

Introduction to biostatistics and machine learning

Olga Dethlefsen, Eva Freyhult, Bengt Sennblad, Payam Emami, Julie Lorent

2020-09-23

Contents

Preface	5
1 Preliminary Mathematics for Statisticians	7
1.1 Mathematical notation, sets, functions, exponents and logarithms	7
1.2 Differentiation	7
1.3 Integration	7
1.4 Vectors and Matrices	7
2 Introduction to R, R Studio and R markdown	9
2.1 R	9
2.2 R Studio	9
2.3 R markdown	9
3 Probability: reasoning under uncertainty	11
3.1 Introduction	11
4 Probability: random variables	13
4.1 Random variables	13
4.2 Discrete random variables	14
5 Summarising and visualising data	15
6 Linear regression	17
6.1 Simple regression	17
6.2 Multiple regression	17

Preface

This “bookdown” book contains teaching and learning materials prepared and used during “Introduction to biostatistics and machine learning” course organised by NBIS, National Bioinformatics Infrastructure Sweden. The course is open for PhD students, postdoctoral researcher and other employees in need of biostatistical skills within Swedish universities. The course is geared towards life scientists wanting to be able to understand and use basic statistical and machine learning methods. It also suits those already applying biostatistical methods but who have never gotten a chance to reflect on or truly grasp the basic statistical concepts, such as the commonly misinterpreted p-value.

More about the course <https://nbisweden.github.io/workshop-mlbiostatistics/>

Chapter 1

Preliminary Mathematics for Statisticians

- 1.1 Mathematical notation, sets, functions, exponents and logarithms
- 1.2 Differentiation
- 1.3 Integration
- 1.4 Vectors and Matrices

Chapter 2

Introduction to R, R Studio and R markdown

2.1 R

2.2 R Studio

2.3 R markdown

Chapter 3

Probability: reasoning under uncertainty

Learning outcomes

- understand the concept of probability
- manipulate probabilities by their rules
- assign probabilities in very simple cases

3.1 Introduction

Some things are more likely to occur than others. Compare:

- the chance of the sun rising tomorrow with the chance that no-one is infected with COVID-19 tomorrow
- the chance of a cold dark winter in Stockholm with the chance of no rainy days over the summer months in Stockholm

We intuitively believe that the chance of sun rising or dark winter occurring are enormously higher than COVID-19 disappearing over night or having no rain over the entire summer. **Probability** gives us a scale for measuring the likeliness of events to occur. **Probability rules** enable us to reason about uncertain events. The probability rules are expressed in terms of sets, a well-defined collection of distinct objects.

Let A , B and S be sets:

- intersection (\cap) and union (\cup) are commutative: $A \cap B = B \cap A$ and $A \cup B = B \cup A$
- test

Chapter 4

Probability: random variables

Learning outcomes

- understand the concept of random discrete and continuous variables
- to be able to use probability density/mass functions and cumulative distribution functions and to understand the relationship between them
- describe properties of binomial, geometric, Poisson, uniform, exponential and normal distributions and identify which distributions to use in practical problems

4.1 Random variables

The outcome of a random experiment can be described by a random variable.

Example random variables:

- The weight of a random newborn baby
- The smoking status of a random mother
- The hemoglobin concentration in blood
- The number of mutations in a gene
- BMI of a random man
- Weight status of a random man (underweight, normal weight, overweight, obese)
- The result of throwing a die

Whenever chance is involved in the outcome of an experiment the outcome is a random variable.

A random variable is usually denoted by a capital letter, X, Y, Z, \dots . Values collected in an experiment are observations of the random variable, usually denoted by lowercase letters x, y, z, \dots .

A random variable can not be predicted exactly, but the probability of all possible outcomes can be described.

The population is the collection of all possible observations of the random variable. Note, the population is not always countable.

A sample is a subset of the population.

4.2 Discrete random variables

A discrete random variable can be described by its *probability mass function*.

→

→

Chapter 5

Summarising and visualising data

Chapter 6

Linear regression

6.1 Simple regression

6.2 Multiple regression