# Wrangle Report

The three datasets viz. archive, image prediction, and API-generated reactions were assessed in two ways:
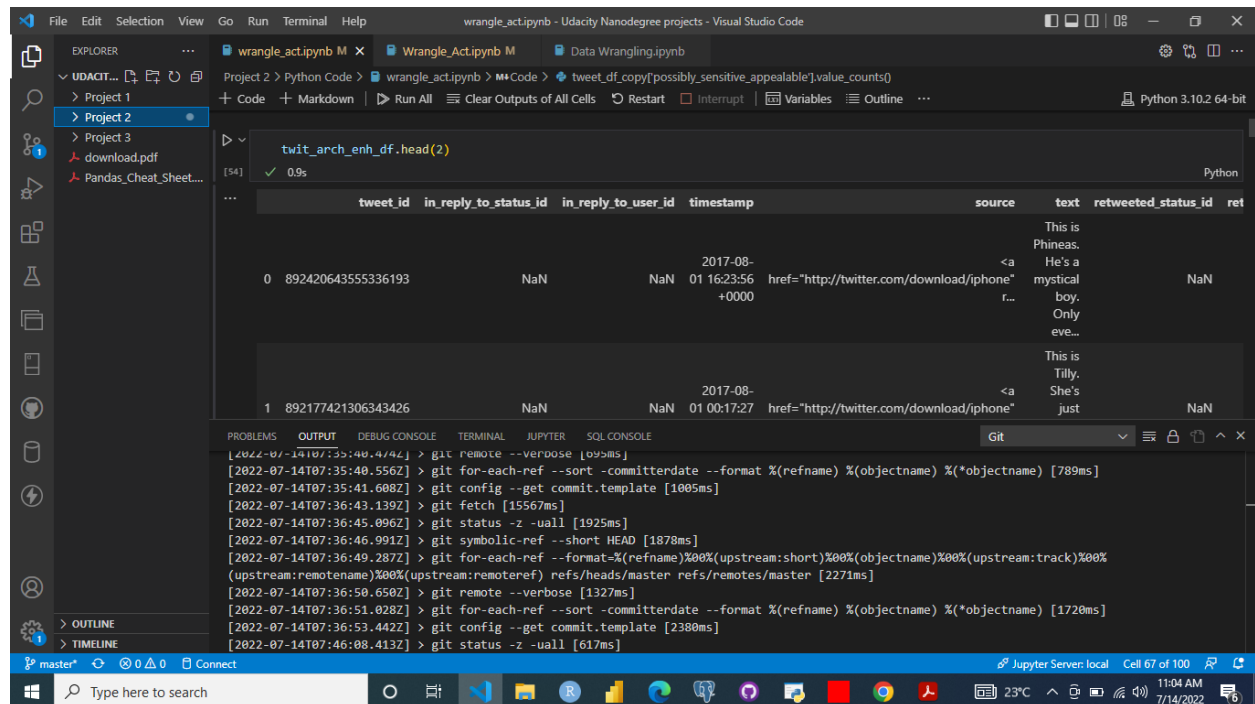
● Visual assessment
● Programmatic assessment

## Visual Assessment

This involves the use of spreadsheet applications such as Google Sheet, Microsoft Excel, and some pandas functions in Python.

In order to visually access these datasets adequately, I tried to view and scroll through the data examining what might be missing, wrongly inputted, or inaccurate structures.

However, I recognize the limitations of using spreadsheets for assessment or examination as in many cases data might be so voluminous or be available in a format not so easily readable by humans. Hence, I utilized the pandas' Python package. Some of the pandas' data frame methods such as

1. Head: This shows the first five observations. The default shows five observations but this may be modified by parsing the count as an argument. I choose to set the argument to 2 to limit the space.



2. Columns: This shows the names of the columns in the data frame.

3. Tail: This shows a specified number of observations. The default shows five observations but the same may be modified by parsing in the count as an argument.
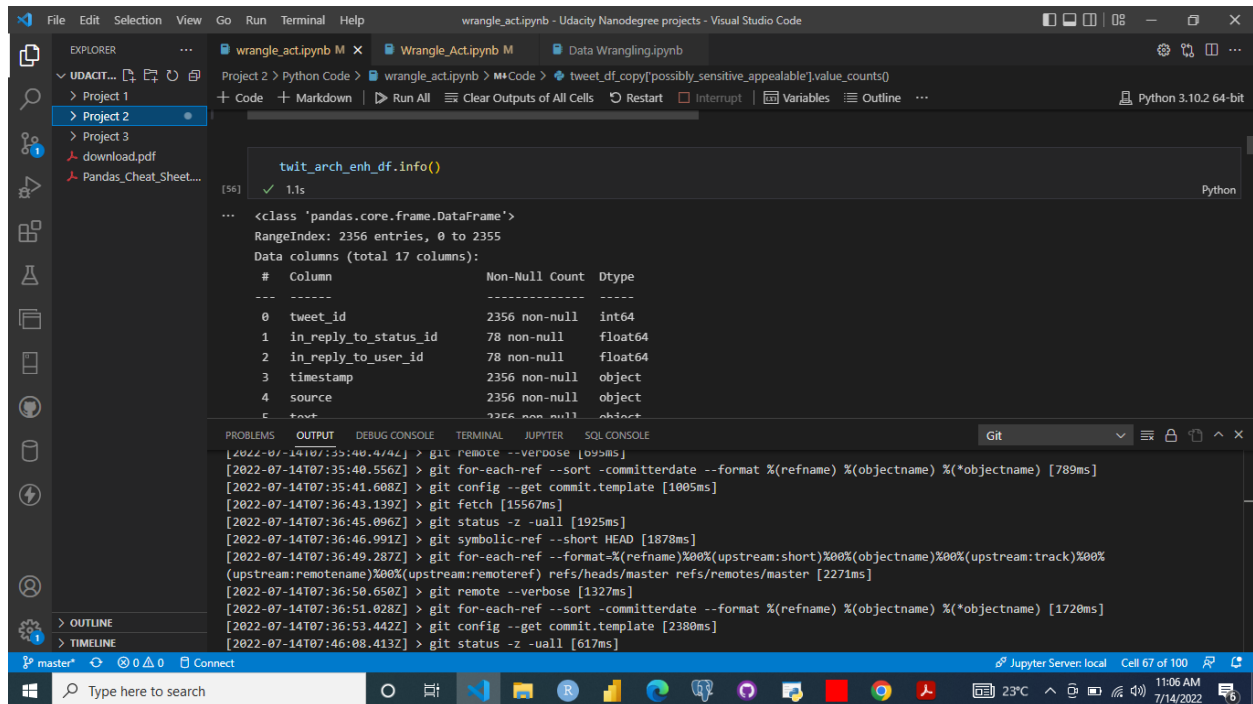


**Programmatic Assessment**

Programmatic assessment involves the use of code to evaluate the possible causes of dirty or untidy datasets.

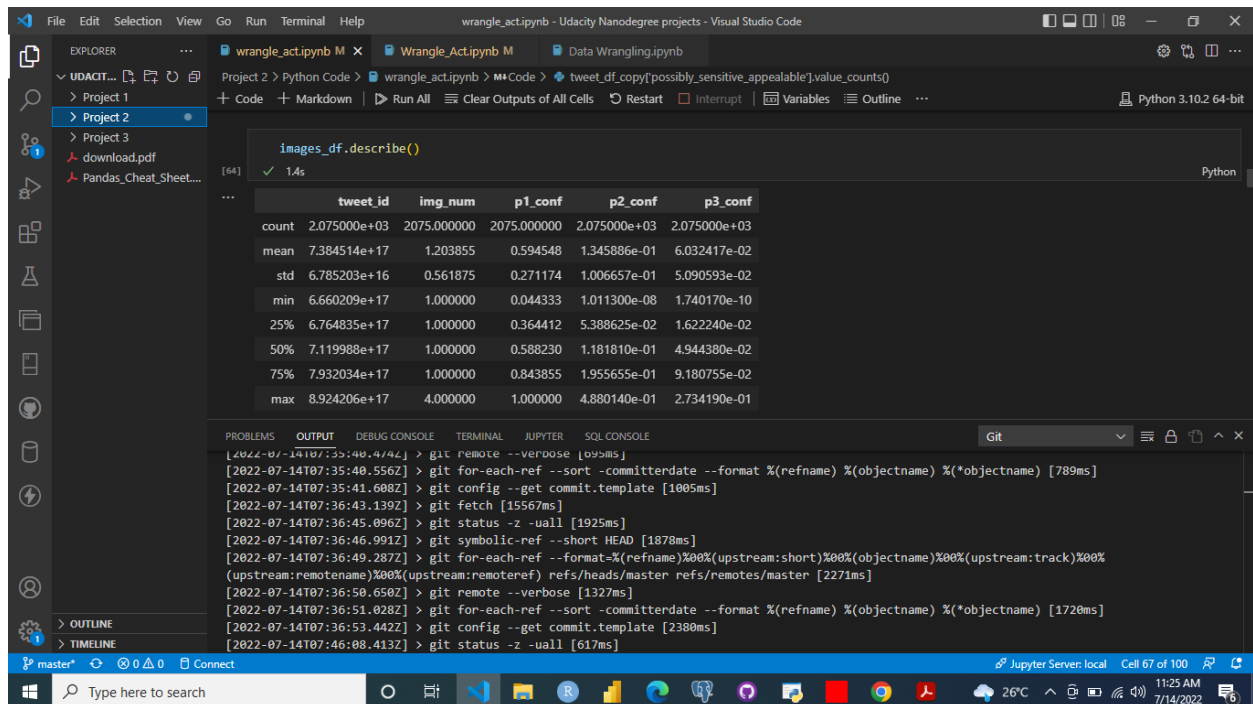The following are the pandas' data frame methods utilized:

1. Info: The info method provides a series of vital information, such as the dimension, columns list, datatypes, data types, the total number of rows, and columns, total nonnull entries, memory usage, etc.



2. Describe: This gives a quick summary of statistics such as mean, min, max, quartiles, etc.

3.  Shape: To get the dimension of the data frame in terms of the number of rows and columns.

Other programmatic assessment functions used include:

4.  IsNull: This checks null values in the dataset. For this dataset, we had a lot of them.
5.  Duplicated: Although we did not have cases of duplicated tweet_id, this function checks for duplicates.
6.  Unique: This finds unique values and combined with the sum will have a reasonable insight into the data
7.  Count
8.  Sum
9.  Isna
10. Value_counts