

Quantitative Methods 2, ZHAW

Jürgen Degenfellner

2025-01-14

Contents

1	Introduction	5
1.1	Books we will heavily borrow from are:	5
2	Introduction	7
2.1	What is statistical modeling and what do we need this for? . . .	7
3	Literature	11
4	Methods	13
4.1	math example	13
5	Applications	15
5.1	Example one	15
5.2	Example two	15
6	Final Words	17

Chapter 1

Introduction

This script is a continuation of the first one for Quantitative Methods 1 at ZHAW.

In the first part, we learned about the basics of probability theory, descriptive statistics, Bayesian statistics, and hypothesis testing.

In this script, we will dive into the basics of statistical modeling - a world of aesthetic wonder and surprises.

This script is a work in progress and will be updated as we go along.

Please feel free to send me suggestions for improvements or corrections.

As the first one, this **should be a collaborative effort** and will (hopefully) never be finished as our insight grows over time.

For the working setup with R, please see this and the following sections in the first script.

1.1 Books we will heavily borrow from are:

- (Free) Statistical Rethinking, YouTube-Playlist: Statistical Rethinking 2023
- (Free) Understanding Regression Analysis: A Conditional Distribution Approach
- Data Analysis Using Regression and Multilevel/Hierarchical Models
- (Free) Doing Bayesian Data Analysis

Chapter 2

Introduction

2.1 What is statistical modeling and what do we need this for?

Typically, one simplifies the complex reality (and loses information) in order to make it better understandable, mathematically treatable and to make predictions.

Underlying our models, there are theories which should be falsifiable and testable. For instance, I would be really surprised if I pull up my multimeter and measure the voltage (V) and electric current (I) at a resistance (R) in a circuit and find that Ohm's law $V = IR$ is not true. This **law** can be tested over and over again and if one would find a single valid counterexample, the law would be falsified. It is also true that the law is probably not 100% accurate, but an extremely good approximation of reality. Real-world measurements carry measurement errors and when plotting the data, one would see that the data points might not lie exactly on a straight line. This is not a problem.

A statistical model is a mathematical framework that represents the relationships between variables, helping us understand, infer, and predict patterns in data. It acts as a bridge between observed data and the real-world processes that generated them. In health research, where variability and uncertainty are inherent, statistical models are valuable tools for making sense of complex phenomena. You can watch this as short intro.

Depending on the task at hand, we would use different models. In any case, logical reasoning and critical thinking comes first, then comes the model. **It makes no sense to estimate statistical models just for the sake of it.**

All models are wrong, but some are useful. Or to quote George Box:

“Since all models are wrong the scientist cannot obtain a ‘correct’ one by excessive elaboration. On the contrary following William of Occam he should seek an economical description of natural phenomena. Just as the ability to devise simple but evocative models is the signature of the great scientist so overelaboration and overparameterization is often the mark of mediocrity.”

In my opinion, statistical model is an art form: difficult and beautiful.

One goal of this course is to improve interpretation and limitations of statistical models. They are not magical turning data into truth. Firstly, the rule garbage in, garbage out (GABA) applies. Secondly, statistical models are based on data and their variability and have intrinsic limitations one cannot overcome even with the most sophisticated models. This is expressed for instance in the so-called bias-variance trade-off. You can’t have it all.

2.1.1 Explanatory vs. Predictive Models

I can recommend reading this article by Shmueli et al. (2010) on this topic.

Statistical models serve different purposes depending on the research question. Two primary goals are **explanation** and **prediction**, and each requires a different approach:

Explanatory Models focus on understanding causal relationships. These models aim to uncover mechanisms and answer “**why**” questions. For example:

- Does smoking increase the risk of lung cancer? **Yes**. (If you want to see what a large effect-size looks like, check out this study.)
- How large is the “effect” of smoking on lung cancer? **Large**.
- Does pain education and graded sensorimotor relearning improve disability (a question we ask in our Resolve Swiss project)?

Explanatory models are **theory-driven**, designed to test hypotheses. Here, one wants to understand the underlying mechanisms and the relationships between variables and hence often uses (parsimonious) models that are more interpretable, like linear regression.

Predictive Models prioritize forecasting future outcomes based on patterns in the data. These models aim to answer “**what will happen?**” For instance:

- Gait analysis using Machine Learning (ML)?
- Skin cancer detection using neural networks?

Predictive models are **data-driven**, often using complex algorithms to achieve high accuracy. Their success is measured using metrics like Root Means Square

2.1. WHAT IS STATISTICAL MODELING AND WHAT DO WE NEED THIS FOR?9

Error (RMSE), Area Under the Curve (AUC), or **prediction error on new, unseen data**. Any amount of model complexity is allowed. One could for instance estimate a neural network (“just” another statistical model) with many hidden layers and neurons in order to improve prediction quality. Interpretability of the model weights is not a priority here.

While explanatory and predictive goals often complement each other, their differences highlight the importance of clearly defining the purpose of your analysis. In applied health research, explanatory models help identify causal mechanisms, while predictive models can guide real-world decisions by providing actionable forecasts. Together, they enhance both our understanding of phenomena and our ability to make informed decisions in complex environments.

2.1.2 Individual vs. Population Prediction

Another important distinction is between **individual vs. population** prediction. In the smoking example above, we can be very sure about the mean effects that smoking has on lung cancer. On an individual level, it is harder to predict the outcome. Nevertheless, individual predictions will be (notably) better than random guessing. We will discuss this in greater detail.

2.1.3 Practical Use of Statistical Models

In my opinion, we should never be afraid to test our statistical models (as honestly as possible) against reality. We could for instance ask ourselves:

- “How much better does this model classify than the arithmetic mean? (i.e., the linear model with just an intercept)”
- “How much better does this model classify than random guessing?”
- Is it worth the effort to collect data and estimate this model by using hundreds of hours of our time?

In some cases, these questions can be answered straightforwardly.

- In advertising (Google, Facebook, ...), a couple of percentage points in prediction quality might make a difference of millions of dollars in revenue offsetting the statisticians salary.
- Increased forecasts of a few percentage points in the stock market or just being slightly better than the average, will make you fabulously rich.
- Increased cancer forecasting might save lives, money and pain and is not only measured in money.

2.1.4 Start at the beginning

What do we actually want to do in general? Very broadly speaking we want to: **describe** the association of variables to each other that carry variability. Hence, the relationship is not deterministic like

$$y = 2x + 3$$

but rather we need to “loosen up” the relationship to account for variability. So, the values 2 and 3 are not fixed but afflicted with uncertainty. Depending on your philosophical view, you might say you want to find the “true” but unknown relationship between variables. This is what we do in simulation studies all the time: We know the true relationship, simulate data with variability and then try to estimate the true relationship we assumed in the first place. For some practical applications, we can get a really nice and complete answer to our question (for instance sample size for proportions).

So we are looking for a function f such that

$$Y = f(X)$$

where Y is the “outcome”, X are the “predictors”. Since this relationship is not deterministic, we have to account for variability. X can be a single Variable x or many variables x_1, x_2, \dots, x_p .

Elaborate on X, Y and F, independent, dependent.....

The simplest statistical model would be the mean model. Let’s use the example in Statistical rethinking starting on page 78.

```
library(rethinking)
data("Howell1")
d <- Howell1
str(d)
```

```
## 'data.frame':    544 obs. of  4 variables:
## $ height: num  152 140 137 157 145 ...
## $ weight: num  47.8 36.5 31.9 53 41.3 ...
## $ age : num  63 63 65 41 51 35 32 27 19 54 ...
## $ male : int  1 0 0 1 0 1 0 1 0 1 ...
```

Chapter 3

Literature

Here is a review of existing methods.

Chapter 4

Methods

We describe our methods in this chapter.

Math can be added in body using usual syntax like this

4.1 math example

p is unknown but expected to be around $1/3$. Standard error will be approximated

$$SE = \sqrt{\frac{p(1-p)}{n}} \approx \sqrt{\frac{1/3(1-1/3)}{300}} = 0.027$$

You can also use math in footnotes like this¹.

We will approximate standard error to 0.027^2

¹where we mention $p = \frac{a}{b}$

² p is unknown but expected to be around $1/3$. Standard error will be approximated

$$SE = \sqrt{\frac{p(1-p)}{n}} \approx \sqrt{\frac{1/3(1-1/3)}{300}} = 0.027$$

Chapter 5

Applications

Some *significant* applications are demonstrated in this chapter.

5.1 Example one

5.2 Example two

Chapter 6

Final Words

We have finished a nice book.