# Summarising datasets in Python (Part 2)

This planning document is intended to support teachers who are delivering the NPA/PDA Data Science or for students who are learning independently. It also aligns with the Data Skills for Work framework.

## Contents

## Lesson Description

| Lesson Overview | Summarising within groups, using summary calculations such as the total, count, min/max and mean. |
|---|---|
| Topic | Data manipulation |
| Book Chapter(s) | "Data Transformation and Manipulation" |

| NPA level | 5, 6 |
|---|---|
| PDA level | 7, 8 |
| Data skills for work level | Core, Analysis |

## Lesson Contents

This lesson consists of:

- A lesson plan (this document)
- A Powerpoint presentation, 'Summarising datasets in Python Part 2'
- Jupyter notebooks:
  - 'summarising_datasets_with_answers_part_2.ipynb' (for teachers), and
  - 'summarising_datasets_part_2.ipynb' (for learners)
- Datasets used in the Jupyter notebooks: the datasets are stored online and imported by the Jupyter notebooks.

## Learning Intentions

We will be learning to summarise datasets in Python, specifically to:

- group rows of data based on logical criteria
- perform summary calculations on grouped data

## Success Criteria

I can *describe* how to create a new variable by extracting data.

I can *create* new variables in Python by extracting data.

I can *describe* how to create a new variable by combining data.

I can *create* new variables in Python by combining data.

## Knowledge Prerequisites

Learners should know:

- Data is held in structured data frames
- Python is a programming language that can be used for data analysis
- How to use a Jupyter notebook to write, edit and run Python code
- How to open a Jupyter notebook to write, edit and run Python code

## Lesson Requirements

| | PDA | NPA | Data Skills for work |
|---|---|---|---|
| **Qualification** | Yes | Yes | Yes |
| **Outcome ID(s)** | WD8.3c, CD8.1g, WD7.3d, WD7.2a, WD7.2b, CD7.3a | DS5.3d, DS6.3d | C2.1, A2.1, A2.3 |

| | | | |
|---|---|---|---|
| **Outcome description(s)** | WD8.3c Transformations<br><br>CD8.1g Preparing data for visualisation<br><br>WD7.3d Data aggregation<br><br>WD7.2a Types of data transformation<br><br>WD7.2b Common transformations including filtering, sorting<br><br>CD7.3a Preparing data for visualisation<br><br>*N.B. out of scope of this lesson,*<br><br>*"WD8.3c … including joins"*<br><br>*"WD7.2b ….combining, separating and grouping"* | DS5.3d Perform analyses including [...] group and summarise,<br><br>DS6.3d Perform descriptive and predictive analyses on the data. | C2.1 Vocabulary used in data science and analytics<br><br>A2.1 Use of tools to analyse data<br><br>A2.3 Data calculation and manipulation |
| **Level** | 7, 8 | 5, 6 | Core, Analysis |
| **Software language** | Python | Python | Python |
| **Required equipment /software for student** | Lesson: PowerPoint<br><br>Python notebook: Jupyter notebook environment | Lesson: PowerPoint<br><br>Python notebook: Jupyter notebook environment | Lesson: PowerPoint<br><br>Python notebook: Jupyter notebook environment |

# Jupyter Notebook

There is a Jupyter notebook for this lesson that provides examples and programming tasks for learners, drawn from the examples in the lesson Powerpoint.

The notebook uses Python 3.x and the following packages:

- pandas - for data manipulation
- s3fs - an API to AWS S3 (Simple Storage Service), used to import datasets

The notebooks can be used with any Jupyter notebook environment. The tasks are described in the table below.

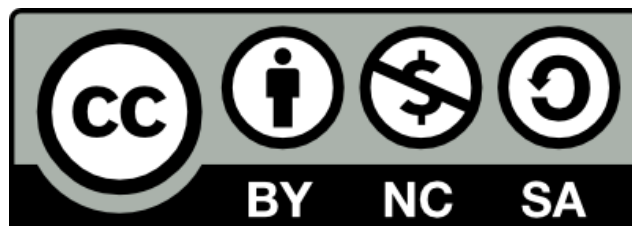| Notebook section | Task | Description |
|---|---|---|
| Group and summarise data | Task 1 - What's the Maximum Height for Each Species? | Calculating the maximum for each group. |
| | Task 2 - Which Gender has the Largest Average Mass? | Calculating the maximum for each group. |
| | Task 3 - Are There More Males Than Females? | Calculating a count of values for each group. |
| | Extension Task 1 - What is the Total Mass of Each Species? | Calculating the sum for each group. |
| | Extension Task 2 - How Many Different Eye Colours Do Each Species Have? | Calculating a count of unique values for each group. |

# Datasets

The following datasets are used in this lesson.

| Dataset name | Description | Link |
|---|---|---|
| Fruit & veg | The sales of fruit and vegetables in a shop. | https://datasets.learn-data.science/fruit_and_veg.csv |
| Star Wars characters | Star Wars character data (from https://swapi.dev/). | https://datasets.learn-data.science/star_wars_characters.csv |

## How you can use this lesson

This lesson has been created by Effini in partnership with Data Education in Schools, The Data Lab and Data Skills for Work, with funding from the Scottish Government.