



[◀ Return to Classroom](#)

# Wrangle and Analyze Data

## REVIEW

## HISTORY

### Meets Specifications

Dear Learner,

I must say this submission was indeed enjoyable to review. I appreciate and commend the efforts and hard work put into this piece. Congratulations 🙌🙌 for making it pass this stage of learning with us and I wish that this spirit is carried forward in subsequent projects. You should be proud of yourself because success is no accident. It is hard work, perseverance, learning, studying, sacrifice and most of all, love of what you are doing or learning to do. I wish you all the best. 💪

### Code Functionality and Readability

All project code is contained in a Jupyter Notebook named `wrangle_act.ipynb` and runs without errors.

Good work! All cells of the notebook run on my end without errors.

### Learning Notes

I am a fan of using shortcuts with Jupyter Notebook. Check out [this medium post](#) on **Jupyter Notebook Shortcuts**.

The Jupyter Notebook has an intuitive, easy-to-follow logical structure. The code uses comments effectively and is interspersed with Jupyter Notebook Markdown cells. The steps of the data wrangling process (i.e. gather, assess, and clean) are clearly identified with comments or Markdown cells, as well.

A remarkable job was done to organize the notebook in an easy-to-follow format.

## Learning Notes

It is also good practice to use functions to avoid any code repetition.

- [Why use functions in programming?](#)

Inline comments could also help with code follow up and ease the work of other programmers working on the same project.

- [Why should we comment code?](#)

## Gathering Data

Data is successfully gathered:

- From at least the three (3) different sources on the Project Details page.
- In at least the three (3) different file formats on the Project Details page.

Each piece of data is imported into a separate pandas DataFrame at first.

Good work gathering data from three different sources and converting all of them into python data frames before performing wrangling.

## Assessing Data

Two types of assessment are used:

- Visual assessment: each piece of gathered data is displayed in the Jupyter Notebook for visual assessment purposes. Once displayed, data can additionally be assessed in an external application (e.g. Excel, text editor).
- Programmatic assessment: pandas' functions and/or methods are used to assess the data.

Both visual and programmatic assessments are used in the notebook and the results are well documented.

At least eight (8) data quality issues and two (2) tidiness issues are detected, and include the issues to clean to satisfy the Project Motivation. Each issue is documented in one to a few sentences each.

Good work identifying the quality and tidiness issues in the dataset. Well done here!

## Learning Notes

Data assessment involves examining data quality and tidiness.

**Quality issues** pertain to the content of data. Low quality data is also known as dirty data. There are four dimensions of quality data:

- **Completeness:** do we have all of the records that we should? Do we have missing records or not? Are there specific rows, columns, or cells missing?
- **Validity:** we have the records, but they're not valid, i.e., they don't conform to a defined schema. A schema is a defined set of rules for data. These rules can be real-world constraints (e.g. negative height is impossible) and table-specific constraints (e.g. unique key constraints in tables).
- **Accuracy:** inaccurate data is wrong data that is valid. It adheres to the defined schema, but it is still incorrect. Example: a patient's weight that is 5 lbs too heavy because the scale was faulty.
- **Consistency:** inconsistent data is both valid and accurate, but there are multiple correct ways of referring to the same thing. Consistency, i.e., a standard format, in columns that represent the same data across tables and/or within tables is desired.

**Tidiness issues** pertain to the structure of data. These structural problems generally prevent easy analysis. Untidy data is also known as messy data. The requirements for tidy data are:

- Each variable forms a column.
- Each observation forms a row.
- Each type of observational unit forms a table.

## Cleaning Data

The define, code, and test steps of the cleaning process are clearly documented.

The different steps of the cleaning process are clearly documented. We have the `define`, `code` and `test` steps which are clearly stated with some explanations of what process you intend to do at each level. Excellent work!

Copies of the original pieces of data are made prior to cleaning.

All issues identified in the assess phase are successfully cleaned (if possible) using Python and pandas, and include the cleaning tasks required to satisfy the Project Motivation.

A tidy master dataset (or datasets, if appropriate) with all pieces of gathered data is created.

Indeed, copies of the original pieces of the data are made before the cleaning process.

## Storing and Acting on Wrangled Data

Students will save their gathered, assessed, and cleaned master dataset(s) to a CSV file or a SQLite database.

The cleaned master datasets are saved to a `CSV` file. Good work!

## Learning Notes

- Check out [this StackOverflow Thread](#) on Pandas writing dataframe to CSV file.
- Also, take a look at the [pandas.DataFrame.to\\_sql](#) and this [Stackoverflow thread](#) for an example of saving pandas dataframe to SQLite.
- [Saving a pandas Dataframe as a CSV](#)

The master dataset is analyzed using pandas or SQL in the Jupyter Notebook and at least three (3) separate insights are produced.

At least one (1) labeled visualization is produced in the Jupyter Notebook using Python's plotting libraries or in Tableau.

Students must make it clear in their wrangling work that they assessed and cleaned (if necessary) the data upon which the analyses and visualizations are based.

Good work analyzing the cleaned data and plotting some visualizations of the data.

## Learning Notes

- Check out the [python visualization documentation](#) for various ways of visualizing data.
- There are several other ways to visualize data including **Box plots**, **Line graphs**, **Pie charts**. Check out [this documentation](#) on Visualization with Seaborn.

## Report

The student's wrangling efforts are briefly described. This document (wrangle\_report.pdf or wrangle\_report.html) is concise and approximately 300-600 words in length.

The write-up ( `wrangle_report.pdf` ) describing the wrangling efforts made in the project.

The three (3) or more insights the student found are communicated. At least one (1) visualization is included.

This document (act\_report.pdf or act\_report.html) is at least 250 words in length.

The insights have been reported and they have been analysed in detail.

## Project Files

The following files (with identical filenames) are included:

- wrangle\_act.ipynb
- wrangle\_report.pdf or wrangle\_report.html
- act\_report.pdf or act\_report.html

All dataset files are included, including the stored master dataset(s), with filenames and extensions as specified on the Project Submission page.

The required files (with identical filenames) are included. Excellent! In addition, all dataset files are included. Well done!

 [DOWNLOAD PROJECT](#)

[RETURN TO PATH](#)

Rate this project