

# Biostatistics

*David Coffey*

*2018-06-29*



# Contents

<b>1</b>	<b>General overview</b>	<b>5</b>
1.1	Introduction . . . . .	5
1.2	Example dataset . . . . .	5
<b>2</b>	<b>Descriptive statistics</b>	<b>7</b>
2.1	Arithmetic mean . . . . .	7
2.2	Median . . . . .	7
2.3	Mode . . . . .	8
2.4	Geometric mean . . . . .	8



# Chapter 1

## General overview

### 1.1 Introduction

This book provides a concise overview of biostatistics and its applications using the R programming language. The textbook *Fundamentals of Biostatistics* (Rosner, 2016) was used extensively in the preparation of this book.

### 1.2 Example dataset

Examples of R functions are performed on a dataset of patients with newly diagnosed multiple myeloma. This dataset contains a variety of categorical and continuous variables.



## Chapter 2

# Descriptive statistics

### 2.1 Arithmetic mean

The arithmetic mean ( $\bar{x}$ ) is a measure of central location. It is calculated from the sum of all the observations ( $n$ ) divided by the number of observations:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

The notation  $\sum_{i=1}^n x_i$  means the sum of all  $x_i$  observations ( $x_1 + x_2 + x_n$ ). One limitation to the arithmetic mean is that it is overly sensitive to extreme values.

```
# Import dataset
load("docs/Example-data.Rda")

# Calculate arithmetic mean
mean(data$Age)
```

```
[1] 59.5
```

### 2.2 Median

If all observations are ordered from smallest to largest, the median is the middle number. More precisely, if  $n$  is odd,  $\frac{n+1}{2}$ , or if  $n$  is even, the average of  $\frac{n}{2}$  and  $\frac{n}{2} + 1$ .

The rationale for using the median is to ensure an equal number of observations on both sides of the sample median. The main weakness of the sample median is that it is less sensitive to the actual numeric values of the data points. If the sample distribution is symmetric, the arithmetic mean is approximately the same as the median. For positively skewed distributions, the arithmetic mean tends to be larger than the median; for negatively skewed distributions, the arithmetic mean tends to be smaller than the median.

```
# Calculate arithmetic mean
median(data$Age)
```

```
[1] 62.5
```

## 2.3 Mode

The mode is the most frequently occurring value among all of the observations in a sample. Some distributions have more than one mode. A distribution with one mode is called unimodal; two modes, bimodal; three modes, trimodal.

```
# Calcuete mode  
library(DescTools)  
Mode(data$Age)
```

```
[1] 64
```

## 2.4 Geometric mean

The geometric mean ( $\log\bar{x}$ ) is the central number in a geometric progression such as exponential growth. The geometric mean is defined as the  $n$ th root of the product of  $n$  numbers:

$$\log\bar{x} = \frac{\sum_{i=1}^n \log x_i}{n}$$

Any base can be used to compute the logarithms for the geometric mean. It is usually preferable to work in the original scale by taking the antilogarithm of  $\log\bar{x}$  to form the geometric mean.

```
# Calcuete geometric arithmetic mean  
library(DescTools)  
Gmean(data$Age)
```

```
[1] 58.61499
```



# Bibliography

Rosner, B. (2016). *Fundamentals of Biostatistics*. Cengage Learning, Boston, Massachusetts, 8th edition. ISBN 978-1-305-26892.