens	Induction
H, then C	H, then C	H, then C
H	not C	C
-----	-----	=====
C	not H	H

While these two inference rules are deductively valid, they do not tell us what to conclude if the implication C is true. There is no valid deductive rule that can be used for the case where H implies C and C is true. We cannot deduce anything from that.

Instead, any rule used here must amplify the information contained in the premises to infer the conclusions. Therefore we must make use of inductive inferences. Inductive inferences are fallible (inductions that fail are common e.g. predicting the weather, stock investing). But fallibility comes in degrees and this degree is affected by the kind and quality of the evidence as well as the inference rule employed. Scientists have attempted to quantify confidence, most prominently by using probabilities. For instance, if an observation O confirms hypothesis H , therefore we say that $P(H|O)$ is greater than $P(H|\neg O)$ where $P(H|O)$ means “the probability of H given O ”.

There is certain debate on this last point. Not everybody agrees that it makes sense to assign probabilities to hypotheses because they differ on the interpretation of the concept of probability. **Frequentists** interpret the probabilities as the frequencies of repeatable observable events. Therefore probabilities cannot be assigned to hypotheses since these are not events, nor observable or repeatable. Another problem is that probabilities are already used to express a property different from confidence. For instance, we may say that the probability of tails when throwing a coin is 50%. But then someone may ask us how confident we are about our claim. Even if we can also answer that second question with a probability, is clear that these two numbers express different things.

Note for data scientists!

Is important to note the relevance of frequentist and Bayesian approaches in artificial intelligence. Both frequentist and Bayesian are statistical approaches to learning from data. But there is a broad distinction between the frequentist and Bayesian. The frequentist learning only depends on the given data, while the Bayesian learning is performed by the prior belief as well as the given data ([Jun, 2016](#)).

The frequentist computes the probability of result or data D given hypothesis H is true, i.e. $P(D|H)$. In comparison, the Bayesian approach focus on the probability of hypothesis H when the result or data D occurs, i.e. $P(H|D)$ ([Orloff and Bloom, 2014](#)).

Understanding that confirmation comes in degrees may help clarify the last step of **H-D method**. Observing C to be true, increases our degree of confidence that H is true. But why is this? A naïve answer to this question is that observing C confirms H because H is compatible with C . But this seems rather weak justification. Indefinitely irrelevant implications could be inferred from a hypothesis. For instance:

I have pancreas cancer, then I have a pancreas

I have a pancreas

I have pancreas cancer

A clear deductive consequence from this example is that indeed I have a pancreas. However, observing that I do have a pancreas should not confirm the claim that I have pancreas cancer. To solve this issue we should introduce a criteria for relevance to make sure that the chosen implications are relevant to the question. This is key part of the scientific process as this often depends on the domain knowledge we have about the matter we are investigating.

An additional problem to the compatibility issue, is that very many hypotheses are compatible with any given observation. This is called the problem of underdetermination.

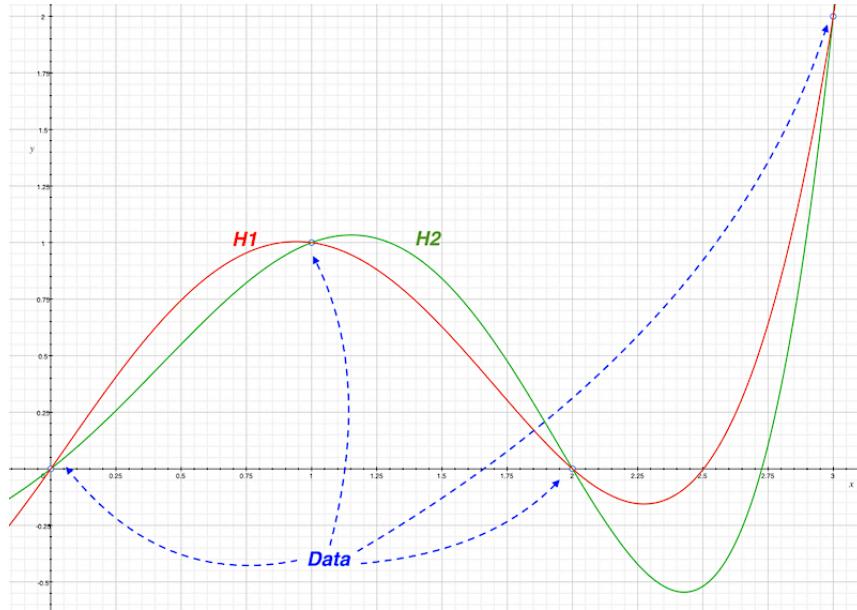


Figure 2.2: The problem of underdetermination illustrated with a chart.

Note for data scientists!

According to **anti-realists**, there will always be multiple **competing** theories about unobservable entities (e.g. atoms) which can account for the data equally well. In other words, such theories are **undetermined** by the empirical data. But then, how do scientists justify choosing one theory over another? **Realists** often reply that aforementioned scenario is only possible in the trivial sense. In fact, scientists often struggle to find even *one* theory that fits the data properly.

But why is this important for data scientists? Often you will find many ML models or solutions that fit your available data or fulfil your requirements, and yet you will have to decide which model/solution is best. If possible, validation with external data and other assessments must be conducted, but sometimes solutions are also chosen based on *non-epistemic* values, such as making society more just or making money.

Is also important to notice how the problem of underdetermination relates to the popular **No Free Lunch Theorem** which is very relevant in the Machine Learning community ([Dotan, 2020](#)). For more on the NFL read ([Domingos, 2015](#)).

2.5 Beyond the Hypothetico-Deductive Method: Bayesian Inference and the Logic of Uncertainty

While the HD method offers a simple and useful framework to think about how science proceeds, from hypothesis formation to testing via deduction, it has important limitations. These become especially clear in data-driven disciplines such as machine learning, statistics, and empirical sciences that rely on probabilistic reasoning. Here, we examine alternative approaches that better accommodate uncertainty, quantify degrees of belief, and reflect how inference is conducted in real-world contexts.

2.5.1 Limitations of the HD method in probabilistic contexts

In the HD method, hypotheses are tested by deducing observable consequences and then checking whether those consequences occur. If any of the deduced implications are false, the hypothesis is rejected (falsified). If they are all true, the hypothesis is not proven true, but may be confirmed. However, **we must distinguish between logical incompatibility (in the rejection of hypothesis) and statistical evidence in hypothesis testing**. In the latter, data is logically compatible with both the null and the alternative hypotheses, i.e., no matter what data is observed, both hypotheses can accommodate, though with varying likelihood extents ([Beall et al., 2024](#)).

Importantly, most scientific theories, especially those in statistical sciences, are not structured as strict universal generalizations like “all swans are white”. Instead, they are probabilistic in nature. For example:

- “If a drug is effective, then 80% of patients will recover.”
- “If a coin is fair, then heads will appear approximately 50% of the time.”

Here, the predictions are not strictly entailed by the hypothesis. They are probabilistic expectations. This weakens the applicability of modus tollens: if only 60% of patients recover, is the hypothesis falsified? Not necessarily. In probabilistic contexts, observing an outcome that deviates from the expected does not conclusively falsify the hypothesis, because such deviations are themselves expected to occur with some probability. **The data may be compatible with both the hypothesis and the null hypothesis**, making it difficult to draw a sharp deductive conclusion. Instead of strict falsification, we assess how well the data supports one hypothesis over another, often in terms of likelihood or posterior probability.

This complexity is further compounded when we consider that hypotheses are rarely

tested alone but often depend on **auxiliary hypotheses** (Grünbaum, 1976), which in practice motivates robustness checks and sensitivity analyses to see how conclusions depend on background assumptions.

Moreover, empirical data is noisy. Measurement error, incomplete information, and stochastic processes mean that consequences of a hypothesis may deviate from what is expected even if the hypothesis is “mostly” correct. Under these circumstances, we need a more flexible form of inference, i.e., probabilistic inference, that allows us to reason under uncertainty.

2.5.2 Bayesian Epistemology and Conditionalization

Bayesian inference offers a model of rational belief revision. In contrast to the HD method, which evaluates hypotheses in an all-or-nothing fashion, the Bayesian approach updates degrees of belief (probabilities) in light of new evidence. Suppose we have a hypothesis H and we observe some data D . The central rule of Bayesian updating is Bayes’ theorem:

$$\text{Bayes' theorem: } P(H | D) = \frac{P(D|H) \cdot P(H)}{P(D)}$$

This formula tells us how to update our belief in H after observing D . Key terms:

- $P(H)$: the prior probability of the hypothesis before seeing the data.
- $P(D | H)$: the likelihood of the data assuming the hypothesis is true.
- $P(H | D)$: the posterior probability of the hypothesis after updating with the data.

This rule allows beliefs to change incrementally as new evidence becomes available, reflecting both prior knowledge and the strength of the data. It avoids the binary logic of the HD method and instead provides a continuous scale of confidence. Where the HD method discards a hypothesis after a single negative test, Bayesian inference merely lowers its probability, unless the evidence is truly decisive (Howson and Urbach, 2006).

2.5.3 Bayesian vs Frequentist Views

There are two dominant interpretations of probability in scientific inference:

- Frequentist: Probabilities refer to long-run relative frequencies of observable events. Hypotheses are fixed and not assigned probabilities. Evidence is used to reject or fail to reject a hypothesis.
- Bayesian: Probabilities reflect subjective or rational degrees of belief. Hypotheses themselves are uncertain, and can be compared based on their posterior

2.5. BEYOND THE HYPOTHETICO-DEDUCTIVE METHOD: BAYESIAN INFERENCE AND T

probabilities.

Both perspectives are widely used in science and data analysis. For example:

Topic	Frequentist Perspective	Bayesian Perspective
Interpretation of Probability	Long-run frequency of events	Degree of belief (subjective or epistemic)
Hypotheses	Fixed and not assigned probabilities	Treated as uncertain; assigned prior and updated posterior
Learning from Data	Estimation via Maximum Likelihood (MLE), confidence intervals	Posterior inference via Bayes' Theorem
Inference	Hypothesis testing using p-values	Probability of hypothesis given data, i.e., $P(H D)$
Updating	No explicit mechanism; fixed once data is observed	Prior updated with data to form posterior
ML applications	Parameter estimation, frequentist confidence	Bayesian networks, uncertainty quantification

Bayesian inference is especially relevant in machine learning, where models are constantly updated as new data is acquired. For example, recommender systems revise their predictions as user preferences evolve, and autonomous vehicles update environmental models in real time.

Note for data scientists!

In Bayesian Machine Learning, priors allow incorporation of domain knowledge into model training. For instance, prior beliefs about plausible parameter ranges can prevent overfitting in small-data regimes. In deep learning, Bayesian neural networks use probability distributions over weights, allowing uncertainty estimation in predictions. These methods are crucial in high-stakes applications such as medical diagnosis, where not knowing the confidence of a prediction can be dangerous.

2.5.4 Bayesian Inference and the Problem of Induction

Bayesian epistemology offers one response to Hume's problem of induction. It does not deny the fallibility of inductive inference, but rather makes it explicit. In Bayesian terms, learning from experience is about adjusting belief in hypotheses as new evidence accumulates.

However, this still does not justify induction in the absolute sense. It formalizes how we might reasonably behave given certain starting beliefs and observations, but it does not explain why those beliefs are justified in the first place. As some philosophers argue, Bayesianism presupposes that some inductive inferences are valid; it does not solve Hume’s problem but reframes it.

2.5.5 Underdetermination and Bayesian Model Selection

As seen with Confirmation in section 2.4.3, many hypotheses can be consistent with the same data. This is known as the problem of underdetermination. Bayesian reasoning provides one approach to this problem by allowing us to compare models in terms of their posterior probabilities. Models with higher likelihood given the data, and more plausible priors, are preferred.

Still, choosing priors is an open issue. Different priors can lead to different posterior conclusions, raising philosophical concerns about **subjectivity**. In practice, however, prior sensitivity can be tested and robustified through sensitivity analysis and hierarchical modeling.

2.6 Other types of inference

So far we have discussed deduction and induction, but we can find additional scientific reasoning methods. In practice, scientists and data scientists often rely on other forms of inference that go beyond strict deduction and enumerative induction. These include abduction (or inference to the best explanation), and logics that depart from the classical, monotonic model of reasoning, such as non-monotonic and defeasible inference.

2.6.1 Abduction and Inference to the Best Explanation (IBE)

Abduction was first introduced by Charles Peirce (1839-1914) to describe the process of inferring the most plausible hypothesis from a set of observations. **Abduction, then, is an explanatory form of reasoning:** given evidence E and a set of possible hypotheses H_1, H_2, \dots, H_n , we infer that H_i is true because it best explains E . This is different from enumerative induction, which generalises from a limited number of instances to a rule.

For example, if we find an empty plate with crumbs and recall hearing scratching sounds in the larder, we infer that the cheese was eaten by a mouse. Similarly, a

doctor observing fever, cough, and shortness of breath will *abduce* (or infer) that the patient might have an infection. Both are plausible explanations, not certain deductions.

Philosopher Peter Lipton (1954-2007) later framed abduction as Inference to the Best Explanation (IBE), stressing that what makes an explanation “best” depends on factors like simplicity, coherence with background knowledge, and explanatory power. While some philosophers distinguish abduction as a psychological process from IBE as its normative version, in practice we treat them as one and the same, both describe how we move from observed data to explanatory hypotheses. For the sake of distinguishing them, abduction, in this sense, is the creative step that proposes a possible explanation, it tells us *how* hypotheses arise. IBE, is the *evaluative* step that selects which of the available hypotheses is most plausible. Of course, these two steps, often happen together.

2.6.2 Monotonic and Non-monotonic Logic

Traditional deductive logic is monotonic: once a conclusion is validly derived from a set of premises, adding more premises can never invalidate it. If $A \vdash p$, then $A \cup B \vdash p$ still holds. This is the logic of mathematics and formal proofs. However, **most real-world reasoning, and especially scientific reasoning, is non-monotonic**. New evidence can force us to withdraw a previous conclusion. For instance, a physician might initially diagnose pneumonia from a chest X-ray, but revise the diagnosis after receiving negative test results or learning that the symptoms were caused by tuberculosis instead.

In its epistemic sense, **monotonicity expresses the fact that adding more premises to an argument allows you to derive all the same conclusions as you could with fewer** (Strasser and Antonelli, 2019). Specifically, under monotonic reasoning, if a conclusion p follows from a set of premises A , (denoted as $A \vdash p$), adding another set of premises B doesn’t alter the conclusion (i.e. $A \wedge B \vdash p$ also holds). Therefore, reasoning is non-monotonic² when a conclusion supported by a set of premises can be retracted in the light of new information. Or in other words, we can infer certain conclusions from a subset of a set S of premises which cannot be inferred from S as a whole. Medical diagnosis fits very well under such definition.

In machine learning terms, this corresponds to models that must update their predictions when new, possibly conflicting, data arrive. This very notion relates to how knowledge is indeed provisional, which in turn justifies practices like MLOps techniques to monitor the performance of deployed machine learning models and act on

²<https://plato.stanford.edu/entries/logic-nonmonotonic/>

events of **concept drift** (when the relationship between input features and the target variable changes over time, e.g. when what “counts” as a fraudulent transaction evolves) and **data drift** (when the distribution of input data itself shifts, e.g. new patterns or populations differing from those seen during training).

2.6.3 Defeasible Reasoning

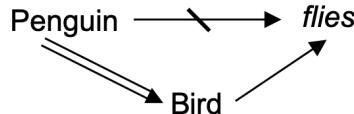


Figure 2.3: Double arrows indicate non-defeasible inferences (hard fact), single arrows depict defeasible inferences, and strikethrough arrows denote a negation. It can be read as: Penguins are birds (no exceptions); Birds usually fly; and Penguins usually don’t fly.

Defeasible reasoning deals with tentative relationships between premises and conclusions, which can be *defeated* by additional information, allowing for the retraction of inferences. For instance, while we may infer that Tweety flies based on the information that Tweety is a bird and the domain knowledge that birds generally fly, we can retract this inference when we learn that Tweety is a penguin. Tweety is indeed a bird but it cannot fly. Defeasible reasoning is not exempt from limitations, requiring from causal information to properly derive conclusions under certain scenarios.

Consider, for example, this problem of Judea Pearl: if the sprinkler is on, then normally the sidewalk is wet, and, if the sidewalk is wet, then normally it is raining. However, we should not infer that it is raining from the fact that the sprinkler is on (Pearl, 2014). Conflicts may arise between hard facts and defeasible conclusions. For instance, both arguments in Figure 2.3 $Penguin \Rightarrow Bird \rightarrow flies$ and $Penguin \rightarrow \neg flies$ finish with a defeasible inference. The transitivity rule $(a \rightarrow b, b \rightarrow c) \Rightarrow a \rightarrow c$ cannot be applied to the first argument. In this case, according to their specificity we can give priority to the argument with more a specific antecedent but is not always as trivial, and complex conflicts can remain unresolved.

Reasoning is defeasible when the corresponding argument is rationally compelling but not deductively valid. The truth of the premises of a good defeasible argument provide support for the conclusion, even though it is possible for the premises to be true and the conclusion false. In other words, the relationship of support between premises and conclusion is a tentative one, potentially defeated by additional information. — (Koons, 2021).

Defeasible reasoning is a particular kind of non-demonstrative reasoning, where the reasoning does not produce a full, complete, or final demonstration of a claim, i.e., where fallibility and corrigibility of a conclusion are acknowledged. In other words, defeasible reasoning produces a contingent statement or claim. Defeasible reasoning is also a kind of ampliative reasoning because its conclusions reach beyond the pure meanings of the premises. Defeasible reasoning finds its fullest expression in jurisprudence, ethics and moral philosophy, epistemology, pragmatics and conversational conventions in linguistics, constructivist decision theories, and in knowledge representation and planning in artificial intelligence.

— Wikipedia

These forms of reasoning are particularly relevant to modern AI and data science. Medical diagnosis, for instance, cannot rely on purely monotonic inference because diseases may overlap, evolve over time, and require multimodal evidence (radiology, lab tests, symptoms). As in the examples of COVID-19 and lung disease, a binary or multiclass classifier that assumes mutually exclusive categories cannot capture such complexity. Reasoning must be able to handle exceptions, conflicting information, and revision of beliefs.

2.7 Summary of the different inference methods

Each of the inference methods plays a distinctive role in both science and data science. **Deduction** ensures internal consistency and allows us to derive testable implications from theories or models. **Induction** supports the generalisation of empirical patterns into broader laws or statistical regularities. **Abduction** drives the generation of explanatory hypotheses when the causes behind observed data are unknown or not fully understood. Finally, **non-monotonic and defeasible reasoning** capture the inherently provisional and revisable nature of scientific knowledge, where conclusions must adapt to new evidence or changing contexts. Together, these inferential modes reflect how science advances: not by accumulating certain truths, but by continuously updating, testing, and improving our best explanations of the world.

Type	Concept	Question	Role
Deduction	From general laws to specific consequences	If this theory is true, what should we observe?	Testing and deriving predictions

Type	Concept	Question	Role
Induction	From specific observations to general regularities	Given what we have seen, what general pattern might hold?	Testing and generalising
Abduction	From surprising facts to plausible explanations	What hypothesis would best explain these facts?	Hypothesis and model generation

2.8 Explanation

At this point, many of you probably have already related data science to two of the goals of science: explanation and prediction. But how do they relate to one another? and what is a scientific explanation? Either to satisfy our natural curiosity or for a further purpose, science has always attempted to understand how the world works. The German philosopher Carl Hempel attempted to answer this question in the 1950s with what is known as the *covering law* model of explanation. He stated that a scientific explanation is an answer given in response to *explanation-seeking why-questions* (e.g. why salt dissolves in water). **In this section, we will first present Hempel's model, then discuss its main problems, and finally explore modern causal and mechanistic views of explanation.**

According to Hempel, explanations are structured like an argument, i.e. a set of premises followed by a conclusion. Therefore, the conclusion of such an argument states that certain phenomenon occurs, e.g. “salt dissolves in water”. On the other hand, the premises indicate why the conclusion is true. Then, the challenge lays in the relationship that should follow between such premises and the conclusion. For Hempel, the premises should all be true and entail the conclusion, i.e. the argument should be deductive and *sound*. Additionally, the premises should contain at least one general law (e.g. all metals conduct electricity), also known as *laws of nature*. The name of the model comes from that fact that the phenomenon to be explained is “covered” by some general law.



Carl Hempel (1905 - 1997).

1. The *explanandum* must be a valid deductive argument.
2. The *explanans* must contain at least one general law actually needed in the

deduction.

3. The *explanans* must be empirically testable.
4. The sentences in the *explanans* must be true.

For instance, Newton explained the elliptical orbits of planets alluding a general rule (his law of universal gravitation) together with some minor assumptions. This example fits Hempel's model very well, but not all scientific explanations do.

General Law (explanans)

Particular Facts (explanans)

Phenomenon to be explained (explanandum)

An interesting consequence of this model lies in the relationship between explanation and prediction. **Hempel argued that explanation and prediction are two sides of the same coin.** Whenever a phenomenon is explained with the help of a covering law, the laws and the particular facts we use could have allowed us to predict the occurrence of the phenomenon. Hempel expressed this by saying that every scientific explanation is potentially a prediction. Hempel argued that the opposite is also true: every prediction is potentially an explanation. **For Hempel, explanation and prediction are structurally symmetric.** For instance, the same information we could use to predict an animal species extinction before it happened will serve to explain that very same fact after it has happened.

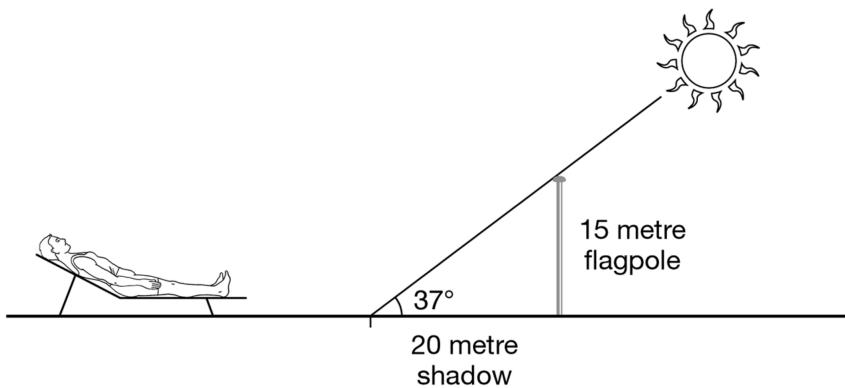


Figure 2.4: A 15-metre flagpole casts a shadow of 20 metres when the sun is 37° overhead. Figure from ([Okasha, 2016](#)).

2.8.1 Limitations of Hempel's model

Hempel's model might be too liberal, as it faces a number of odd counterexamples. For example, consider Figure 2.4. In order to explain why the shadow is 20 metres

long. Indeed, this is a *explanation-seeking why-question* and a possible answer could be the following: The light rays from the sun hit the flagpole (15 metres high), the sun's elevation angle is 37° and light travels in straight lines. The trigonometric calculation $\tan(37^\circ) = 15/20$ demonstrates that the flagpole will cast a shadow of 20 metres long. This example can be found in both ([Rosenberg and McIntyre, 2019](#)) and ([Okasha, 2016](#)).

General Law (Light travels in straight lines)

General Law (Trigonometric laws)

Particular Fact (Flagpole is 15 metres high)

Particular Fact (Sun's angle of elevation is 37°)

Phenomenon to be explained (Shadow is 20 metres long)

This explanation can be structured according to Hempel's schema and indeed fits Hempel's covering law model. However, when we swap the *explanandum* with the particular fact that the flagpole is 15 metres high, a problem arises. The explanation still complies with the covering law pattern, but it would be rather odd to regard it as an explanation of why the flagpole is 15 metres high. In this case, we know that the height of the flagpole is not conditioned by the sun's angle of elevation but rather because it was manufactured with such height. We can *calculate* or *predict* its height but this height will not change upon the other variables, so they do not *explain* the flagpole height.

General Law (Light travels in straight lines)

General Law (Trigonometric laws)

Particular Fact (Shadow is 20 metres long)

Particular Fact (Sun's angle of elevation is 37°)

Phenomenon to be explained (Flagpole is 15 metres high)

The moral of this example is that the concept of explanation showcases an important **asymmetry**. The length of the shadow can be explained by the height of the flagpole, given aforementioned general laws. But this does not happen in the other direction. In general, if x explains y , then it will not be true that y explains x given the same laws and facts. **Explanation is then an asymmetric relation and Hempel's covering law model does not respect such asymmetry.** Information that allow us to predict a fact before we know it does not serve to explain that very same fact after we know it, which **contradicts Hempel's thesis**.

The general conclusion is that a good explanation of a phenomenon should contain information that is relevant to the phenomenon's occurrence. Hempel's model left

unanswered why some explanations feel more informative than others. Philosophers and scientists increasingly turned toward causation as the missing ingredient: what truly makes an explanation explanatory.

2.8.2 Explanation and causality

There are alternatives to the covering law model that help us understand what counts as scientific explanation. For many, explaining a phenomenon is simply to say what caused it. Obviously, causality is also an asymmetric relation. If a faulty appliance caused a fire, then it is clear that the fire did not cause the appliance's failure. The asymmetry of explanation derives from the asymmetry of causality.

However, the criticism against Hempel covering law is a bit unfair as he was an empiricist. Empiricists are sceptical about the concept of causality and argue that all our knowledge comes from experience. David Hume argued that is impossible to experience causal relations, and that causality is just what we humans project to understand the world.

There are however some examples where explanation and causality do not align. For example, to say that an object's temperature is the average kinetic energy of its molecules is to explain what temperature *is*, but this does not yield the cause of such temperature.

The law $PV = nRT$ explains the temperature of a gas at equilibrium by appeal to its pressure and the volume it takes up. But volume and pressure cannot be *causes* of temperature since all of them — the temperature, the volume, and the pressure — vary, in the way the law describes, instantaneously. The changes in volume at one time do not cause changes in temperature at a later time; instead, the change in temperature occurs during exactly the same interval that pressure is changing (Rosenberg and McIntyre, 2019). The nature of causation is still an open debate, but most philosophers have agreed that causes somehow necessitate their effects and that mere regularity cannot express this necessity.

Another example, this time from (Pearl and Mackenzie, 2018), show us the gap between causal vocabulary and ordinary scientific vocabulary. Consider the problem of expressing the following causal relationship: The barometer reading B tracks the atmospheric pressure P . We can write down the relationship as $B = kP$, where k is a constant of proportionality. Thanks to algebra we can rewrite the equation in multiple ways such as: $P = B/k$, $k = B/P$, or $B - kP = 0$. They mean the same and given two of the variables we can calculate the third. But in these equations there is no account of directionality. We cannot express that *is* the pressure which *causes* the barometer to change and not the other way around. Similarly, we cannot

express the fact that the singing of the rooster *does not cause* the sun to rise. Why have scientists not captured such facts in formulas as is done in other areas like mechanics, geometry or optics? For a better understanding of causation and its role in data science I recommend you the Book of Why ([Pearl and Mackenzie, 2018](#)).

Data can tell you that the people who took a medicine recovered faster than those who did not take it, but they can't tell you why. — Book of Why

2.8.3 Summary and modern views

We can summarise Hempel's view as follows:

- Hempel proposed that scientific explanations answer why-questions using a logical argument.
- The structure of an explanation resembles a deductive argument, consisting of explanans (premises including at least one general law and specific facts) and an explanandum (the phenomenon to be explained).
- For Hempel, a valid scientific explanation must be a deductive and sound argument; contain at least one general law of nature; have empirically testable premises; use true statement.
- Hempel held that explanation and prediction are structurally symmetric.

In short, Hempel's model helped formalize scientific explanation, but it faces key limitations:

- It considers “explanations” many arguments that don’t really explain (e.g., the flagpole-shadow example).
- The model treats explanation and prediction as equivalent, but in reality, they are asymmetric (knowing the shadow length doesn’t explain the pole’s height).
- Ignores causality: The model focuses on logical form, not on causal relationships.
- Relevance problem: Explanations must contain information relevant to the phenomenon (explanandum); Hempel’s model doesn’t ensure this.
- No directionality: It doesn’t distinguish between causes and effects; causality is one-directional, but deduction is symmetric.
- Empiricist skepticism of causality: Hempel, influenced by Hume, downplayed causal reasoning, yet most explanations are causal in nature.

The covering-law model was foundational, but incomplete: it motivated the causal-mechanistic turn in philosophy of science. A good explanation should provide relevant information about why the event occurred, not just a correlation or logical

entailment.

Because the covering-law model treated explanation as mere deduction from general laws, it failed to capture how scientists actually explain phenomena through causal or mechanistic reasoning. Contemporary models emphasize that explaining is not just showing that an event fits under a law, but showing why it occurs, by revealing the underlying causal mechanisms, counterfactual dependencies, and contextual relevance that make the phenomenon intelligible.

- Causal-Mechanistic Turn:
 - Explanation moved from logical form to causal relevance and counterfactuals.
 - Focus on causal relevance and counterfactual dependence.
 - To explain = to show how changing X would change Y.
 - Causation tied to manipulability and intervention (e.g., Pearl's causal diagrams).
- Mechanistic Models (esp. in biology & medicine):
 - “To explain is to describe a mechanism”.
 - Explanation shows how entities and activities are organised to produce a phenomenon.
 - Mechanisms replace general laws in complex domains.
 - In complex sciences, general laws are rare but mechanisms are abundant.
- Unificationist and Pragmatic Views:
 - Unification: explanation means showing how diverse phenomena fit under a few general patterns (e.g., Newton's laws unify celestial motion).
 - Pragmatic: explanation is context-dependent, addressing why A happened rather than B, for a given audience or contrast.
- Pluralism in Explanation:
 - There is no single canonical model today.
 - Different sciences adopt different explanatory ideals:
 - * Physics: laws and mathematical deduction.
 - * Biology/Medicine: mechanisms and probabilities.
- Modern consensus
 - Explanation implies revealing the causal, probabilistic, or structural dependencies that make a phenomenon understandable within scientific practice.

In sum, while Hempel's covering-law model provided the first systematic account of scientific explanation, its neglect of causality, asymmetry, and context opened the path for contemporary approaches that view explanation as revealing the mechanisms, dependencies, and conditions that make phenomena intelligible rather than

merely deducible.

In data science, explanatory and predictive goals often overlap but differ in purpose: prediction seeks accuracy about future observations, whereas explanation seeks understanding of why they occur. Both, however, depend on our ability to uncover the underlying structures—logical, causal, or mechanistic—that connect data to the world

Note for data scientists!

In the biological and social sciences, instead of strict laws one finds statements of probabilities, or statistical regularities, and explanations that appeal to them. In the medical contexts, explanations often employ relations that are reported in statistical form in order to express causal relationships. For instance, it is accepted that smoking causes lung cancer because it is associated with a big increase in the probability of contracting lung cancer. Nonetheless, we know that **statistical correlation does not warrant causal connection**. There are some problems with the statement that smoking causes cancer. Some smokers never contract cancer, while some lung cancer victims never smoked. The latter issue is easy, smoking is not the only cause of lung cancer. However, the first problem is harder to tackle.

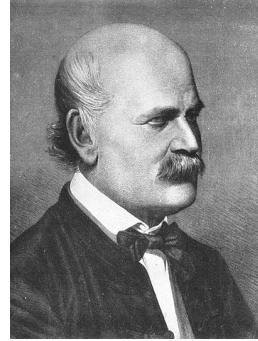
Smoking can be said to cause cancer if and only if, among all the different background conditions we know about (heredity, diet, exercise, air pollution, etc.), there is no correlation between smoking and a lower than average incidence of lung cancer, and in one or more of these background conditions, smoking is correlated with a higher incidence in lung cancer rates. — ([Rosenberg and McIntyre, 2019](#))

2.9 Examples

2.9.1 The problem is in your hands!

Ignaz Semmelweis, a Hungarian physician, was a member of the First Maternity Division at the Vienna General Hospital from 1844 to 1848. Semmelweis was distressed to find a big proportion of the women who delivered their babies contracted a serious and often fatal illness known as childbed fever. In 1844, 8.2% of mothers died from the disease, 6.8% in 1845 and 11.4% in 1846. However, in the adjacent Second Maternity Division which had as many women as the first, the death toll was much lower (2.3%, 2% and 2.7% respectively).

From this moment on, various explanations were considered, subjected to test and then rejected.



Dr. Ignaz Semmelweis in 1860.

The first explanation attributed the issue to “epidemic influences” described as “atmospheric-cosmic-telluric changes” spreading over districts and causing childbed fever. This hypothesis did not explain why the first division had more cases than the second. Neither explained the lack of cases in the city of Vienna. Epidemics such as cholera are not so selective. Finally, Semmelweis notes that women who had to give birth in the street on their way to the hospital had a lower death rate than the average for the first division.

Year	First Division Mortality (%)	Second Division Mortality (%)
1844	8.20	2.30
1845	6.80	2.00
1846	11.40	2.70

On a different view, overcrowding of the first division was proposed as a cause but Semmelweis pointed out that the second division was much crowded. Moreover, there were no differences regarding diet or general care of the patients.

In 1846, a commission was appointed to investigate the issue, which attributed the prevalence to injuries in the first division resulting from rough examination by medical students. Semmelweis refuted this view since: a) the injuries of birth itself are more extensive than those from the examination. b) midwives’ examinations from the second division were similar. c) as a consequence of the commission the number of students was halved and the examinations were reduced to a minimum. The mortality increased.

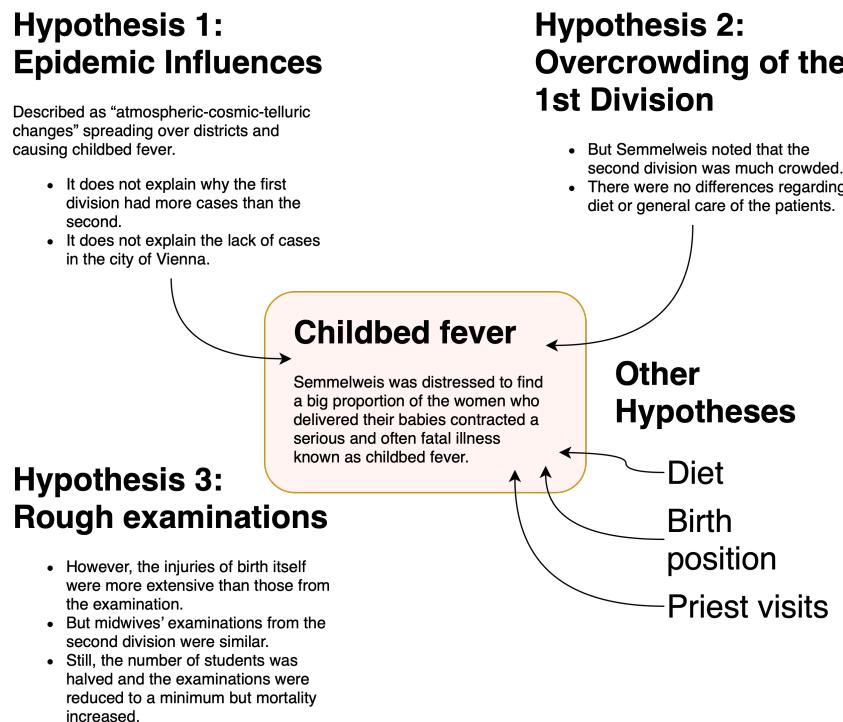


Figure 2.5: Summary of the hypotheses considered to explain the deaths in the first division.

After considering peculiar conjectures (e.g. delivery position, priest visits), an accident gave Semmelweis the decisive clue. In 1847, a colleague of his received a puncture wound in the finger, from the scalpel of a student while performing an autopsy. His colleague died after an illness with similar symptoms to those observed in the victims of childbed fever. Note, that the role of micro-organisms had not yet been recognized at the time. Semmelweis ordered all medical students to wash their hands with a chlorinated lime solution before making examinations, especially after performing dissections in the autopsy room.

Mortality fell to 1.27% in the First Division compared to 1.33% in the second. In further support of his hypothesis, Semmelweis notes that midwives from the Second Division did not dissect cadavers. This also explained the “street births” low mortality since women were rarely examined as they already gave birth. Semmelweis concluded that the cause was infection by cadaveric material and putrid matter.

At the time, Semmelweis’s findings lacked a scientific explanation. This understanding only emerged in the 1860s and 1870s, thanks to the work of Louis Pasteur, Joseph Lister, and others who advanced the germ theory of disease. Semmelweis’s idea that a single cause existed, centered solely on cleanliness, was considered radical at that time and was mostly ignored, rejected, or mocked. He was politically dismissed from the hospital and faced harassment from the medical community in Vienna, ultimately leading him to relocate to Budapest. People close to him, including his wife, believed he was losing his mind, and in 1865, he was admitted to a lunatic asylum, where he died just 14 days later. Semmelweis’s practices gained widespread acceptance only years after his death, when Louis Pasteur further refined the germ theory of disease, providing a theoretical basis for Semmelweis’s observations.

Note for data scientists!

Notice how this story relates to data science in many aspects. First, Semmelweis considered multiple factors (diet, division sizes, etc.) and collected information about them. Similarly, in data science, the quality and variety of data can significantly help in the data analysis. Semmelweis formulated a hypothesis based on his observations and sought to test it through intervention (handwashing). In data science, hypotheses are also formulated and contrasted using data to compare pre-intervention and post-intervention observations. The medical community initially rejected Semmelweis's hypothesis despite evidence. This reflects a common bias in scientific inquiry where existing beliefs can overshadow new findings. In data science, **confirmation bias** can lead researchers to favour results that confirm pre-existing hypotheses rather than exploring contrary evidence. Semmelweis's approach demonstrates the iterative process of refining hypotheses based on new data and observations. Importantly, data showed that Semmelweis was probably right, but only after the advent of the germ theory an explanation could be provided. This shows the important role of domain knowledge in the interpretation of data-driven results.

2.9.1.1 How a hypothesis is tested

Some conjectures (e.g. differences in diet, crowding or care) were trivial to test as their assumptions conflict with readily observable facts (e.g., there were no differences in diet across divisions). Others were not as straightforward and required certain interventions. For example, changing the routine of the priest or the birth position. If the hypothesis H is true, then certain observable events I should occur (e.g. drop in mortality) under specified circumstances (e.g. lateral delivery position). Semmelweis experiment showed the test implication to be false (he changed the routine of the priest or the delivery position and did not find different results in mortality), rejecting the hypothesis in consequence.

If H is true, then so is I .

But (as the evidence shows) I is not true.

H is not true.

This is a good example of *modus tollens* (see 2.2.2). However, let us consider now the case where observation or experiment confirms the test implication I . From the hypothesis that childbed fever is blood poisoning produced by cadaveric matter, Semmelweis infers that antiseptic measures will reduce mortality rates. Now, the experiment shows the test implication to be true. But this favourable outcome does not prove the hypothesis true.

If H is true, then so is I .
 (as the evidence shows) I is true.

H is true.

This reasoning is deductively invalid and referred to as the *fallacy of affirming the consequent* (see 2.2.2). The conclusion may be false even if its premises are true. Thus, even if many implications of a hypothesis have been confirmed by tests, the hypothesis may still be false.

If H is true, then so are I_1 , I_2 , ...
 (as the evidence shows) I_1 , I_2 , ... are all true.

H is true.

Above's argument still commits the fallacy. Note that although the many tests do not provide conclusive proof for a hypothesis, they provide at least some support or confirmation for it.

Chapter 4 of ([Hempel, 1966](#)) continues on this.

In the absence of unfavorable evidence, the confirmation of a hypothesis will normally be regarded as increasing with the number of favorable test findings. [...] the increase in confirmation effected by one new favorable instance will generally become smaller as the number of previously established favorable instances grows. If thousands of confirmatory cases are already available, the addition of one more favorable finding will raise the confirmation but little.

Note for data scientists!

Notice how ML models can also be affected by the previous statement. Many researchers blindly rely on the dogma *the more data, the merrier* but is not just the amount of data that matters but also its variety. The greater the variety, the stronger the resulting support for the trained model.

2.9.2 Wason selection task



Figure 2.6: Wason selection task or four-card problem.

Consider the following hypothetico-deductive reasoning problem created by Peter Cathcart Wason employing the logical rule of implication:

You are shown a set of four cards placed on a table, each of which has a number on one side and a colored patch on the other side. The visible faces of the cards show 3, 8, red and brown. Which card(s) must you turn over in order to test the truth of the proposition that if a card shows an even number on one face, then its opposite face is red?

Hypothesis H: “If a card shows an even number on one face, then its opposite face is red”

Test whether H is false. Which consequences of H do you need to consider - i.e. which cards do you need to turn over? Under what conditions would this statement be false?

These are the possible situations:

- If the 3 card is red (or brown), that doesn't violate the rule. The rule makes no claims about odd numbers. (Denying the antecedent)
- If the 8 card is not red, it violates the rule. (Modus ponens)
- If the red card is odd (or even), that doesn't violate the rule. The red color is not exclusive to even numbers. (Affirming the consequent)
- If the brown card is even, it violates the rule. (Modus tollens)

Table 2.4: Truth table for $p \rightarrow q$. (*) In instances of *modus ponens* we assume as premises that $p \rightarrow q$ is true and p is true. Only one line of the truth table — the first — satisfies these two conditions (p and $p \rightarrow q$). On this line, q is also true. Therefore, whenever $p \rightarrow q$ is true and p is true, q must also be true. (**) In instances of *modus tollens* we assume as premises that $p \rightarrow q$ is true and q is false. There is only one line of the truth table — the fourth line — which satisfies these two conditions. In this line, p is false. Therefore, in every instance in which $p \rightarrow q$ is true and q is false, p must also be false.

p	q	$p \rightarrow q$
T	T	*T
T	F	F
F	T	T
F	F	**T

There are two ways to face the problem and reach the solution. First, we can choose the cards based on *modus ponens* and *modus tollens* as follows: From *modus ponens*

we need to check the cards that are even. If even cards are not red, then the claim is false.

If even, then red. (claim)

even (obs)

Therefore, red. (conclusion)

From *modus tollens* we need to check the cards that are not red i.e. brown. If brown cards are even, then the claim is false.

If even, then red. (claim)

not red (obs)

Therefore, not even (conclusion)

Another approach is to take the truth table of $p \rightarrow q$ and take the case where $p \rightarrow q$ is false - second line - i.e. when p is true and q is false or for our case, when a card is even and its back not red (so brown). From this we need to take the cards that are p and $\neg q$ i.e. even cards and brown cards.

2.9.3 U.S.A. Presidents



Figure 2.7: Presidents of the United States of America as of 2021.

Suppose we aim to predict whether the next president of the United States of America will be a woman or not. If we rely solely on the gender of previous presidents, by induction we will predict a zero chance. But by understanding how a person becomes a presidential candidate, and how previously became a candidate for their party, we can take into account the network of people involved in the process and recalculate our forecast with higher precision. In this case the rules are clearly defined in the law. Pouring these bits of domain knowledge into our model will show that chances are increasing over time. Encoding the rules behind the data heavily increased the robustness and precision of our model. Thanks to these rules our inference became deductive rather than inductive, since the conclusion necessarily follows from the premises; and as long as the premises are true the conclusion will also be true.

We can identify two issues in the first approach of our example: First, partial data can misrepresent the underlying phenomena that shapes the data, producing a model that does not resemble the real world. This is especially notable in the case of bias and confounders which are further aggravated by the lack of domain knowledge in designing the solutions. The second issue relates to induction. Contrary to deduction, where the truth of the premises guarantees the truth of the conclusion, inductive inferences are *ampliative* — since whose conclusions go beyond what is contained in their premises — and their conclusions could be totally wrong even if infinitely many examples confirm them (Bergadano, 1991). This *ampliative* factor has also an amplifying effect over the partial data from which we infer a conclusion. In this case, considering only the final results of the elections amplified the bias derived from a partial collection of the data, reducing the chances of women being predicted as president to zero.

From (Vega, 2021).

2.9.4 *Yersinia pestis*

This excerpt from Plague and Cholera is a great example of how laboratory conditions can act as an unintended auxiliary hypothesis that must be taken into account during research. It was 1894 in Hong Kong and all was set for an intellectual duel between Alexandre Yersin and Kitasato Shibasaburō that eventually unveiled the cause of the disease plague.

From the moment of his disembarkation in torrential rain, Yersin sees the bodies of plague victims lying in the street, in pools of standing water, in parks, aboard moored junks. British soldiers, acting on authority, remove the sick and empty their houses, pile everything up and set fire to it. [...]

‘I notice many dead rats lying on the ground.’ The first note scribbled by Yersin that evening concerns sewers spewing out decomposed bodies of rats. Since Camus, that has seemed obvious, but not then. [...] By telegram, and as a concession to diplomacy, British governor Sir William Robinson gives Yersin explicit authority to come and study plague in Hong Kong. However, bad faith on the British side is clear to see, and it is even worse with the Japanese team under Shibasaburo Kitasato, who intends to reserve all autopsies for himself. [...]

Never again, in the history of humanity, will there be such an opportunity to become the person who vanquished plague. A few more weeks of devastation will mean a few thousand more bodies to study. [...] Kitasato, though, has a handicap advantage. Not a single cadaver will be placed at Yersin’s disposal. [...]

For Yersin’s benefit he [Father Vigano] arranges, in just two days, to have a bamboo-framed, straw-covered hut erected near the Alice Memorial Hospital. With the matter of his living quarters and laboratory settled, Yersin installs a camp bed, unlocks the cabin trunk, and sets out microscope and test tubes. Vigano then greases the palms of the British sailors in charge of the hospital mortuary, where the bodies are stacked prior to being cremated or buried, and buys several from them. Yersin proceeds to ply his scalpel. [...] ‘The bubo is quite distinct. In less than a minute I have it out and take it up to my laboratory. I make a quick preparation and place it under the microscope. One glance reveals



Alexandre Yersin.

a veritable mess of microbes, all similar. They are small stubby rods with rounded ends.' [...] Yersin becomes the first human being to observe the plague bacillus, as Pasteur was the first to observe those of silkworm pebrine, ovine anthrax, chicken cholera and canine rabies.

What Kitasato describes, having sampled organs and blood and disregarded the bubo, is the pneumococcus of a collateral infection, which he mistakes for the plague bacillus. Without luck, without chance, genius is nothing. The agnostic Yersin is blessed by the gods. Subsequent studies will show that one reason for Kitasato's failure is that he enjoyed the benefits of a proper hospital laboratory, including an incubator set at the temperature of the human body, a temperature at which pneumococcus proliferates, whereas the plague bacillus develops best at approximately twenty-eight degrees centigrade, the mean temperature in Hong Kong at that time of year and the temperature at which Yersin, with no incubator, conducts his observations.

From Plague and Cholera, by Patrick Deville. ([Deville, 2014](#))

I absolutely recommend this book about Alexandre Yersin life. A Swiss-French physician and bacteriologist, pupil of Louis Pasteur, that trying to run away from himself became an agronomist and an explorer of the highlands of Vietnam and Cambodia.

Note for data scientists!

In scientific and data science research, having too much control over experimental conditions can limit our understanding of complex phenomena. The text illustrates this through the contrasting approaches of Alexandre Yersin and Kitasato Shibasaburō during their investigation of the plague. Yersin's observations were conducted in a hut laboratory without the ideal conditions, which allowed him to discover the plague bacillus by directly examining the natural state of the environment. In contrast, Kitasato's access to a fully equipped hospital laboratory may have created an environment that favoured the growth of contaminants, leading to his misidentification of the pneumococcus as the plague bacillus. This situation highlights the importance of balancing control (like in a laboratory setting) and real-world conditions in scientific inquiry; while controlled experiments are essential for reproducibility, they can also obscure critical variables and insights that emerge from more natural settings. Thus, understanding how much control to exert is crucial for ensuring the validity and applicability of scientific findings.

2.9.5 Risks of induction and non-epistemic values in ML

I recommend the following blog post³ from Simon Fischer. I copy a fragment here but the whole article is very interesting.

For example, when we think of the problem of *filter bubbles* we are less and less confronted with opposing world views. Moreover, the idea that the future resembles the past, gives us examples of how Amazon has developed an algorithm for recruiting new staff which only hired males (Dastin, 2018). Even though the model might be correct from an epistemological point of view, such as accuracy or simplicity, it questions non-epistemic values, such as fairness. [...]

Another problem arises with regard to Popper's falsification approach. We cannot be sure what we have falsified: the hypothesis, the auxiliary assumptions, or even both? Consequently, under these considerations, it appears that the risks of drawing conclusions from machine learning outweigh the benefits. [...]

In the case of Amazon the false hypothesis and background assumptions were found rather quickly. But there could be more subtle biases around us which we are not yet aware. This again shows the twofold consequences in terms of inductive risk: The danger of scientists implementing these biases into the algorithms and the benefit of amplifying these biases, and thus making them visible to us. — (Fischer, 2020)

³<https://simonfischer.me/the-necessity-of-non-epistemic-values-in-machine-learning-modelling/>

Chapter 3

Empirical Practices and Models

3.1 Overview

Empirical: based on, concerned with, or verifiable by observation or experience rather than theory or pure logic.

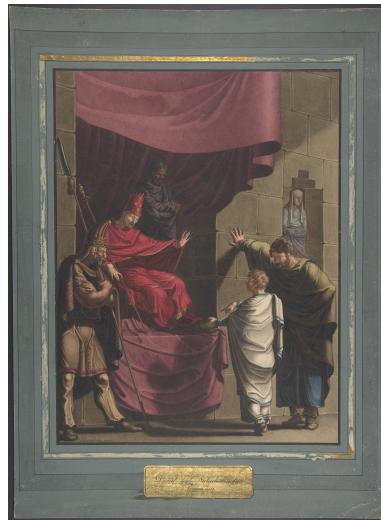


Figure 3.1: Daniel interprets Nebuchadnezzar’s Dream.

I would like to introduce this chapter in the same way the Book of Why ([Pearl and Mackenzie, 2018](#)) introduces its fourth chapter “Slaying the lurking variable”. During the times of Babylonian King Nebuchadnezzar (642 BC - 562 BC), one captive – Daniel – refused to eat royal meat offered by the King as part of their education and service in the court since it did not comply with his religious beliefs. Instead, Daniel asked to be fed on a vegetable diet. The overseer was reluctant as he thought the

servants would lose weight and become weaker. Daniel proposed an experiment to convince his overseer. For ten days, one group of servants would be given a vegetable diet, while another group of servants would eat the king's meat. Then, the overseer would compare both groups and see that the vegetable diet did not reduce their strength. Of course, the experiment was a success, and the king was so impressed that he granted Daniel a favoured place in the court.

This example synthesizes the process of controlled experiments employed nowadays in experimental science. The overseer poses a question, *will the vegetarian diet cause my servants to lose weight?*. There it is our hypothesis. To address the question, Daniel proposed a methodology. Divide the servants in two identical groups. Give one group a new treatment (e.g. diet or a drug), while another group (control) remains under no special treatment. Of course, the two groups should be comparable and representative of some population in order to transfer the conclusions to the population at large. This process allowed Daniel to show the *causal effect* (beware, we will tackle this in Chapter 4) of the diet. Moreover, Daniel's experiment was prospective (in contrast to retrospective studies) as the groups were chosen in advance. Prospective controlled trials are a common characteristic of sound science. Still, Daniel did not think of everything, but we will see that in Chapter 4.

Course objectives

This chapter addresses objectives 5 to 8 by illustrating the varied roles of data across experiments, observational studies, and models. Aimed at students unfamiliar with scientific studies and causal inference, it highlights their limits through examples that raise questions to be tackled in the chapter dedicated to experimental control and statistical abuse.

3.2 Experiments

Many data scientists believe their role should be limited to data analysis, but experiment design is fundamental for data collection, which conditions how the data must be analysed. Conclusions drawn from data can be biased or determined by decisions and errors taken during experiment design. Understanding this can help you spot issues during the data analysis and ask the right questions to your colleagues in charge of the experiments.

An experiment is an observation process in which we control background variables through manipulation, intervene on target variable (through manipulation) and observe the difference produced by such intervention thanks to measurements.

Experiment is the kind of scientific experience in which some change is deliberately provoked, and its outcome observed, recorded and interpreted with a cognitive aim. — (Bunge, 2017)

3.2.1 Observational studies

However, there are whole research areas where scientists cannot make experiments. For instance, astrophysics is mainly observational and theoretical as it is not possible to manipulate the observed entities (e.g. stars). It aims to find out measurable implications of physical models. Sometimes it is not feasible, legal or ethical to conduct certain types of experiments, conducting observational studies instead. So, in **observational studies** there is no manipulation, no intervention on the target variable, neither control of background variables.

3.2.1.1 Natural experiments

Natural experiments¹ on the other side share the first two characteristics but it is possible to control background variables (but not through manipulation though). See § 3.5.2 for an example. A major limitation of natural experiments when inferring causation is the presence of unmeasured confounding factors. Natural experiments are appealing for public health research because they enable the evaluation of events or interventions that are difficult or impossible to manipulate experimentally, such as many policy and health system reforms (de Vocht et al., 2021).

For example the Canterbury earthquakes in 2010-2011 could be used to study the impact of such disasters because about half of a well-studied birth cohort lived in the affected area with the remainder living outside. [...] More recently, the use of the term ‘natural’ has been understood more broadly as an event which did not involve the deliberate manipulation of exposure for research purposes, even if human agency was involved. [...] Natural experiments describing the study of an event which did not involve the deliberate manipulation of an exposure but involved human agency, such as the impact of a new policy, are the mainstay of ‘natural experimental research’ in public health. — (de Vocht et al., 2021)

See Figure 3.2 for a schema depicting the conceptualisation of natural and quasi-experiments. Some authors differentiate between natural experiments and *quasi-experiments*. In a quasi-experiment, the criterion for group assignment of the study units (e.g. study participants) is selected by the researchers, whereas, in a natu-

¹https://en.wikipedia.org/wiki/Natural_experiment

ral experiment, the assignment occurs *naturally*, without the intervention of the researchers.

Quasi-experiment: A quasi-experiment is an empirical interventional study used to estimate the causal impact of an intervention on target population without random assignment. Quasi-experimental research shares similarities with the traditional experimental design or randomized controlled trial, but it specifically lacks the element of random assignment to treatment or control. Instead, quasi-experimental designs typically allow the researcher to control the assignment to the treatment condition, but using some criterion other than random assignment. Quasi-experiments are subject to concerns regarding internal validity, because the treatment and control groups may not be comparable at baseline. In other words, it may not be possible to convincingly demonstrate a causal link between the treatment condition and observed outcomes. This is particularly true if there are confounding variables that cannot be controlled or accounted for. — Wikipedia on ([Rossi et al., 1985](#)) and ([DiNardo, 2010](#)).

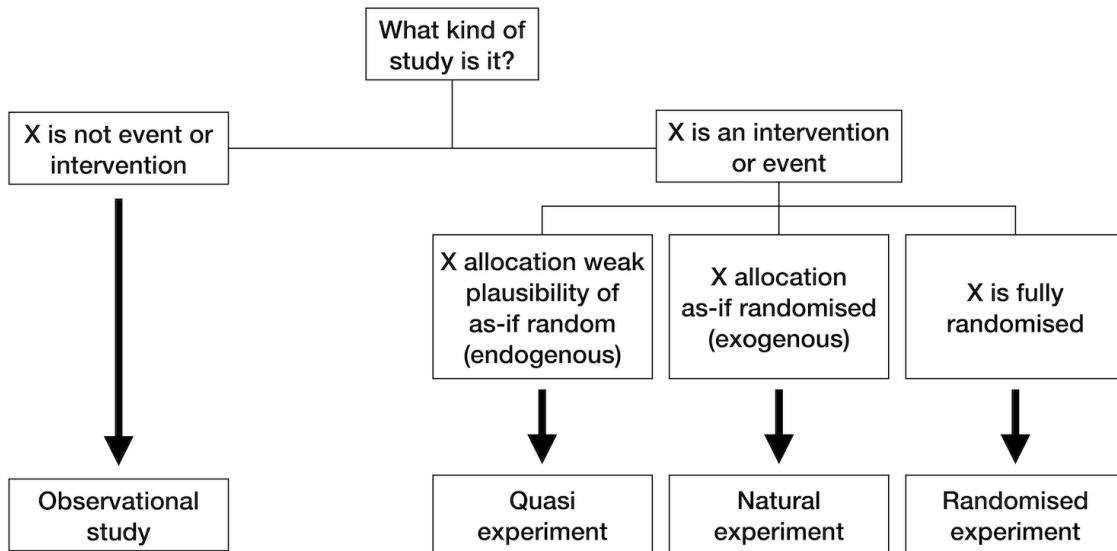


Figure 3.2: Diagram depicting the conceptualisation of natural and quasi-experiments within the evaluation framework of Thad Dunning. Re-drawn from ([de Vocht et al., 2021](#)). Note that the same article provides three additional conceptualisations from different frameworks. For example, a different conceptualisation makes a distinction between quasi and natural experiments, arguing that natural experiments describe unplanned events whereas quasi-experiments describe events that are planned (but not controlled by the researcher).

Dunning takes this concept further and defines a ‘natural experiment’

as a quasi-experiment where knowledge about the exposure allocation process provides a strong argument that allocation, although not deliberately manipulated by the researcher, is essentially random, referred to as ‘as-if randomization’ — ([de Vocht et al., 2021](#))

Definitions:

Target variables: The target variable of a dataset is the feature of a dataset about which you want to gain a deeper understanding. They also receive the name “dependent variables” because, in an experiment, their values are studied under the supposition or demand that they depend, by some law or rule (e.g., by a mathematical function), on the values of other variables. The dependent variable is the *effect*. Its value depends on changes in the independent variable.

Independent variables: It is a variable that stands alone and isn’t changed by the other variables you are trying to measure. The independent variable is the *cause*. Its value is independent of other variables in your study.

Background variables: An explanatory variable that can affect other (dependent) variables but cannot be affected by them. For example, one’s schooling may affect one’s subsequent career, but the reverse is unlikely to be true.

We can recognise five elements in the observation process: the *object* of observation; the *subject* (or observer) and its perceptions; the *circumstances* of observation (e.g. environment of object and subject); the observation *media* (e.g. senses, instruments, procedures); and the body of *knowledge* used to relate all the previous elements. The last two can be grouped into *tools* (concrete and conceptual). So, an observation statement has the form “*w* observes *x* under *y* with the help of *z*”. ([Bunge, 2017](#))

3.2.1.2 Observability

We can try to define observability by stating that a fact is *observable* “only if there exists at least one subject, one set of circumstances, and one set of observation tools, such that the fact can appear to the subject armed with those tools under those circumstances” ([Bunge, 2017](#)). This definition is rather unsatisfactory since someone could claim the existence of ghosts or aliens. We should define what is objectively observable. Then, *x* is observable only if there exist at least one recording instrument *w*, one set of circumstances *Y*, and one set of observation tools *Z*, such that *w* can register *x* under *y* helped by *z*. Here we have eliminated the possibility of the subject’s perceptual delusions, but devices (e.g. a camera) have limitations too.

Observations are often expressed in the form of a rule so that other researchers can reproduce their results under similar conditions. Some facts cannot be repeated,

such as the eruption of a volcano or a supernova. So very often, we expect results of the same kind to be reproducible by observers. Exact duplication is desirable but not always achievable. Even independent observers may make the same wrong observations due to faulty equipment or false hypotheses.

3.2.1.3 Indicators

Most facts we know about are indirectly observable, i.e. we infer them through an intermediary. For instance, the wind is not directly observable but inferred from bodies apparently moved by it. We *objectify* an unobservable fact by establishing its relationship to some perceptible fact(s) that serve us as an *indicator* of the fact. In other words, hypotheses are made concerning unperceived facts and tested through evidence consisting of data about other directly observable facts, assuming that the latter are **collaterally connected with** or **effects** of the former. Of course, that such relationship should hold is as well a hypothesis (see Figure 3.3).

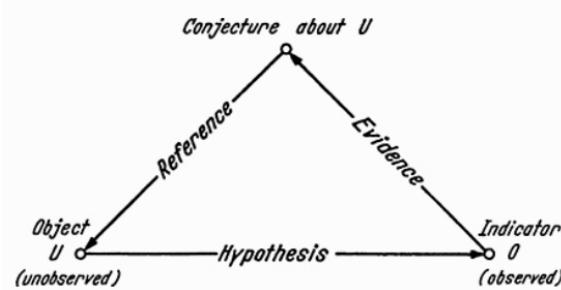


Figure 3.3: The physical object-indicator relation, is expressed by a hypothesis enabling us to infer the object from observations made on its indicator. Figure extracted from ([Bunge, 2017](#)).

3.2.1.4 Data and Evidence

Every evidence is a *datum* but not every datum constitutes *evidence*. What turns a datum into evidence is that is relevant to some idea, that it makes sense under some theory or body of knowledge. In particular, we believe a datum constitutes an evidence in favour of a theory and assign the theory some *credence* because it justifies or predicts that evidence. The evidence must be related to a specific hypothesis, and this relationship is justified because of a body of theoretical knowledge. In fact, no evidence is absolute. Consider the following example from ([Bunge, 2017](#)):

The observed deviation of a magnetic needle in the vicinity of an electric circuit (datum e) supports the hypothesis h_x that electricity is flowing

through the circuit, on the theory T_1 that electric currents produce magnetic fields which in turn interact with the fields of magnetic needles. But exactly the same datum e might be taken as an evidence in favour of the rival hypothesis h_2 that a big magnet somewhere nearby has been switched on, on the theory T_2 that magnets can interact directly with one another. The given datum is then *ambiguous* and only an independent checking of h_x and h_2 , i.e. a test independent of e , will enable us to reach a decision between the two rivals.

Importantly, the characteristics that make data count as evidence must be agreed prior to observation and on the basis of theory. Sometimes a scientist may obtain data that seems incompatible with a theory. Instead of getting rid of such data (or the theory), the scientist will attempt to reproduce the data and assess whether is anomalous data (e.g. due to a faulty instrument) or not. The *raw* data may contain any information, but *refined* data should express only relevant and useful information for the problem at hand. Of course, some information is always lost in the process. In consequence, the refinement process is irreversible. Data are *means* rather than *ends* and we aim to systematise data in order to disclose patterns on it. For this reason *noise* must be removed. The systematization of refined data may involve displaying information in graphs or tables as well as arranging information in data structures such as matrices.

3.2.2 Field, laboratory and simulation experiments

3.2.2.1 Field experiments

In contrast to observational experiments, **field experiments** randomly assign the sampling units (e.g. study participants) into two groups (treatment and control) to test causal relationships. The same conditions are maintained for both groups only varying the intervention on the factor of interest (e.g. two parts of soil (fertilized/unfertilized)). The background variables are considered as given and not manipulated.

- No manipulation.
- No intervention on the target variable.
- Control for the background variable (but not through manipulation).

In the example (see Figure 3.4), the background variables are controlled, we do not alter the soil, the number of hours of sun light received by the two groups of plants, nor the watering conditions. The only intervention is giving fertiliser to one side of the field. In this case, the seeds can be randomly assigned the treatment (fertiliser) or control groups.



Figure 3.4: Fertiliser experiment.

Potential threats to internal validity

- Excludability: The assumption of excludability states that the randomization does not affect outcomes through other variables than the reception of the treatment. If this assumption is violated, the causal effect identified in a study is a combination of the treatment and other variables ([Hansen and Tummers, 2020](#)). For instance, that the two fields do not receive the same amount of light.
- Interference: Interference occurs when experimental units alter each other's outcomes. This generates a bias that precludes the proper estimation of causal effects by the researchers.
- Attrition: Attrition occurs when outcome data are missing. Attrition becomes a problem for causal inference when two conditions are present: (1) units with missing outcomes differ systematically on the outcome from those that are not missing and (2) attrition is different in experimental groups. There is greater potential for attrition in field experiments than in laboratory experiments because field experiments confer less control ([Hansen and Tummers, 2020](#)). For instance, some participants may leave a study if they do not get any improvement.

3.2.2.2 Laboratory experiments

On the other side, **laboratory experiments** construct the same background conditions in both groups manipulating the environment (lab settings) and varying the intervention on the factor of interest. Background conditions are controlled through manipulation. For instance, temperature, pressure, humidity can be controlled for a fertiliser trial. Laboratory experiments tend to have higher internal validity, but at

the cost of lower external validity (generalisation), owing to the artificial setting in which the study is conducted may not reflect the real world.

3.2.2.3 Simulation experiments

Finally, **simulation experiments** are constructions representing a real system on a computer to perform interventions. This type of experiments are done when it is not feasible to experiment on the real entities (e.g. climate simulations or geological simulations). The important consideration is that all interventions and manipulations are performed on the computer representation instead of the real target itself.

3.2.3 Summary of the different experiments

Therefore, an experiment is a controlled observation in which the observer manipulates the real variables (independent variables) that are believed to influence the outcome (dependent variable), both for the purpose of intervention and control. The following article provides a good description of the basics of experiments².

In Chapter 4 we will see some examples of experimental errors (e.g. confirmation bias, selection bias, etc) as well as examples of statistical abuse. All in all, the experiment process is also a craft which entails learning from previous experiments (ours and others), as well as applying all available knowledge (theoretical and experimental) for the design of experiments.

Experiments differ in how much control and intervention they allow, shaping what kind of knowledge they produce. **Observational studies** only reveal correlations, useful for description but weak for causal inference. **Natural** and **quasi-experiments** bridge the gap exploiting either naturally occurring or researcher-driven interventions without full randomisation, allowing limited causal claims under specific assumptions. **Field experiments** introduce real interventions with random assignment in natural settings, increasing realism but reducing control. Finally, **laboratory experiments** maximise internal validity through artificial control, while **simulation experiments** model reality computationally to test counterfactuals when real interventions are impossible.

For data science, understanding these designs matters because **datasets inherit the assumptions and biases of their underlying study**. A data scientist's analytical choices must respect these epistemic limits: **no method can fix a flawed design**.

- All empirical designs involve assumptions.
- The strength of causal inference lies on control and intervention degrees.

²<https://opentextbc.ca/researchmethods/chapter/experiment-basics/>

- Underdetermination remains, i.e., different theories can explain the same evidence, so empirical adequacy does necessarily equal truth.
- Data are not neutral: they reflect study design choices, measurement biases, and contextual limits.
- Good analysis starts with good design.

Type			Back.			Purpose	Limit.	e.g.
	Man./ Interv.	Var. Ctrl.	Alloc.	Setting	Purpose			
Obs. Study	None	None	None	Real-world	Describe patterns or correlations when intervention is impossible, unethical, or impractical	No causal inference; strong risk of confounding	Astrophysics, geology, epidemiological observation	
Natural Exp.	No direct manipulation	<i>As-if ran-domised</i> exogenous factors	“As-if” random (by nature)	Real-world	Estimate causal effects from naturally occurring variation	Hidden founders may persist; ran-domisation only assumed	John Snow’s cholera study; Canterbury earth-quakes	

Type	Back.		Alloc.	Setting	Purpose	Limit.	e.g.
	Man./ Interv.	Var. Ctrl.					
Quasi- Exp.	Int. may be present	Groups chosen by criteria	Non-random assignment	Real-world / policy contexts	Estimate causal impact when ran-dom)	Internal validity issues: treatment or feasible (ethical or practical limits)	Evaluation of new school policy
Field Exp.	Real intervention	Back. conditions given, not manipulated	Random assignment of participants	Real-world (schools, hospitals, farms)	Tests causal relationships in natural settings; higher external validity	Lower internal validity due to exclusivity, violations, interference, attrition	Fertilised vs. un-fertilised soil; clinical field trials
Lab. Exp.	Man. of target variable and environment	Full control over background	Random assignment	Lab	High internal validity; lab isolates causal mechanisms	Low external validity; lab results may not generalise	Physics or chemistry experiments; cognitive psychology tasks

Type	Back.		Setting	Purpose	Limit.	e.g.
	Man./ Interv.	Var. Ctrl.				
Sim.	Man.	Full	Virtual	Comp.	Explore	Dependent Climate
Exp.	within a computer model	control within simu- lated environ- ment	ran- domisa- tion	causal hy- pothe- ses when real exper- imenta- tion is infeasi- ble or danger- ous;	on realism and assump- tions of the model; must be vali- dated	models; epi- demic or eco- nomic simula- tions

Definitions:

Repetition: An experiment is repeatable if enough information is provided about the used data and the experiment methods and conditions. With such information, it should be possible to repeat the experiment.

Reproduction: An experiment is considered as reproduced if the repetition of the experiment yields the same result. For instance, in computer science, reproducing involves using the original data and code.

Replication: An independent experiment, in the spirit of the original experiment produces the same result. For example, in computer science replication entails collecting new data and use similar methods to reach similar conclusions in answer to the same scientific question. Or implementing a new software following similar design principles and reaching similar results.

3.2.4 How to evaluate experiment success

Very often, success is not defined by a single goal or metric. For instance, the best car is not always the fastest car. In fact, there are many other values to bear in mind, such as gasoline consumption, pollution, ease of manufacture, etc. Similarly an experiment success is rarely assessed with a single metric in mind.

Moreover, some metrics must not be degraded, often called **guardrail metrics**. This type of metrics can include security, speed, robustness, etc. But very often include *non-epistemic values* too. In this context, non-epistemic values are metrics not directly related to the instance to be designed, such as fairness, justice, or making money (or saving it), in contrast to metrics that make the instance at issue *internally* or *intrinsically* better (e.g. speed). For instance, a car is not necessarily a better car depending on its price if what is judged is the *car itself* in isolation, but a low price might make it easier to sell. In another example, the fastest data processing system might not necessarily be the best choice since other requirements must be considered too (e.g. ease of use).

A non-epistemic value that is always at stake is money, or in a different shape, OPEX (operational expenditure) and CAPEX (capital expenditure). Very often, they condition other metrics, such as performance (e.g. use less/worse resources) or safety (e.g. employ less/worse materials). For example, I had the opportunity to work on the design of enterprise log processing systems. In this case, we wanted to maximise speed while reducing resources, as mid-sized companies often wish to reduce the number of servers deployed, which ultimately affects their operational costs (e.g. space and electricity). Most commercial solutions scale horizontally, requiring the use of on-site server clusters to handle large amounts of data (at prohibitively high prices) or cloud-hosted clusters (impractical due to data protection). Our proposal optimised vertical scalability and coped with tens of millions of events per second with a single server. But of course, such an approach was specifically designed for a particular task, in contrast to the flexibility offered by commercial alternatives.

In data science, success should be defined by how well the analysis answers the research questions. For this reason, setting the research questions at the very beginning of the process remains crucial. They not only determine the data analysis but, more importantly, the data collection design. However, very frequently, the data science process starts with a given dataset. Still, it is essential to assess if the collected data can answer the posed questions.

3.3 Scientific models

Before starting, we need to remind ourselves that data alone does not yield explanations. Data describes regularities, but only models tell us reasons. Still, reasoning alone is also insufficient. Theories require empirical confrontation to validate assumptions. Progress comes from models failures, not just their success. Humility stems as a key value for scientific progress. Otherwise, precise data and a wrong model lead to misleading certainty. **Models bridge reasoning and observation.**

Scientific models are widespread and varied. These include scale models (e.g. plane models for aerodynamic studies), laboratory animals (e.g. mice for drug trials), simulation models (e.g. computational model for weather forecast). In all the previous cases, a model is made to replace or to stand in for what we are ultimately interested in. Models are characterised by being **representations**, containing **idealisations**, being **purpose dependent**, and **ready to be manipulated**.

Typically models are representations in which details, that appear inessential for intended uses, are omitted. A model is intended to represent the real thing in certain significant aspects. — Guy Orcutt

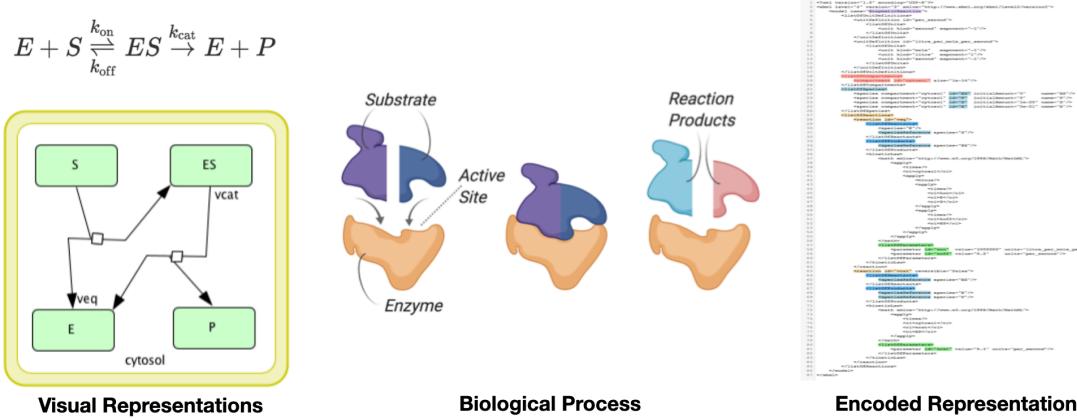


Figure 3.5: Different representations of a biological process. On the right side, the process is encoded in a XML file in Systems Biology Markup Language (SBML) format.

We shall be concerned with model objects and theoretical models as hypothetical sketches of supposedly real, though possibly fictitious, things or facts. Thus a fluid may be modeled as a continuum endowed with certain properties, such as compressibility and viscosity. Such a model object may be grafted onto any of a number of general theories, say classical mechanics, or general relativistic mechanics. Likewise, a learning organism may be modeled as a black box equipped with certain input and output terminals, and this model object may then be expanded into a hypothetico-deductive system. In either case a specific theory, or theoretical model, of a concrete (or supposedly concrete) object, results. What can be subjected to empirical tests are such theoretical models: on the other hand general theories, being unconcerned with particulars, remain empirically untestable unless enriched with models of their referents. — (Bunge, 2012)

For instance, a physical model of DNA on a table represents the real DNA. Obviously, such a model is not a piece of real DNA. It is made of something else (e.g. plastic) and at a different scale. In this case, such a model is useful for pedagogic purposes. Although there are clear differences between models and targets, the key relationship is that **a model represents (in some way) the target**. From the methodological point of view, we must justify why to represent targets with models instead of investigating the targets themselves?. Possible answers include physical impossibility, costs, ethical and legal reasons, etc. However, very often the main justification to employ models is that targets are very complex. Therefore, employing a model that simplifies the target complexity might allow us to get a better understanding of the main factors operating in the target system. Our cognitive limits very often determine how we investigate complex systems, starting with a simpler model and increasing its complexity as we gain understanding.

A schematic representation of an object may be called a model object. If the represented object (or referent) is concrete or physical, then its model is an idealization of it. The representation may be pictorial, as in the case of a drawing, or conceptual, as in the case of a mathematical formula. It may be figurative, like the ball-and-spoke model of a molecule, or semisymbolic, as in the case of the contour map of the same molecule; or finally symbolic like the hamiltonian operator for that same object. —
(Bunge, 2012)

Therefore, we are intentionally choosing or building a model that differs from the target in some properties. Precisely because of this condition, we cannot assume that whatever is the case in the model is also the case in the target. **Models come with idealisations.** Not bearing in mind this key condition of the models can lead us to produce false claims about the target.

- **Idealized models.** Idealized models are models that involve a deliberate simplification or distortion of something complicated with the objective of making it more tractable or understandable. Frictionless planes, point masses, completely isolated systems, omniscient and fully rational agents, and markets in perfect equilibrium are well-known examples. Idealizations are a crucial means for science to cope with systems that are too difficult to study in their full complexity.
- **Scale models.** Some models are down-sized or enlarged copies of their target systems (Black 1962). A typical example is a small wooden car that is put into a wind tunnel to explore the actual car's aerodynamic properties.

- **Phenomenological models.** Phenomenological models have been defined in different, although related, ways. A common definition takes them to be models that only represent observable properties of their targets and refrain from postulating hidden mechanisms and the like.
- **Exploratory models.** Exploratory models are models which are not proposed in the first place to learn something about a specific target system or a particular experimentally established phenomenon. Exploratory models function as the starting point of further explorations in which the model is modified and refined.
- **Models of data.** A model of data (sometimes also “data model”) is a corrected, rectified, regimented, and in many instances idealized version of the data we gain from immediate observation, the so-called raw data. Characteristically, one first eliminates errors (e.g., removes points from the record that are due to faulty observation) and then presents the data in a “neat” way, for instance by drawing a smooth curve through a set of points. These two steps are commonly referred to as “data reduction” and “curve fitting”. — ([Frigg and Hartmann, 2020](#))

3.3.1 The models of the atom

For example, Bohr’s model of the atom assumes that electrons orbit the atomic nucleus in circles. The success of such a model relied that the Bohr assumptions reproduced the series that fitted the hydrogen emission spectra. In 1913 it predicted the correct frequencies of the specific colours of light absorbed and emitted by ionised helium. One could say that Bohr was very lucky as despite his model is wrong in some ways, it also has some bits of truth, enough for his predictions about ionised helium to work out.

However, other predictions about the properties of the atom were wrong, and its implications were not observed in experiments. In the Schrödinger model, the electron of a one-electron atom, rather than travelling in fixed orbits around the nucleus, has a probability distribution allowing the electron to be at almost all locations in space, some being much more likely than others. Bohr theory (1913) was rejected in 1925 after the advent of quantum mechanics, but its model remains because despite its flaws and idealisations, Bohr’s model is useful for education⁴.

⁴<https://blogs.scientificamerican.com/guest-blog/why-it-s-okay-to-teach-wrong-ideas-in-physics/>

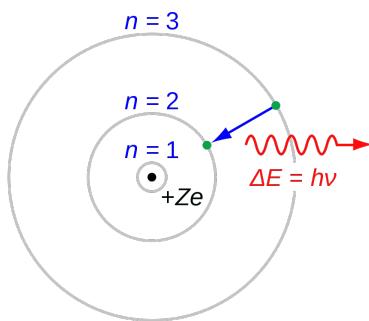


Figure 3.6: Illustration of bound-bound transition in the Bohr atomic model. Source: Wikipedia Commons³.

3.3.2 The models of benzene

Models are as well purpose-dependent. Suppose the next question. Which benzene model is better? A quantum mechanic model, or a structural formula?. On one side, the quantum mechanic model is more precise about the potential position of electrons. Additionally, is more similar to the target as it represents better its relevant properties. The structural model is simpler and easier to work with. In this case, theoretically tractable models such as structural models allow for functional group analysis in chemistry.

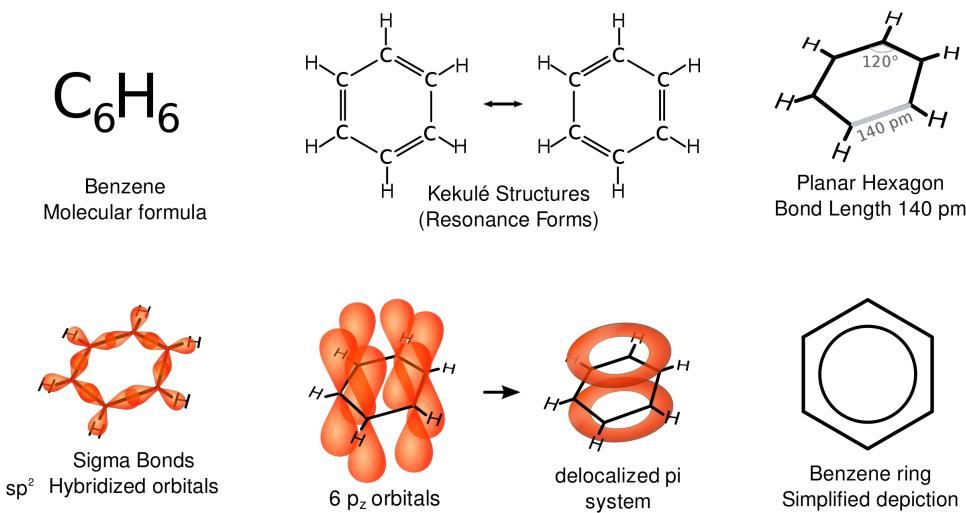


Figure 3.7: Various representations of Benzene. Source: Wikipedia Commons⁵.

3.3.3 Models as analogies

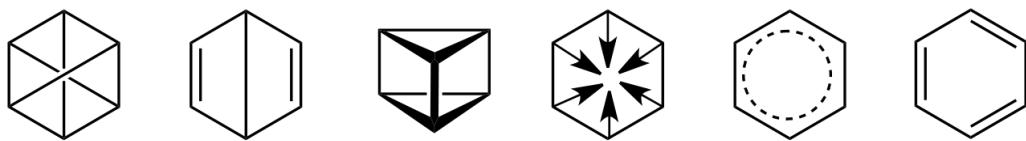


Figure 3.8: Historic benzene structures (from left to right) by Claus (1867), Dewar (1867), Ladenburg (1869), Armstrong (1887), Thiele (1899) and Kekulé (1865). Dewar benzene and prismane are distinct molecules that have Dewar's and Ladenburg's structures. Thiele and Kekulé's structures are used today.

Philosopher Mary Hesse (1924-2016) argued that models act as analogies rather than descriptions of the targets. She distinguished between 3 kinds of analogies. In the first place, she considered that **positive analogies** hold between the aspects of a model and its target for which we have reasons to believe they are similar. An example of positive analogies can be found between mice and humans, which have similar hormone systems and physiology. On the other hand, the idealisations constitute **negative analogies**, such as the differences in size or lifespan between mice and humans. These negative analogies cover the properties in which model and target differ. Finally, **neutral analogies** concern the properties that we cannot investigate directly in the target, requiring the model for their study.

Without analogy there might be no knowledge: the perception of analogies is a first step towards classification and generalization. ([Bunge, 2012](#))

For instance, the reaction to a certain drug or treatment of interest. At the initial stage, is not possible to tell whether the model-target relationships concerning these properties constitute positive or negative analogies because we do not know yet how the relevant target properties are affected. Instead, such properties are investigated in the model, and researchers hypothesise that the model is analogous to the target in such properties. For instance, we assume that the effects of a drug in mice will give us knowledge about its effects in humans.

The positive analogy between two items consists in the properties or relations they share (both gas molecules and billiard balls have mass); the negative analogy consists in the properties they do not share (billiard balls are colored, gas molecules are not); the neutral analogy comprises



Prof. Mary Hesse (1924-2016),
by Peter Mennim.

the properties of which it is not known (yet) whether they belong to the positive or the negative analogy (do billiard balls and molecules have the same cross section in scattering processes?). Neutral analogies play an important role in scientific research because they give rise to questions and suggest new hypotheses. — (Frigg and Hartmann, 2020)

Consider again the Michelson and Morley experiment. Before the XX century, most physicists considered light as a wave. Their beliefs were justified on the many positive analogies between light, water and sound waves. For instance, light produces a diffraction pattern when encountering an obstacle, just as water and sound waves do. With this model in mind, physicists inferred a neutral analogy, namely, that light needs a medium to travel, as other waves require. They called this medium: luminiferous aether. The experiment from Michelson and Morley is a consequence of such model. However, the experiment revealed that such analogy was indeed a negative analogy, an idealisation. This discovery led people to replace the model of light for more precise models. Therefore, model manipulation allows discovering the effects of neutral analogies.

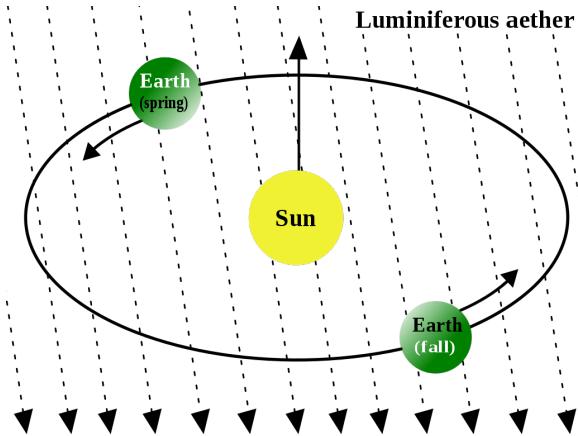


Figure 3.9: The luminiferous aether: it was hypothesised that the Earth moves through a "medium" of aether that carries light. Source: Wikipedia Commons⁶.

3.3.4 Differences between Models and Experiments

There are several commonalities between models and experiments. Setting parameters and variables in models resembles experimental control. In the same way, model manipulation is akin to experimental manipulation. Moreover, model manipulation yields results that are observed, just as in experimental observations.

Notwithstanding, there are also important differences to bear in mind, which mainly concern the source of errors. For instance, the most troubling experimental errors

concern internal validity questions, i.e. the degree of confidence that the causal relationship put to test is not influenced by other factors or variables. For these reasons, researchers require careful design and control of experiments. Nonetheless, models are generally less sensitive to these issues since the modeller is aware of the idealisations and mechanisms of its models.

For modelling, the key concern is whether the relevant analogies between model and target hold. Such concern is usually not a problem for experiments, especially those conducted directly on the target. Once internal validity is assessed, researchers are confident that the inferences drawn from the experiment refer to the target. However, inferences drawn from model manipulation constitute neutral analogies considered as hypotheses regarding the target, requiring further testing and justification. This last step is error-prone.

3.3.5 What makes a good model?

Since models are purpose dependent, there is no exhaustive set of sufficient and necessary conditions to define what a good model is. Nonetheless, there are some common criteria (e.g. robustness, simplicity, tractability) that can be balanced, but is often impossible to optimise all criteria at the same because some criteria are complements of each other.

The **similarity** criterion can for example assess physical resemblance. More generally, we could say that a model M is similar to target X if and only if M is similar to X with respect to the set of properties P to a certain degree. However, this definition does not tell us which properties should be optimised. For instance, for a scale model of an air plane aimed at aerodynamic studies, it might be more justified to maximise the similarities of geometric properties over the interior design (e.g. the number of seats might not be relevant since the cabin is closed). Therefore, the purpose of the model justifies maximising one set of properties over another, in particular, the properties that are relevant for the research purpose.

Robustness expresses how model results are affected by condition changes. Therefore, a model result is robust (w.r.t. some condition) if changing such condition does not alter the result. For example, all properties except one (e.g. plane hull colour) can be kept fixed to test whether painting the plane with a different colour might affect its aerodynamic properties. If the result remains equal, we can say the model is robust with respect to the hull colour. Perhaps such property is not relevant to the research purpose, but that is not enough to justify removing the property.

Another model criterion to consider is **precision** (w.r.t. parameters). We say that model M_1 is more precise than M_2 if the parameter specifications of M_1 *imply* those

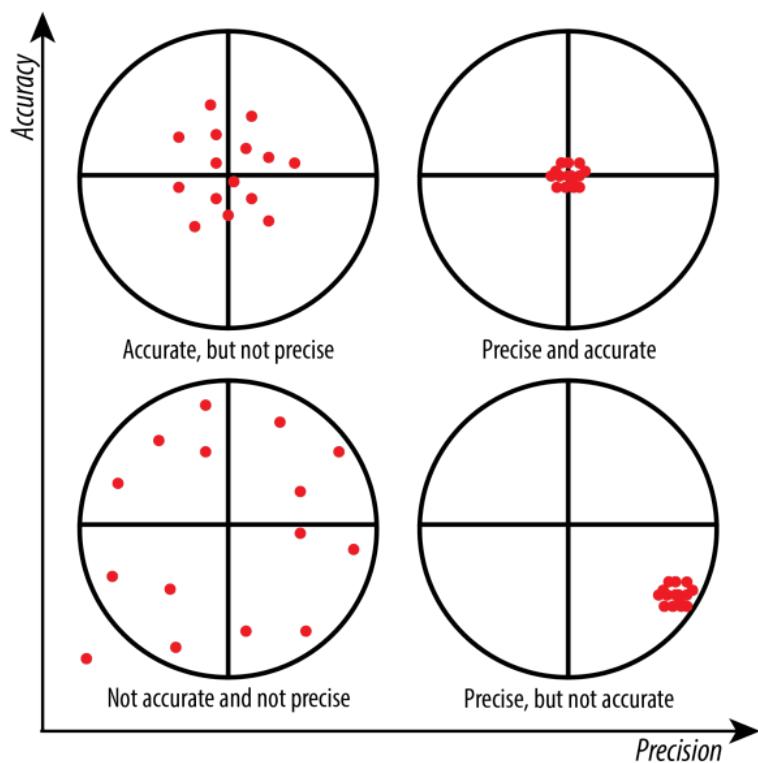


Figure 3.10: Accuracy consists of trueness (proximity of measurement results to the true value) and precision (repeatability or reproducibility of the measurement). Source: St. Olaf College⁷.

from M_2 . This definition is better understood through an example. Consider the following models M_1 , M_2 , and M_3 and below definitions. The first model describes a rate of changes as a function of X . M_2 is more precise as describes the rate of changes as a linear function of X . The description of M_2 implies the description of M_1 . Linear functions of X are a subset of functions of X . Finally, the third model is yet more precise as it indicates an absolute value of the parameter a , reducing the subset of linear functions from the definition of M_2 to a particular linear function. Importantly, parameter precision is a property of the model alone, not of the relationship between model and target. Although precision offers potential for high accuracy, it is no warrant for it. For instance, if the actual rate of linear change would be other than $1.2X$, then the less precise model M_2 would be more accurate than M_3 . Similarly, if the rate of change would not change linearly, the more general model M_1 would be more accurate than the alternatives.

- $M_1 : dX/dt = f(X)$
- $M_2 : dX/dt = aX$
- $M_3 : dX/dt = 1.2X$

Note for data scientists!

Questions regarding scientific models also concern Machine Learning models to a great extent. For example, consider the precision and accuracy criteria. The following paragraph is extracted from the article “A Few Useful Things to Know About Machine Learning” by Pedro Domingos.^a

Everyone in machine learning knows about overfitting, but it comes in many forms that are not immediately obvious. One way to understand overfitting is by decomposing generalization error into bias and variance. Bias is a learner’s tendency to consistently learn the same wrong thing. Variance is the tendency to learn random things irrespective of the real signal. Figure illustrates this by an analogy with throwing darts at a board. A linear learner has high bias, because when the frontier between two classes is not a hyperplane the learner is unable to induce it. Decision trees do not have this problem because they can represent any Boolean function, but on the other hand they can suffer from high variance: decision trees learned on different training sets generated by the same phenomenon are often very different, when in fact they should be the same. Similar reasoning applies to the choice of optimization method: beam search has lower bias than greedy search, but higher variance, because it tries more hypotheses. Thus, contrary to intuition, a more powerful learner is not necessarily better than a less powerful one. — (Domingos, 2012)

^a<https://dl.acm.org/doi/pdf/10.1145/2347736.2347755>

Simplicity is often another criteria that affects models. A simpler model might fit very well its purpose. For example, underground maps often misrepresent distances or omit unnecessary details such as roads, monuments, etc. Such simplification is suited for the particular purpose of travelling in the underground but is not useful for other purposes. We can say a model is simpler if it contains less parameters, considers less variables and uses less operations than another model. Therefore, simplicity is a virtue with respect to the models and not to targets. Is usually a practical criterion that facilitates the model use.

Related to simplicity, we can find **tractability**. We say a model is tractable (w.r.t. a set of rules) if the relevant model result may be obtained by applying certain principles to the model. For instance, models solved through analytical methods (e.g. mathematical proofs) are called analytically tractable in contrast to models for which results can only be approximated through numerical simulations methods. Tractability implies the existence of methods to analyse and solve such models.

In this sense, **theoretical tractability** considers theoretical principles to assess the suitability of the model for certain operations. For instance, a structural representation of a chemical compound allows for the application of functional group classification (this is an example of theoretical principle to fulfil). In contrast, a quantum mechanical model is more accurate but does not allow for such operation. Therefore we consider it a less tractable model.

Finally, **transparency** is an epistemic value that assesses the degree to which the model user can cognitively understand how the model result is produced. This criterion is known in artificial intelligence as interpretability and/or explainability. For example, a decision tree is often human-readable while the nature of neural networks creates obfuscated models difficult to interpret. Transparent models allow to back-track the result and understand how it was produced from a given input. A transparent model enables the scientist to check the correctness of the results, which is especially important when employing models developed by third parties.

Again, most of the previous virtues will require certain trade-off. Increasing an epistemic value often entails decreasing another one. Therefore, building or choosing a model requires finding the best trade-off for the model's purpose.

3.3.6 Models as mirrors

A common way to consider models is as mirrors of the real world. Very often, models are built to be as similar as possible to the target. This is common in highly complex projects such as brain simulations of the neural networks that represent human brain activity or epidemic simulations in which all the available demographic information (e.g. transportation, habits) is considered. The aim of such complex projects is to build a model as a replacement of the actual target system. In the case of epidemic simulations, it is clear that is not feasible to conduct real-world experiments, but the simulation can serve as a way to try different vaccination strategies that might, for instance, prioritise the vaccination of potential super spreaders, i.e. people who are in contact with many people during their daily routine (e.g. supermarket cashiers, waiters). In this sense, we say that the simulation is *mirroring* the world. Another example from engineering includes finite element analysis (FEA) which is used to divide a complex problem into smaller elements to facilitate calculations.

This type of model is very similar to the target and require high precision modelling. However, these advantages come at a cost. For instance, they compromise the simplicity, transparency and sometimes the tractability of the models. Despite the high similarity of the models to their target systems, it is not enough to avoid external validity issues. For example, FEA employs a mesh consisting of millions

of small elements that mould the physical shape of the analysed structure. Then, calculations are made for each element. Such approximations are usually polynomial, which means that the structural elements have been interpolated, and their precision is bounded to the mesh size. Therefore, the accuracy depends on the purpose of the analysis (e.g. car, hardware tools, camera).

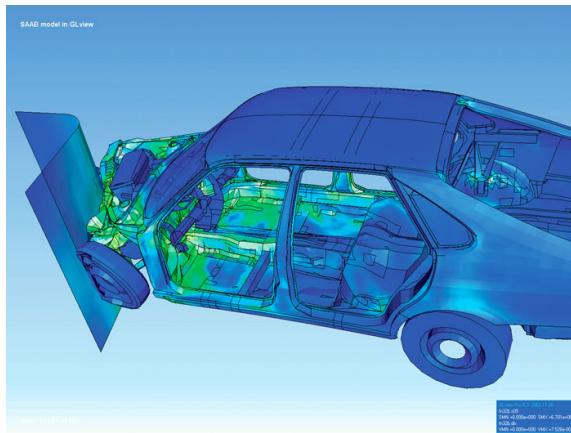


Figure 3.11: A visualization of an asymmetrical collision analysis using the finite element analysis method.

3.3.7 Models as isolations

An alternative perspective is to consider models as isolations of particular features of the complex world. This consideration is motivated by the following question. Can models be similar to their target systems and still be simple? We have previously seen that these two factors are very often inversely related. Isolated models choose a particular aspect of the target, disregarding all the rest. Then, a model is built to represent the behaviour of such factors as accurately as possible.

During the development of the airplane. George Cayley (1773 - 1857) proposed to separate the airplane system into 3 subsystems: Lift, Propulsion, Control. After this, the problem is divided into three problems, i.e. how to obtain lift, how to provide propulsion, and how to offer control. The Wright Brothers developed a separate model for each of this subsystems. As can be observed in Figure 3.12, the lift balance model does not resemble an airplane. Similarly, their propeller model was not attached to the airplane. Likewise, they employed gliders to test their control systems. This contrast with the approach of other inventors, such as Hiram Maxim, who attempted to build a full scale from scratch. An overview of the Wright Brothers Invention Process can be found at NASA's website⁸.

⁸<https://wright.nasa.gov/overview.htm>

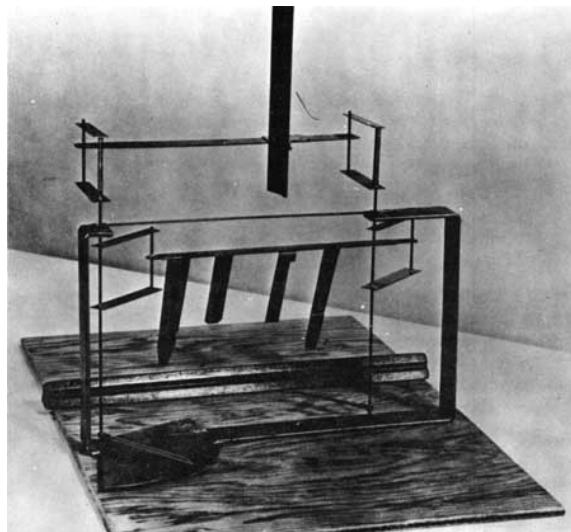


Figure 3.12: Reproduction of lift balance used in 1901 wind tunnel; model airfoil in testing position. Source: NPS.

However, a limitation of this approach is that the target system must be dividable in different subsystems. Similarly, the results of valid isolating models might not look like the real world phenomena simply because the latter is a combination of effects. On the other side, the model represents a single effect in isolation. Therefore, the validation of isolated models is problematic.

3.3.8 Summary of the different model views

The previous views on scientific models reflect different epistemic goals:

- Analogy: Understanding and conceptual exploration.
- Isolation: Causal inference and mechanism testing.
- Mirror: Prediction and faithful representation.

Real scientific practice usually combines all three: a model resembles (analogy), simplifies (isolation), and represents (mirror) reality, but never completely. Analogical models appear in AI when we use biological metaphors (neural networks, genetic programming). Isolating models appear in causal inference, A/B testing, and simulations. Mirror models dominate in predictive ML like fitting data to replicate the world statistically, or in digital twins. Still, as in philosophy of science, **no data-science model truly mirrors reality: it abstracts, compresses, and interprets through assumptions and training data.**

- All models are useful distortions, not exact copies.
- The purpose of a model (understanding, explanation, prediction) determines

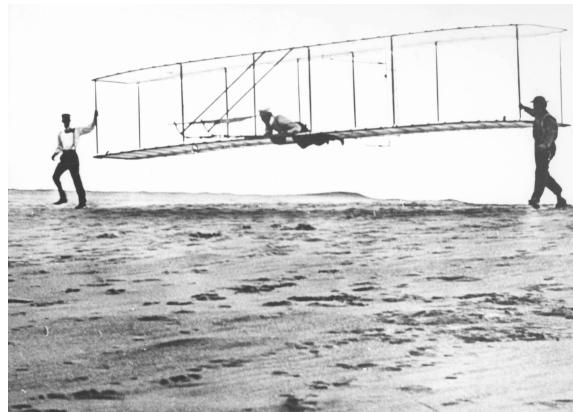


Figure 3.13: Historic photo of the Wright brothers' third test glider being launched at Kill Devil Hills, North Carolina, on October 10, 1902. Wilbur Wright is at the controls, Orville Wright is at left, and Dan Tate (a local resident and friend of the Wright brothers) is at right.

how it should be evaluated.

- In data science, confusing isolation (causal explanation) with mirror-like prediction leads to common errors, e.g., assuming accurate predictions imply causal truth.
- Awareness of these distinctions helps data scientists to design, interpret, and communicate models responsibly.

View	Main Idea	Meaning	Examples	In Data Science	Take-away Message
As Analogies	Models work by resemblance: they are like the target system in relevant aspects, but not identical.	Models help understand unknown systems by comparing them with something familiar.	Bohr's atom system (electrons orbiting like planets).	Neural networks inspired by brain architecture; genetic algorithms mimic evolution.	Analogies guide hypotheses, but may also mislead if taken too literally.

View	Main Idea	Meaning	Examples	In Data Science	Take-away Message
As Isolations	A model is a simplified, idealized system that isolates mechanisms or variables while ignoring others.	Reality is too complex to study directly; we deliberately remove factors to expose causal relationships.	Planes in physics; predator-prey models in ecology.	Isolate variables in simulations; A/B tests; causal DAGs to infer mechanisms.	Useful simplifications show causal structure but risk missing contextual effects (external validity).
As Mirrors	A model aims to mirror the real world faithfully to capture its structure accurately.	Aims for the ideal of reflecting reality as it is, through accurate description and measurement.	Detailed climate or population models; molecular dynamics simulations; digital twins.	Predictive ML models aiming to replicate observed data patterns.	Pure mirroring is an illusion: all models are partial and theory-laden . What matters is fitness for purpose , not perfection.

3.4 A tool for scientific reasoning

Course Note:

The following content is important for the mid-term exam. More examples will come soon.

This tool should help you spotting what assumptions were made, how evidence was connected to theory, whether the conclusions follow logically. This program or tool, comes from the book “Understanding scientific reasoning” from Ronald N. Giere,

John Bickle and Robert F. Mauldin ([Giere et al., 2006](#)). The following schema is a tool for evaluating whether a scientific model seems to fit the world, doesn't fit, or remains uncertain. It is a systematic way to identify the main components:

- What is being studied.
- What is the proposed model/theory.
- What prediction is made.
- What data was collected.
- Does data and prediction agree?
- What the agreement/disagreement does it mean for the model/theory?

The tool reduces the process of evaluating reports involving theoretical models to six steps. It is divided in 2 parts. (1) Understanding the episode (Steps 1–4) and (2) Evaluating the hypothesis (Steps 5–6).

- **Step 1: Real-world.** Identify what part of reality the study is about. Describe the phenomenon in ordinary or broadly scientific terms, not using words that belong to a specific model. For instance: “Turtles migrating between Brazil and Ascension Island.”
- **Step 2: Model.** Identify the theoretical model being tested. Explain it briefly, how it represents or explains the real-world phenomenon. For example: “Turtles began migrating 40 million years ago when the island was near the coast and kept going as it drifted away.”
- **Step 3: Prediction.** State what the model predicts should be observed if it truly fits the real world. This links theory to observation or experiment, e.g., “Genetic data should show the turtle population is about 40 million years old”.
- **Step 4: Data.** Identify the actual observations or experimental results reported, e.g., “DNA analysis shows the turtles have been distinct for only tens of thousands of years”.
- **Step 5: Negative Evidence?** Do the data and prediction agree? If no then, negative evidence, the model likely does not fit the real world. Otherwise, step 6.
- **Step 6: Positive Evidence?** Would the prediction likely be true even if the model were wrong? If no, then positive evidence: the data support the model. For instance, “Subjects guessed 90 % correctly. This is unlikely by chance. Then, supports model”. If yes, then other plausible models could also explain the data (inconclusive). For instance, “Warm decade could occur naturally. then, our data is inconclusive”.

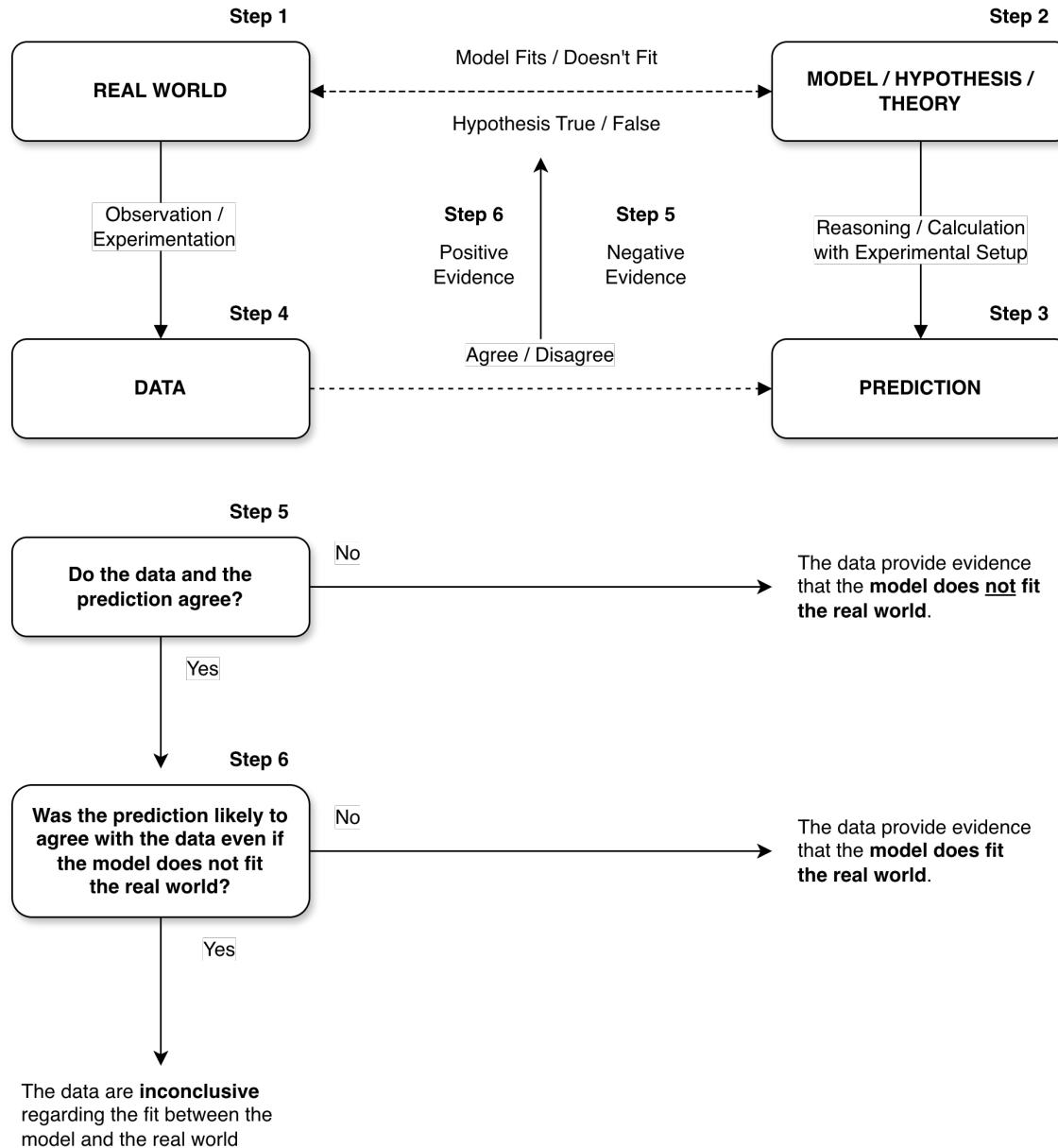


Figure 3.14: Reasoning tool extracted from ([Giere et al., 2006](#)).

This yields 3 possible conclusions.

- **Case 1. Negative evidence:** (Disagreement) When data and prediction disagree, it usually means the model fails to represent reality accurately. Step 3 ensures the prediction really follows from the model. Step 5 compares prediction and data. If they conflict, something is wrong with the model. Therefore we reject the model.
- **Case 2. Positive evidence:** (Agreement with no plausible alternatives) When data and prediction agree, Step 6 asks whether that agreement could happen even if the model were false. If no other plausible model would produce the same result, the agreement counts as good evidence that the model fits the real world. Then, we support the model.
- **Case 3. Inconclusive evidence:** (Agreement but plausible alternatives exist) If data and prediction agree, but other plausible models could also explain the same data, then the agreement does not tell us which model is correct. The model is not disproven, but also not confirmed. Therefore, our data is inconclusive regarding the fit of the model and the world.

Below, three examples are given that cover each of the cases aforementioned.

Case	Topic	Data vs. Prediction	Result	Evidence Type
1	Turtle migration theory	Disagree	Model rejected	Negative
2	Unconscious perception	Agree, unlikely by chance	Model supported	Positive
3	Greenhouse effect	Agree, but alternative plausible	Model uncertain	Inconclusive

3.4.1 A Case of Negative Evidence. “Gene Analysis Upsets Turtle Theory”

Researchers have concluded from genetic analyses that a widely accepted theory explaining why green turtles migrate 1250 miles to an island in the middle of the Atlantic Ocean to lay their eggs and then swim back to Brazil is invalid.

According to the popular theory, the turtles started coming to the island more than 40 million years ago, when it was close to the shore of South

America, and just kept coming as the island moved farther and farther away. But the new research concluded that the turtles had been using the island for only a few tens of thousands of years.

The now-invalid explanation was advanced 15 years ago by Archie Carr and Patrick Coleman. It was a bold idea, based on knowledge that the Atlantic Ocean was born 70 million years ago and began spreading. This gradually increased the isolation of Ascension Island, formed by volcanic activity on the ocean's centerline.

As long as the Atlantic was narrow, according to this hypothesis, remote ancestors of today's green turtles had no trouble reaching the island and laying their eggs in its beaches. They then returned to the shallow waters along the Brazil coast to feed on its marine grasses. After 40 million years, the ocean had grown to substantial width, but, according to the hypothesis, the turtles continued to reach the island by some mysterious form of navigation that still enables them to find it.

Now three scientists have examined the extent to which genetic material in the Ascension Island turtles has changed from that of the same species elsewhere. The difference, they believe, is far too small to have evolved over 40 million years. But it is sufficient to show that the Ascension Island turtles are a distinctive group from those nesting at other Atlantic sites, having probably used the island for not more than a few tens of thousands of years.

A comparison was made of a specific locus (the mitochondrial DNA) in the genetic material of turtles from four widely separated regions of Earth. Earlier work had indicated how fast, in a particular population, subtle changes in this material occur over centuries or millions of years. In the 1987 nesting season, eggs, or turtle hatchlings, were taken from twelve nests on French Frigate Shoal in the Hawaiian Islands, ten on Hutchinson Island off Florida, eight on Aves Island off Venezuela, and sixteen on Ascension Island in the mid-Atlantic.

The turtles in Hawaii had presumably been isolated from the other locations since the Isthmus of Panama formed about 3 million years ago. Their DNA, as expected, was most distinctive and provided an index of how fast turtle DNA becomes modified.

— Extracted from ([Giere et al., 2006](#))

- **Step 1. Real World.** Green turtles migrate between the coast of Brazil

and Ascension Island, 1,250 miles away. Scientists want to know how this long-distance migration began.

- **Step 2. Model.** The accepted theory held that turtles began nesting on Ascension Island about 40 million years ago when it was close to South America and continued as the island drifted outward.
- **Step 3. Prediction.** If that model fits the world, genetic analysis should show the Ascension turtle population to be roughly 40 million years old.
- **Step 4. Data.** DNA comparisons showed that the population has been distinct for only tens of thousands of years.
- **Step 5. Negative Evidence?** Yes. The data (40 thousand years) disagree sharply with the prediction (40 million years).
- **Step 6. Positive Evidence?** Not reached. The disagreement already shows that the model does not fit the real world.

In this example, because the data and prediction clearly do not match, the evaluation ends at Step 5. We already have enough to conclude the model do not fit reality. It is also important to note that the prediction is not directly “the turtles are 40 million years old,” but rather that the DNA data should indicate an age consistent with that. The prediction is always about what kind of observable data would appear if the model were correct.

3.4.2 A Case of Positive Evidence. “New View of the Mind Gives Unconscious an Expanded Role”

For decades, mainstream research psychologists suppressed the notion that crucial mental activity could take place unconsciously. But now, in the wake of exciting new studies, experimental psychologists are taking the unconscious more seriously. Among the most influential of the new studies are investigations into the role of the unconscious in the visual perception of objects and words.

One of the main researchers in this new area is Dr. Anthony Marcel of Cambridge University. He has developed a model of unconscious perception in which the unconscious mind perceives and remembers things of which the conscious mind is unaware. One of the most impressive tests of this model involves what is called “unconscious reading”.

In these experiments, Dr. Marcel flashes a word on a screen for a very short time. In addition, the word of interest is “masked” by being surrounded with other nonsense words such as esnesnon. When asked directly, the subjects were unable to say what real word appeared on the

screen. Dr. Marcel then asked his subjects to guess which of two words looks like the masked word. For example, the masked word might be blood and the two choices for look-alikes might be flood and week. The subjects were correct in their guesses an astonishing 90 percent of the time.

— Extracted from ([Giere et al., 2006](#))

- **Step 1. Real World.** Human perception and memory. Namely, whether the unconscious mind can perceive things the conscious mind misses.
- **Step 2. Model.** Dr. Anthony Marcel's model proposes an unconscious component that perceives and remembers stimuli invisible to consciousness.
- **Step 3. Prediction.** If this model is right, subjects should correctly identify or “guess” masked words more often than chance allows.
- **Step 4. Data.** In “unconscious reading” experiments, subjects chose the correct word 90 percent of the time even though they could not consciously report it.
- **Step 5. Negative Evidence?** No. The prediction and the data agree.
- **Step 6. Positive Evidence?** Yes. It is highly unlikely that random guessing could yield 90 percent accuracy. Therefore, the data provide good evidence that the model fits the real world.

The model in this case is not described in much detail, and the connection between theory and prediction is somewhat loose. Normally that would make the outcome uncertain. But because the data are so striking (subjects guessing correctly 90% of the time), it is hard to imagine any other reasonable explanation besides unconscious perception. That is why this example still counts as positive evidence, even though the model itself is only roughly defined.

3.4.3 Inconclusive Data. “Was That a Greenhouse Effect? It Depends on Your Theory”

The memorably uncomfortable summer of 1988 left many Americans with a suspicion that nature is at last getting even for mankind's wanton pollution of the atmosphere. From California to the Carolinas, the summer's heat wave and drought took a sobering toll. Electric power faltered, vast forests went up in flames, river navigation was throttled, and crops failed.

The “greenhouse effect”—the trapping of solar heat by pollutant gases in the atmosphere—became a household phrase. Some climatologists

warned that unless we quickly mend our ways, the world's grain belts will turn to dust bowls, coastal regions will be flooded, forests will die, and countless species will become permanently extinct. On June 23, Dr. James A. Hansen of the National Aeronautics and Space Administration caught the nation's attention when he told a Senate committee that the warming trend almost certainly stems from the greenhouse effect. A crisis, he warned, may not be long in coming.

But forecasting climate has never been as straightforward as scientists could wish. Many are not even sure that the summer's weather was really symptomatic of any trend at all. A. James Wagner, an analyst at the Weather Service, acknowledges that during this decade the world has seen the four warmest years of the past century—1980, 1983, 1986, and 1987. "But I do not feel that the evidence is overpowering that this is anything more than a normal fluctuation," he said.

Climatologists have invented several models in an attempt to understand fluctuations in the weather. One such model, which seems to mimic the real climate quite realistically, was devised by Dr. Edward Lorenz of the Massachusetts Institute of Technology. This model, which does not take carbon dioxide into account but does reckon on the interactions of the atmosphere with the ocean, exhibits large variations.

"The Lorenz model was run backward on a computer for the equivalent of about 400 years," Mr. Wagner said, "and the large fluctuations it sometimes produced, which were not entirely random but were not cyclical either, were quite startling." The swings, he said, were as much as $\pm 3.6^{\circ}\text{F}$ ($\pm 2.0^{\circ}\text{C}$) in global temperature from one year to the next. The model sometimes produced clusters in which several years close together were unusually hot—a pattern imitating the real climate of the 1980s.

— Extracted from ([Giere et al., 2006](#))

- **Step 1. Real World.** Earth's climate, especially the unusually hot, dry summer of 1988 and the generally warm 1980s.
- **Step 2. Model.** The greenhouse model: rising CO₂ traps heat in the atmosphere, raising global temperatures.
- **Step 3. Prediction.** If the model fits, average global temperatures should be increasing.
- **Step 4. Data.** Recent decades included the four warmest years of the century; 1988 was exceptionally hot and dry.
- **Step 5. Negative Evidence?** No. The data and the prediction appear to

agree.

- **Step 6. Positive Evidence?** Unclear. An alternative, the Lorenz climate model, can also produce large, random-like temperature swings without extra CO₂. Therefore, the data are inconclusive regarding the greenhouse model's fit to the real world.

Here, the data seem to support the greenhouse model, but since another credible model (Lorenz's) explains the same warming trend, the results do not clearly prove anything. The evidence is inconclusive about which model best fits the real world. Also, it is worth remembering that not all numbers count as "data". The ±3.6°F (±2.0°C) change came from computer equations — a prediction, not an observation. So it does not serve as direct evidence.

3.5 Examples

3.5.1 Willow tree experiment

Jan Baptist van Helmont (1580 - 1644) was a chemist, physiologist, and physician from the Spanish Netherlands (current Belgium). In 1634, Jan Baptist was arrested by agents of the Spanish Inquisition for the crime of studying plants and other natural phenomena. During his house arrest, he studied how plants grew. The prevailing theory at the time, stated that **plants grew by eating soil**, and Jan Baptist conceived an experiment to test this idea. Such prevailing theory has its origins in the ancient Greeks, from the Aristotelian view that all matter is composed of four terrestrial elements: earth, water, air, and fire.

Jan Baptist started by weighing a small willow tree (2.28 kg) and then weighed the dry soil (90 kg) in which he planted the tree. To prevent the dust from the surrounding air from mixing with the earth, the rim of the pot was protected, covered with a sheet of iron covered with tin and pierced by many holes. The tree was watered with rainwater or (when necessary) with distilled water. Jan Baptist left the tree for five years. After the five years had passed, Jan Baptist re-weighed the tree, which weighed 169 pounds (about 77 kg). He also re-weighed the dried soil and found the same 200 pounds (90 kg) minus about 2 ounces (56 gm).



Portrait of Jan Baptist van Helmont by Mary Beale, c.1674.

He wrongly concluded that the mass gain of the tree was produced by the water, which was his only intervention on the system. Although the

experiment was carefully conducted, the conclusions derived from the experiment were wrong because the theory on which it was based was incorrect. Importantly, the fact that Helmont used soil contradicted his hypothesis that only water was needed for plant growth ([Hershey, 1991](#)).

Jan Baptist did not know anything about photosynthesis. During the photosynthesis, carbon from the air and minerals from the soil are used to generate new plant tissue. Ironically, Helmont has been credited with discovering carbon dioxide ([Hershey, 1991](#)).

The employment of the balance during van Helmont's experiment was an important improvement; Jan Baptist believed that the mass of materials had to be accounted for during the study of chemical processes. **This experiment is considered the first quantitative experiment in plant nutrition.** It is also a great example of how firm conclusions can be misled by lack of knowledge of the studied system. Jan Baptist failed to control for an important background factor.

Therefore, **even rigorous observation can lead to false conclusions if the conceptual framework or background theory is wrong.** Observation is theory-laden: what one sees and how one interprets data depends on one's conceptual assumptions. But scientific progress is often iterative and self-correcting. Scientific change involves conceptual transformation, not just new data (see Thomas Kuhn).

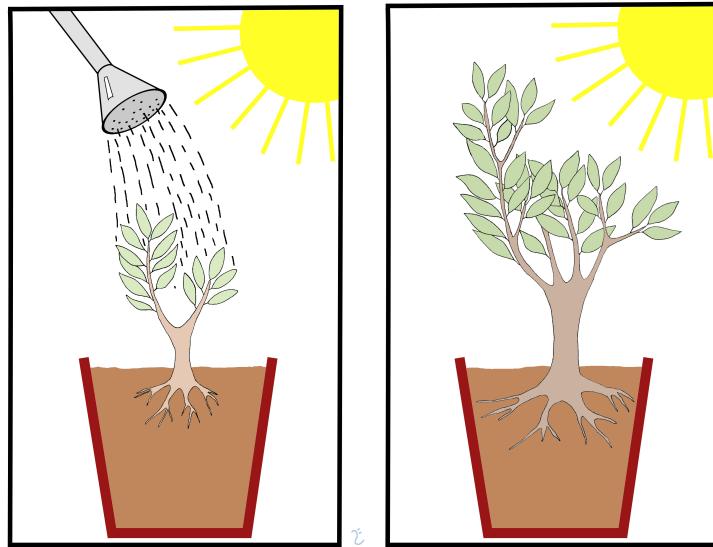


Figure 3.15: Illustration of the willow tree experiment. By Lars Ebbersmeyer⁹.

Importantly, the same data (soil mass unchanged, tree mass increased) could fit multiple theories, such as “Water is transformed into wood” (Van Helmont’s view) or

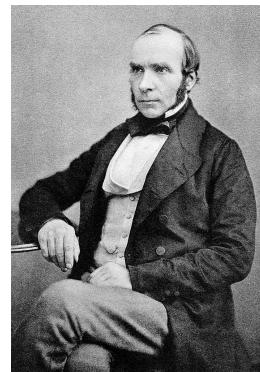
“Substance comes from air (CO_2)” (modern view). This is a case of underdetermination, data alone cannot uniquely determine theory. **In science, distinguishing empirical findings from theoretical inferences is crucial.** Van Helmont observed correctly that the soil’s mass did not change but his interpretation (all matter came from water) exceeded what the data justified. To properly differentiate between competing theories, we would need an independent experiment that controls for the previously untested variable (air). In other words, finding discriminating evidence.

This example contrasts with Semmelweis story, where he correctly identifying the causal factor (cadaveric contamination) despite lacking a background theory (germs theory) to justify it. Although his findings correctly solved the death issues in the maternity ward, his reasoning was empirically correct but theoretically unsupported at the time. Similarly, Van Helmont lacked an adequate theoretical framework, resorting to a naive inference. Still, **both cases anticipated a paradigm shift** in different areas such as germ theory and a new view of air composition chemistry and plant physiology.

3.5.2 1854 Broad Street cholera outbreak

The birth of epidemiology and public health is often attributed to the Natural Experiment described by Dr. John Snow in the mid-1800s when he investigated the relationship between drinking contaminated water and the incidence of cholera (Montelpare, 2021). The case of Dr. Snow is very popular in public health science and can be found in several books and posts on-line but I recommend the explanation given in Chapter 7 from The Book of Why (Pearl and Mackenzie, 2018) as the authors also re-formulate the case in causal terms, using concepts that were not available in the mid-1800s. Below, I give a summarised account of this study.

Miasma: “A vaporous exhalation formerly believed to cause disease.”



Dr. John Snow
(1813-1858), British physician.

Since the 1830s, various epidemics spread across Europe but were often attributed to social unrest and political upheaval. During the 1850s, there were several theories and misconceptions about the causes of cholera outbreaks. The miasma theory attempted to explain outbreaks of bubonic plague and cholera, stating that they were caused by a form of “bad air”. According to the competing theory, i.e. the germ theory of disease, the cause of the outbreak was a yet unknown germ. However, in 1853, disease-causing germs had not yet been observed under a micro-

scope and the germ theory of disease was not established yet. Louis Pasteur did not demonstrate the relationship between germ and disease until the 1860s. Therefore, the prevailing theory was that a miasma of unhealthy air caused cholera.

In August 1854, a major outbreak of cholera occurred in Soho, London (United Kingdom). In just three days, 127 people died in Broad Street. By September, 500 people had died. John Snow was skeptical of the prevailing miasma theory and theorized that the cause was the presence of an agent in the contaminated water source from certain water supplying companies ([Montelpare, 2021](#)). Namely, the Southwark and Vauxhall Company and the Lambeth Waterworks Company. The main difference between the two companies was that the former drew its water from London Bridge, which was downstream from London's sewers. Years earlier, Lambeth had moved its water intake so that it would be upstream of the sewers.

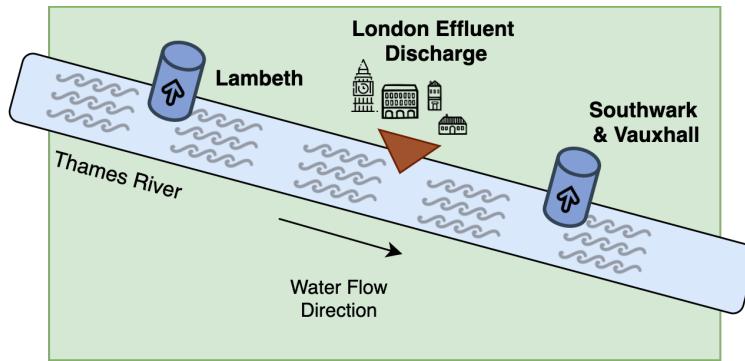


Figure 3.16: Water distribution by the Lambeth Water and the Southwark & Vauxhall Companies.

Therefore, the customers of the Southwark and Vauxhall Company were getting water contaminated by the excrements of cholera victims, whereas Lambeth users were drinking uncontaminated water.

In consequence, districts supplied by the Southwark and Vauxhall Company had a death rate eight times higher than other districts. At this point, the evidence supporting the hypothesis of water contamination is just circumstantial. The causal diagram from Figure 3.17 depicts the situation. A proponent of the miasma theory could argue that the effect of miasma was strongest in those districts ([Pearl and Mackenzie, 2018](#)).

Then, Snow noted that in some districts water was served by both companies (see Fig. 3.19), and even there, the death rate was still higher in the houses where water was supplied by the Southwark and Vauxhall Company. Those households did not showcase any difference in terms of poverty or miasma. Snow wrote: “each company supplies to both rich and poor, both large houses and small; there is no

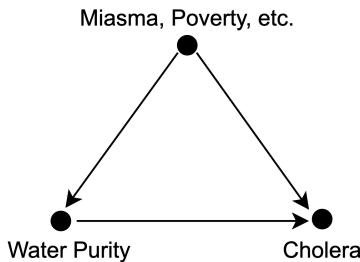


Figure 3.17: Causal diagram for cholera for the case of John Snow.

difference either in the condition or occupation of the persons receiving the water of the different companies.” ([Pearl and Mackenzie, 2018](#)).

Table 3.4: Relation of the household water source and deaths

Source of water	Deaths
Southwark and Vauxhall company	286
Lambeth company	14
Direct from the river	22
Pump wells	4
Ditches	4
Unknown	4

It is precisely at this point where the natural experiment takes all its strength. Around 300 people of both sexes, every age, and socio-economic class were naturally divided into two groups without them to know. One group received pure water, whereas the other received water mixed with sewage.

The observations of John Snow introduced a new variable into the causal diagram (see Fig. 3.18), the **Water Company**. In this new diagram we can see that there is no arrow between Miasma and Water Company, therefore both variables are independent. Moreover, we can note the presence of an arrow between Water Company and Water Purity. Finally, the diagram depicts a third assumption. There is no direct arrow from Water Company to Cholera, i.e. water companies do not deliver cholera to their customers in any other way.

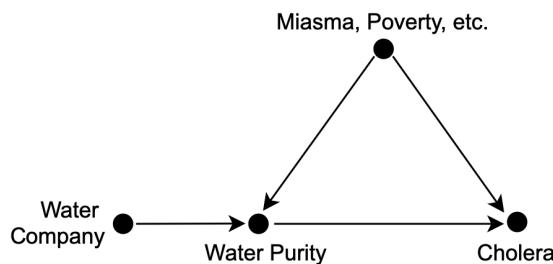


Figure 3.18: Causal diagram after the introduction of an instrumental variable.

A variable that satisfies these properties is called an **instrumental variable**. See

relation between Water Company and Cholera, any observed association must be causal. Similarly, the association between Water Purity and Cholera is also causal.

Specifically, an instrumental variable Z is an additional variable used to estimate the causal effect of variable X on Y. The traditional definition qualifies a variable Z as an instrumental (relative to the pair (X, Y)) if (i) Z is independent of all variables (including error terms) that have an influence on Y that is not mediated by X and (ii) Z is not independent of X. — (Pearl, 2000)

Although today miasma theory has been discredited, poverty and location are clear confounders. In (Pearl and Mackenzie, 2018), the authors show how instrumental variables can be used to determine the number of lives that could have been saved by purifying the water supply. The instrumental variable **Water Company** allow us to find the effect of **Water Purity** on **Cholera** even without being able to control, or collect data on, the confounder variables (poverty, location, etc.).

One of the main innovations of John Snow approach was to focus on death rates in districts served by two water companies rather than on data from victims of the Broad Street pump which drew water from a well.

A transitional period began in the late 1850s with the work of Louis Pasteur. This work was later extended by Robert Koch in the 1880s. By the end of that decade, the miasma theory was struggling to compete with the germ theory of disease. Viruses were initially discovered in the 1890s. Eventually, a “golden era” of bacteriology ensued, during which the germ theory quickly led to the identification of the actual organisms that cause many diseases. — Wikipedia, Germ theory of disease¹¹.

3.5.3 Causal models: Estimating treatment effect in the presence of a confounder

This example illustrates the use of a causal model and causal inference methods to estimate the effect of **study hours** (per week) in the final exam **scores**. Our outcome again is students’ scores. We aim to estimate the causal effect that the “*treatment*” **study hours** has on students’ **scores** (effect). Our beliefs and assumptions are encoded into a causal model. We assume that **study hours** (per week) is directly related to **scores** but also that the **prior knowledge** affects both the effect and the cause (i.e., it is a confounder). For this exercise we will employ Python

¹¹https://en.wikipedia.org/wiki/Germ_theory_of_disease

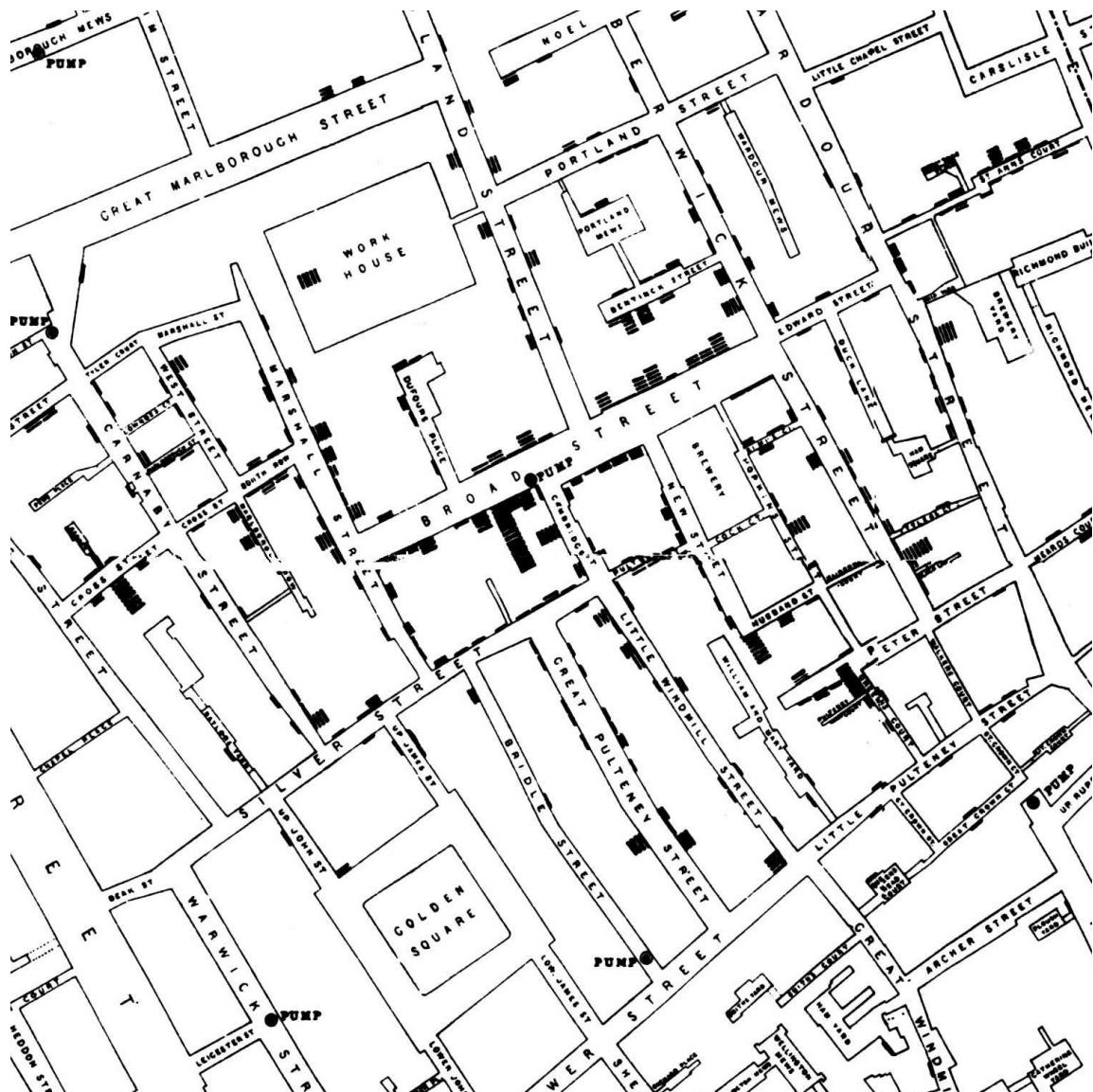


Figure 3.19: Detail of John Snow's map of Cholera in the Broad Street outbreak in 1854. Each bar represents one death in a topography that attempted to relate the water source ("Pump") to pattern of cases in the neighborhood outbreak.

language and the library doWhy¹².

This example shows how to estimate the effect while controlling for the confounder. Remember that a confounder Z affects both the cause X and the effect/outcome Y.

We import the necessary packages to work.

```
import matplotlib
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import dowhy
from dowhy import CausalModel
```

Then, we begin by creating a synthetic dataset in which study hours and prior knowledge are random variables normally distributed. Then, exam score is a function of both study hours and prior knowledge. In a real world scenario, we would normally not know about this function because is what we are trying to estimate. Note the values in the function, specifically, the `2 * data["study_hours"]`.

```
np.random.seed(42)
data = pd.DataFrame({
    "study_hours": np.random.normal(10, 2, 1000),
    "prior_knowledge": np.random.normal(3, 1, 1000),
    "exam_score": np.zeros(1000) # Placeholder for exam score
})

# Define the relationship:
# Exam score is influenced by both study hours and prior knowledge
data["exam_score"] = (5 + 2 * data["study_hours"]
                      + 3 * data["prior_knowledge"]
                      + np.random.normal(0, 1, 1000))

data.head(5)

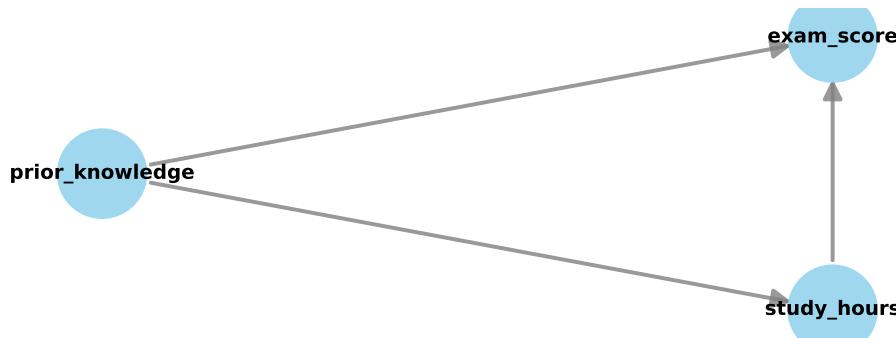
##   study_hours  prior_knowledge  exam_score
## 0      10.993428        4.399355  39.509745
## 1       9.723471        3.924634  36.076325
## 2      11.295377        3.059630  35.977225
## 3      13.046060        2.353063  37.843348
## 4       9.531693        3.698223  33.264442
```

¹²<https://www.pywhy.org>

Now, we define the model and visualize it.

```
model = CausalModel(
    data=data,
    treatment="study_hours",
    outcome="exam_score",
    common_causes=["prior_knowledge"])

model.view_model(size=(8,3));
```



```
identified_estimand = model.identify_effect()
estimate = model.estimate_effect(identified_estimand,
    method_name="backdoor.linear_regression")
```

The result below means that, according to our causal model, an additional hour of study per week is estimated to increase the exam score by approximately 2 points. This effect is derived after accounting for the confounder (in this case, `prior_knowledge`) using the backdoor criterion. If this estimate is accurate, it suggests a positive causal relationship between study hours and exam score, where more study hours lead to higher scores.

As can be seen, this estimation is very close to the degree stated in the equation of our synthetic data in `2 * data["study_hours"]`.

```
print("Causal Estimate of Study Hours on Exam Score:",
      estimate.value)

## Causal Estimate of Study Hours on Exam Score: 2.010902911319434
```

A placebo test in causal inference is a refutation technique designed to test the robustness of our causal estimate by using a “fake” treatment variable. The idea is to see if our model incorrectly detects a causal effect when there should be none,

which would indicate potential bias or flaws in the model. The placebo test helps us check if this estimated causal effect could be due to biases in the model rather than a true causal relationship.

```
refute_result = model.refute_estimate(identified_estimand,
estimate, method_name="placebo_treatment_refuter")
```

This test replaces the actual treatment (study hours) with a “placebo” or random variable that should have no causal effect on the outcome (exam score).

```
print(refute_result)
```

```
## Refute: Use a Placebo Treatment
## Estimated effect:2.010902911319434
## New effect:0.0
## p value:1.0
```

- New effect (0.0): This shows that the placebo treatment has no effect on the outcome, as expected (since it’s a random variable unrelated to exam score).
- Original effect (2.01) vs. New effect (0.0): The fact that the original effect remains 2.01 while the placebo effect is 0 suggests that the observed causal effect of study hours on exam score is not due to random bias. Instead, it’s likely a true effect.
- p-value (2.0): A high p-value indicates that there is no statistically significant effect with the placebo treatment. This further supports that any observed effect in the placebo scenario is likely due to chance.

The placebo test confirms that our causal model appears robust, as it detects no effect when a placebo treatment is used. Therefore, we can have more confidence that the estimated effect of 2.01 is a genuine causal effect of study hours on exam score, rather than an artifact of model bias or unmeasured confounding.

These results suggest that our initial model is likely valid, meaning that the relationship we estimated is trustworthy given the data and assumptions. But things could still be wrong... as we will see in the next example.

3.5.3.1 Backdoor criterion

The backdoor criterion is a method used in causal inference to control for confounders—variables that influence both the treatment and the outcome. By satisfying the backdoor criterion, we aim to isolate the causal effect of the treatment on the outcome by “blocking” other paths of influence. Suppose we have a causal

structure where a variable Z (like prior knowledge) affects both our treatment X (study hours) and outcome Y (exam score). The backdoor criterion tells us to condition on (i.e., adjust for) variables like Z which lie on “backdoor paths” that might introduce bias in estimating the effect of X on Y . By conditioning on Z , we effectively “block” any non-causal paths from X to Y , allowing us to observe the effect of X on Y without interference from confounding.

When we adjust for a variable like Z in causal inference, we’re essentially isolating the relationship between the treatment X (e.g., study hours) and the outcome Y (e.g., exam score) by removing the influence of Z (e.g., prior knowledge) on both X and Y . This adjustment is typically done in one of a few statistical ways, such as conditioning, stratification, or regression. Here a summary of them:

- Conditioning on Z : Conditioning involves looking at the relationship between X and Y within subgroups defined by different values of Z . By examining these subgroups, we can remove the influence of Z on X and Y . For instance, if Z represents prior knowledge levels (low, medium, high), we can separately look at the effect of study hours on exam scores within each of these prior knowledge groups. This way, any relationship between X and Y is not confounded by Z .
- Stratification by Z : Stratification divides the data into strata (groups) based on the values of Z , and then estimates the effect of X on Y within each stratum. The overall effect is a weighted average of the effect within each stratum. In our example, if Z (prior knowledge) has three levels, we could stratify our analysis into three groups and then calculate the effect of study hours on exam scores within each level of prior knowledge. By averaging these effects, we get a more accurate estimate that accounts for the confounding effect of Z .
- Regression Adjustment: Regression adjustment uses statistical models (usually linear or logistic regression) to “hold constant” the effects of Z on both X and Y . By including Z as a variable in the regression model, we can estimate the relationship between X and Y while “controlling” for Z .

Adjusting for Z removes any indirect associations between X and Y that arise because Z influences both. By blocking this confounding pathway, we effectively “close” the backdoor path and allow the model to estimate a more accurate causal effect between X and Y . In causal inference, this adjustment step is crucial because it attempts to mimic a randomized experiment where all confounders would ideally be balanced, allowing us to make stronger causal claims about X and Y .

3.5.4 Causal models: Detecting an overestimation

This example illustrates the use of a first causal model that is **WRONG** because it leaves out an important variable (either because we don't know about it or because we **believe** it does not have an effect on the outcome).

- Our **outcome** again is students' grades.
 - We aim to estimate the causal effect that the “*treatment*” private tutoring has on students' grades (**effect**).
 - We **believe** that **study hours** (per week) is the only variable that affects both tutoring and grades (i.e., it is a **confounder**)

We will see how to **detect that the model is wrong** and how once corrected it estimates the effect properly.

We generate a synthetic data with different random values. As you can see **parental involvement** is biased with a 70% probability to be 0 and 30% of being 1. In the last two lines we fill tutoring and grade as we have synthetically made them to depend on the study hours and parental involvement. Imagine that instead of being synthetically created this data was given to you from data collection (e.g. via surveys).

```
np.random.seed(0)

#generate synthetic data
data = pd.DataFrame({
    "study_hours": np.random.normal(5, 1, 1000),
    "parental_involvement": np.random.choice([0, 1],
                                              size=1000,
                                              p=[0.7, 0.3]),
    "tutoring": np.zeros(1000), # Placeholder for tutoring
    "grade": np.zeros(1000) # Placeholder for grade
})

# Assume tutoring is influenced by
# both study hours and parental involvement
data["tutoring"] = (data["study_hours"]
                    + 10 * data["parental_involvement"]
                    + np.random.normal(0, 1, 1000) > 6).astype(int)

data["grade"] = (50 + 10 * data["tutoring"]
                 + 5 * data["study_hours"]
                 + 10 * data["parental_involvement"])
```

```
+ np.random.normal(0, 2, 1000))

data.tail(10)
```

	study_hours	parental_involvement	tutoring	grade
## 990	5.441033		0	1 86.784749
## 991	5.178793		0	0 75.005558
## 992	4.200578		1	1 92.496516
## 993	5.240788		0	0 74.929059
## 994	5.289121		1	1 98.189120
## 995	5.412871		0	0 77.813355
## 996	4.801601		0	0 76.318058
## 997	5.094192		1	1 96.811745
## 998	3.852389		0	0 67.153054
## 999	4.641886		0	0 71.478162

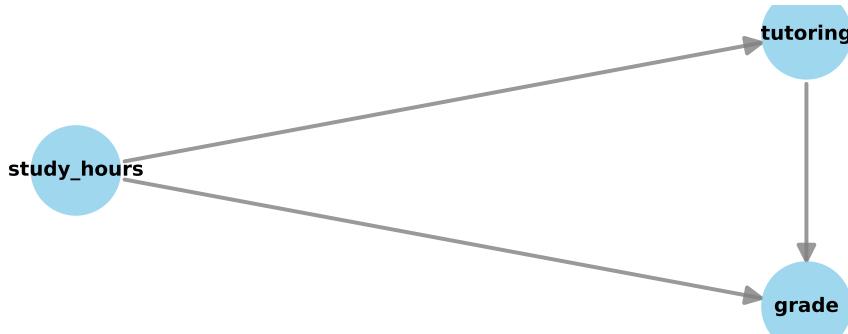
Next we create our **initial causal model** in which we **wrongly** assume that **parental involvement** does not affect anything else. Again, the reasons for this error could be two:

- We don't know about this variable. So we didn't even capture it to begin with (e.g. we didn't ask in the survey).
- We captured this variable (but we asked so many things in our survey !) but we believe not to affect the other variables.

Note that, whether unknown to us or wrongly assumed by us, variables can be affected by **parental involvement** in reality. The model shows that **study hours** is a **confounder** of both **tutoring** and **grade**.

```
# Create a causal model excluding parental involvement initially
model = CausalModel(
    data=data,
    treatment="tutoring",
    outcome="grade",
    common_causes=["study_hours"])

model.view_model(size=(8,3))
```



Our initial results show an estimated effect of 14.98 meaning that having `private_tutoring` is estimated to increase a student's final `grade` by approximately 17 points. This effect estimate was obtained after controlling for `study_hours` but without considering the additional confounder `parental_involvement`. At this point we could be naively happy with this result. But what if we are overestimating or underestimating the effect of `tutoring`?

```

identified_estimand = model.identify_effect()
estimate = model.estimate_effect(identified_estimand,
                                 method_name="backdoor.linear_regression")

print("Estimated Effect (initial model):",
      estimate.value)
  
```

```
## Estimated Effect (initial model): 17.585373682989484
```

Again, the placebo test below replaces the actual treatment (in this case, `tutoring`) with a “placebo” or random variable that has no relationship to the outcome (final `grade`).

```

refutation = model.refute_estimate(identified_estimand, estimate,
                                    method_name="placebo_treatment_refuter")
  
```

The placebo test here suggests that the original model's structure is robust, meaning that the observed effect of `tutoring` on `grades` is not due to random bias. However, this does not yield out the possibility that the obtained estimate could be biased due to an unmeasured confounder.

```

print("Placebo Test Result:", refutation)

## Placebo Test Result: Refute: Use a Placebo Treatment
## Estimated effect:17.585373682989484
  
```

```
## New effect:-0.05521038628415212
## p value:0.92
```

We could ask ourselves too whether the estimated effect changes significantly when we replace the given dataset with bootstrapped samples from the same dataset? (Hint: It should not).

```
bootstrap_refutation = model.refute_estimate(identified_estimand,
                                              estimate,
                                              method_name="bootstrap_refuter")
```

As seen below, the effect does not deviate much from the original estimated effect. Still, something else could be going on.

```
print(bootstrap_refutation)

## Refute: Bootstrap Sample Dataset
## Estimated effect:17.585373682989484
## New effect:17.60096969638309
## p value:0.8400000000000001
```

If we suspect about the presence of unknown confounders, we could test how sensitive is the effect estimate when we add an additional common cause (confounder) to the dataset that is correlated with the treatment and the outcome. (Hint: It should not be too sensitive).

```
sensitivity_analysis = model.refute_estimate(identified_estimand,
                                              estimate,
                                              method_name="add_unobserved_common_cause")
```

In this case we see a huge drop. This may indicate that we over estimated the effect of `tutoring` on the `grades` and that we are missing an important confounder.

```
print(sensitivity_analysis)

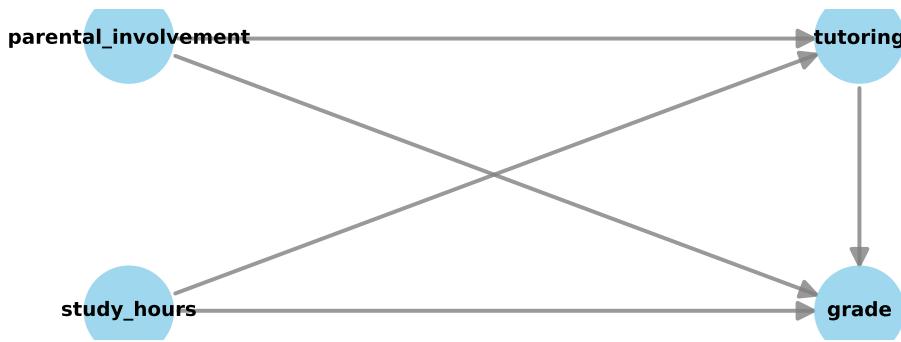
## Refute: Add an Unobserved Common Cause
## Estimated effect:17.585373682989484
## New effect:1.6491017911987171
```

At this point we should get back to the design table and see that we missed an important variable. In our case we knew already: `Parental involvement`. In real life, we would either ask the students again (through new surveys) or if we already had this variable in our surveys we would add it to the model as a confounder affecting the `grades` and the `tutoring`. Note that in “our reality” `parental involvement`

did not affect study hours directly.

```
# Realize missing confounder: Parental involvement
model_v2 = CausalModel(
    data=data,
    treatment="tutoring",
    outcome="grade",
    common_causes=["study_hours", "parental_involvement"])

# Visualize the causal graph
model_v2.view_model(size=(8,3))
```



We then re-run our estimation.

```
identified_estimand_with_confounder = model_v2.identify_effect()
estimate_with_confounder = model_v2.estimate_effect(
    identified_estimand_with_confounder,
    method_name="backdoor.linear_regression")
```

As can be seen, this value is closer to the degree of our synthetic data $10 * \text{data}["tutoring"]$. Still, in real world scenarios we do not necessarily know the functions that describe the behaviour of real world phenomena.

```
print("Estimated Effect (revised model):",
      estimate_with_confounder.value)

## Estimated Effect (revised model): 10.095828768814414
```


Chapter 4

Experimental Control and Statistical Abuse

4.1 Overview

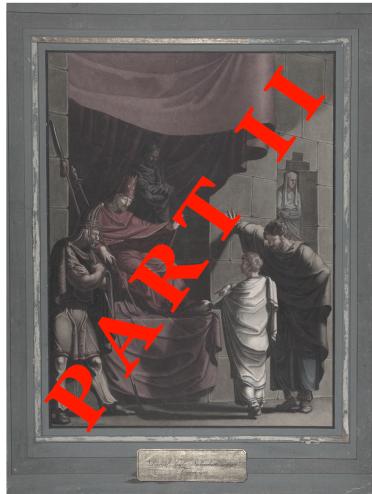


Figure 4.1: Daniel interprets Nebuchadnezzar's Dream.

Chapter 3 started with an account of an experiment from the times of Babylonian King Nebuchadnezzar, in which Daniel proposed to feed a group with a vegetable diet whilst another group continued eating the King's royal meat. However, Daniel did not think of everything when designing his experiment. He did not take confounding bias into account. For instance, Daniel and his friends could have been healthier than the control group. Under such supposition, their strong appearance after ten days of a vegetarian diet may have nothing to do with the diet itself. Perhaps, they

would have become even stronger if they had eaten the meat from the king. As we have seen in previous examples presented in this book, confounding bias happens when a variable influences both who is chosen for the treatment group as well as the experiment outcome. These variables might be known variables or act as a lurking third variable we are not aware of. Such variables are easy to spot in causal diagrams.

The term “*confounding*” means “*to pour, mingle, or mix together*”, and Figure 4.2 illustrates why such name was chosen to denominate these situations. The true causal effect $X \rightarrow Y$ is *mixed* with the spurious correlation between X and Y induced by the fork $X < -Z -> Y$ (Pearl and Mackenzie, 2018).

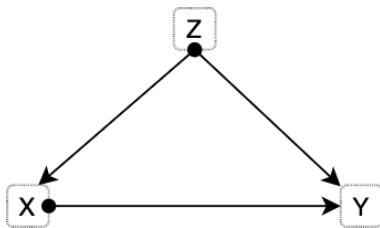


Figure 4.2: Most basic version of confounding situation: Z is a confounder of the proposed causal relationship between X and Y .

Course objectives

At the end of this chapter, most objectives will have been strengthened, with this chapter focused on wrapping up objectives 5 to 8. Examples from natural and social sciences help illustrate control techniques, which are then connected to data science. This fosters discussion with the students, motivating them to recognise similar issues in data science and understand techniques to mitigate them.

4.2 The smoke debate - Part II

Many famous cases present these confounding situations. For instance, this book previously tackled the debate around smoking and lung cancer. Austin B. Hill and Richard Doll noticed that hidden biases could be present in their previous case-control studies, and the replication of the studies would not be enough to overcome them. In consequence, they began a prospective study (1951) in which they considered 60.000 physicians from United Kingdom consisting of questionnaires tackling their smoking habits. These physicians were followed over time. In just five years, heavy smokers showed a death rate from lung cancer 24 times higher than non-smokers. A similar study conducted in the United States showed that smokers died from lung cancer 29 times more often than non-smokers while heavy smokers died



Figure 4.3: 1946 cigarette advertisement launched by R.J. Reynolds Tobacco Company. Source: Tobacco Ads¹.

90 times more frequently. However, former smokers reduced their risk by a factor of two. This behaviour is often called the “dose-response effect”, indicating that a prolonged dose of a drug causes a stronger response.

Still, R. A. Fischer and Jacob Yerushalmy remained sceptical, stating that such prospective studies failed to compare smokers to non-smokers, arguing that they were not identical groups. The rationale of the critic is that smokers in the study are self-chosen. Moreover, there might be a constitutional difference between smokers and non-smokers. For instance, smokers might be more risk-taking, or more prone to be alcoholics, which might cause adverse health effects which are then wrongly attributed to smoking by Hill and Doll studies. Another possibility they appealed to is the existence of a smoking gene that caused people to become smokers and made them more likely to develop lung cancer.

The *constitutional hypothesis* was almost impossible to test. In 2000, the sequencing of the human genome became real and with it the possibility to study links between genes and lung cancer. Actually, such genes do exist, as with breast cancer that make people more prone to develop certain types of cancer. In 1959, a couple of researchers published a rebuttal of Fischer’s arguments that settled the debate. One of the researchers, Cornfield, was not a statistician, nor a biologist, but instead, a historian with statistical knowledge who worked in the department of agriculture (this is of course not a cause of his family name). Cornfield aimed to debunk such constitutional hypothesis with the following reasoning: suppose the possibility of a confounding factor (e.g. smoking gene) that would fully explain the cancer risk of smokers. If smokers have 9 times the risk of developing lung cancer, the supposed

confounding factor ought to be at least nine times more common in smokers to account for such risk difference. Let's exemplify this. If 11% of non-smokers have such a gene, then 99% (since they have 9 times more risk: 11×9) of smokers would have to have the gene. But if 12% of non-smokers would have the smoking gene, then it is not mathematically possible for the cancer gene to fully explain the association between smoking and lung cancer. This is known as Cornfield's inequality, and led to the development of sensitivity analysis.

The above's explanation shows that the association between smoking and lung cancer is too strong to be explained by appeal to a smoking gene (or any other constitutional hypothesis). In essence, Cornfield's rationale gives us a way to choose between both causal diagrams. Once it becomes evident that such constitutional hypothesis is not able to fully explain the association, the relationship between smoking and lung cancer (left diagram) becomes apparent.

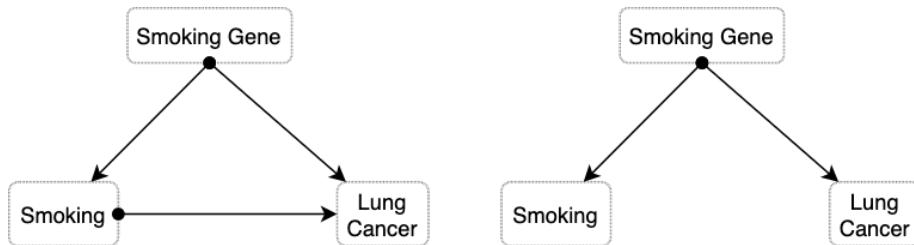


Figure 4.4: The causal diagram on the left presents the situation in which the constitutional hypothesis is insufficient to explain the association between smoking and lung cancer. The diagram on the right side depicts the alternative situation in which the smoking gene fully explains the observed association.

The tobacco industry magnified any bit of controversy they could find on the scientific studies. Such organised denialism explains why the link between smoking and cancer remained so controversial in the public long after the debate was settled among epidemiologists.

Remarkably, even researchers at the tobacco companies were convinced—a fact that stayed deeply hidden until the 1990s, when litigation and whistle-blowers forced tobacco companies to release many thousands of previously secret documents. In 1953, for example, a chemist at R.J. Reynolds, Claude Teague, had written to the company's upper management that tobacco was “an important etiologic factor in the induction of primary cancer of the lung,” nearly a word-for-word repetition of Hill and Doll’s conclusion. — The Book of Why (Pearl and Mackenzie, 2018).

4.3 Experimental Control

Experimental control entails a series of procedures for experiment and observation design aimed at minimising the effects of extraneous variables (i.e. confounding factors) other than the manipulated variables (i.e. independent variable) to ensure that the measured variable (i.e. dependent variable) is only affected by the independent variables. To evaluate the effects of manipulating the independent variables, some control system is needed in which no such deliberate changes are introduced. As we have seen, sampling units (e.g. study participants) are often divided into two groups (the experimental group and the control group) in a way that the only noticeable (or significant) difference between them lies in the stimuli exerted by the experiment. Therefore, the control and experimental groups must be *homogeneous* in all relevant factors.

In general, there are two techniques for the formation of such homogeneous groups: individual and collective control (Bunge, 2017). **Individual control** requires simultaneous pairing of individuals in both groups, i.e. every member of the experimental group has a corresponding equivalent member in the control group. For instance, for every thirty years old Asian man in the control group another thirty years old Asian man is assigned to the experimental group. Simultaneous pairing is complex and expensive. **Statistical control** has two main types. On one side, the *control of distributions* should be performed to equate certain parameters such as averages, spreads (i.e. std. dev.) and other collective properties (e.g. medians). This technique is more flexible as only some properties are kept under control. In this case, we would take two samples of people with the same age and height distributions. Both simultaneous pairing and distribution control share a common disadvantage regarding the formation of the groups, which could be unintentionally biased. For instance, we could assign the strongest people to the treatment (or experimental) group to make sure they bear the treatment. To prevent this issue the two groups are usually formed at random. Thanks to **randomisation**, all variables (including most unknown factors) that were not previously controlled become randomly distributed, minimising their effect on the dependent variables. However, randomisation is not an alternative to other techniques, but rather a complement.

4.3.1 Other experimental control techniques

There are multiple strategies for experimental control. We have previously seen the method of division into treatment and control groups. The control and treatment groups can entail two moments in **time**, with the initial setting being the control scenario which is later on manipulated through the intervention of certain variables

(e.g. measure noise from bats in a dark chamber before and after turning a light). Another technique requires **holding certain factors constant** or finding scenarios (like in a field experiment) with the same background conditions. Nonetheless, constructing such conditions in a laboratory can also achieve this goal. In an **elimination** strategy some factors are removed to simplify study conditions, such as air resistance in a vacuum chamber or drop tower², radio waves in a Faraday cage, or gravity in space experiments. A common case of elimination is **blinding**, where subjects do not know which group they are assigned to (single blinding). Moreover, double-blinding implies hiding this information from the experimenter and/or the data analyst. Finally, we can **separate factors** by measuring their effect and correcting for it. For example, the measurements of time dilation require taking into account the Doppler effect caused by the changing distance between the observer and the moving clock. GPS systems perform adjustments due to the effects of time dilation and gravitational frequency shifts. Another example, missile trajectories are often adjusted for the effect of Coriolis force³.

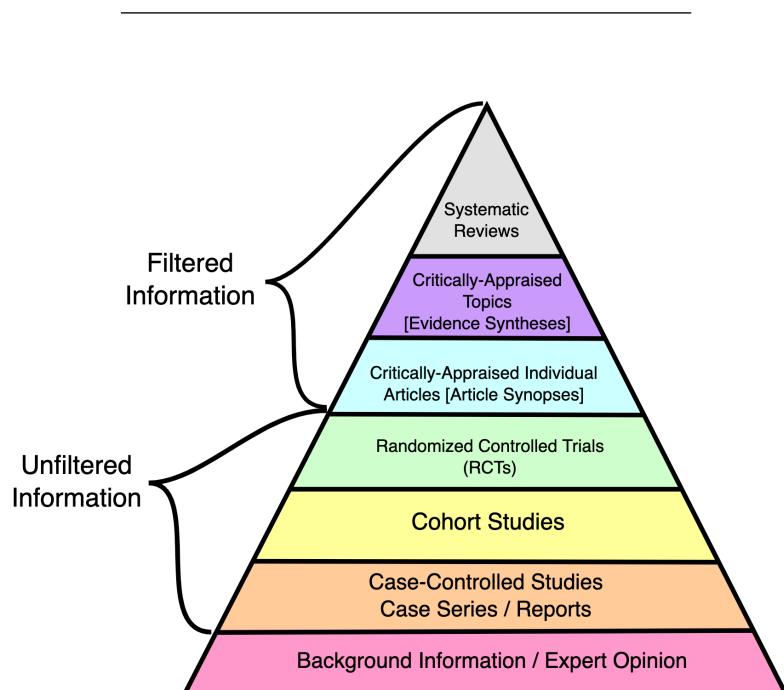


Figure 4.5: Figure from Wikimedia⁴ by CFCF. Keep in mind that this hierarchy is not free from criticism⁵ and take it just as a useful simplification.

²https://en.wikipedia.org/wiki/Drop_tube

³https://en.wikipedia.org/wiki/Coriolis_force

4.4 Randomised Control Trials

The fundamental problem of causal inference tell us that it is impossible, by definition, to observe the effect of more than one treatment on a subject over a specific time period. A study participant cannot both take the pill and not take the pill at the same time. Directly observing causal effects is impossible. Nonetheless, this does not make causal inference impossible. There are certain techniques and assumptions that allow to circumvent the fundamental problem. In this context, randomized experiments allow for the estimation of population-level causal effects.

Randomisation offers a systematic solution for the division of participants (or sampling units) into two groups. In particular, RCTs are frequently regarded as a gold standard for clinical trials and among the highest quality evidence available (see Figure 4.5). However, as with every method, it will only yield fruitful results if applied correctly, and its sole employment does not warrant against other errors.

There are different types of randomisation. In **simple randomisation**, subjects are assigned into two groups purely randomly but in small samples, we risk creating uneven groups. **Block randomisation** works by randomising participants within blocks such that an equal number are assigned to each treatment. For example, given a block size of 4, there are 6 possible ways to equally assign participants to a block (AABB, ABAB, ABBA, BAAB, BABA, BBAA). Allocation proceeds by randomly selecting one of the orderings and assigning the next block of participants to study groups according to the specified sequence. A major disadvantage of this method is that it might be possible to predict the next sequence. **Stratified randomisation** is crucial whenever all other properties (except for the factors of interest) need to be assigned equally. The study population is first stratified into subgroups (i.e. *stratas*) sharing attributes, then followed by simple or block random sampling from the subgroups.

One of the main advantages of RCTs is the reduction of selection bias or allocation bias. In Chapter 4 we will see biases in more detail. The randomisation process reduces mistrust towards a potential rigged distribution of the participants. Another common advantage is that it facilitates blinding the groups from investigators and participants.

Terminology Note:

Very often terms are used interchangeably in many domain but they can also mean different things depending on the are.

By “allocation bias” we understand the bias caused by allocating patients with better prognosis to either the experimental or the control group. In the context of a randomized trial the term “selection bias” is sometimes used instead of allocation bias to indicate selection of patients into treatment arms. We avoid the term “selection bias” as it has a different meaning in epidemiology more broadly: selection of non-representative persons into a study. — ([Paludan-Müller et al., 2016](#))

However, RCTs do not necessarily ensure that background factors are equally distributed in the treatment and control groups. For small samples randomisation can provide unequal distributions. The average number after rolling a dice an infinite amount of times will converge to 3.5, but we should not be surprised if we roll a dice 10 or 20 times obtaining considerably more occurrences of the number 6 than the other numbers. The danger of relying on pure randomisation to balance covariates has been described in ([Krause and Howard, 2003](#)) ([Morgan and Rubin, 2012](#)). For this reason is essential to check for imbalances in known factors after randomisation. Stratified randomisation also helps balancing known factors. Nonetheless, randomisation does not necessarily guarantee full control of unknown factors but *on average* their effect should be significantly smaller than the treatment applied ([Deaton and Cartwright, 2018](#)).

When we use an RCT to evaluate an intervention, we do so with respect to one or more endpoints (or outcomes) that will be measured in the future, after the period of intervention. It could be blood pressure, death, quality of life, etc. We want to understand the causal effect of the intervention on that outcome, but this is tricky. That’s because to really understand the effect of the intervention, we would need to give it to someone and measure the outcome to see what happened. Then we would need to reset the universe back to the exact point when the intervention was given, withhold it this time, and see what happened when they were left untreated. The difference in the outcomes between the two scenarios would be our estimate of the causal effect of the intervention. This is clearly a fantasy, but hope is not lost. Thankfully we can mimic this counterfactual situation by randomizing people into groups, and since we are now talking about groups, we have to start

talking about distributions of future outcomes. — Darren Dahly, PhD⁶

Although RCTs are still preferred to observational studies, there are scenarios in which intervention is not possible. For instance, we cannot assign participants to be obese or not in order to study the effect of obesity on heart diseases.

4.4.1 Origins of RCTs

R.A. Fisher (1890-1962) conceived the RCTs in the 1930s for its employment in agriculture experiments. Fisher designed intricate approaches to disentangle the effects of fertiliser from other variables. Using the *Latin Square*, he would divide the field into a grid of subplots to test each fertiliser with each combination of soil type and plant. However, in this scenario the experimenter would observe the effects of the fertiliser *mixed* (i.e. *confounded*) with a variety of other things (e.g. soil fertility, drainage, microflora). Fisher realised that the only design that would “*trick nature*” is one where the fertilisers are assigned randomly to the subplots. Of course, sometimes you might be unlucky and assign a certain fertiliser to the least fertile subplots, but other times you might get the opposite assignment. A new random allocation is generated each time the experiment is conducted. By running the experiment multiple times the luck of each allocation is *averaged*.



Ronald Aylmer Fisher in 1913

But Fisher realized that an uncertain answer to the right question is much better than a highly certain answer to the wrong question. [...] If you ask the right question, getting an answer that is occasionally wrong is much less of a problem. You can still estimate the amount of uncertainty in your answer, because the uncertainty comes from the randomization procedure (which is known) rather than the characteristics of the soil (which are unknown). — Section “Why RCTs work” in Chapter 4 from ([Pearl and Mackenzie, 2018](#))

The Book of Why describes the aforementioned experiment in causal terms ([Pearl and Mackenzie, 2018](#)). The causal diagram from Figure 4.6 depicts a model describing how the yield of each plot is determined by both the fertiliser and other variables, but the effect of the fertiliser is also affected by the same variables (red arrows). The experimenter aims to know about the effect of the fertiliser controlling for the latter effects. In other words, a model in which the effects represented by the red arrows

⁶<https://statsepi.substack.com/p/out-of-balance>

are controlled. In this second scenario, the relation between Fertilizer and Yield is *unconfounded* since there is no common cause of Fertiliser and Yield.

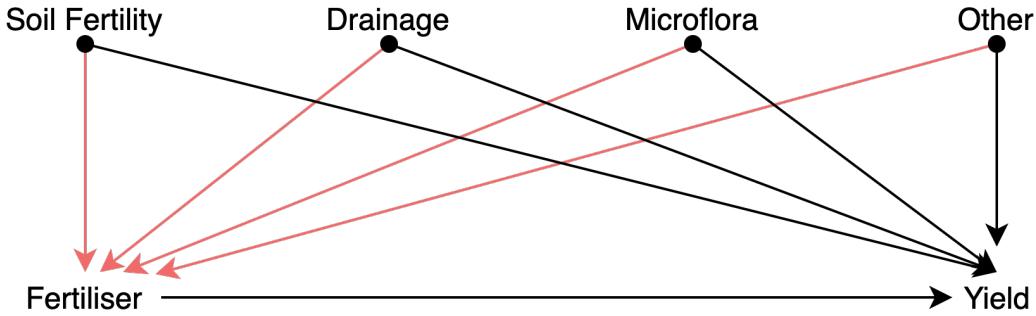


Figure 4.6: Causal diagram depicting an improperly controlled experiment.

4.4.2 Validity

When a hypothesis is designed to explain certain observed phenomena, it will of course be so constructed that it implies their occurrence; hence, the fact to be explained will then constitute confirmatory evidence for it. But it is highly desirable for a scientific hypothesis to be confirmed also by “new” evidence — by facts that were not known or not taken into account when the hypothesis was formulated. Many hypotheses and theories in natural science have indeed received support from such “new” phenomena, with the result that their confirmation was considerably strengthened. — (Hempel, 1966)

Transferring RCTs results to other scenarios is not trivial. All in all, RCTs results concern a particular sample used during the study. The study sample is of course drawn from a larger group, i.e. the population, but the RCT results cannot be simply applied to another sample drawn from the population. Randomisation is not the same as random sampling from the population. In fact, there are many RCT studies that misrepresented certain population groups. An example of women inclusion issues in clinical studies includes the under-representation of women in stroke randomized controlled trials, which leads to misleading conclusions that affect stroke care delivery (Tsivgoulis et al., 2017). A similar bias exists in animal research, including lab mice⁷.

Most rodents used in biomedical studies — the ones that suss out the effects of treatments before they make it to humans — have boy parts and boy biological functions. And that particular kind of gender imbalance

⁷<https://www.wired.com/2016/07/science-huge-diversity-problem-lab-rats/>

has cascading effects. A growing body of evidence indicates that females process pain differently than males. But many lab scientists who study ways of treating pain *still* use all-male cohorts of lab mice. They say it's because male mice and rats aren't as hormonal as females—because isn't that what they always say—and are therefore more reliable in terms of getting data. And that means the scientific community is ignoring research that might help women manage pain better. — **Science Has a Huge Diversity Problem... in Lab Mice - Wired**

Of 2,347 articles reviewed, 618 included animals and/or cells. For animal research, 22% of the publications did not specify the sex of the animals. Of the reports that did specify the sex, 80% of publications included only males, 17% only females, and 3% both sexes. A greater disparity existed in the number of animals studied: 16,152 (84%) male and 3,173 (16%) female. — ([Yoon et al., 2014](#))

Therefore, RCTs must be internally valid, — i.e. the design must eliminate the possibility of bias — but to be clinically useful the result must also be relevant to a well-defined group of patients in a particular clinical setting (i.e. external validity). Differences between trial protocol and routine practice also affect the external validity of RCTs. In ([Rothwell, 2006](#)), the authors list some of the most important potential determinants of external validity.

4.5 Cross-validation in Machine Learning

As data scientists, you may wonder why the previous practices are relevant to your job. In this section I want to show how similar control measures must be considered regarding machine learning (ML). When applying supervised ML methods, is important to prevent over-fitting and under-fitting situations. In particular, over-fitting occurs when a model begins to *memorize* training data rather than *learning* to generalize from a trend (see Figure 4.7). One of the techniques to detect or lessen the effect of over-fitting includes cross-validation. The basis of this technique is to test the generalization power of the model by evaluating its performance on a set of data not used during the training stage.

Importantly, the purpose of cross-validation is to assess how well the model will generalize to an independent dataset. To do that, CV tests (K-folds) the model on K different partitions of the data. If the model we chose is strong and appropriate for our task, it should showcase a performance consistency across folds. The purpose

of cross-validation **is not** to come up with our final model. **The purpose of CV is model checking, not model building.** Afterwards, the model is trained on all the data.

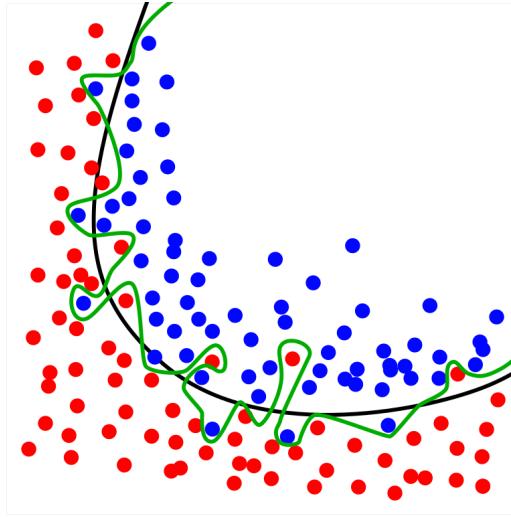


Figure 4.7: Source: Wikimedia⁸. The green line represents an overfitted model and the black line represents a regularized model. While the green line best follows the training data, it is too dependent on that data and it is likely to have a higher error rate on new unseen data, compared to the black line.

To do this, the simplest approach is the **hold out method** which entails splitting the dataset into a train and test sets. However, yet another part of the dataset is often held out (validation set) so that the model training proceeds on the training set, the model evaluation on the validation set, and once the hyperparameters are successfully tweaked, the final evaluation is conducted on the test set. This process reduces the amount of data available for training. Cross-validation (CV) alleviates this issue.

The following procedure (see Figure 4.8) is followed for each of the k “folds”:

- A model is trained using $k - 1$ of the folds as training data.
- The resulting model is validated on the remaining part of the data.

The performance measure reported by K-fold CV is then the average of the values computed in the loop. This approach can be computationally expensive, but does not waste too much data (Pedregosa et al., 2011).

However, the vanilla approach to K-fold CV does not consider certain properties of the dataset. In particular, K-fold CV is not affected by classes or groups. For instance, the training set of the first CV iteration in Figure 4.9 does not contain one of the classes.

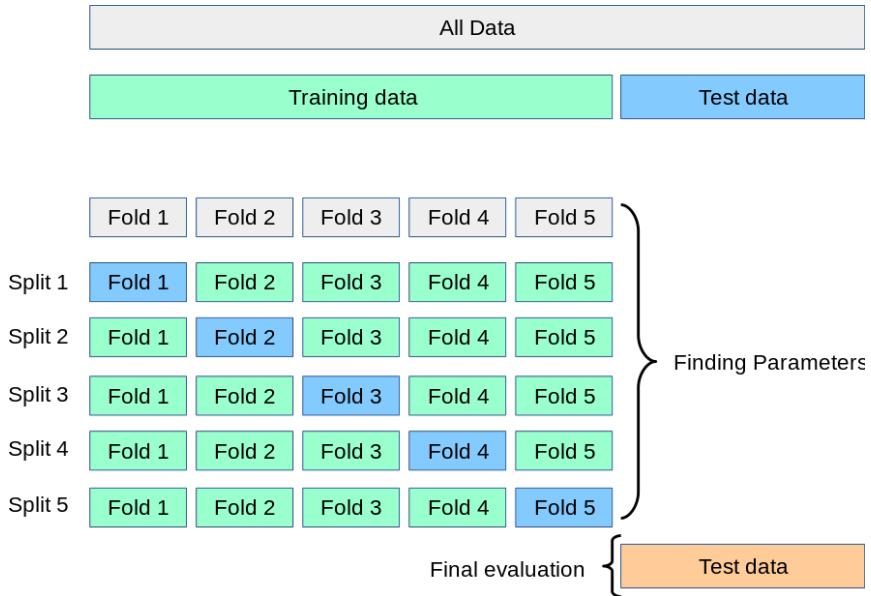


Figure 4.8: Source: Scikit-Learn⁹. A test set should still be held out for final evaluation, but the validation set is no longer needed when doing CV. In the basic approach, called K-fold CV, the training set is split into k smaller sets (Pedregosa et al., 2011).

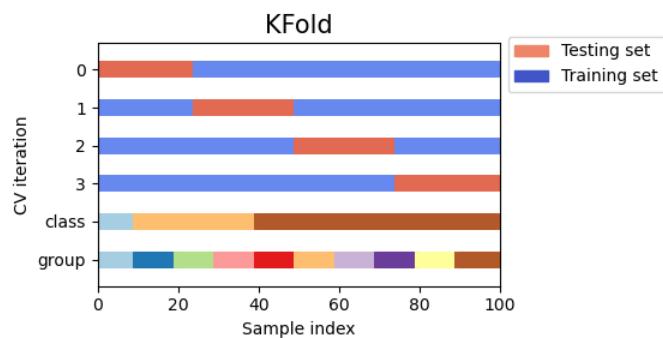


Figure 4.9: Source: Scikit-Learn. K-fold CV is not affected by classes or groups.

Issues similar to the ones previously studied regarding RCTs can arise when conducting cross-validation. Some problems exhibit a large imbalance in the distribution of the target classes. For example, the negative class can be more representative than the positive class. In such cases, stratified sampling is recommended (see Figure 4.10) to preserve relative class frequencies in each train and validation fold.

One strong assumption of machine learning theory is that data is Independent and Identically Distributed (i.i.d.), i.e. that all samples stem from the same generative process and that such process is assumed to have no memory regarding past samples. For example, a succession of throws of a fair coin is i.i.d. since the coin has no memory, so all the throws are independent. In this sense, if we know that the generative process has a group structure (e.g. samples collected from different subjects, experiments, measurement devices) we should use group-wise CV. The grouping of data depends on the context. For instance, in medical data, we can find multiple samples for each patient, so it makes sense to group the samples by patient to prevent any *data leakage*¹⁰. Similarly, problems where the samples have been generated using a time-dependent process call for time-series aware CV schemes¹¹.

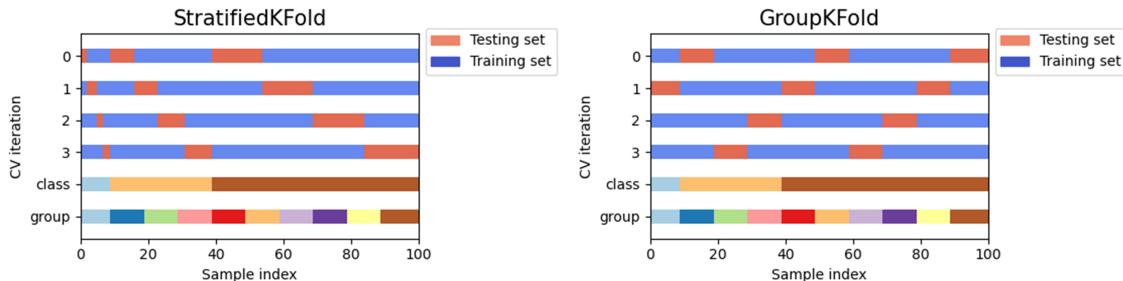


Figure 4.10: Source: Scikit-Learn. Other K-fold CV strategies. **GroupKFold**¹² is a variation of K-fold which ensures that the same group is not represented in both testing and training sets. **StratifiedKFold**¹³ is a variation of K-fold which returns stratified folds: each set contains approximately the same percentage of samples of each target class as the complete set.

Similar to RTC internal validity, cross-validation does not ensure transferability to other scenarios. External validation must be performed with independent datasets to ensure robustness against new scenarios. Consider a deep-learning algorithm trained to predict the number of years a patient will survive based on its characteristics and the medication administrated. This system could be then transferred to a different hospital, in another country, region, or city where the population characteristics (diet,

¹⁰https://scikit-learn.org/0.24/common_pitfalls.html#data-leakage

¹¹https://scikit-learn.org/stable/modules/cross_validation.html#timeseries-cv

hygiene, professions) are different. The model will require undertaking a certain recalibration process to learn the new conditions.

In this chapter we tackle experimental control, and we contrast methods from natural sciences to methods of data science. This does not mean they are equal or equivalent. Despite their differences in nature and aim we can find similarities between a method used for experimental control and methods used in data science (like CV). Unlike CV, RCT aims to determine causality, but both methods aim to reduce or detect spurious associations. Both use randomisation to reduce selection bias and employ a kind of blinding. In a RCT participants and practitioners ignore the assigned group; for CV group-wise CV prevents data leakage (that could break i.i.d.) by grouping samples of the same individuals so that they fall in the same set.

4.6 Surrogates, Proxies, Confounders and Colliders

In this section we address certain induced biases that may affect many data scientist tasks, including data analysis and machine learning solutions. Induced biases often arise during the design and data collection process, but they can also arise during the data analysis step.

4.6.1 Surrogates and Proxies

During this course we may have seen the term surrogate and proxy many times in the context of particular examples. These terms can be used interchangeably. In general, surrogates and proxies are variables that can be measured and employed in place of some other variable that cannot be measured. This impediment may be related to legal, ethical restrictions or technical difficulties. Of course, the employment of such surrogate metric or variable depends upon the assumption that such relationship with the true variable holds for our purposes. For instance, it may not hold for the general population and just for the acquired sample, therefore this surrogate may present external validity issues.

The employment of surrogate metrics may seem attractive for certain purposes but is a source of problems when they are inadvertently exploited in machine learning solutions, algorithms or any other modelling approach.

For instance, the ZIP code is a widely accepted proxy that can often reveal race, ethnicity, or age. From an ethical point of view, such proxies can be employed for the good or the bad. In the United States of America, it is forbidden by law

from considering race in admissions. However, consider that a campus may want to ensure certain diversity or lack of it. In Texas, the “Top 10 Percent Law” is the common name for the state law passed in 1997 that guarantees Texas students who graduated in the top ten percent of their high school class automatic admission to all state-funded universities. The goal was to find a way to increase the enrollment of black and Latino students after a federal court banned race-based affirmative action in the state. Without entering into the fairness and impact of the law, which is still open to debate, it is clear that in scenarios where most high schools have a majority of a certain race, such a race will be favoured over the rest. Therefore, such rules may have a racial or discriminatory impact on university admissions. However, in this case, there may be other factors in place, since being granted access does not mean that you can, for instance, assume the costs of university tuition and so on.

Similarly, say lenders would like to exclude borrowers of certain race. Despite the law prohibiting from doing this, the lenders could find a proxy or surrogate such as neighborhood to continue excluding by race by benefiting from knowing that certain districts or areas are populated by racialised citizens.

Importantly, should the underlying association change, the proxy or surrogate would have to be re-calibrated.

4.6.2 Confounding factors

Related to this, we can find confounding variables. Confounders (see age in Fig. 4.11) are variables that affect both the potential predictor variable (physical activity) and the outcome (cardiac problems). When the presence of confounders is unknown and in lack of experiments specifically designed to minimise them (e.g., randomised controlled trials), we cannot control for them. Uncontrolled confounders may lead to wrongly conclude that a given feature is a strong predictor of the outcome when in reality the association is spurious.

Importantly, such associations may not hold anymore when the sample comes from a different setting where the confounder is differently expressed, e.g. income level may play a different role on diabetes treatment depending on the country. In the context of machine learning, when a model learns spurious associations between predictors and outcomes, an undetected overfitted model is produced, resulting in poor generalisation capabilities that eventually unveil during its translation into real-world settings ([Garcia et al., 2022](#)).

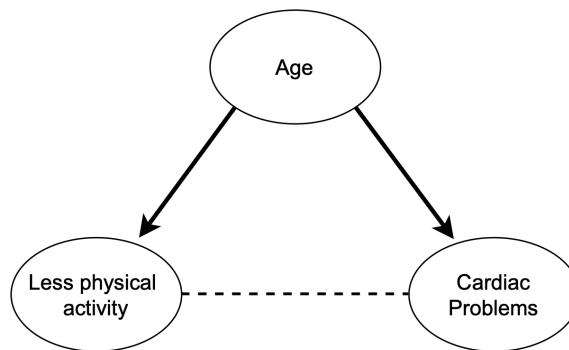


Figure 4.11: DAG depicting a scenario with a confounding factor (age) acting (solid lines) in both the predictor and outcome (relationship depicted with a dashed line).

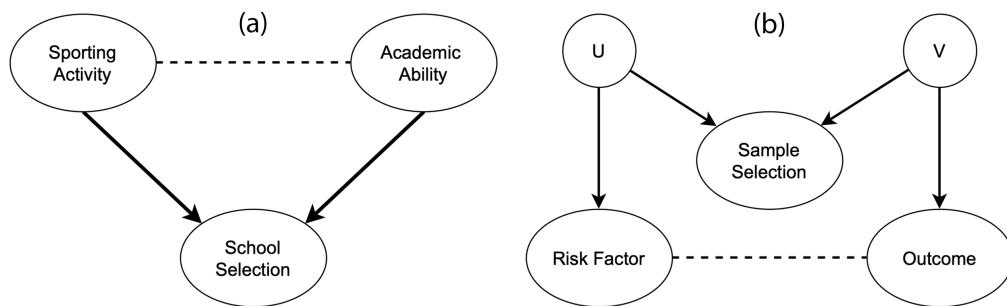


Figure 4.12: Directed arrows indicate causal effects and dotted lines indicate induced associations. (a) shows a scenario in which collider bias could distort the estimate of the causal effect of sporting activity on the academic ability. As shown in (b), the relation between the two associated variables can be indirect, with the risk factor and the outcome being indirectly associated with sample selection through unmeasured confounding variables (U and V).

4.6.3 Collider bias and M-bias

A collider entails a variable that is influenced by two other variables, i.e. collider bias occurs when an exposure and outcome (or factors causing these) each influence a common third variable. The associations induced by collider bias are properties of the sample, rather than the individuals that comprise such a sample. Thus, such associations fail to generalise beyond the sample and may be inaccurate even within the sample, threatening validity (Garcia et al., 2022). Figure 4.12 (a) shows how academic and sporting abilities can influence selection into a prestigious school. These two factors are barely correlated in the general population, but they become strongly correlated in the sample because the school enrolment depends on them (see Figure 4.13). Figure 4.12 (b) shows that the association of interest can be distorted without their variables being directly influencing the collider. Factors affecting the sample selection can themselves influence the variables of interest, distorting the relationship between them (Griffith et al., 2020). This effect is known as **M-bias**.

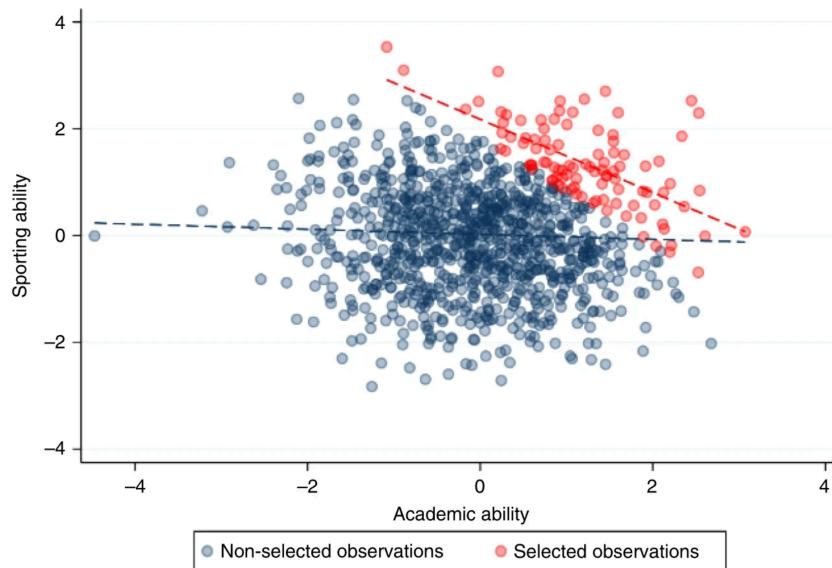


Figure 4.13: These traits are negligibly correlated in the general population (blue), but because they are selected for enrolment they become strongly correlated when analysing only the selected individuals (red). Extracted from (Griffith et al., 2020).

“If we assume that a given **covariate** influences both the hypothesised risk factor and the outcome (a confounder), it is appropriate to condition on that covariate to remove bias induced by the confounding structure. However, if the covariate is a consequence of either or both the exposure and the outcome (a collider), rather than a common cause (a confounder), then conditioning on the covariate can induce, rather than reduce bias.

That is, **collider bias can also be introduced when making statistical adjustments for variables that lie on the causal pathway between risk factor and outcome.** A priori knowledge of the underlying causal structure of variables and whether they function as a common cause or common consequence of risk factor and outcome in the data generating process can be hard to infer. Therefore, it is appropriate to treat collider bias with a similar level of caution to confounding bias.” — ([Griffith et al., 2020](#))

Note for data scientists!

Covariate: An independent variable that can influence the outcome of a given statistical trial, but which is not of direct interest. A covariate is a continuous variable that is expected to change (“vary”) with (“co”) the outcome variable of a study.

4.7 Data alone is not enough

The confirmation of a hypothesis is often considered to increase as the number of favourable test findings grows, but the increase in confirmation, produced by one new favourable instance, will generally become smaller as the number of previously established favourable instances grows ([Hempel, 1966](#)). Many researchers and data scientists blindly rely on the dogma *the more data, the merrier* but the addition of one more favourable finding raises the hypothesis confirmation but little. The confirmation of a hypothesis depends not only on the quantity of the favourable evidence available but also on its variety.

As we have seen during this course, data alone is not enough. Note that this is especially a problem for solutions based on Machine Learning, since domain knowledge or *context* should be introduced somehow to *direct* the model in the desired direction.

“There is no learning without bias, there is no learning without knowledge” — ([Skansi, 2020](#)) ([Domingos, 2015](#)).

A relevant example on how data depends on its context is user ratings or opinions. For instance, the meaning of *fashionable clothes* changes over time, as do political terms. This issue is known as *concept drift* ([Kubat, 2017](#)). Similarly, a text-mining engine to tag biology terms with the corresponding ontology terms may confuse elements between species, as several entities appear in multiple animals or organisms. Context is crucial for external validation and translation of solutions into real-world settings. A system for clothes recommendation should adapt to countries, cultures or ages. Similarly, a health system to predict patient risk based on disease comor-

bilities must be *calibrated* for each country or region (e.g. Diabetes treatment is often affordable in the EU, but an expensive treatment in the USA, which increases its mortality rate).

Figure 4.14 represents an ideal causal inference engine for scientific questions. Today, causal models for scientific applications are based on a similar design. It is important to notice how this diagram showcases the importance of extra-observational information (i.e. information other than data) such as **assumptions**, which derive from the available **knowledge**. With them, a **causal model** is built in any of its different forms, e.g. logical statements, structural equations, causal diagrams, etc. Causation (or a causation assumption) can be defined from the following analogy, X is a cause of Y if Y *listens* to X and determines its value in response to what it hears. For instance, the patient's lifespan L is determined by the intake of drug D . In this case, D acts as a cause of L (although it might not be the only cause), which is represented by an arrow from D to L in a causal diagram (see Figure 4.15). For simplicity, the other causes of L can be grouped in an additional variable Z .

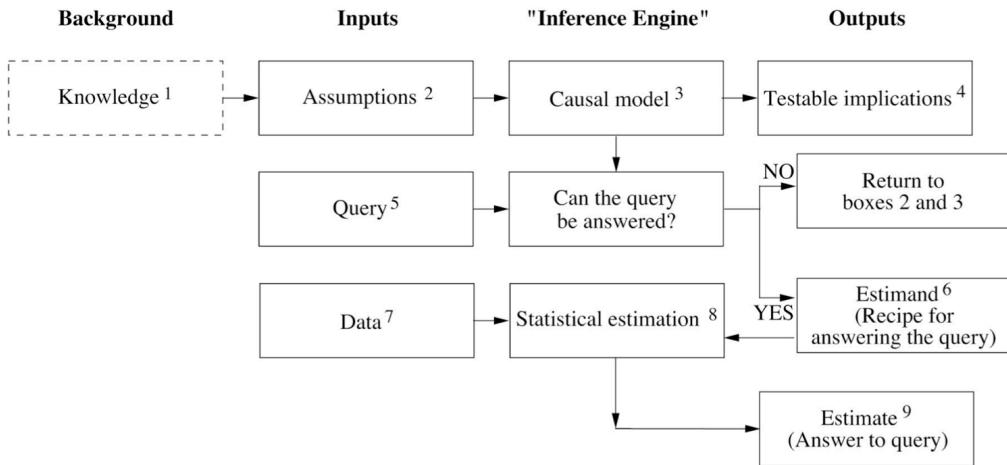


Figure 4.14: Diagram extracted from The Book of Why. The diagram depicts a hypothetical inference engine that combines data and causal knowledge to produce answers to questions of interest. Knowledge (dashed) is not part of the engine but required for its construction. Boxes 4 (testable implications) and 9 (estimate) could also feedback such knowledge to incrementally improve the engine.

In box 4 the patterns encoded in the paths of the causal model yield a series of observable consequences (or data dependencies), that we know as **testable implications** (remember the hypothetico-deductive method?). These implications can be used to test the model. For instance, the lack of path between D and L would imply that D and L are independent, meaning that a variation of D will not alter L . If such implication is contradicted by the data, the model should be revised bearing in

mind this new knowledge. The box 5 is in charge of the scientific **query** which must be encoded in causal vocabulary, e.g. $P(L|do(D))$, i.e. what is the probability that a typical patient would live L years given that it takes the drug D ?

The $do()$ operator represents an intervention in the system, in contrast to an observation $P(L|D)$. An instance of the latter would entail letting the patient decide between taking or not the drug (see left side of Figure 4.15). Such decision might be affected by other variables we are not aware, like the patient's education, family, etc. However, when we make an intervention and assume that we are giving the drug to the patient, the arrow illustrating the patient's decision disappears (right side of Figure 4.15).

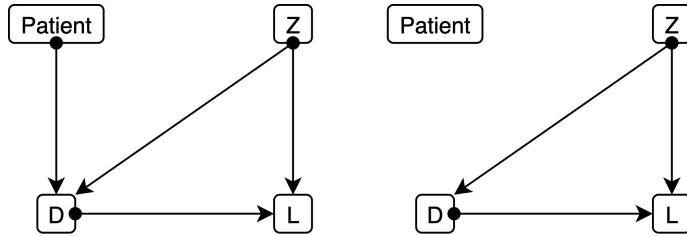


Figure 4.15: Diagram depicting two different scenarios: before/after an intervention.

The ***estimand*** is the recipe to answer the scientific query, written as a probability formula, such as $P(L|D, Z) \times P(Z)$. Once the **data** is introduced, an **estimate** can be calculated. Importantly, some queries may not be answerable regardless of the amount of data collected. For instance, our causal model could indicate that both D and L depend upon a third variable Z . If there would not be any way to measure Z , the query $P(L|do(D))$ would be unanswerable. Collecting data for this question would be worthless. Under such a scenario, the causal model needs to be reviewed. Either to introduce new knowledge to enable estimating Z , or to simplify the previous assumptions, potentially increasing the risk of wrong answers, e.g. stating that Z has a negligible effect on D .

Following the analogy, the data acts as the ingredients of the recipe provided by the ***estimand***. Our estimate (box 9) represents an approximate answer to the query. This answer is approximate because data always represents a finite sample from a theoretically infinite population (Pearl and Mackenzie, 2018). An example answer in this case could be that drug D increases lifespan L of diabetic patients by $30\% \pm 10\%$.

The most important fact about the diagram in Figure 4.14 is that data and causal model are two independent pieces of the puzzle that later work together. Data is collected after the causal model and stating that the scientific query can be answered. The ***estimand*** computation does not require any data. Comparing this to conventional machine learning (ML) systems, a ML solution would have to be re-trained

when moved from one hospital to another since such model just fitted a function to data, without levering from any causal model.

4.8 Examples

4.8.1 Covid-19: How can efficacy versus severe disease be strong when 60% of hospitalized are vaccinated?

There are three kinds of lies: Lies, damned lies, and statistics

In this blog post¹⁴, biostatistics Professor Jeffrey Morris demonstrates how without properly controlling for age, efficacy against severe disease in Israel may appear weak when in fact within each age-group it is extremely strong. Consider the table from Figure 4.16 and the following data from the the Israeli government. As of August 15, 2021 nearly 60% of all patients currently hospitalized for COVID-19 are vaccinated. Out of 515 patients currently hospitalized with severe cases in Israel, 301 (58.4%) of these cases were fully vaccinated.

Age	Population (%)		Severe cases		Efficacy vs. severe disease
	Not Vax %	Fully Vax %	Not Vax	Fully Vax	
All ages			214	301	Vax don't work!

Figure 4.16: Misleading table. This kind of tables have been used to claim that vaccines do not work or that its efectiveness reduces over time.

The numbers are true, but we need more than that to draw a proper conclusion about vaccine efficacy. Consider the following extreme scenarios. If the number of vaccinated people would be 0 we would expect all severe cases to be not vaccinated (obviously). On the other hand, if 100% of people would have been vaccinated, we would expect all severe cases to proceed from vaccinated people and 0 from non vaccinated. In this case, we have an in-between situation where 80% of residents (older than 12 years) have been vaccinated. Therefore, since the group of vaccinated people is larger than the non-vaccinated, we can expect more severe cases in absolute numbers. However, once we adjust for vaccination rates and normalise the counts, the story changes. The rate of severe cases is three times higher in unvaccinated individuals.

¹⁴<https://www.covid-datascience.com/post/israeli-data-how-can-efficacy-vs-severe-disease-be-strong-when-60-of-hospitalized-are-vaccinated>

Age	Population (%)		Severe cases		Efficacy vs. severe disease
	Not Vax %	Fully Vax %	Not Vax per 100k	Fully Vax per 100k	
All ages	1,302,912 18.2%	5,634,634 78.7%	214 16.4	301 5.3	67.5%

Figure 4.17: Table adjusted for vaccination rates.

Vaccine Efficacy vs. Severe disease = $1 - 5.3/16.4 = 67.5\%$. The interpretation of this number is that the vaccines are preventing $>2/3$ of the serious infections leading to hospitalization that would have occurred sans vaccination.

Still, the obtained efficacy is lower than what we would expect. There are other factors that contribute to this confusion, including: age disparity in vaccinations, old people is more likely to be hospitalized than young people, etc.

I recommend going through the blog post to see how the author continues to apply adjustments and stratifications to find the true efficacy of the vaccines. Moreover, this is a good example of the Simpson's paradox, where misleading results can be obtained in the presence of confounding factors.

In conclusion, as long as there is a major age disparity in vaccination rates, with older individuals being more highly vaccinated, then the fact that older people have an inherently higher risk of hospitalization when infected with a respiratory virus means that it is always important to stratify results by age; Even more fundamentally, it is important to use infection and disease rates (per 100k, e.g.) and not raw counts to compare unvaccinated and vaccinated groups to adjust for the proportion vaccinated. Use of raw counts exaggerates the vaccine efficacy when vaccinated proportion is low and attenuates the vaccine efficacy when, like in Israel, vaccines proportions are high.

4.8.2 Misinterpretations of hurricane forecast maps

The following article¹⁵ by Alberto Cairo published in The New York Times explains how hurricane cone forecast maps can mislead the public and produce real-world consequences.

¹⁵<https://www.nytimes.com/interactive/2019/08/29/opinion/hurricane-dorian-forecast-map.html>



Figure 4.18: Example of hurricane forecast cone graphic in TV.

Studies¹⁶ show that people can misinterpret this type of map as indicating that the hurricane will get bigger over time. Others think it shows areas under threat. Recent research¹⁷ suggests that 40% of people would not feel threatened if they lived just outside of the cone. Moreover, people who live inside the cone, but far from the center, take less precautions than those closer to the central line. These misunderstandings have real-world consequences. Actually, the cone represents a range of possible positions and paths for the storm’s center. The dots in the middle of the cone correspond to the forecast of where the hurricane’s center could be in the following five days. But there’s a good chance that the actual center of the storm will not end up being at those positions.

To create the cone, the National Hurricane Center (N.H.C.) surrounds each estimated position of the storm center with circles of increasing size. **These circles represent uncertainty**, meaning that the storm center may end up being anywhere inside the circles — or even outside of them. The uncertainty circles grow over time because it is more difficult to predict what will happen in five days from now than in one day. Finally, a curve connects the circles, resulting in what is popularly known as the ‘cone of uncertainty’.

N.H.C. says cones will contain the path of the storm center only 60 to 70 % of the time. In other words, one out of three times we experience a storm like this, its center will be outside the boundaries of the cone. Hurricanes are also hundreds of miles wide, and the cone shows only the possible path of the storm’s center. Heavy rain, storm surges, flooding, wind and other hazards may affect areas outside the cone. The cone,

¹⁶<https://www.semanticscholar.org/paper/Misinterpretations-of-the-%E2%80%9CCone-of-Uncertainty%E2%80%9D-in-Broad-Leiserowitz/f7c04b6eb883cf7d7fdee007cda056ed18182829>

¹⁷<https://interactive.miami.edu/hurakan/>

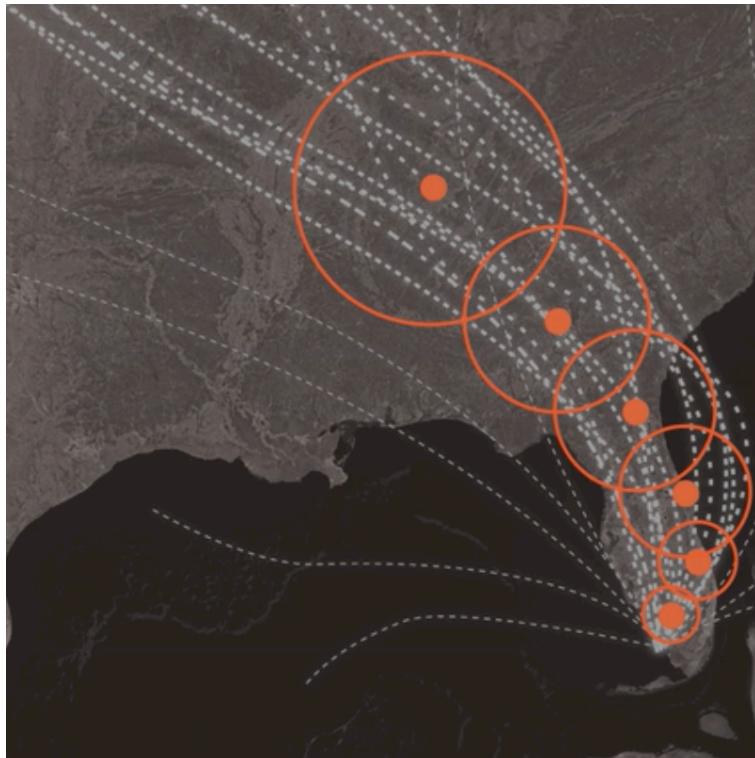


Figure 4.19: Cone of uncertainty.

when presented on its own, doesn't give us much information about a hurricane's dangers. The N.H.C. designs other graphics, including this one showing areas that may be affected by strong winds. But these don't receive nearly as much attention as the cone. The cone graphic is deceptively simple. That becomes a liability if people believe they're out of harm's way when they aren't. As with many charts, it's risky to assume we can interpret a hurricane map correctly with just a glance. Graphics like these need to be read closely and carefully. Only then can we grasp what they're really saying.

From a NYT article¹⁸ by Alberto Cairo

¹⁸<https://www.nytimes.com/interactive/2019/08/29/opinion/hurricane-dorian-forecast-map.html>

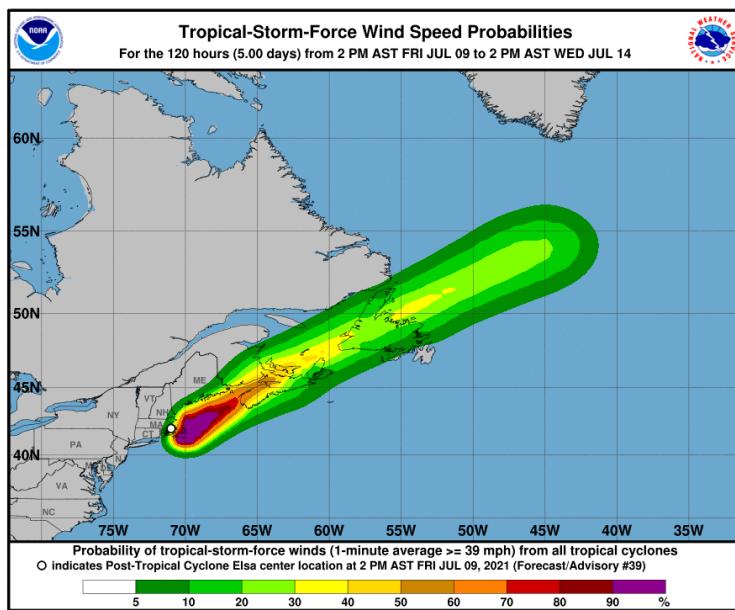


Figure 4.20: Other graphics designed by USA National Hurricane Center.

Chapter 5

Ethics and Responsibility

5.1 Overview

Very often, progress, ethics and laws go hand in hand. The industrial revolution increased the living standards but also brought new challenges, problems and ethical questions. Labour conditions were harsh, and child labour was common long before the industrial revolution. Progress and technology raise new ethical questions. Bit by bit certain practices and conditions became culturally unacceptable. Then a growing social consensus in the same direction eventually forced laws to reflect those changes. Nowadays, working rights such as office restrooms or the lunch break are widespread.



Figure 5.1: A young drawer pulling a coal tub along a mine gallery.

Similarly, we live in the information revolution, which has transformed the way societies communicate in a way not seen since the Gutenberg printer revolution, enabling the world to transmit information worldwide almost instantaneously through the interconnection of computers. The information era brings ethical challenges to a new dimension. Gutenberg's print helped to spread literature and knowledge, but it also made it easier to spread hate and fake news (see this article for some examples¹). Today, tools like Facebook have been used to incite the genocide of Rohingya in Myanmar². Every data scientist must understand the impact of their work for the

¹<https://www.politico.com/magazine/story/2016/12/fake-news-history-long-violent-214535/>

²<https://www.nytimes.com/2018/10/15/technology/myanmar-facebook-genocide.html>

good and the bad. Something as simple as a good visualization can be enlightening and convince people to get a vaccine or find the root of an outbreak (see John Snow’s map). But bad data science can have unexpected real-world consequences too.

In this chapter, we will not address all the philosophical questions regarding ethics, which are lengthy³ but will try to focus on the ethical practice of data science. Many ethical questions affecting data ethics have already been tackled during this book. Such questions have been interleaved during this book together with other topics, e.g., algorithmic fairness, validity issues and their consequences. Therefore, this book provides a utilitarian framework of ethics. This chapter describes how ethical values, regardless of which they may be, affect the practice of data science, either because of the lack of commitment to them or to how data science is conditioned when abiding by them.

This chapter begins by delving into the definitions of ethics, some ethical frameworks and the values of science. Next, it will describe how ethics are meant to be articulated. Afterwards we will be ready to tackle data ethics and the topics affecting it. Finally, a short summary of the GDPR and examples affecting data ethics are provided to illustrate the impact and significance of this new field.

Course objectives

This final chapter supports objectives 9 and 10, introducing the role of non-epistemic values into every stage of the data science workflow. Historical and current examples showcase the role of regulations, ethical concerns and the consequences of neglecting them. After all, data scientists will have to comply with codes of conduct, legal constraints and consider ethical debates.

5.2 Morality and Ethics

Ethics comes from the Greek word *Ethos*, meaning habit or custom. Ethics are theories that offer normatively valid reasons to rationally endorsing a code of behaviour. In general, the cost of following an ethical rule is less than the benefit we obtain from others following the same rule (e.g. “not stealing does not require an effort from myself, and I get a huge benefit from others following the same rule”). In this sense, like many other social phenomena, ethical rules can be regarded as a network which follows Metcalfe’s Law, which implies a critical mass point in network size, after which network value begins to exceed its cost (Metcalfe, 2013). This behaviour explains why ethical or societal changes may often be adopted very quickly after a

³<https://en.wikipedia.org/wiki/Ethics>

long *plateau*. However, this may happen bothways, either after reaching a certain mass point or receding from it. For example, feminist ethics aims to understand, criticise and correct how gender operates within our moral beliefs and practices ([Norlock, 2019](#)). Similarly, we can find environmental ethics and, of course, data ethics. In science, we often need to answer ethical questions like:

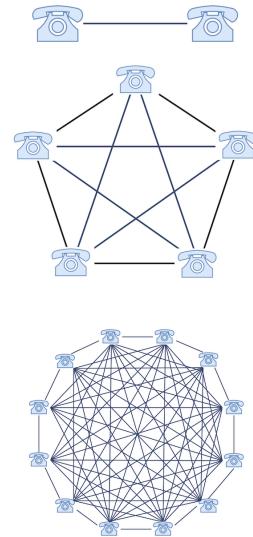
- Is it right to use RCT in this specific scenario?
- Under what conditions is justified to conduct animal experiments?

All ethics groups aim to understand, criticise and correct how certain actions impact a particular sphere. Similarly, **data ethics aim to understand, criticise and correct how our data practices (from collection, to processing, analysis) operate within our moral beliefs and impact our society.**

Descriptively, *morality* refers to a code of conduct that is put forward by a society or group (e.g. a religion), or accepted by an individual. In a normative sense, *morality* refers to a code of conduct that would be accepted by anyone who meets certain intellectual conditions (e.g. being rational) ([Gert and Gert, 2020](#)). For example, a company is not a *moral agent* but it can encourage a certain *code of conduct* to their workers. Even though *morality* is the subject matter of ethics, it is most often used interchangeably with *ethics*. Finally, ethics are not laws but laws are often used to enforce certain shared social values.

In essence, ethics tell us about what is right and wrong. They are the cornerstone of civilisation and no human practice is foreign to ethical values, although these may vary from region to region or over time. Therefore, ethics are context dependent. For many, their moral principles stem from religious teachings, ideologies, philosophies, etc. These can promote ethical behaviour, but it does not mean that ethics should necessarily stem from a religion or belong to a particular ideology or philosophy.

Finally, ethics are not laws but very often laws reflect the societal ethical consensus. Progress brings new ethical questions, resulting in a changing societal consensus that eventually translates into new laws. A notable example is the advent of the industrial era which radically reshaped the working methods and workers conditions. Workers fought to claim certain rights and conditions improvement (e.g. have a place to have lunch, time to have lunch, toilets, etc.). However, ethical changes may not



2 phones can establish just one connection, 5 may establish 10 connections, and 12 can set 66.

always imply *replacing* a set of rules by another set, or *replacing* a code of conduct for another code of conduct. Back then when western countries forbid women from showing their bare legs, many women defended their right to wear whatever clothes they wanted and as short as they wanted. Clothes like the mini-skirt became a symbol because it challenged the established code of conduct. In the same way, trousers also became a symbol for women empowerment in the 1850s (see bloomers), also to show that they could also be part of the workforce. Nevertheless, **what is important about this** kind of changes is not that they were replacing their clothes for another set of clothes. They wanted to have the choice, to have control over their clothes. All in all, just because people could wear whatever they wanted it does not mean that they will wear it. But people now have the right to choose.

Exactly the same can be applied to data science in many scenarios. The most notable case affects browsing the internet and the recent regulations for cookies. Thanks to the new regulations, any website you browse should give you the choice to either accept analytics cookies or not. **Data ethics is about having control as data subjects over your data.** You have the choice to withdraw from a study or remove your data from Facebook or Twitter. **Data ethics extend your rights and those from which you may collect data.**



Figure 5.2: Ethics guide the creation of laws. Progress brings new ethical questions.

5.3 Ethical Frameworks

In this book we will not delve into the different ethical frameworks. However, it is important to regard at least three of them superficially.

5.3.1 Consequentialism

In a consequential ethical framework, our choices are morally assessed only by their consequences (e.g. utilitarians may evaluate their choices based on the obtained pleasure, happiness, welfare, safety). People may obtain pleasure or satisfaction

from a car based on different attributes. Speed may satisfy your need to arrive earlier to your destination. Public transport may be preferred in terms of safety or pollution. Similarly, there are consequences to consider in science and data science. Is your computer program faster although it requires 100 computers?

Sometimes we need to trade off such consequences and for example trade off speed in exchange of a less costly solution (using 100 computers needs a lot of energy). Thus, science also deals with epistemic consequences, such as whether the model predictions are precise or accurate; or if the explanations of our scientific theories are successful. Other actions may affect institutions. If the scientific method and the appropriate means are not properly followed, errors may happen, damaging the trust of an institution, the impartiality of a jury, and so on.

The consequentialism framework is not exempt from issues. Consequentialism is very demanding since all acts are either required or forbidden because every choice has consequences. Importantly, consequences are not always clear and evident. For instance, sharing my medical records may improve general health care at the cost of losing my individual privacy.

5.3.2 Deontology

“Act as if the maxim of your action should become, by your will, the universal law of nature.” — Categorical Imperative of Immanuel Kant

For deontology, to choose morally is to satisfy relevant rules and duties. Deontology comes from the greek word *deon* (duty). Under this framework, some choices are not justified because of their effects. People’s rights are a good example (e.g. right to demonstrate, expression).

Some examples include the journalist deontology code which considers the right of people not to answer a question. Nurses have a deontology code under which they are required to obtain consent from the patient or its relatives. Note that the actions of a nursing team could have beneficial consequences for the patient but they need to wait for the consent. Of course, deontology codes always provide exceptions (e.g. an unconscious person in the sidewalk must be helped).

“Some choices cannot be justified by their effects, no matter how good their consequences are, some choices are forbidden and others mandated.”
— (Alexander and Moore, 2021)

The advantages of deontology over consequentialism lie in the ability to create checklists and conduct codes without the need to consider all consequences. Some examples of such checklists exist in data science.

- FAIR⁴: Findability, Accessibility, Interoperability, and Reuse of digital assets ([Wilkinson et al., 2016](#)).
- DOME: recommendations for supervised machine learning validation in biology ([Walsh et al., 2021](#)).

Deontology issues relate to the blind compliance with deontological norms, which may bring disastrous consequences. This problem may lead to deontology fundamentalism that forbids questioning the rules. Categorical rules may lack degrees of wrongness.

There have been various attempts to reconcile deontology with consequentialism. The proposed “threshold deontology” indicates that rules must govern up to a certain point regardless of the adverse consequences. But once the consequences become so dire that they cross a stipulated threshold, consequentialism takes over.

5.3.3 Virtues

Finally, we have virtue ethics. Defenders of this view argue that morality consists of having proper character traits, i.e. virtues, and to exemplify them through their use. Some argue that certain traits define humans, what we usually refer as humanity (hospitality, etc.). Philosophers like Aristotle argued that there are some traits that distinguishes us as humans (e.g. compassion, help, generosity).

In science, there are many virtues that we need to consider every day. Scientists need to defend their colleagues or find the courage to criticise them. Scientists need to be honest and sincere. The problem, as can be seen, is that virtue ethics is less demanding than a defined set of rules or codes of conduct. Virtues do not state any set of rules or quantities regarding the amount of your honesty or help. However, virtues are usually accompanied by justifications from a moral perspective. For instance, “I do help because I believe helping is good”.

Importantly, deontology codes are better understood when they are properly justified. However, many prescriptions from deontology codes cannot be justified by virtues right away. For instance, helping is a strong virtue, but providing treatment without consent is not acceptable. Moral conflicts can be difficult to solve by reference to virtues alone. These conflicts are common in science and data science, e.g., should we employ animals for research?.

⁴<https://www.go-fair.org/>

5.4 Values in Science

“One cannot deduce an ‘ought’ from an ‘is’ ” — David Hume

In his book “Philosophy of Science for Scientists” (Springer), Lars-Göran Johansson dedicates the 13th Chapter to discuss the role of values in science and to different views regarding the meaning of value-free science. I would like to superficially tackle this topic as I believe it to impact data science as it does for general science. As Johansson argues, a strong interpretation of value-free science would not contain any expression of values. However, as the author points out, this view ignores potential effects on some groups (Johansson et al., 2016). For instance, such a view of data science would perhaps ignore fairness values, hence ignoring an appropriate data collection, including a diversity of individuals or a proper evaluation of the outcome performance (either of an analysis or a model) on different groups. Another alternative, perhaps more reasonable, is that researchers should not convey their values. However, a researcher could provide an accurate description of the effects of some measures and yet recommend a particular action based on its subjective analysis of the analysed effects. It is then up to the reader to distinguish between the first part and the prescription.

A third interpretation, the author adds, is that “one ought not to deduce value statements from premises that only contain descriptive statements” (see Hume’s quote). However, very often, political argumentation breaks this rule, usually by jumping from descriptive statements to normative statements while leaving out (sometimes deliberately) an intermediate premise. This last issue stems from the is-ought problem⁵ and fact-value distinction⁶.

Note for data scientists!

A **descriptive statement** gives an account of how the world is without stating whether that is good or bad. Thus, such statements describe in an objective or non-judgmental way. Conversely, a **normative statement** expresses an evaluation, stating that something is good or bad, better or worse, relative to some standard or alternative (i.e. norm).

Regardless of the interpretation debate⁷ of Hume’s words (Cohon, 2018), it is never easy not to deduce normative conclusions from descriptive statements, specially in fields such as psychology, economics or political science. And very often, words such as prosperity, democracy are used to describe, but they are nonetheless normatively

⁵https://en.wikipedia.org/wiki/Is%E2%80%93ought_problem

⁶https://en.wikipedia.org/wiki/Fact%E2%80%93value_distinction

⁷<https://plato.stanford.edu/entries/hume-moral/#io>

charged words. The important conclusion is to explicitly indicate the norm and definitions of such terms in our research, specially considering that definitions and norms may evolve over time.

Consider the following quotes regarding value-free and value-laden science.

Everybody would agree that scientific knowledge has sometimes been used for unethical ends — to make nuclear and chemical weapons, for example. But such cases do not show that there is something ethically objectionable about scientific knowledge itself. It is the *use* to which that knowledge is put that is unethical. Indeed many philosophers would say that it makes no sense to talk about science or scientific knowledge being ethical or unethical *per se*. For science is concerned with facts, and facts in themselves have no ethical significance. It is what we do with those facts that is right or wrong, moral or immoral. On this view, science is essentially a *value-free* activity. — (Okasha, 2016)

Research aimed at finding better treatment for diseases are certainly value-laden; we would give very high value to positive results in this area. But, of course, reports of such results, if they are found, should be value-free; we want objective knowledge about what to do for curing and/or preventing severe diseases. That is to say, scientific activity can be value-laden and value-free at the same time. — (Johansson et al., 2016)

One takeaway message from this section is that science is driven by values. Science is, therefore, value-laden, but its results can, and should, be value-free. Whether a researcher comes to a conclusion that stands in line or that undermines the researcher's non-scientific interests, one must, under all circumstances, evaluate the matter itself, as such considerations are not relevant arguments for or against the scientific work at issue.

A relevant example is the feminist critique of science (Crasnow, 2020) and the existence of certain unconscious assumptions from which researchers might be unaware until recently. A notable example of this includes the research on myocardial infarction⁸, which was mainly conducted on men. Today we understand symptoms vary between women and men. This issue resulted in underdiagnosis and mistreatment of infarction in women. Similarly, some assumptions derived from the modern state of society, have been employed to describe social patterns in prehistory in research. For instance, the division of labour between sexes with men dedicated to outdoor activities and women dedicated to home activities has conditioned the research about

⁸https://en.wikipedia.org/wiki/Cardiovascular_disease_in_women#History

prehistoric humans. Recent studies challenge the common assumption that prehistoric men hunted while women gathered⁹. The relevant conclusion is not whether the previous research is false but that unconscious norms regarding what we consider ‘normal’ may affect our criteria to focus on some circumstances and ignore others. What we consider ‘normal’ or ‘common’ may vary across regions and over time. All in all, humans are not neutral observers of the external world.

Sorting phenomena into categories is an interest-related task (Johansson et al., 2016) but is also unavoidable. Nevertheless, is important to make such interests explicit. An example of such categories defined by humans is diseases, which are progressively updated, hierarchized or redefined according to new knowledge. Diseases are not found in nature as entities *per se*; they refer to a definable deviation from a normal phenotype evident via symptoms and/or signs. The different sets of symptoms, pathologies and signs are grouped into diseases, and likewise diseases are grouped into categories, all of them organised into disease taxonomies (e.g. International Classification of Diseases). Thus, one disease can have more than one etiology, and one etiology can lead to more than one disease (Vega, 2021).

5.5 Ethics in action

“There are years, centuries, in which nothing happens, and there are days, like yesterday, into which a whole lifetime is compressed.” — The devil
(Van Westrum, 1908)

Some events in history completely changed how we see the world. Section 5.8 will present some of these events and how they relate to data ethics. However, the consequences of these events share common aspects. In the case of data ethics and science, some events have reshaped the societal consensus and the conception of certain rights. For instance, the scandals and ban of Dichloro-diphenyl-trichloroethane (DDT¹⁰) as insecticide changed the cultural awareness regarding chemicals. In the USA, the Tuskegee tragedy¹¹ raised ethical questions regarding voluntary informed consent, which led to the establishment of the National Commission for the Protection of Human Subjects of Biomedical and Behavioral Research¹² and the National Research Act¹³. Both were intended to shape bioethics policies in the USA and the

⁹<https://www.nationalgeographic.com/science/article/prehistoric-female-hunter-discovery-upends-gender-role-assumptions>

¹⁰<https://www.nytimes.com/2004/04/11/magazine/what-the-world-needs-now-is-ddt.html>

¹¹https://en.wikipedia.org/wiki/Tuskegee_Syphilis_Study

¹²https://en.wikipedia.org/wiki/National_Commission_for_the_Protection_of_Human_Subjects_of_Biomedical_and_Behavioral_Research

¹³https://en.wikipedia.org/wiki/National_Research_Act

establishment of institutional review boards for the studies.

More recently, IBM faced an scandal¹⁴ in 2019 regarding the scrapping of millions of on-line pictures without consent for the development of facial recognition technology. Also, in 2019, Project Nightingale¹⁵, was a data storage and processing project by Google and Ascension (one of the largest private healthcare systems in the United States), which involved the secret transfer of medical data of up to 50 million US Americans. The data was transferred without consent from doctors and patients containing full personal details, including name and medical history and could be accessed by Google staff (The Guardian, 2019¹⁶). The raw data included lab results, diagnoses, hospital records such as medical history, and personal details of patients. In the 2010s, the Facebook-Cambridge Analytica scandal¹⁷ involved the unconsented personal data collection from millions of Facebook users by British consulting firm Cambridge Analytica. This data was predominantly employed for political advertising. This event led to the development of the #DeleteFacebook movement, among other cultural and political consequences. The Federal Trade Commission fined Facebook with \$5 billion. The cultural impact of these scandals was huge and led to massive media coverage, TV shows and films. Regulators and law makers also acted to further protect users data ownership.

In this sense, **ethics are intended to be actionable**. For instance, environmental ethics (Brennan and Lo, 2022) make considerations about the moral relationship of human beings to the environment and its non-human contents, e.g., “is it morally acceptable to pollute a river for the production of energy?”. Once we answer such questions, actions need to *make a difference*, e.g., translating the answer into rules, laws or regulations to either allow, forbid or limit such practice. Answering ethical questions or showcasing a position to them is not enough. For example, if the societal consensus (either at a company level or country level) agrees on the need for toilets in the workplace (or kitchens), the next step is to provide such spaces to *make a change*.

However, it is not always clear how to articulate certain ethical positions. Consider the following ethical goal: “artificial intelligence models should be equally fair for everybody disregarding of their demographic characteristics”. Then consider the following issue: “does preventing any data collection of demographic information avoid the development of biased machine learning models?”. How do we test if the

¹⁴<https://www.nbcnews.com/tech/internet/facial-recognition-s-dirty-little-secret-millions-online-photos-scraped-n981921>

¹⁵https://en.wikipedia.org/wiki/Project_Nightingale

¹⁶<https://www.theguardian.com/technology/2019/nov/12/google-medical-data-project-nightingale-secret-transfer-us-health-information>

¹⁷https://en.wikipedia.org/wiki/Facebook%20%93Cambridge_Analytica_data_scandal

model is fair against different population groups (e.g. minorities) if we did not collect information that let us conduct such tests?

Sometimes it may be impossible to act according to some ethical values. For instance, you may want to buy a t-shirt manufactured under some specific conditions that you find ethical, but you may find none. Therefore you are forced to take an unethical decision according to your values. Similarly, you may need your data to meet some conditions or requirements to avoid producing a biased output (such as an analysis or an ML model). But sometimes, such data might not exist or proving such an assumption becomes challenging. Garcia et al. showcase an example of this issue in their work addressing public Covid-19 X-Ray datasets. The authors highlight that the most popular datasets were found to have a very high risk of inducing bias in models. Importantly, some solutions employed datasets containing adult individuals for the unhealthy group and a paediatric population as the control group ([Santa Cruz et al., 2021](#)). Such an approach raises ethical concerns regarding the risk of harmful outcomes deriving from working with poor data. In this case, new data acquisition would be needed to improve data quality and reduce the risk of producing a biased solution.

Preventing issues such as **data drift** and **concept drift** require constant monitoring of the models. Unsurprisingly, any intelligent entity interacting with a changing world requires frequent re-training and re-education. Any job requires a similar knowledge update process, from medical doctors to computer scientists or pilots. Such jobs require learning new diseases, treatments, regulations, protocols, devices, technologies, etc. Similarly, constant model evaluation is paramount to maintain the appropriate performance and fairness. Regulations and guidelines help to operate ethical considerations. For instance, the recent document “Ethics Guidelines for Trustworthy AI” published by the “high-level expert group on artificial intelligence” group from the European Commission presents transparency as a requirement for trustworthy AI. This document recommends transparency of all elements relevant to an AI system: the data, the system and the business models ([AI, 2019](#)). Abiding to this requirement entails designing solutions with traceability and explainability in mind (among other features) so that models can be properly audited and their decisions thoroughly traced.

Other actions that may help to prevent bias and/or external validation issues include accompanying the data with metadata describing acquisition details, causal models describing the assumptions taken during the design stages, and documentation, among others ([Garcia et al., 2022](#)).

“In many product design (and other) endeavors, there often comes a

moment when the question is asked, “Are we doing the right thing?” or “Is this the right solution?” In this context, the word right can mean many things. It can mean: Are we meeting the customer’s expectations? Is the design solution appropriate to the problem? Are we honoring the scope of the work? Is this a profitable feature to add to our product? Will people buy this? It can also mean: Do we agree that this action is acceptable to perform based on our values?” — ([Davis, 2012](#))

5.6 Data Ethics

We often face situations where we care about the privacy of our data, e.g., we do not want our location to be shared with certain companies, but we wish to benefit from the same data to see if the road is jammed on our way to work. Similarly, we want our medical records to remain private but also wish to benefit from the analysis of medical records to improve the care and treatment that we receive. Moreover, we aim to automate and ease decisions thanks to data-driven algorithms but we also worry about unintended bias. Getting to know the consequences of our actions make us aware of what is right or wrong. Without that, we lack any motivation to fix potential problems. Data science is a new field and we are still defining what is right and wrong. We are starting to experience certain consequences regarding privacy, fairness, equity, etc. In this sense, a social consensus is growing, with laws and regulations enacted to enforce such values.

In this section we will tackle the origins, and the specifics of data ethics, including topics such as informed consent, data ownership and privacy. They all affect anonymity of individuals, the validity of experiments as well as the algorithmic fairness of the developed solutions.

5.6.1 Origins

It is important not to confuse data ethics with computer ethics and information ethics. In the mid 1940s, developments in science and philosophy led to the creation of a new branch of ethics that would later be called “computer ethics” or “information ethics”. The founder of this new philosophical field was the American scholar Norbert Wiener, a professor of mathematics and engineering at MIT ([Bynum, 2018](#)).

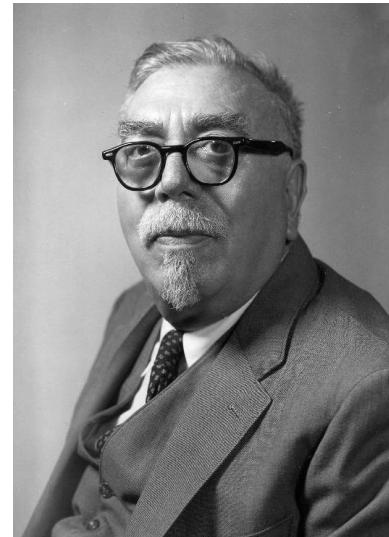
Moreover, other terms such as “cyberethics” and “Internet ethics” have been used to refer to computer ethics issues associated with the Internet.

“Because of the breadth of Wiener’s concerns and the applicability of his ideas and methods to every kind of information technology, the term

“information ethics” is an apt name for the new field of ethics that he founded. As a result, the term “computer ethics”, as it is typically used today, names only a subfield of Wiener’s much broader concerns.” — (Bynum, 2018)

These are variously termed “computer ethics”, “information ethics”, or “ICT ethics” (information and communication technology ethics). “Information ethics” is perhaps the most appropriate name for Wiener’s field of ethical research, because it concerned all means of storing, transmitting and processing information, including, for example, perception, memory, printing, telephones, telegraph, recorders, phonographs, television, radio, computers, and so on. “Computer ethics” involves ethical questions and problems that are altered by the use of computers or that would not have existed if computers had not been invented.

Although these two branches of ethics are not handled in this course, they are tightly related to data ethics.



Norbert Wiener (1894-1964).
Copyright MIT Museum.

Data ethics builds on the foundation provided by computer and information ethics but, at the same time, it refines the approach endorsed so far in this research field, by shifting the level of abstraction of ethical enquiries, from being information-centric to being data-centric. This shift brings into focus the different moral dimensions of all kinds of data, even data that never translate directly into information but can be used to support actions or generate behaviours, for example. It highlights the need for ethical analyses to concentrate on the content and nature of computational operations—the interactions among hardware, software and data—rather than on the variety of digital technologies that enable them. — (Luciano and Mariarosaria, 2016)

5.6.2 What are data ethics?

As seen in the previous section, data ethics are tightly related to computer ethics and information ethics. Data ethics may be as well be found termed as big data ethics¹⁸ which tackles data ownership, consent, privacy, and openness among other

¹⁸https://en.wikipedia.org/wiki/Big_data_ethics

principles. To answer this question, we need to clarify what is data science. I like the answer given by Rachel Schutt and Cathy O’Neil in their book “Doing data science: Straight talk from the frontline” ([O’Neil and Schutt, 2013](#)) in which a list of data science tasks is given (see some below). However, it is important to remember some lessons from previous chapters when we address this question. **Data science is not just to solve problems using data**, but rather to solve problems *through* data by employing a series of guidelines, practices, methods and knowledge. Consider cooking, which may look as a series of ingredients and instructions. We need the ingredients, the tools, the recipe, but in order to come up with new recipes, or to understand previous recipes, we need knowledge.

- Exploratory data analysis
- Visualization
- Dashboards and metrics
- Find business insights
- Data-driven decision making
- Data engineering / Big Data
- Get the data themselves
- Build data pipelines
- Build products instead of describing existing product usage
- Patent writing
- Detective work, ask good questions, make hypothesis, research
- Predict future behaviour or performance
- Write up findings in reports, presentations, and journals
- Programming (proficiency in R, Python, C, Java, etc.)
- Conditional probability, optimization, algorithms, statistical models, and machine learning
- Make inferences from data
- Build data products
- Find ways to do data processing, munging, and analysis at scale
- Sanity checking
- Interact with domain experts (or be a domain expert)
- Design and analyse experiments
- Find correlation in data, and try to establish causality

Data science is not limited to big technological companies, it affects other domains such as neuroscience, medicine, environmental sciences, and so on. This broad impact comes with ethical challenges that need to be addressed.

Data ethics can be defined as the branch of ethics that studies and evaluates moral problems related to data (including generation, recording,

curation, processing, dissemination, sharing and use), algorithms (including artificial intelligence, artificial agents, machine learning and robots) and corresponding practices (including responsible innovation, programming, hacking and professional codes), in order to formulate and support morally good solutions (e.g. right conducts or right values). This means that the ethical challenges posed by data science can be mapped within the conceptual space delineated by three axes of research: the ethics of data, the ethics of algorithms and the ethics of practices. — ([Luciano and Mariarosaria, 2016](#))

Therefore, data ethics is a branch of ethics that tackles data practices involving the collection, generation, analysis and dissemination of data which have the potential to impact people and society. Data ethics provides guidelines regarding the concepts of right and wrong conduct within the scope of data management. **All in all, data lies at the core of most solutions and research, serving as evidence of the studied phenomena and it is employed accordingly to draw conclusions about hypotheses and build theories.**

Whereas data is ethically neutral, the *use of data is not* ([Davis, 2012](#)). The perception regarding a particular data practice, or in other words, the employment of data for a certain task, changes over time. What once was acceptable may become unacceptable. A good example of this is the unsolicited email or SPAM. This was considered a great idea in the 1990s for marketing purposes. Overtime it became a problem and socially unacceptable. The Can't SPAM act¹⁹ was created to regulate this practice (e.g. offering a method to unsubscribe).

“Organizations that fail to explicitly and transparently evaluate the ethical impacts of the data they collect from their customers risk diminishing the quality of their relationships with those customers, exposing their business to the risks of unintended consequences.” — ([Davis, 2012](#))

For data science, data ethics asks questions such as:

- Should I write the questions differently depending on the group of people they are addressed to?
- Should I reject data if I suspect it to be falsified?
- May I collect data from participants?
- Which data should I collect exactly?
- From what sample of people?

¹⁹https://en.wikipedia.org/wiki/CAN-SPAM_Act_of_2003

5.7 General Data Protection Regulation (GDPR)

The General Data Protection Regulation (GDPR) represents an inflection point in the way organizations handle personal data. This regulation represents an intersection between legal frameworks and the current societal consensus regarding moral values in the digital age. At its core, GDPR represents a commitment to safeguarding individual privacy and promoting transparency, principles that align closely with the ethical values of autonomy, dignity, and personal rights. This regulation is a powerful reminder that in our increasingly data-driven society, ethics and data protection are inseparably linked, highlighting the need to consider the rights of data subjects in the digital realm.

The GDPR emphasizes individuals' rights to control their personal data, sets principles for lawful data processing, demands strong data security and accountability, regulates data transfers across borders, specifies the conditions to appoint Data Protection Officers, mandates reporting of data breaches, encourages privacy integration into system design, enforces accountability and governance through documentation and audits, imposes penalties for non-compliance, and emphasises the importance of obtaining clear consent for data processing. These pillars collectively form a robust framework aimed at safeguarding personal data and upholding individuals' privacy rights.

Here we will summarise some important aspects.

5.7.1 Personal data

Any data containing any identifiable datum. It is also referred as **personal identifiable information** (PII). Personal data is any data able to directly or indirectly (combined with other information) point to the individual. For instance, individually, information like university, date of birth or hair colour are not enough to directly identify anyone by themselves. But used in combination it is easily possible to find someone on social media, i.e., in combination with other databases, in this case, social media. It is difficult to provide a complete list and GDPR avoids doing so. However, some common examples include: on-line ids (e.g., IP address), names, addresses, e-mail, location information, etc. For more information, see definitions under Art. 4 of the GDPR²⁰.

²⁰<https://gdpr-info.eu/art-4-gdpr/>

5.7.1.1 Special Category Data

Some types of information require extra protection under GDPR. For instance, racial information, ethnic origin, political opinions, religious and philosophical beliefs, trade union membership, genetic and biometric data, health data, etc.

Apart from the lawful bases, 1 condition out of the following 10 must be met to process this type of data. (a) Explicit consent. (b) Employment, social security, and social protection. (c) vital interests. (d) Not for profit bodies. (e) Made public by the data subject. (f) Legal claims or judicial acts. (g) Reasons of substantial public interest. (h) Health or social care. (i) Public health. (j) Archiving, research and statistics. For example, we may need to collect ethnic data to ensure that an artificial intelligence solution has a similarly good performance across different subgroups (stratified performance assessment). This data will not be used to train the model, but rather for its evaluation. For more information see Art. 9 of the GDPR²¹.

5.7.1.2 Children's data

This kind of data is common in schools, hospitals, etc. Collection of children's data generally entails writing a privacy notice which can be understood by a child. Some additional steps may be needed to prove a person's age. For more information see Art. 8 of the GDPR²².

5.7.2 Principles GDPR

The following section is a very short summary from Chapter 2 of the GDPR²³.

5.7.2.1 Lawfulness, Fairness and Transparency

All data must be processed lawfully, fairly and transparently. **Lawfulness** means following the law (this may include other local or institutional laws apart from GDPR) and making sure that you have a lawful basis to process data. **Fairness** means that data was used as intended and expected by the data subject (an unfair use is to do SPAM without consent). **Transparency** entails honesty with the data subject by informing them (via Privacy Policy) who the data collector is, why data is needed, how it is processed, where it is stored, how it is kept safe, whether other organizations have access to the data and the rights of the data subjects.

²¹<https://gdpr-info.eu/art-9-gdpr/>

²²<https://gdpr-info.eu/art-8-gdpr/>

²³<https://gdpr-info.eu/chapter-2/>

The six lawful reasons included in the GDPR are: Consent. Necessary for contract. Legal obligation. Vital interest. Public Interest. Legitimate Interest.

5.7.2.2 Purpose Limitation

The data controller must be clear and open about the purpose and not use it for anything else. It is key to plan and understand why people's data is required. Under GDPR there is no "just in case" data collection. This helps avoiding function creep, i.e., when small changes add up over time until the way data is used changes from one purpose to another. Regular reviews. New purposes must be compatible with old ones. Three of them are automatically compatible: Archiving, scientific or historical research and statistical purposes.

5.7.2.3 Data Minimization

Adequate, relevant and limited. Only collect data necessarily for the purpose. E.g. Collect health data for job positions with both positions with added health risks and jobs without, goes against data minimization. Accuracy: Data must be accurate and kept up to date. (The most recent address, or conditions such as allergies etc). But only if it helps to ensure the rights of the data subject. Otherwise pooling the data subject too frequently may result in annoyance.

5.7.2.4 Storage Limitation

Only keep data as long as needed for a purpose. For instance, the length of a study. For this we need to consider what the data is used for, under what lawful basis, its sensitivity and the safeguards taken to protect the data in the long term. Other reasons like accountability or legal requirements may justify storing data for longer. If data is kept for longer with the purpose of archiving, research or statistics, then this data cannot be used for other purposes than those. Once the limit is reached, data must be removed or anonymised.

5.7.2.5 Integrity and Confidentiality

Data should be processed in a way that ensures appropriate security of the data (see AI examples). It is key to make sure that only authorized people can access data. Documentation regarding who has access to the systems with PII is advised. Similarly, procedures for removing users, updating passwords or expiration of permissions must be in place. Importantly, this also applies for physical security.

5.7.2.6 Accountability

This principle relates to demonstrating compliance with the other principles. It is also a way to avoid high penalties in case something goes wrong as it serves to demonstrate that best efforts were in place. Policies, guidelines, dedicated jobs (data stewards), training and tools like data catalogues help achieve high accountability. See the list of sanctions from the French organisation “Commission Nationale de l’Informatique et des Libertés” CNIL²⁴ for examples on penalties and their reasons.

5.7.3 Controllers, Processors and Subjects.

If a company is the main data decision making, when to collect data, what data to collect and how it should be processed, then this company is acting as **data controller**. Two companies sharing the same data for the same purposes are called **joint controllers** (e.g. two auditing companies). If a company processes data on behalf of another organization and acts under their instructions, then the company is acting as a **data processor**. Last, **data subjects** are physical living people, i.e., the person identified by personal data.

Data controllers bear the higher responsibility. They must comply with the regulation, demonstrate and document their compliance and ensure compliance of the data processors. Data processors must follow data controllers instructions and help the controller comply with the regulation and provide them with all necessary information. Ensure all sub-processors comply with GDPR.

5.7.4 Rights of the data subject

The following is a very short summary of Chapter 3 from the GDPR²⁵.

- **Right to be informed:** The data subjects have the right to know what happens with their data, how is collected, etc. This should be found in the privacy policy.
- **Right to access:** Data subjects have the right to request a copy of their data held by the company. For exceptions, see Recital 73 of GDPR tackling restrictions of rights and principles.
- **Rectification:** This allows data subjects to correct their data and improves data accuracy. However, data subjects may be required to provide justification and proof before making changes so that mistakes are not made upon rectification.

²⁴<https://www.cnil.fr/en/sanctions-issued-cnil>

²⁵<https://gdpr-info.eu/chapter-3/>

- **Erasure:** This right is also known as the right to be forgotten. This right has limitations. For instance, removing a data subject from all backups retrospectively. Should this be done they would not be backups any more. Therefore, a best effort must be made and in the case of backups, we should keep a list of removed data subjects to erase their records upon recovery should we need to restore a backup.
- **Right to restrict processing:** Data subjects can temporarily limit how their data is used.
- **Data portability:** Data subjects have the right to receive their data in a structured, common and machine readable format. They can also request to move their data from one data controller to another.
- **Right to object:** Allows data subject to object the processing of their data on the grounds that it is being done unlawfully.
- **Rights related to automated decision making:** Data subjects can contest automated decisions that did not involve humans with some exceptions.

5.7.5 International transfers

The regulations regarding international transfers can be regarded as a 3 stage hierarchy listing the ways to safely transfer data. As a previous step, we need to consider whether we actually need to transfer data, e.g., can we transfer the computation instead? (1) The first stage is to find if the destination country is covered by an **adequacy decision** (check EU list²⁶). If there is no adequacy decision, then stage (2) asks whether there are **appropriate safeguards** in place such as legally binding instruments, clauses, etc. (see GDPR Art 46). Stage 3 asks whether there is an **exemption** (e.g., explicit consent from data subject) in place which could allow you to transfer data but this is not suitable for regular data transfers. For more information see Chapter 5 of the GDPR²⁷.

5.7.6 Data breaches

A breach does not only occur in the event of lost or stolen data. If data is lost, stolen, corrupted, destroyed, accessed without permission, sent to the wrong person or it becomes unavailable or unusable in any way, then we can say there has been a data breach. Under GDPR in case of a data breach, organisations must **notify** their supervisory authority (data protection officer or authority) within 72 hours. They should **communicate** to the data subjects in case the breach may result in a high

²⁶https://commission.europa.eu/law/law-topic/data-protection/international-dimension-data-protection/adequacy-decisions_en

²⁷<https://gdpr-info.eu/chapter-5/>

risk to individuals' rights and freedoms. Finally, they should **keep record** of the breaches. Of course, they must **investigate and remediate** the root cause of the data breach.

5.8 Examples

5.8.1 A Genocide Incited on Facebook



Figure 5.3: Headline of the article published in The New York Times regarding the genocide of Rohingya people in Myanmar.

In 2018, a planned campaign from Myanmar military personnel turned the social network Facebook into a tool for encouraging ethnic cleansing and systematic targeting

of the Muslim Rohingya minority group. The military exploited the wide reach of Facebook among Myanmar's internet users to propagate anti-Rohingya propaganda, inciting murders, rapes. In turn, this translated into the largest forced human migration events in recent history ([Mozur, 2018](#)). 700.000 Rohingya fled the country in a year.

The chairman of the U.N. Independent International Fact-Finding Mission on Myanmar stated that Facebook played a “determining role” in the Rohingya genocide, and Facebook has been accused of enabling the spread of Islamophobic content which targets the Rohingya people ([Miles, 2018](#)). In 2021, Facebook removed accounts which were owned by the Myanmar Armed Forces and banned the account from Myanmar military from their platform. In December 2021, a group of around one hundred Rohingya refugees launched a lawsuit against Facebook.

5.8.2 The Theranos scandal - A drop of blood for hundreds of different assays



Figure 5.4: Examples of news coverage and media praising.

Theranos' goal was to quickly conduct several disease tests with a single drop of blood in a way that could be done at a fraction of the costs and with unrivalled ubiquity levering the spread of pharmacies and supermarkets. Elizabeth Holmes and Ramesh Balwani raised more than US\$700 million through “elaborate, years-long fraud in which they exaggerated or made false statements about the company’s technology, business, and financial performance”, as the US Securities and Exchange Commission put it in March this year ([Topol, 2018](#)).

Holmes studied chemical engineering at Stanford University in California but dropped out in her second year to start the company that later became Theranos. The company's board was crowded with political figures and relevant names. It included former secretaries of state George Shultz and Henry Kissinger, secretaries

of defence, and former senators. It also managed to partner with companies such as the supermarket giant Safeway, and the pharmacy chain Walgreens. Channing Robertson, a professor of Holmes's at Stanford, was also a board member and adviser to the company. She was praised by then-president Barack Obama and then-vice-president Joe Biden.

The main mistakes made by the company and its founders included a complete oversell of the product. The idea was of course great, but good ideas (or desires) are not enough. For instance, it is rather easy to say “wouldn't be nice to breathe underwater with cheap and light equipment that provided several days of autonomy?” or “wouldn't be nice to produce energy at half the cost?”. Yes, it would be amazing. But the idea is not enough, we need to show *how* we reach the goal, idea or desire we propose. That is where the strength of the idea lies. Moreover, several investors made the mistake of not sufficiently checking the technology readiness of Theranos. Hidden behind non-disclosure agreements and secrecy, Theranos managed to grow and escape inspections. The key difference with other technological start-ups is that healthcare-related companies must undertake an even greater degree of quality control and validation, since lives may be at risk.

In particular, the solution of Theranos provided both false positives and false negatives, which put at risk the lives of many patients at. In the end, Theranos required more than a single drop of blood to conduct the tests, but still less amount than the competition. However, the company diluted blood and ran the tests in machines from a different manufacturer (hence not employing even their product). Their accuracy numbers were forged by removing anomalous data points to increase the apparent performance.

[...] nearly 1 million lab tests had been run in California and Arizona. A significant proportion of these were erroneous; all had to be voided. An untold number of people were harmed by the erroneous results: some underwent unnecessary procedures, received misdiagnoses of serious conditions and experienced emotional turmoil. — (Topol, 2018)

For example, a test for Lyme disease claimed a 97% “true positive rate” – a phrase with no clear public health meaning – but actually yielded mostly inaccurate results, according to a 2015 FDA review. A false positive for Lyme disease could mean unnecessary medications and a delayed diagnosis of true underlying conditions. — The Conversation²⁸

²⁸<https://theconversation.com/how-theranos-faulty-blood-tests-got-to-market-and-what-that-shows-about-gaps-in-fda-regulation-168050>

The Netflix show, The Dropout²⁹, addresses the rise and fall of Theranos and its founder Elizabeth Holmes and provides some hints on the issues that led to the development of a scam that put the lives of real people at risk. The show is based on the podcast from ABC News of the same name³⁰.

During the rise of Theranos, Elizabeth Holmes became famous and a symbol of entrepreneurship in the USA, with the technology and economic newspapers and magazines dedicating articles and covers to her. For instance, an article from Wired titled [This Woman Invented a Way to Run 30 Lab Tests on Only One Drop of Blood]](<https://www.wired.com/2014/02/elizabeth-holmes-theranos/>). In 2015, she was appointed a member of the Harvard Medical School Board of Fellows and was named one of Time magazine's "Time 100 most influential people".

On January 3, 2022, **Holmes was found guilty on four counts of defrauding investors** – three counts of wire fraud, and one of conspiracy to commit wire fraud. She was found not guilty on four counts of defrauding patients – three counts of wire fraud and one of conspiracy to commit wire fraud. The jury returned a “no verdict” on three counts of wire fraud against investors – the judge declared a mistrial on those counts and the government soon after agreed to dismiss them. Holmes is awaiting sentencing while remaining ‘at liberty’ on \$500,000 bail, secured with property. **She faces a maximum sentence of twenty years in prison, and a fine of \$250,000, plus restitution**, for each count of wire fraud and for each conspiracy count. The sentences would likely be served concurrently thus an effective maximum of 20 years total.

Sentencing is scheduled for October 17, 2022. — Wikipedia³¹

The company and its founders broke several ethical principles and standards such as honesty, integrity, and credibility. False statements were provided to investors, employees, partners and customers (See more at SF Magazine³²). Many Medical Doctors worked at Theranos, and one of their principles from the medical Hippocratic Oath is: *primum non nocere*, i.e. *first do not to harm* (or “I will abstain from all intentional wrong-doing and harm”). Wrong methodologies were employed, poor methods were used, and finally, wrong results were knowingly introduced on purpose into the system, reaching real people and jeopardising their lives.

²⁹<https://www.theguardian.com/tv-and-radio/2022/mar/09/how-theranos-drama-the-dropout-gets-scam-and-tech-culture-right>

³⁰<https://abcaudio.com/podcasts/the-dropout/>

³¹https://en.wikipedia.org/wiki/Elizabeth_Holmes#U.S._v._Holmes,_et_al.

³²<https://sfmagazine.com/post-entry/june-2022-theranos-cautionary-tale-of-ethical-failings/>

Chapter 6

Extra Material

This chapter discusses related questions and topics not tackled in this course. As in any course, being exhaustive in the content taught to students is not feasible due to time constraints and the need for curricular coherence. Topics cannot be included merely for the sake of completeness, they must be introduced with clear context, purpose, and alignment with the program's learning objectives. For this reason, several rich but complex philosophical debates are excluded from the core syllabus. These include Thomas Kuhn's theory of scientific revolutions, the debate between realism and anti-realism, the Gettier problem, the pessimistic meta-induction, discussions on instrumentalism, and detailed philosophical debates over the aims of science.

While some of these ideas, such as instrumentalism or the plurality of scientific aims, are occasionally mentioned informally to enrich class discussion or highlight philosophical context, the course does not dedicate time to formally exploring them in depth.

For instance, Kuhn's notion of paradigm shifts offers a compelling sociological view of scientific change, but properly integrating it would need extended dedication to the history and sociology of science. Similarly, the realism debate and the pessimistic meta-induction ask deep questions about whether scientific theories aim at truth or are merely useful fictions. The Gettier problem, while central to analytic epistemology, has little relevance for applied reasoning in data contexts. Likewise, instrumentalism, which treats theories as predictive tools rather than truth claims, could be a fruitful lens but would need to be introduced alongside contrasting positions to avoid confusion.

6.1 History of science

Although many historical examples are given to the students to illustrate the different chapters and topics addressed in this course, the history of science is not part of the course curriculum. However, we aim to include a rich but brief summary of the history of science as optional reading in further semesters of this course.

6.2 Theory-relatedness of observations

In the philosophy of science, observations are said to be “theory-laden” when they are affected by the theoretical presuppositions held by the investigator. This thesis of theory-ladenness is associated with the works of Thomas Kuhn and perhaps first put forth by Pierre Duheem ([Boyd and Bogen, 2021](#)).

“A related topic is the theory-relatedness of observations; some have claimed that there are no such things as fully theory-independent observations. If true, it would undermine the possibility of objectivity of science and force us to accept strong relativism. I believe that this disastrous consequence can be avoided and that there really is a basis of theory-neutral data, also in the humanities.” – ([Johansson et al., 2016](#)).

The question then arises: is it just as easy to distinguish between theoretical statements and observational statements? The answer is no, as can be seen from the previous examples regarding how unconscious background beliefs can affect what is observed and reported even in a very simple tasks such as time measurements. – ([Johansson et al., 2016](#)).

6.3 Thomas Kuhn and the idea of scientific revolutions

The Duhem problem points out that neither the confirmation nor falsification of a single hypothesis is as clear and unequivocal as it might be supposed to. Accounts of the nature of science, such as Thomas Kuhn’s view, diverge from inductivism and falsificationism, suggesting that the progress of science is not a continuous accumulation of knowledge imposed by observational evidence. They deny as well that scientists are purely critical rationalists prepared to renounce to their theoretical commitments when experiments contradict them.

Thomas Kuhn describes the history of science as periods of conservative scientific activity disrupted by revolutions, highlight the role of social factors in this process.

Many scholars see Kuhn's view as motivating "social constructivism", which understands scientific knowledge as a cultural product instead of the pure discovery of better and better approximation to the truth. According to the view of Thomas Kuhn, scientific theories are socially negotiated instead of purely determined by nature and experiments.

The strength of social constructivism depends, therefore, on how much freedom is allowed for the construction of theories by the social factors affecting science. Social constructivists appeal to the underdetermination argument, as it shows that evidence is always compatible with multiple theories.

6.4 Gettier problems

The definition of knowledge is an ongoing debate among epistemologists. Although the three criteria from Plato are necessary conditions, they are not sufficient as there are situations that satisfy all these conditions and yet don't constitute knowledge (see Gettier cases¹) but such cases are rather philosophical and will not be discussed during this course.

6.5 Realism and anti-realism

For now this will not be included as part of the course curriculum. For a short account of this topic, read Chapter 4 from ([Okasha, 2016](#)).

6.6 Pessimistic meta-induction

One of the most compelling arguments against scientific realisms is the 'Pessimistic meta-induction' argument, by Larry Laudan. Instead of appealing to the underdetermination argument, it appeals to history. Recalling Induction to the Best Explanation (IBE), there is a connection between the success of scientific theories and their truth, for which scientific realism offers the only, or the best, explanation of the progress of science. However, Laudan turns this argument around and argues that we have good reasons, by induction, for not believing in the existence of the theoretical entities described by our current theories ([Ladyman, 2012](#)). Laudan then proceeds to enumerate a number of now abandoned theories that once had predictive and explanatory success (e.g. phlogiston, ether...).

¹https://en.wikipedia.org/wiki/Gettier_case

“Therefore, we should not believe in the approximate truth or the successful reference of the theoretical terms of our best current theories”.

My personal take on this, is as follows: It is argued that from a future perspective we will criticise and debunk our current theories as we now do with past theories. However, for this to be true, we must necessarily assume that we will gain knowledge over time. Such knowledge will be used to debunk, improve or replace the previous theories. Notably, this knowledge will mostly be obtained thanks to prior theories. For instance, the rejection of luminiferous aether theory happened thanks to an experiment that was designed with the help of previous theories. Scientists follow their theories to their limits, and when found, they require theory reformulation or replacement. In the same way the Americas were found on the assumption that they would reach East Indies. Failures lead to new discoveries and theories are always under development, is not just theories what realists use as truth but rather all the knowledge that transcends in the process of developing theories, including, of course, the theories themselves.

Bibliography

- AI, H. (2019). High-level expert group on artificial intelligence.
- Alexander, L. and Moore, M. (2021). Deontological Ethics. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Winter 2021 edition.
- Barker, G. and Kitcher, P. (2014). *Philosophy of science: A new introduction*. Oxford University Press New York.
- Beall, J., Restall, G., and Sagi, G. (2024). Logical Consequence. In Zalta, E. N. and Nodelman, U., editors, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Summer 2024 edition.
- Bergadano, F. (1991). The problem of induction and machine learning. In *IJCAI*.
- Boyd, N. M. and Bogen, J. (2021). Theory and Observation in Science. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Winter 2021 edition.
- Brennan, A. and Lo, N. Y. S. (2022). Environmental Ethics. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Summer 2022 edition.
- Bunge, M. (2012). *Method, model and matter*, volume 44. Springer Science & Business Media.
- Bunge, M. (2017). *Philosophy of science: volume 2, from explanation to justification*. Routledge.
- Bynum, T. (2018). Computer and Information Ethics. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Summer 2018 edition.
- Cadeddu, M., Farrokhyar, F., Thoma, A., Haines, T., Garnett, A., Goldsmith, C. H.,

- Group, E.-B. S. W., et al. (2008). Users' guide to the surgical literature: how to assess power and sample size. *Canadian Journal of Surgery*, 51(6):476.
- Cairo, A. (2019). *How charts lie: Getting smarter about visual information*. WW Norton & Company.
- Cohon, R. (2018). Hume's Moral Philosophy. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Fall 2018 edition.
- Coles, P. (2019). Einstein, eddington and the 1919 eclipse. *Nature*, 568(7752):306–308.
- Crasnow, S. (2020). Feminist Perspectives on Science. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Winter 2020 edition.
- Davis, K. (2012). *Ethics of Big Data: Balancing risk and innovation*. " O'Reilly Media, Inc.".
- de Vocht, F., Katikireddi, S. V., McQuire, C., Tilling, K., Hickman, M., and Craig, P. (2021). Conceptualising natural and quasi experiments in public health. *BMC medical research methodology*, 21(1):1–8.
- Deaton, A. and Cartwright, N. (2018). Understanding and misunderstanding randomized controlled trials. *Social Science & Medicine*, 210:2–21.
- Deville, P. (2014). *Plague and cholera*. Hachette UK.
- Díez, J. A. and Moulines, C. U. (1997). Fundamentos de filosofía de la ciencia.
- DiNardo, J. (2010). Natural experiments and quasi-natural experiments. In *Microeconometrics*, pages 139–153. Springer.
- Domingos, P. (2012). A few useful things to know about machine learning. *Communications of the ACM*, 55(10):78–87.
- Domingos, P. (2015). *The master algorithm: How the quest for the ultimate learning machine will remake our world*. Basic Books.
- Dotan, R. (2020). Theory choice, non-epistemic values, and machine learning. *Synthese*, pages 1–21.
- Douven, I. (2021). Abduction. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Summer 2021 edition.

- Eubanks, V. (2018). *Automating inequality: How high-tech tools profile, police, and punish the poor*. St. Martin's Press.
- Fischer, S. (2020). The necessitiy of non-epistemic values in machine learning modelling.
- Frigg, R. and Hartmann, S. (2020). Models in Science. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Spring 2020 edition.
- Garcia, B., Vega, C., and Hertel, F. (2022). The need of standardised metadata to encode causal relationships: Towards safer data-driven machine learning biological solutions. In *International Meeting on Computational Intelligence Methods for Bioinformatics and Biostatistics*, pages 200–216. Springer.
- Gert, B. and Gert, J. (2020). The Definition of Morality. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Fall 2020 edition.
- Giere, R. N., Bickle, J., and Mauldin, R. F. (2006). Understanding scientific reasoning.
- Goodman, S. and Greenland, S. (2007). Why most published research findings are false: problems in the analysis. *PLoS medicine*, 4(4):e168.
- Griffith, G. J., Morris, T. T., Tudball, M. J., Herbert, A., Mancano, G., Pike, L., Sharp, G. C., Sterne, J., Palmer, T. M., Davey Smith, G., et al. (2020). Collider bias undermines our understanding of covid-19 disease risk and severity. *Nature communications*, 11(1):1–12.
- Grüne-Yanoff, T. (2014). Teaching philosophy of science to scientists: why, what and how. *European Journal for Philosophy of Science*, 4(1):115–134.
- Grünbaum, A. (1976). Ad hoc auxiliary hypotheses and falsificationism. *The British Journal for the Philosophy of Science*, 27(4):329–362.
- Hansen, J. A. and Tummers, L. (2020). A systematic review of field experiments in public administration. *Public Administration Review*, 80(6):921–931.
- Hansson, S. O. (2021). Science and Pseudo-Science. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Summer 2021 edition.
- Hempel, C. G. (1966). *Philosophy of natural science*. Prentice-Hall Englewood Cliffs, N.J.

- Henderson, L. (2024). The Problem of Induction. In Zalta, E. N. and Nodelman, U., editors, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Winter 2024 edition.
- Hershey, D. R. (1991). Digging deeper into helmont's famous willow tree experiment. *The American Biology Teacher*, 53(8):458–460.
- Howson, C. and Urbach, P. (2006). *Scientific reasoning: the Bayesian approach*. Open Court Publishing.
- Hume, D. (1739). *A Treatise Upon Human Nature*. Oxford University Press, Oxford.
- Ioannidis, J. P. (2005). Why most published research findings are false. *PLoS medicine*, 2(8):e124.
- Jager, L. R. and Leek, J. T. (2014). An estimate of the science-wise false discovery rate and application to the top medical literature. *Biostatistics*, 15(1):1–12.
- Johansson, L.-G. et al. (2016). *Philosophy of science for scientists*. Springer.
- Jun, S. (2016). Frequentist and bayesian learning approaches to artificial intelligence. *International Journal of Fuzzy Logic and Intelligent Systems*, 16(2):111–118.
- Koons, R. (2021). Defeasible Reasoning. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Fall 2021 edition.
- Krause, M. S. and Howard, K. I. (2003). What random assignment does and does not do. *Journal of Clinical Psychology*, 59(7):751–766.
- Kubat, M. (2017). *An introduction to machine learning*. Springer.
- Ladyman, J. (2012). *Understanding philosophy of science*. Routledge.
- Loukides, M., Mason, H., and Patil, D. (2018). *Ethics and data science*. O'Reilly Media.
- Luciano, F. and Mariarosaria, T. (2016). What is data ethics. *Phil. Trans. R. Soc. A*. 37420160360.
- Mayo, D. G. (1996). *Error and the growth of experimental knowledge*. University of Chicago Press.
- Mayo-Wilson, C. (2021). Statistical inference as severe testing: How to get beyond the statistics.
- Metcalfe, B. (2013). Metcalfe's law after 40 years of ethernet. *Computer*, 46(12):26–31.

- Miles, T. (2018). U.N. investigators cite Facebook role in Myanmar crisis.
- Montelpare, W. J. (2021). *Applied Statistics in Healthcare Research*. University of Prince Edward Island.
- Morgan, K. L. and Rubin, D. B. (2012). Rerandomization to improve covariate balance in experiments. *The Annals of Statistics*, 40(2):1263–1282.
- Mozur, P. (2018). A Genocide Incited on Facebook, With Posts From Myanmar's Military.
- Nidditch, P. H. (1968). *The philosophy of science; edited by P. H. Nidditch*. Oxford U.P London.
- Norlock, K. (2019). Feminist Ethics. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Summer 2019 edition.
- Okasha, S. (2016). *Philosophy of science: very short introduction*. Oxford University Press.
- Olsson, E. (2025). Coherentist Theories of Epistemic Justification. In Zalta, E. N. and Nodelman, U., editors, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Fall 2025 edition.
- O’Neil, C. and Schutt, R. (2013). *Doing data science: Straight talk from the frontline*. ” O'Reilly Media, Inc.”.
- Orloff, J. and Bloom, J. (2014). Comparison of frequentist and bayesian inference. class 20, 18.05, spring 2014.
- Ortner, R. (2023). Adaptive algorithms for meta-induction. *Journal for General Philosophy of Science*, 54(3):433–450.
- Paludan-Müller, A., Laursen, D. R. T., and Hróbjartsson, A. (2016). Mechanisms and direction of allocation bias in randomised clinical trials. *BMC medical research methodology*, 16(1):1–10.
- Pearl, J. (2000). Causality: models, reasoning and inference cambridge university press. *Cambridge, MA, USA*, 9:10–11.
- Pearl, J. (2014). *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier.
- Pearl, J. and Mackenzie, D. (2018). *The Book of Why: The New Science of Cause and Effect*. Basic Books, Inc., USA, 1st edition.

- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Percival, R. S. (2015). Confirmation versus falsificationism.
- Pokropek, A. (2016). Introduction to instrumental variables and their application to large-scale assessment data. *Large-scale Assessments in Education*, 4(1):1–20.
- Potochnik, A. (2015). The diverse aims of science. *Studies in History and Philosophy of Science Part A*, 53:71–80.
- Rosenberg, A. and McIntyre, L. (2019). *Philosophy of science: A contemporary introduction*. Routledge.
- Rossi, P. H., Freeman, H. E., and Wright, S. R. (1985). Evaluation: a systematic approach. beverly hills.
- Rothwell, P. M. (2006). Factors that can affect the external validity of randomised controlled trials. *PLoS clinical trials*, 1(1):e9.
- Russell, B. (1912). *The problems of philosophy*.
- Santa Cruz, B. G., Bossa, M. N., Sölter, J., and Husch, A. D. (2021). Public covid-19 x-ray datasets and their impact on model bias—a systematic review of a significant problem. *Medical image analysis*, 74:102225.
- Skansi, S. (2020). *Guide to Deep Learning Basics*. Springer.
- Strassner, C. and Antonelli, G. A. (2019). Non-monotonic Logic. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Summer 2019 edition.
- Thornton, S. (2021). Karl Popper. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Spring 2021 edition.
- Topol, E. (2018). Blood, sweat and tears in biotech—the theranos story. *Nature*, 557(7706):306–308.
- Tsivgoulis, G., Katsanos, A. H., and Caso, V. (2017). Under-representation of women in stroke randomized controlled trials: inadvertent selection bias leading to suboptimal conclusions. *Therapeutic advances in neurological disorders*, 10(5):241–244.
- Van Westrum, A. S. (1908). *The Devil*. M.A. Donohue, Chicago. Founded on Ferenc Molnar's play, as produced by Harrison Grey Fiske at the Belasco Theatre, New York.

- Vega, C. (2021). From hume to wuhan: an epistemological journey on the problem of induction in covid-19 machine learning models and its impact upon medical research. *IEEE Access*, 9:97243–97250.
- Walsh, I., Fishman, D., Garcia-Gasulla, D., Titma, T., Pollastri, G., Harrow, J., Psomopoulos, F. E., and Tosatto, S. C. (2021). Dome: recommendations for supervised machine learning validation in biology. *Nature methods*, 18(10):1122–1127.
- Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.-W., da Silva Santos, L. B., Bourne, P. E., et al. (2016). The fair guiding principles for scientific data management and stewardship. *Scientific data*, 3(1):1–9.
- Yoon, D. Y., Mansukhani, N. A., Stubbs, V. C., Helenowski, I. B., Woodruff, T. K., and Kibbe, M. R. (2014). Sex bias exists in basic science and translational surgical research. *Surgery*, 156(3):508–516.