

R Lab - Day 2 (part 2)

Descriptive statistics

MF9130E V24

2024.04.09

Chi Zhang

Oslo Center for Biostatistics and Epidemiology

chi.zhang@medisin.uio.no

Overview

Plan for this lecture

Review key concepts from descriptive statistics, exploratory data analysis

Practice: load a dataset, produce some **summary statistics**, make some **plots**

Descriptive statistics, EDA

EDA: Exploratory Data Analysis

In contrast to Confirmatory analysis (e.g. hypothesis tests)

The goal of EDA is to get a first impression of your data

Descriptive statistics is part of the process of exploration

For example, what is the average of 'height' in my data?

In this session, we learn how to explore a dataset with

- Review **descriptive (summary) statistics**
- Some **simple data manipulation** techniques
- **Visualisation** with histogram, boxplot, scatterplot

Descriptive statistics

Central measures

Mean (average)
 $(x_1 + x_2 + \dots + x_n)/n$

Median
Half values smaller than this value; half greater

Mean is sensitive to extreme values (outliers)

Variation measures

Range

Interquartile range (percentiles, quartiles)

Variance

Standard deviation

Descriptive statistics

Mean

Median

Minimum, maximum

Quantiles (top 5% = 0.95
quantile)

Quartiles (0.25, 0.5, 0.75)

Variance, standard deviation

```
# x is a continuous variable
```

```
mean(x)
```

```
median(x)
```

```
min(x), max(x)
```

```
summary(x)
```

```
quantile(x, 0.95)
```

```
quantile(x, 0.25)
```

```
var(x), sd(x)
```

Simple data manipulation

When you get a dataset, the first thing to do is to get an overview of your dataset:

How many observations?

How many variables are measured?

What data types exist?

```
# df is a data.frame

# first 6 rows
head(df)

# number of observations
nrow(df)

# column names (variables)
colnames(df)

# what data types?
str(df)
class(df$var1)
```

Descriptive statistics with plots

Data visualization is a very effective way to explore, and present your data.

We focus on **base R**
(rather than more complex solutions: ggplot2)

```
# x is a continuous variable  
hist(x)  
boxplot(x)
```

Demo: birth data

Now we are going to practice what we have introduced just now, using **birth** dataset.

You can check the lab notes after class: **Descriptive statistics**

- load the dataset
- print out the first few rows of the data, how many rows? Column names?
- take a **numeric** variable, produce some statistics (mean, variance, min, max...)
- make a plot to describe the variable visually (histogram, box plot)
- Take a **categorical** variable, count the number in each category

Exercise 1 (weight)

The solution to the exercises are at the bottom.

1a) Generate a variable named **weight**, with the following measurements

50	75	70	74	95	83	65	94	66	65
65	75	84	55	73	68	72	67	53	65

Exercise 1 (weight)

1b) Make a simple descriptive analysis of the variable. What are the mean, median, maximum, minimum and quantiles?

Exercise 1 (weight)

1c) Make a histogram of the variable.

Exercise 1 (weight)

1d) Make a boxplot. What do the two dots on the top represent?

Exercise 2 (lung function)

2a) Download and open **PEFH98-english** data into R

(Use the file **PEFH98-english.csv** or **.rda** format)

Exercise 2 (lung function)

2b) How many observations are there? (Number of subjects)

How do you get a list of variables from your dataset?

Exercise 2 (lung function)

2b) Make a histogram of the following variables. Compute means, and interpret the results.

Height, weight, age, pefsitm, pefstam

(Illustrate height)

Exercise 2 (lung function)

2c) Make histograms for the variable **height** and **pefmean** for **men** and **women** separately.

Also make boxplots.

What conclusion can you draw?

(Illustrate height for men)

Exercise 2 (lung function)

2d) Make three scatterplots to compare

Pefmean with **height**

Pefmean with **weight**

Pefmean with **age**

(Illustrate pefmean with height)