Scientific methods and data ethics

Bertil Wegmann

Division of Statistics and Machine Learning Dept. of Computer and Information Science Linköping University

2024-08-28

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Outline

Summary

What is science? What is knowledge?

Hypothetico-deductive method

Models

Scientific Revolutions

Ethics in Science and Statistics

Summary

◆□▶ ◆□▶ ◆ □▶ ◆ □ ● ● ● ●

Summary

Take home message:

- Be critical! Question it!
 - data, plots, graphs, tables
- Science is hard
- Correlation does not imply causation
- Who is behind the results? What are the incentives?

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● ○ ○ ○

- money, power, prestige, reputation
- Don't do bad things with data: Ethics matter
 - "With big data comes big responsibility"

Intro

This lecture

science = "all science", not just "natural science"

(ロ)、(型)、(E)、(E)、 E) のQ()

A smorgasbord of different topics

Intro

Apartment prices in Linköping



Figure: From Svensk Mäklarstatistik

▲□▶ ▲□▶ ▲□▶ ▲□▶ = 三 のへで

Intro

Apartment prices in Linköping



Figure: From Svensk Mäklarstatistik

▲□▶ ▲□▶ ▲□▶ ▲□▶ = 三 のへで

Example

Consider the observed dataset

$$x = \left(\begin{array}{ccccc} 0 & 1 & 2 & 3\end{array}\right) \qquad y = \left(\begin{array}{cccccc} 2 & 3 & 4 & 5\end{array}\right)$$

Problem: We want to understand the relation between x and y. The assumption y = f(x) is given. Which function f should be used? Why?

 \rightarrow Science to the rescue!

- "Large amount of relatively secured knowledge"
- Results description of facts or explanation of dependencies
 - Results are often published in scientific journals, conferences or books

- Peer review
- ▶ A process the methods and activities that lead to the results
 - E.g. Experiments, computations, theories

What is science?

- Methods are established within the scientific community
- Scientific methods have clear and explicit rules and procedures
- Replication is important: details matter
 - Methods that are arbitrary and cannot be repeated are not scientific
- Science should be as objective as possible
 - Researcher bias should be reduced
- Scientific knowledge is created within the scientific community
- All new results should be related to existing knowledge in the field

Categories

- Nomothetic (general) studies
 - General laws
- Idiographic (specific) studies
 - Describes specific objects and processes
 - What happened at the Battle of Hastings?
- Formal science
 - The study of constructed and formal systems
 - Logic, mathematics, statistics
- Empirical science
 - empirical evidence gained through experimentation or collection of data
 - The objects and processes are studied in the "real world"
 - Medicine, history, economics

Classic definition of knowledge

- Classic definition of knowledge (Plato): Justified true belief
- A knows that a proposition, P, holds if and only if the three following conditions are met:
 - P is true,
 - A believes that P is the case
 - A is justified in believing that P is the case.

Epistemology

Rationalism: knowledge through logic and deductive reasoning

Empiricism: knowledge through empirical evidence

Knowledge

Correspondence theory of truth

- Truth or falsity of a statement is determined only by how it relates to the world
- Coherence theory of truth
 - A statement is true if it is coherent with other statements in the system under study

Pragmatism

A statement is true if the predicted consequences are correct

Scientific explanations

Deductive explanations

- Based on a number of premises
- With help of the premises, conclusions are deduced with the help of logic

If the premises are true, then the conclusions are also true

- Probabilistic explanations
 - No general law. Premises that have a (high) probability to be true (or to happen)
 - The "probabilistic conclusions" are not true in a formal way, but may be probable

Probability theory and statistical inference are used to formalize the probabilistic explanations

Philosophy of Statistics

Probability:

- Relative frequencies in series of events
- Degree of belief, epistemic
- Statistical inference
 - Frequentist statistics
 - Bayesian statistics
- Philosophy of Statistics:

https://plato.stanford.edu/entries/statistics/

Induction

Induction:

- From the observed (specific) to the general
- Based on verification of theories
- Can we be absolutely sure that...
 - The sun will rise tomorrow?
 - All swans are white?
- The Problem of Induction
 - We cannot logically verify cause and effect based on induction

Circular reasoning:

- the future will resemble the past, because we know that
- the future will resemble the past, because we know that

Hypothetico-deductive method

- Axiomatic method: Summarizes and justifies a large number of statements by proving that they can be deduced from a small number of axioms
- Problems:
 - How to justify the axioms?
 - How to relate to the "real world"?
- Karl Poppers suggestion: Hypothetico-deductive method
 - Focus on hypotheses
 - Falsification
 - Popper: Psychoanalysis cannot be considered a science

It is not falsifiable

Hypothetico-deductive method

- Hypothesis: a proposed explanation for a phenomenon
 - Scientific hypothesis: should be testable (falsifiable)
- Hypothesis: a statement that fulfills the requirements:
 - We are not sure if it is true
 - We can deduce logical consequences from it, with the purpose of:

- Testing it
- Making predictions
- Explaining observed facts

Hypothetico-deductive method

- Infinite number of observations cannot verify the hypothesis
- Degree of belief
 - Hypotheses are *confirmed* to different degrees, when we fail to falsify them
 - The best option with current knowledge, but we do not know if they are true

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Auxiliary hypotheses

- Theoretically we only need one incorrect observation to falsify a hypothesis
- Auxiliary hypothesis: Hypothesis necessary for the deduction of a prediction, but which is not tested in the study
- Reality is complex: occasional incorrect observations are often explained as measurement error or are assumed to depend on some of the auxiliary hypotheses.

Hypothetico-deductive system

- Axiomatic system: consequences are absolutely sure based on the premises
- Hypothetico-deductive system = a deductive system with hypotheses as starting point
 - All principles and statements that are not entirely secure are called hypotheses
 - Justification is based on that the consequences deduced from the hypotheses conforms with our experience

We are never entirely sure about the truth value of our hypotheses

Creation of hypotheses: Example

Consider the observed dataset

$$x = \left(\begin{array}{ccccc} 0 & 1 & 2 & 3\end{array}\right) \qquad y = \left(\begin{array}{cccccc} 2 & 3 & 4 & 5\end{array}\right)$$

Problem: We want to understand the relation between x and y. The assumption y = f(x) is given. Which function f should be use?

Suggestions:

$$H_1: y = x + 2$$

$$H_2: y = x^4 - 6x^3 + 11x^2 - 5x + 2$$

$$H_3: y = x^5 - 4x^4 - x^3 + 16x^2 - 11x + 2$$

Ref: "Argumentationsteori, språk och vetenskapsfilosofi"

Creation of hypotheses:

- Hypotheses are created as "guesses"
 - Deduce the consequences of hypotheses with logic

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで

Choice of hypotheses: Simplicity

- Several (infinite number) hypotheses can be used to explain the same observations.
 - We must choose!
 - Ex: Regression, what is a good fit to data?
- Criteria:
 - Security: Should we be aggressive or defensive when we create hypotheses?

- Strength: How much can the hypothesis explain?
- Simplicity: Most important criteria

Choice of hypotheses: Simplicity

In practice researchers often choose the simplest hypothesis.

- Aesthetic
- Makes understanding easier
- Occam's razor: If H_1 and H_2 do the same prediction

Choose the one that has fewest assumptions

Poppers recipe:

- Boldness in the guesses/hypotheses and stringency in rejection
- You should declare precisely under which circumstances you are willing to abandon your hypothesis

Models

- A model can be compared with a map
- What kind of map is important for the following persons?

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- 🕨 taxi driver
- orienteerer
- epidemiologist

Models

- "All models are wrong. Some are useful." George E.P. Box
- Science often speaks about models
- Model: a representation of a process or a system
 - Important features are a part of it (often emphasized)
 - Other features are not included
- Historically: mechanical models were important
- Nowadays: Theoretical or mathematical models are much more important.

Models

Purpose of scientific models:

- Understand, define, quantify, visualize or simulate a process or a system
- Common approach when working with complex problems
- Calculations and predictions become easier
- A theory is always a model: makes a phenomenon understandable
 - All models are not theories, e.g. mechanical model
- Causal models
 - Describes the causal mechanisms of a system.
 - Causal diagram is a directed graph that displays causal relationships

- Confounding factor: important in statistics
 - Correlation does not imply causation

Scientific Revolutions

Research is often a cumulative process

- Continuous revision of old knowledge
- Sometimes there are revolutions: Old theories are rejected and replaced with new ones
 - Chemical revolution: Lavoisier
 - Scientific theory of evolution: Darwin
 - Theory of relativity: Einstein
 - Quantum mechanics
 - Convolutional neural network and deep learning within image classification (2012)

▲□▶ ▲□▶ ▲三▶ ▲三▶ - 三 - のへの

Scientific Revolution

Thomas Kuhn: "The structure of Scientific Revolutions", 1962

- ► Normal science → revolution and crisis → Normal science
- Paradigm = central hypotheses
 - researchers are laying a puzzle
 - After a while: to many pieces do not fit
 - A revolution happens when the central hypotheses are rejected: A new puzzle start

- It can be hard in practice to define and observe scientific revolutions.
 - When does a hypothesis become a paradigm?
- Paradigms can be subjects of "religious" belief

Ethics in science

Ethics or moral philosophy:

- Deals with what is right or wrong
- How to act?
- Research ethics: How to handle moral issues that arise during or as a result of research activities

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Ethics in science

Bad examples from history:

- Nazi human experimentation
 - A large number of prisoners were forced to participate, the experiments typically lead to death, trauma, permanent disability etc.

Lead to the Nuremberg Code after the Nuremberg trials

- Tuskegee syphilis experiment (US, 1932-1972)
 - African-American men were used to see the effect of untreated syphilis infection, without consent of the participants
- Vipeholm experiments Vipeholm experiments (Sweden, 1945-1955):
 - Intellectually disabled were fed with sweets in order to provoke dental caries, the aim was to determine the role of carbohydrates

Ethics in science

- A researcher's work is regulated by rules and regulations
- Researcher's own ethical responsibility that
 - Research has good quality
 - Is morally acceptable
- Professional Ethics
 - Research activity is driven by a number of implicit and explicit norms that decide what good science is, e.g. Helsinki Declaration about ethical principles regarding human experimentation

- Follow national and local rules: issues like discrimination, harassment and humiliation, gifts to the researcher
- Field specific codes of ethics: Ethical code

Ethics in Statistics

- American Statistical Association (USA): "Ethical guidelines for statistical practice"
- Royal Statistical Society (UK): "Code of conduct"
- International Statistical Institute: "Declaration of Professional Ethics"

Swedish Statistical Society: "Svenska statistikfrämjandets etiska kod för statistiker och statistisk verksamhet"

Ethics in Statistics

"Declaration of Professional Ethics"

- Pursuing Objectivity
- Clarifying Obligations and Roles
- Assessing Alternatives Impartially
- Conflicting Interests
- Avoiding Preempted Outcomes
- Guarding Privileged Information
- Exhibiting Professional Competence
- Maintaining Confidence in Statistics

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● ○ ○ ○

Ethics in Statistics

"Ethical Guidelines for Statistical Practice"

- Professional Integrity and Accountability
- Integrity of data and methods
- Responsibilities to Science/Public/Funder/Client
- Responsibilities to Research Subjects
- Responsibilities to Research Team Colleagues
- Responsibilities to Other Statisticians or Statistics Practitioner
- Responsibilities Regarding Allegations of Misconduct
- Responsibilities of Employers, Including Organizations, Individuals, Attorneys, or Other Clients Employing Statistical Practitioners

Handling of data has a special role in statistics and data science.

Special care must be taken when data is collected, stored and used for statistics and machine learning.

General Data Protection Regulation (GDPR)

GDPR

Look here and here.

Scope: "The General Data Protection Regulation exists to protect individuals' fundamental rights and freedoms, in particular their right to protection of their personal data."

- ► GDPR:
 - EU law on data protection
 - Regulate the use of personal data
 - Called "Dataskyddsförordningen" in Sweden

"In today's most common digital business model, consumers pay for 'free' products with their personal data."

from: Big data, artificial intelligence, machine learning and data protection

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Ethics in Big Data

A few starting principles

- 1. Ownership: Individuals own their own data
- 2. Transaction Transparency: The use of the data should be transparent
- 3. Consent: informed and explicitly expressed consent is needed to use the data
- 4. Privacy
- 5. Currency: Individuals should be aware of financial transactions resulting from the use of their personal data and the scale of these transactions

6. Openness: Aggregated data sets should be freely available

Ethics in Big Data

5 Principles for Big Data Ethics from "Towards Data Science":

- Private customer data and identity should remain private: private data obtained from a person with their consent should not be exposed for use by other businesses or individuals with any traces to their identity.
- Shared private information should be treated confidentially: Third party companies share sensitive data — medical, financial or locational — and need to have restrictions on whether and how that information can be shared further.
- Customers should have a transparent view of how our data is being used or sold, and the ability to manage the flow of their private information across massive, third-party analytical systems.
- Ref [here]

5 Principles for Big Data Ethics from "Towards Data Science":

Big Data should not interfere with human will

- Big Data should not institutionalize unfair biases like racism or sexism. Machine learning algorithms can absorb unconscious biases in a population and amplify them via training samples.
 - Example: Microsoft's Twitter chatbot "Tay": the robot began releasing racist and sexually-charged messages

Ref [here]

▲□▶ ▲圖▶ ▲国▶ ▲国▶ - 国 - のへで

- Facebook-Cambridge Analytica data scandal: Big political scandal in early 2018
- The company Cambridge Analytica had collected personal data from millions of peoples' Facebook profiles

- without consent
- used it for political advertising purposes.
- Whistle-blower: Christopher Wylie

- Aleksandr Kogan researcher at Cambridge University created an app
 - "This Is Your Digital Life"
- Several hundred thousands of Facebook users gave consent to be part of the survey only for academic use.
- Facebook's design allowed data to be collected from the social network of the participants
 - This allowed Cambridge Analytica to collect data from up to 30/50/87 million users

Cambridge Analytica used the data to

- Create psychographic profiles of Facebook users
- Profiles used to choose advertisement that most effectively persuade specific groups of persons
- Used in political campaigns with the aim to affect elections, examples
 - 2016 United States presidential election
 - 2016 United Kingdom European Union membership referendum

Many other countries and elections

Discussion

- What is ethical to do with user data on social media platforms?
- "Personal data as gold": companies using data as main source of profit, what is good ethics in such business? Do the users understand what their data are used for?
- What responsibilities does a machine learner or data scientist working for a social media company have?
 - What to do if your boss asks you to do something that maybe feels wrong? Eg. collect or analyze personal data when it is unclear if consent is given
- Is it right to produce a machine learning system (in a democratic country) and then sell the system to totalitarian regime, who wants to use the system control its citizens?

Summary

Take home message:

- Be critical! Question it!
 - data, plots, graphs, tables
- Science is hard
- Correlation does not imply causation
- Who is behind the results? What are the incentives?

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● ○ ○ ○

- money, power, prestige, reputation
- Don't do bad things with data: Ethics matter
 - "With big data comes big responsibility"

References |

Books:

- Ladyman, James, Understanding Philosophy of Science, Routledge, London, 2002
- Dagfinn Föllesdal, Lars Wallöe, Jon Elster, Argumentationsteori, språk och vetenskapsfilosofi, Thales, Stockholm, 2001

Data Ethics – The New Competitive Advantage, [link]

References II

Links

- Stanford Encyclopedia of Philosophy [here]
 - Philosophy of Statistics [here]
 - Scientific Method [here]
 - Science and Pseudo-Science [here]
 - Scientific Progress [here], Scientific Revolutions [here]
 - The Problem of Induction [here]
 - Bayes' Theorem [here]
- Probability and Induction [here]
- CODEX website rules and guidelines for research [here]
- Big data, artificial intelligence, machine learning and data protection [here]

DataEthics: [dataethics.eu]