Chapter 1

Introduction, descriptive statistics, Python and data visualization

Exercises

Contents

| 1 | oduction, descriptive statistics, Python and data visualization rcises | | |
|---|---|---------------------|--|
| | 1.1 | Infant birth weight | |
| | 1.2 | Course Grades | |
| | 1.3 | Cholesterol | |
| | 1.4 | Project start | |

Chapter 1 CONTENTS

Initilize Python packages

import numpy as np import pandas as pd import scipy.stats as stats import matplotlib.pyplot as plt import statsmodels.api as sm import statsmodels.stats.proportion as smprop

1.1 Infant birth weight

In a study of different occupational groups the infant birth weight was recorded for randomly selected babies born by hairdressers, who had their first child. The following table shows the weight in grams (observations specified in sorted order) for 10 female births and 10 male births:

| Females (<i>x</i>) | 2474 | 2547 | 2830 | 3219 | 3429 | 3448 | 3677 | 3872 | 4001 | 4116 |
|----------------------|------|------|------|------|------|------|------|------|------|------|
| Males (y) | 2844 | 2863 | 2963 | 3239 | 3379 | 3449 | 3582 | 3926 | 4151 | 4356 |

Solve at least the following questions a)-c) first "manually" and then by the inbuilt functions in Python. It is OK to use Python as alternative to your pocket calculator for the "manual" part, but avoid the inbuilt functions that will produce the results without forcing you to think about how to compute it during the manual part.

- a) What is the sample mean, variance and standard deviation of the female births? Express in your own words the story told by these numbers. The idea is to force you to interpret what can be learned from these numbers.
- b) Compute the same summary statistics of the male births. Compare and explain differences with the results for the female births.
- c) Find the five quartiles for each sample and draw the two box plots with pen and paper (i.e. not using Python.)
- d) Are there any "extreme" observations in the two samples (use the *modified box plot* definition of extremeness)?
- e) What are the coefficient of variations in the two groups?

1.2 Course Grades

Exercise 1.1 Course grades

To compare the difficulty of 2 different courses at a university the following grades distributions (given as number of pupils who achieved the grades) were registered:

| | Course 1 | Course 2 | Total |
|----------|----------|----------|-------|
| Grade 12 | 20 | 14 | 34 |
| Grade 10 | 14 | 14 | 28 |
| Grade 7 | 16 | 27 | 43 |
| Grade 4 | 20 | 22 | 42 |
| Grade 2 | 12 | 27 | 39 |
| Grade 0 | 16 | 17 | 33 |
| Grade -3 | 10 | 22 | 32 |
| Total | 108 | 143 | 251 |

a) What is the median of the 251 achieved grades?

b) What are the quartiles and the IQR (Inter Quartile Range)?

1.3 Cholesterol

Exercise 1.2 Cholesterol

In a clinical trial of a cholesterol-lowering agent, 15 patient cholesterol (in mmol L^{-1}) was measured before treatment and 3 weeks after starting treatment. Data is listed in the following table:

| Patient | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---------|-----|-----|-----|------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Before | 9.1 | 8.0 | 7.7 | 10.0 | 9.6 | 7.9 | 9.0 | 7.1 | 8.3 | 9.6 | 8.2 | 9.2 | 7.3 | 8.5 | 9.5 |
| After | 8.2 | 6.4 | 6.6 | 8.5 | 8.0 | 5.8 | 7.8 | 7.2 | 6.7 | 9.8 | 7.1 | 7.7 | 6.0 | 6.6 | 8.4 |

- a) What is the median of the cholesterol measurements for the patients before treatment, and similarly after treatment?
- b) Find the standard deviations of the cholesterol measurements of the patients before and after treatment.
- c) Find the sample covariance between cholesterol measurements of the patients before and after treatment.
- d) Find the sample correlation between cholesterol measurements of the patients before and after treatment.
- e) Compute the 15 differences (Dif = Before After) and do various summary statistics and plotting of these: sample mean, sample variance, sample standard deviation, boxplot etc.

f) Observing such data the big question is whether an average decrease in cholesterol level can be "shown statistically". How to formally answer this question is presented in Chapter 3, but consider now which summary statistics and/or plots would you look at to have some idea of what the answer will be?

1.4 Project start

Exercise 1.3 Project start

a) Go to Learn or the website and take a look at the first project. Read the project page on the website for more information (02323.compute.dtu.dk/projects or 02402.compute.dtu.dk/projects). Choose a project and read the project description. Follow the steps to import the data into Python and get started with the explorative data analysis.