# Statistical analysis of experimental data
## Introduction to Python

Artur Kalinowski

**FACULTY OF PHYSICS**
UNIVERSITY
OF WARSAW

**Lecture 02**
October 12, 2023

**Lecture concept**

The goal of the lecture is not only to present the theoretical concepts
but also to demonstrate how they can be applied in practical terms to data analysis.

Different concepts as well as analysis methods will be discussed using code examples.

We will address many (relatively simple) problems and try to look at their numerical solutions.

All examples presented are based on Python and multiple additional packages,
which will be shortly introduced today.

**Web resources**

**Kampus platform** will be used for lecture slides, Python notebooks, as well as home exercises (and final exam):

https://kampus-student2.ckc.uw.edu.pl/course/view.php?id=14456

All files will also be available from the dedicated web page:

http://www.fuw.edu.pl/~zarnecki/SAED/

accessible without USOS account...

**Containers**

To avoid the burden of installing and configuring all the required packages on your personal computers, we will use virtual environments known as containers. Containers are lightweight versions of old time virtual machines. Two popular environments used to run the containers:

- Apptainer (previously singularity; still difficult to find in Google)
- Docker

Both Docker and Apptainer can use the same input files.
Docker is available for any operating system, Apptainer only for Linux ones.

Notebooks presented at this lecture are developed with akalinow/root-fedora35 container prepared for the other course some time ago. This is a container based on Fedora Linux ditribution, and conains all packages we will use during the classes.

Detailed instructions on installing and starting Docker are given here

**Git and Google Colaboratory**

All scripts discussed during the course will be available for download at the Kampus platform and the dedicated lecture web page. They will also be uploaded to github repository:

https://github.com/zarnecki/SAED

You also open them in your browser using the Google Colab platform.
To do so, use  link and select the notebook you want to load.

Please remember that you need to save the file to (your) Google Disc before you make any modifications. Otherwise, all your changes will be lost!..

**Python tools**

One of the main advantages of Python is the large number of diverse packages developed for various applications...

During the course we will use the following packages for computation and plotting:

- numpy - The fundamental package for scientific computing with Python
- matplotlib - Visualization with Python
- SciPy - Fundamental algorithms for scientific computing in Python
- scikit-learn - Machine Learning in Python, simple and efficient
- pandas - Handling more complex data structures

See lecture notebook for more details and examples:  CO Open in Colab

pandas | Getting started  User Guide  API reference  Development  Release notes | 2.1.1 ▾

# pandas documentation

**Date**: Sep 20, 2023 **Version**: 2.1.1

**Download documentation**: Zipped HTML

**Previous versions**: Documentation of previous pandas versions is available at pandas.pydata.org.

**Useful links**: Binary Installers | Source Repository | Issues & Ideas | Q&A Support | Mailing List

`pandas` is an open source, BSD-licensed library providing high-performance, easy-to-use data structures and data analysis tools for the Python programming language.

## Getting started

New to *pandas*? Check out the getting started guides.

## User guide

The user guide provides in-depth information on the

# Python packages

## Graphical presentation of the data

Look for the recent data describing age and gender structure of the population of Poland. One of the possibile sources is the Central Statistical Office (Polish: Główny Urząd Statystyczny) or Poland's Data Portal.

1. Prepare the plot showing the gender balance (men to women ratio) as a function of age.
2. Can the observed dependence be used to draw any conclusions concerning the life expectancy? Why do you think so?

Solutions should be uploaded until October 26.

Solutions (dedicated files or screenshots of the notebook final output) should be uploaded to Kampus in readable format (PDF, JPG, PNG). You should add the link to your notebook in Google Colab as a comment.
Please share your Google Colab space with azarnecki@uw.edu.pl, so I can have a look at your notebook, if needed.