



Mastering Pandas: An in-depth Guide in Data Science Techniques for Researchers

Lecture 1 – 5 April 2023

Overview of Pandas data structures, data manipulation, filtering, and aggregation

George Leonard defines mastery as a continual journey, and there's no end or "perfection" of the skill—rather, the practice is the point in and of itself.

Agenda

- About Me
- Admin
- Mini-School Overview
- Data
- EDA
- Pandas Part 1
- Useful Resources

About Me

- Benjamin (pronounced Binyamin) is the Lead Software Engineer at the Centre for High Performance Computing (CHPC), a research centre that forms part of the CSIR
- BEng in Mechatronic Engineering from Stellenbosch University and a MEng degree in Radar Digital Signal Processing from UCT
- Worked in industry focusing on the fields of electronic engineering, data science, and software development
- Currently lead the software engineering architecture, and the design and development of HPC support software
- Lead the annual Coding Summer School which focuses on training researchers in fundamental programming skills to assist with their research.
- Research interest: DSP for IoT, and NLP for African languages – see <https://www.masakhane.io/>

- Form - <https://forms.office.com/r/QJs4SbGXBB>
- All code and lessons – github repo:
https://github.com/kode2go/nithecs/tree/main/lecture_01

Mini-School Overview

Timetable:

- Lecture 1 (5 April) - Overview of Pandas data structures, data manipulation, filtering, and aggregation*
- Lecture 2 (12 April) – Aggregation*, data wrangling, merging, joining, and time series analysis
- Lecture 3 (19 April) - Data visualization, statistical modelling, and managing varied datasets types
- Lecture 4 (26 April) - Advanced data analysis techniques, and real-world applications

General Topics*:

- EDA
- Big Data
- Data Pipelines

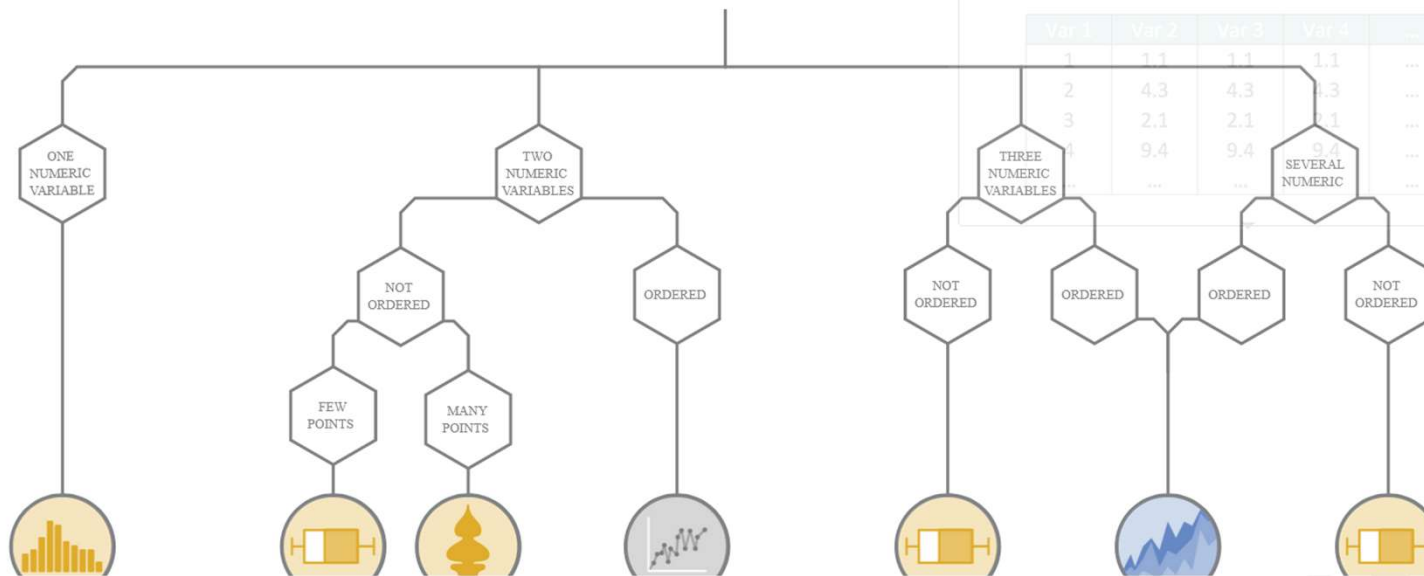
Overview:

- This mini-school is aimed at those who are familiar with Python and Pandas.
- By the end of the mini-school, attendees should have learnt skills to analyse their datasets more effectively and derive meaningful insights.



What kind of data do you have? Pick the main type using the buttons below. Then let the decision tree guide you toward your graphic possibilities.

Numeric [Categoric](#) [Num & Cat](#) [Maps](#) [Network](#) [Time series](#)



- Exploratory data analysis on the Iris dataset
- 1936 by the statistician and biologist Ronald Fisher

iris setosa



petal sepal

iris versicolor

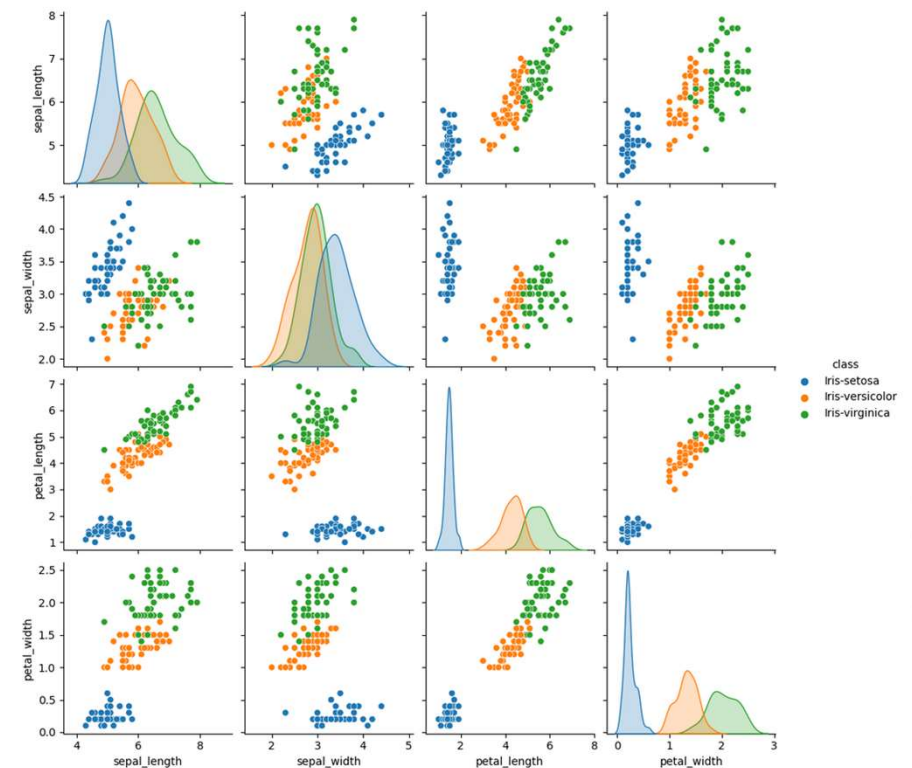


petal sepal

iris virginica

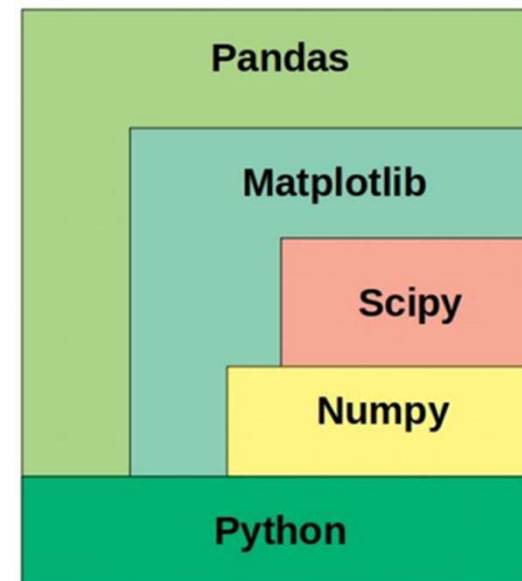


petal sepal



Pandas Part 1

- Pandas – “Panel Data” – econometrics term for data sets
- Numpy -> Pandas
- Pandas Data Structures and Functionality
- Pandas Profiling



Mastering Pandas – Mini-School 2023

Useful Resources Fundamentals

- DataCamp (<https://www.datacamp.com/>)
- Kaggle (<https://www.kaggle.com/>)
- Towards Data Science (<https://towardsdatascience.com/>)
- Analytics Vidhya (<https://www.analyticsvidhya.com/>)
- Real Python (<https://realpython.com/>)
- Dataquest (<https://www.dataquest.io/>)
- Python Data Science Handbook (<https://jakevdp.github.io/PythonDataScienceHandbook/>)