

# SIT220/731 2022.T3: Task 6HD

## Data Mining Challenge

Last updated: 2022-10-30

### Contents

<b>1</b>	<b>Task</b>	<b>1</b>
<b>2</b>	<b>Additional Tasks for Postgraduate (SIT731) Students (*)</b>	<b>1</b>
<b>3</b>	<b>Artefacts</b>	<b>1</b>
<b>4</b>	<b>Intended Learning Outcomes</b>	<b>2</b>

## 1 Task

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

1. Download at least five different datasets that are part of the NHANES 2017–2020 study, see <https://wwwn.cdc.gov/nchs/nhanes/continuousnhanes/default.aspx?Cycle=2017-2020>. Merge them into a single data frame.
2. Create at least five nontrivial data visualisations and/or tables.
3. Draw insightful and interesting conclusions. Do not forget to reflect on the potential data privacy and ethics issues that arise during the data analysis process.

The PDF version of the report must be at least 10 pages long. Make it aesthetic and interesting to read.

*This HD-level task is purposely under-defined – you will not be told precisely what to do. Your aim is to discover, visualise, and explain some **interesting** relationships between data features.*

In the course of the report preparation, you should apply a wide range of data frame wrangling techniques, including filtering, aggregation in groups, missing value handling, column transformation, etc.

## 2 Additional Tasks for Postgraduate (SIT731) Students (\*)

There are no specific additional tasks, because the whole exercise has an open-ended formulation.

## 3 Artefacts

The solution to the task must be included in a single Jupyter/IPython notebook (an .ipynb file) running against a Python 3 kernel.

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**, **email address**, and whether you are an **undergraduate (SIT220)** or **postgraduate (SIT731)** student.

Make sure that your notebook has a **readable structure**; in particular, that it is divided into sections. Use rich Markdown formatting.

Imagine it is a report 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.

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

---

Submit one file via OnTrack:

1. the version of the Jupyter/IPython notebook converted to a PDF file (e.g., via *File* → *Export Notebook As* → *PDF* or convert to HTML and from that to PDF with your web browser; any method will do).

You do not need to submit the .ipynb file via OnTrack, but you must store it for further reference – a marking tutor might ask for it later, e.g., at the end of the trimester.

## 4 Intended Learning Outcomes

ULO	Is Related?
ULO1 (Data Processing/Wrangling)	YES
ULO2 (Data Discovery/Extraction)	YES
ULO3 (Requirement Analysis/Data Sources)	YES
ULO4 (Exploratory Data Analysis)	YES
ULO5 (Data Privacy and Ethics)	YES