



[Return to Classroom](#)

# Investigate a Dataset

## REVIEW

## HISTORY

### Requires Changes

5 specifications require changes

Great submission!

You met most of the requirements, just a few minor issues that need to improve.

Keep working on it!

### Code Functionality

- All code is functional and produces no errors when run.
- The code given is sufficient to reproduce the results described.

All code works well without any errors, great job!

#### Further Consideration

Working through some examples of [Automate the Boring Stuff with Python](#) will level up your skills in Python.

- The project uses NumPy arrays and Pandas Series and DataFrames where appropriate rather than Python lists and dictionaries.

- Where possible, vectorized operations and built-in functions are used instead of loops.

Numpy and Pandas have been used, well done!

Below is an article about the Pandas performance, you may have a look for your reference:[Pandas Performance](#)

- The code makes use of at least 1 function to avoid repetitive code.
- The code contains good comments and meaningful variable names, making it easy to read.

Although you did not define any function, I did not see the significant repetitive code. Good job!

#### Required Improvement:

Try to leverage the comments to explain your code, so the whole analysis will be easier to follow.

You may refer to the below article about [the importance of good comments]  
(<https://realpython.com/lessons/importance-writing-good-code-comments/>)

## Quality of Analysis

- The project clearly states one or more questions, then addresses those questions in the rest of the analysis.

Questions have been listed at the beginning and addressed in the later analysis, well done!

## Data Wrangling Phase

- The project documents the steps that were taken to clean the data, such as merging multiple files, handling missing values, etc.

You did a perfect data wrangling, you have looked into the numeric statistic, data type, missing value, duplicates and outliers. Each step has been well documented.

## Exploration Phase

- The project investigates the stated question(s) from multiple angles.

- The project explores at least three variables in relation to the primary question. This can be an exploratory relationship between three variables of interest, or looking at how two independent variables relate to a single dependent variable of interest.
- The project performs both single-variable (1d) and multiple-variable (2d) explorations.

Multiple-variable explorations have been done, fantastic.

### Required Improvement:

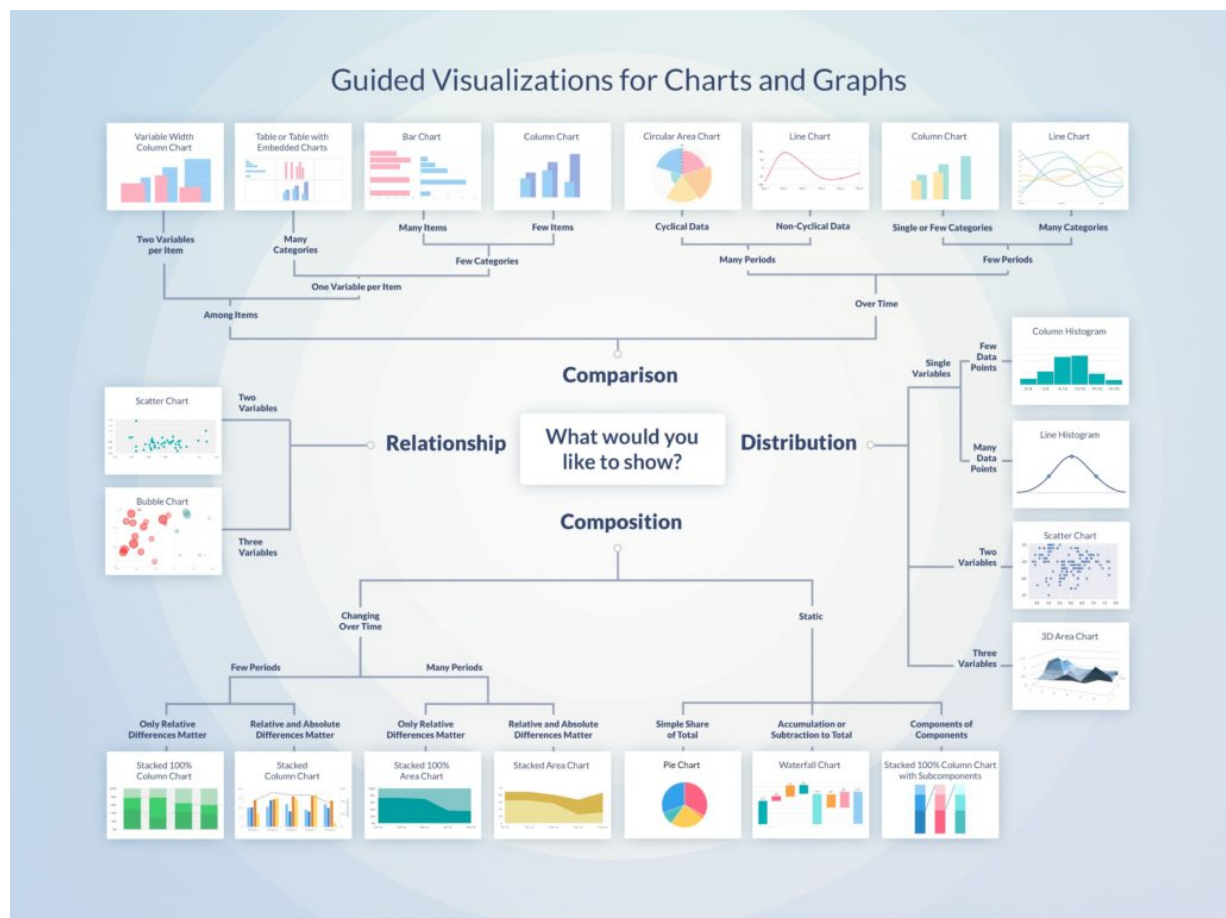
In order to meet the rubric, you still need to do single-variable explorations for at least **THREE** variables. Histogram or boxplot will be a good choice.

Below is a very good article to introduce these two explorations: [Data Exploration](#)

- The project's visualizations are varied and show multiple comparisons and trends.
- At least two kinds of plots should be created as part of the explorations.
- Relevant statistics are computed throughout the analysis when an inference is made about the data.

Multiple plots have been properly used, great job!

Below is a guideline for graphing, just for your future reference.



## Conclusions Phase

- The Conclusions have reflected on the steps taken during the data exploration.
- The Conclusions have summarized the main findings in relation to the question(s) provided at the beginning of the analysis accurately.
- The project has pointed out where additional research can be done or where additional information could be useful.
- The conclusion should have at least 1 limitation explained clearly.
- The analysis does not state or imply that one change causes another based solely on a correlation.

You did conclude the findings, fantastic!

#### Required Improvement:

To meet the rubric, you still need to address the limitations.

Limitations can exist due to constraints on analysis design or data, and these factors may impact the findings of your conclusion.

You should clearly acknowledge any limitations in your conclusion ( with subsection ) in order to show the audience that you are aware of these limitations and to explain how they affect the conclusions.

- Is the data provided by the dataset sufficient to answer your question?
- Is the size of the dataset is sufficient to give a good judgment about the questions you asked?

You may refer to the below article.

[What are the limitations of a study and how to write them?](#)

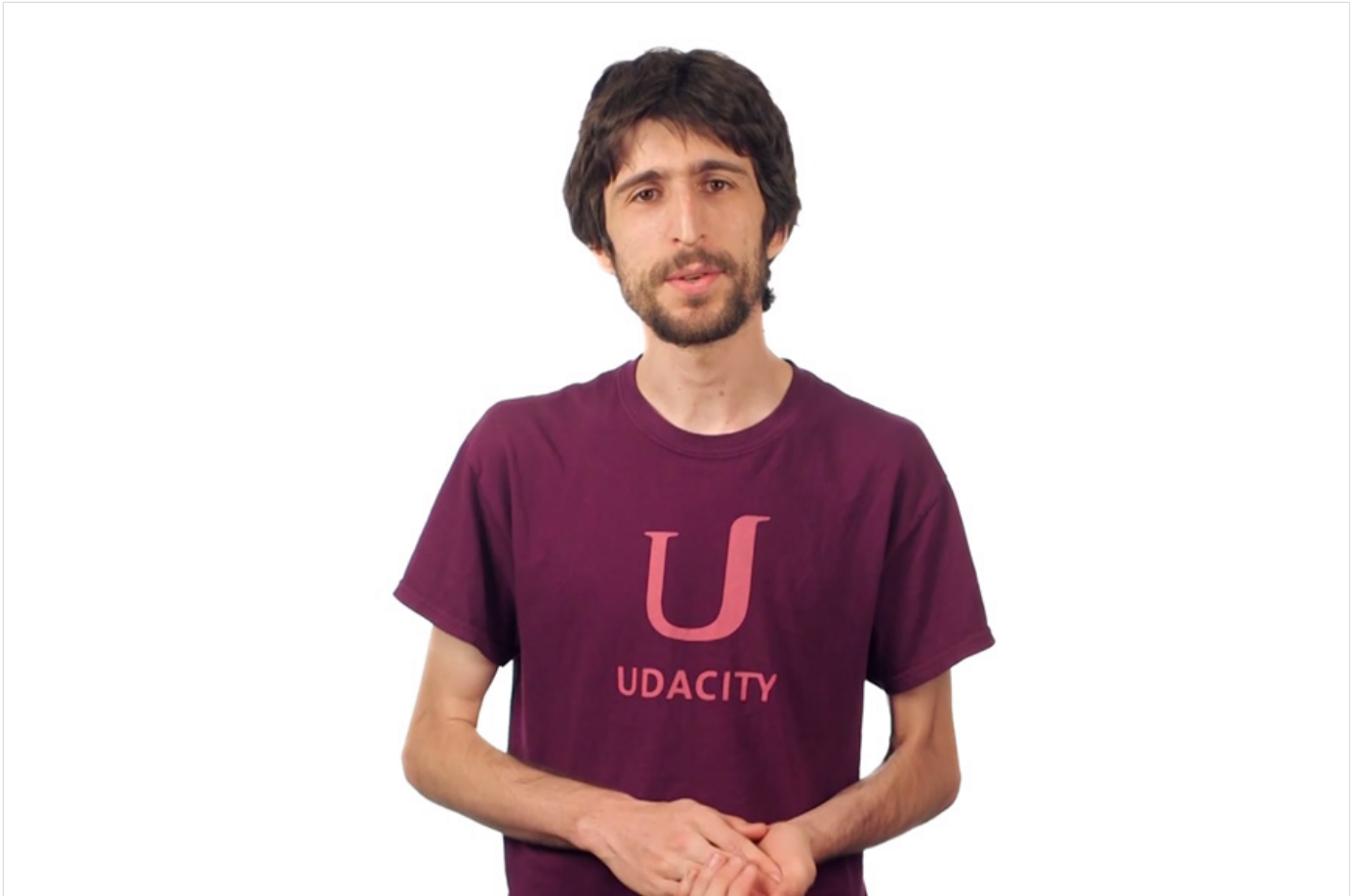
## Communication

- The code should have ideally the following sections: Introduction; Questions; Data Wrangling; Exploratory Data Analysis; Conclusions, Limitation.
- Reasoning is provided for each analysis decision, plot, and statistical summary.
- Interpretation of plots and application of statistical tests should be correct and without error.
- Comments are used within the code cells.
- Documented the flow of analysis in the mark-down cells.

Validate this point once single variable explorations added

- Visualizations made in the project depict the data in an appropriate manner (i.e., has appropriate labels, scale, legends, and plot type) that allows plots to be readily interpreted.

Validate this point once single variable explorations added

[RESUBMIT](#)[DOWNLOAD PROJECT](#)

## Best practices for your project resubmission

Ben shares 5 helpful tips to get you through revising and resubmitting your project.

[Watch Video \(3:01\)](#)

[RETURN TO PATH](#)