

Introduction to biostatistics and machine learning

Olga Dethlefsen, Eva Freyhult, Bengt Sennblad, Payam Emami

2020-09-24

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
3.2 Basic set definitions	11
3.3 Basic set operations	12
3.4 Exercises	13
3.5 Answers to exercises	13
4 Probability: random variables	15
4.1 Random variables	15
4.2 Discrete random variables	16
5 Summarising and visualising data	17
6 Linear regression	19
6.1 Simple regression	19
6.2 Multiple regression	19

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.

3.2 Basic set definitions

- **set**: a well-defined collection of distinct objects, e.g. $A = \{2, 4, 6\}$

- **subset**, \subseteq : if every element of set A is also in B, then A is said to be a subset of B, written as $A \subseteq B$ and pronounced A is contained in B, e.g. $A \subseteq B$, when $B = \{2, 4, 6, 8, 10\}$. Every set is a subset of itself.
- **empty set**, \emptyset : is a unique set with no members, denoted by $E = \emptyset$ or $E = \{\}$. The empty set is a subset of every set.

3.3 Basic set operations

- **union of two sets**, \cup : two sets can be “added” together, the union of A and B, written as $A \cup B$, e.g. $\{1, 2\} \cup \{2, 3\} = \{1, 2, 3\}$ or $\{1, 2, 3\} \cup \{1, 4, 5, 6\} = \{1, 2, 3, 4, 5, 6\}$
- **intersection of two sets**, \cap : a new set can be constructed by taking members of two sets that are “in common”, written as $A \cap B$, e.g. $\{1, 2, 3, 4, 5, 6\} \cap \{2, 3, 7\} = \{2, 3\}$ or $\{1, 2, 3\} \cap \{7\} = \emptyset$
- **complement of a set**, A^c , A^C : are the elements not in A
- **difference of two sets**, \setminus : two sets can be “subtracted”, denoted by $A \setminus B$, by taking all elements that are members of A but are not members of B, e.g. $\{1, 2, 3, 4\} \setminus \{1, 3\} = \{2, 4\}$. This is also in other words a relative complement of A with respect to B.
- **partition of a set**: a partition of a set S is a set of nonempty subset of S, such that every element x in S is in exactly one of these subsets. That is, the subset are pairwise *disjoint*, meaning no two sets of the partition contain elements in common, and the union of all the subset of the partition is S, e.g. Set $\{1, 2, 3\}$ has five partitions: i) $\{1\}, \{2\}, \{3\}$, ii) $\{1, 2\}, \{3\}$, iii) $\{1, 3\}, \{2\}$, iv) $\{1\}, \{2, 3\}$ and v) $\{1, 2, 3\}$

Some further properties

- intersection (\cap) and union (\cup) are commutative: $A \cap B = B \cap A$ and $A \cup B = B \cup A$
- intersection and union are associative: $(A \cap B) \cap C = A \cap (B \cap C)$ and $(A \cup B) \cup C = A \cup (B \cup C)$
- union is distributive over intersection and vice versa: $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ and $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$
- de Morgan’s laws: $(A \cup B)^c = A^c \cap B^c$ and $(A \cap B)^c = A^c \cup B^c$

Example

Example 3.1. Here comes example

a = 1

3.4 Exercises

3.5 Answers to exercises

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