

# Practise data understanding in Python

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

Version Control .....	1
Lesson Description.....	2
Lesson Contents .....	2
Learning Intentions.....	2
Success Criteria .....	3
Knowledge Prerequisites.....	3
Lesson Requirements.....	4
Jupyter Notebook.....	4
Datasets.....	5
How you can use this lesson .....	5
Alternative format.....	7

## Version Control

Version number	Purpose/Change	By	Date
1.0	Published by Effini	John Bell	12 April 2022

## Lesson Description

<b>Lesson Overview</b>	<p>This lesson is intended to follow the <b>Data Understanding in Python Part 1</b> and <b>Part 2</b> lessons.</p> <p>This is a consolidation activity to give learners the chance to apply the data understanding skills they have learned in <b>Data Understanding in Python Part 1</b> and <b>Part 2</b>. Learners will be provided with a single, unfamiliar dataset and undertake a set of activities to get a better understanding of it.</p> <p>Note that learners will use what they find out about this dataset in a later lesson, <b>Practise Data Cleansing in Python</b>, where they will take what they have learned and apply it to the task of cleansing the same dataset as used in this lesson. It is recommended that learners record their findings from this lesson so that they can access this during <b>Practise Data Cleansing in Python</b>.</p>
<b>Topic</b>	Data Manipulation and Data Analysis
<b>Book Chapter(s)</b>	Analysing data

<b>NPA level</b>	5, 6
<b>PDA level</b>	7, 8
<b>Data skills for work level</b>	Core, Analysis

## Lesson Contents

This lesson consists of:

- A lesson plan (this document)
- A PowerPoint presentation, 'Practise Data Understanding in Python'
- 2 Jupyter notebooks:
  - 'practise\_data\_understanding.ipynb' (for learners)
  - 'practise\_data\_understanding\_with\_answers.ipynb' (for teachers)

## Learning Intentions

We will be learning how to apply data understanding techniques to **understand an unfamiliar dataset using Python**, specifically,

- how to **import** a dataset without importing **metadata**
- how to use a **data dictionary** to find out about a dataset
- how to find the **shape, size** and **format** of datasets, using Python
- how to find the **data types** of variables in a dataset, using Python
- how to **identify outliers** and **missing values** in Python

## Success Criteria

I can *import* a dataset without importing metadata in Python

I can *use* a data dictionary to find out about a dataset

I can *find* the shape, size and format of datasets, using Python

I can *find* the data types of variables in a dataset, using Python

I can *identify* outliers and missing values in Python

## Knowledge Prerequisites

Learners should know:

- Python programming to at least the level defined in SQA Computer Programming Level 5 (HY2C 45)
- How to use a Jupyter notebook to write, edit and run Python code
- Data understanding is part of the analysis steps
- The fundamentals of data understanding, as covered in **Data Understanding in Python**

## Lesson Requirements

	<b>PDA</b>	<b>NPA</b>	<b>Data Skills for work</b>
<b>Qualification</b>	Yes	Yes	Yes
<b>Outcome ID(s)</b>	CD7.1c, CD7.1f, WD8.1e, WD8.1f	DC5.2b, DC6.2b	A1.2, A1.3, C2.1
<b>Outcome description(s)</b>	CD7.1c Types of data CD7.1f Data quality WD8.1e Data quality WD8.1f Stages in the data analysis process	DC5.2b Explain how data can be analysed, DC6.2b Explain how data can be analysed	A1.2 Data quality A1.3 Interpretation and insight C2.1 Vocabulary used in data science and analytics
<b>Level</b>	7, 8	5, 6	Core, Analysis
<b>Software language</b>	Python	Python	Python
<b>Required equipment /software for student</b>	Lesson: PowerPoint Python notebook: Jupyter notebook environment	Lesson: PowerPoint Python notebook: Jupyter notebook environment	Lesson: PowerPoint 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:

- [numpy](#) – for scientific computing
- [pandas](#) - for data manipulation
- [s3fs](#) - an API to AWS S3 (Simple Storage Service), used to import datasets

The tasks are described in the table below.

Notebook section	Task	Description
Remove the metadata	Task 1 - Import the dataset	Import a dataset that contains metadata, without importing the metadata.
Use the metadata	Task 2 - Use the data dictionary	Use a data dictionary to find out some basic information about an unfamiliar dataset.
Find the shape, size and data types of the dataset	Task 3 - Find the shape, size and data types of the dataset	Find out some basic attributes of an unfamiliar dataset.
Identify the format of the dataset	Task 4 - Identify the format of the dataset	Through visual inspection, determine if a dataset is wide or long.
Identify outliers and missing data	Task 5 - Spot the outliers using describe()	Identify outliers in a dataset.
	Task 6 - Spot the outliers using sort_values()	Identify outliers in a dataset.
	Task 7 - Find the rows where <i>all</i> the values are missing	Find the rows where <i>all</i> the values are missing.

## Datasets

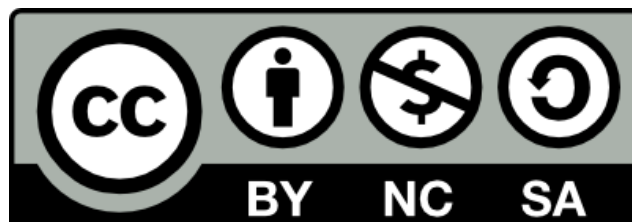
The following dataset is used in this lesson.

Dataset name	Description	Link
music_messy	<p>This is the only dataset used in this lesson.</p> <p>Contains information about popular songs from 1945-2019. The data is originally from <a href="#">Spotify</a>.</p> <p>The dataset has been created by taking a clean dataset and introducing specific changes to it that require cleansing.</p>	<a href="https://datasets.learn-data.science/music_messy.csv">https://datasets.learn-data.science/music_messy.csv</a>

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.

© 2021. This work is licensed under a [CC BY-NC-SA 4.0 license](#).



You are free to:

- **Share** – copy and redistribute the material in any medium or format
- **Adapt** – remix, transform and build upon the material

Under the following terms:

- **Attribution** — You must give [appropriate credit](#), provide a link to the license, and [indicate if changes were made](#). You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
- **NonCommercial** — You may not use the material for [commercial purposes](#).
- **ShareAlike** — If you remix, transform, or build upon the material, you must distribute your contributions under the [same license](#) as the original.

Alternative format

**If you require this document in an alternative format, such as large print or a coloured background, please contact**

**hello@effini.com**

**or**

**4th Floor, The Bayes Centre**

**47 Potterrow**

**Edinburgh**

**EH8 9BT**