

SIG731 2023: Task 3P

Working with **numpy** Matrices (Multidimensional Data)

Last updated: 2023-11-24

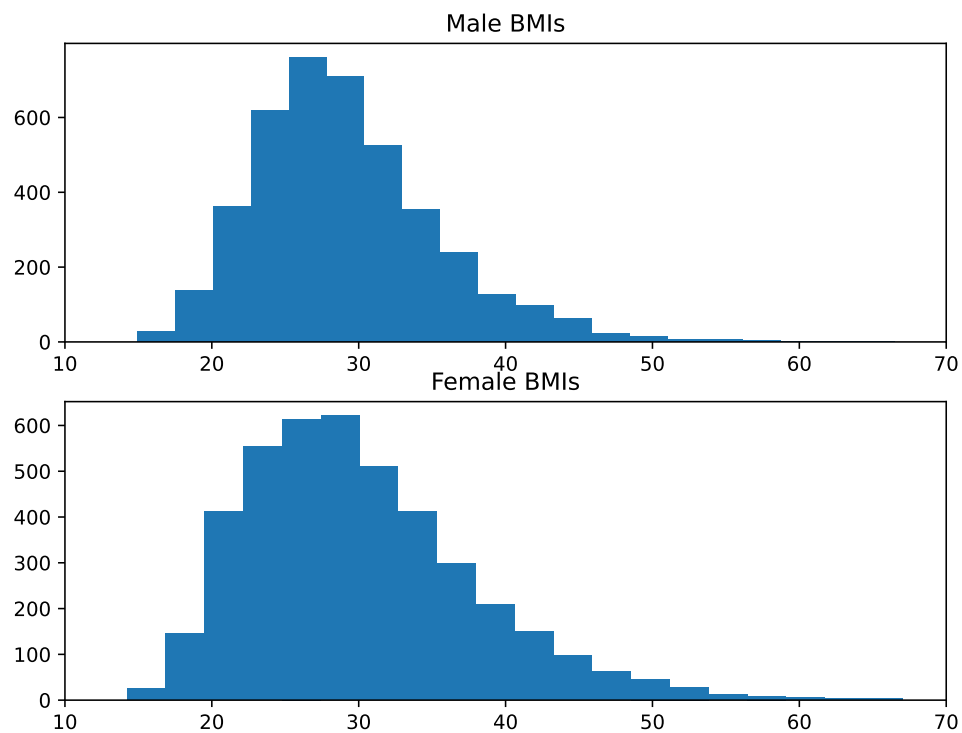
Contents

1	Task	1
2	Artefacts	4

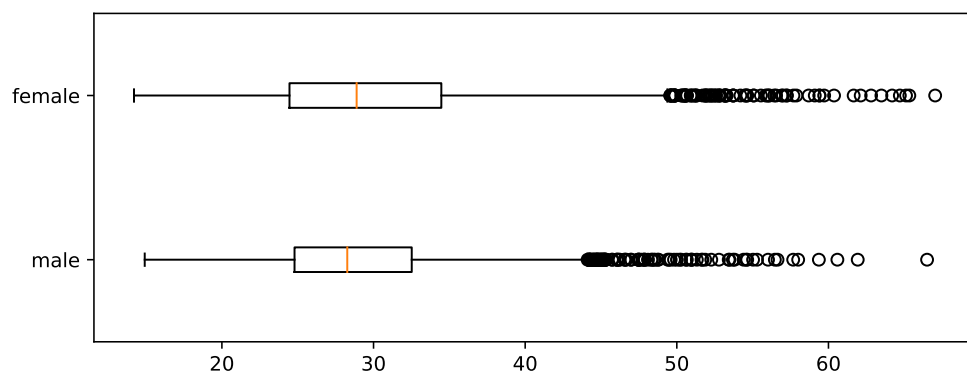
1 Task

Create a single Jupyter/IPython notebook (see the *Artefacts* section below for all the requirements), where you perform what follows.

1. From <https://github.com/gagolews/teaching-data/tree/master/marek>, download the two following excerpts from the National Health and Nutrition Examination Survey (**NHANES** dataset) that give body measurements of adult males and females.
 - `nhanes_adult_male_bmx_2020.csv`,
 - `nhanes_adult_female_bmx_2020.csv`.
2. Read them as **numpy** matrices named `male` and `female` using `numpy.genfromtxt`. Each matrix consists of seven columns:
 1. weight (kg),
 2. standing height (cm),
 3. upper arm length (cm),
 4. upper leg length (cm),
 5. arm circumference (cm),
 6. hip circumference (cm),
 7. waist circumference (cm).
3. In both cases, add the eighth column which stores the **body mass indices** of the participants.
4. On a **single** plot, draw two histograms: for male BMIs (top subfigure) and for female BMIs (bottom subfigure) **one below another**. Set the number of histogram bins to 20. Use `matplotlib.pyplot.subplot` to create two subplots in one figure. Call `matplotlib.pyplot.xlim` to make the **x-axis limits identical for both subfigures** (work out the appropriate limits yourself). For example:



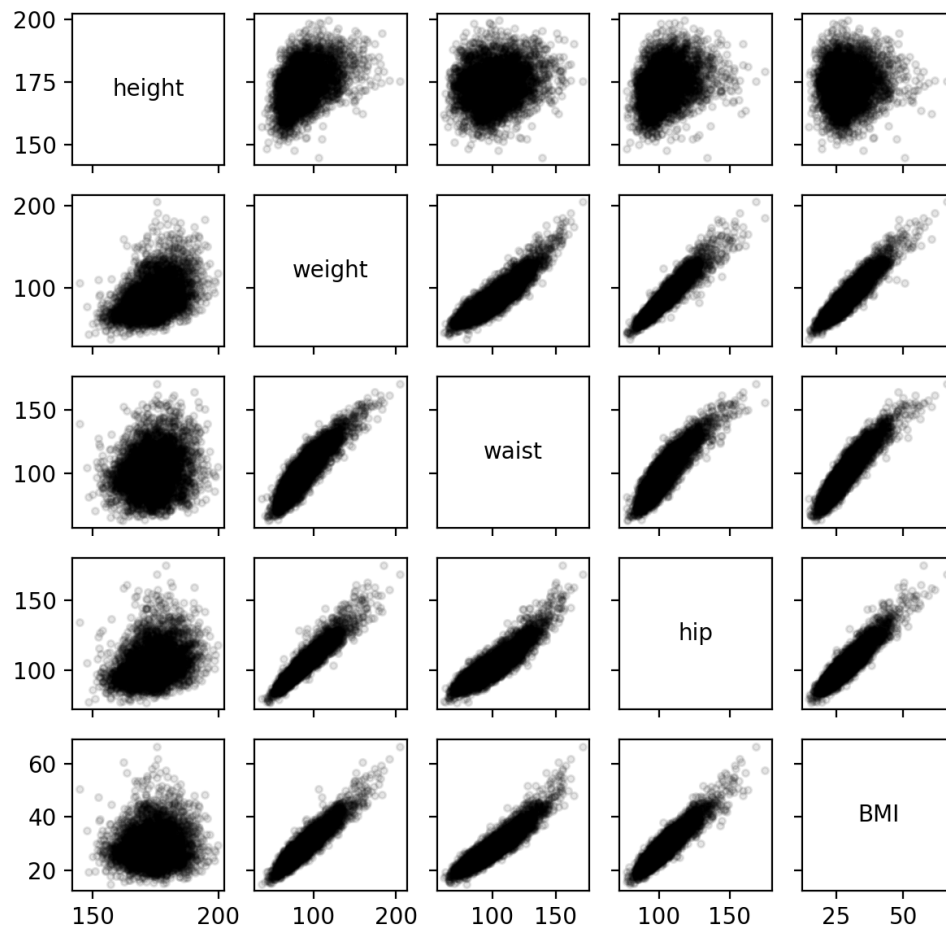
- Using a *single* call to `matplotlib.pyplot.boxplot`, draw a box-and-whisker plot giving the male and female BMIs, with **two boxes one below another** (on one plot) so that they can be compared to each other. Note that the boxplot function can be fed with a list of two vectors like `[male_BMIs, female_BMIs]`. For example:



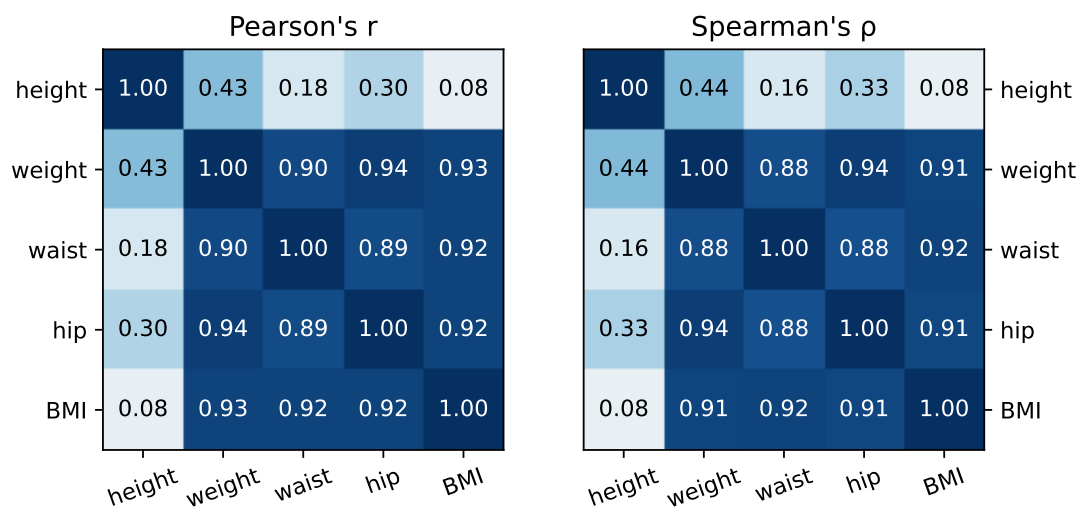
- Compute the basic numerical aggregates of the male and female BMIs (measures of location, dispersion, and shape). Report them in a readable format. Example formatting of the aggregates:

```
##          female    male
## BMI mean    30.10  29.14
##   median    28.89  28.27
##    min     14.20  14.91
##    max     67.04  66.50
##   std       7.76   6.31
##   IQR      10.01   7.73
##   skew       0.92   0.97
```

7. In your own words, describe the two distributions based on the results obtained in subtasks 4, 5, and 6 above (e.g., are they left-skewed, how they differ, which one has more dispersion, and so forth).
8. Draw a scatterplot matrix (pairplot) for the male heights, weights, waist circumferences, hip circumferences, and BMIs (these five columns only); see the `pairplot` function in section 7.4.3 of our book. Example output (yours can be more aesthetic):



9. Compute Pearson's *and* Spearman's correlation coefficients for all pairs of variables mentioned in subtask 8. Present/visualise these coefficients on two correlation heatmaps (with correlation coefficients printed inside the coloured cells); see the `corrheatmap` function in Section 9.1.2 of our book. Example outputs:



10. Discuss the findings from subtasks 8 and 9.

Important. Remember that this is an exercise where you demonstrate the mastery of **numpy** matrices, and not **pandas** data frames. The use of **pandas** is forbidden. You can use **scipy**, though.

All packages must be imported and data must be loaded at the beginning of the file (only once!).

2 Artefacts

Make sure that your notebook has a **readable structure**; in particular, that it is divided into sections. Use rich Markdown formatting (text in dedicated Markdown chunks – not just Python comments).

Do not include the questions/tasks from the task specification. Your notebook should read nicely and smoothly – like a report from data analysis that you designed yourself. Make the flow read natural (e.g., *First, let us load the data on... Then, let us determine... etc.*). Imagine it is a piece of work that you would like to

show to your manager or clients — you certainly want to make a good impression. Check your spelling and grammar. Also, use formal language.

At the start of the notebook, you need to provide: the **title** of the report (e.g., *Task 42: How Much I Love This Unit*), your **name**, **student number**, and **email address**.

Then, add 1–2 introductory paragraphs (an introduction/abstract – what the task is about).

Before each nontrivial code chunk, briefly **explain** what its purpose is. After each code chunk, **summarise and discuss the obtained results** (in a few sentences).

Conclude the report with 1–2 paragraphs (summary/discussion/possible extensions of the analysis etc.).

Checklist:

1. Header, introduction, conclusion (Markdown chunks).
2. Text divided into sections, all major code chunks commented and discussed in your own words (Markdown chunks).
3. Every subtask addressed/solved. In particular, all reference results that are part of the task specification have been reproduced (plots, computed aggregates, etc.).
4. The report is readable and neat. In particular:
 - all code lines are visible in their entirety (they are not too long),
 - code chunks use consecutive numbering (select *Kernel - Restart and Run All* from the Jupyter menu),
 - rich Markdown formatting is used (# Section Title, * bullet list, 1. enumerated list, | table |, *italic*, etc.),
 - the printing of unnecessary/intermediate objects is minimised (focus on reporting the results specifically requested in the task specification).

Submissions which do not *fully* (100%) conform to the task specification *on* the cut-off date will be marked as FAIL.

Good luck!