

[Return to Classroom](#)

Investigate a Dataset

REVIEW

HISTORY

Requires Changes

1 specification requires changes

Hello Mohamed

Nice job and your project look great 😊

small easy step and your project will meet specifications. KEEP UP THE GOOD WORK!

Summary of the Changes Required

🟡 You have to implement a user-defined function to avoid writing repetitive code.

Great work in the project

- ✅ I really like the usage of comments throughout the project. Kept me interested while going through the project
- ✅ The reasoning beneath every visualization is appropriate and well explained.
- ✅ I appreciate that you have included the introduction of the project which is quite detailed and gives reader an overview of the project as well.
- ✅ I must say I enjoyed seeing the visualizations that you created within the project

Don't get upset or disappointed, you did a great job which deserves a big compliment, think that those changes are a great opportunity to learn more and perfect your skills.

I wish you good luck with your Nanodegree!

Stay Safe and Stay Udacious 🟡!

Suggestion

If you have a question or query, please visit the [Knowledge Platform](#), where you can ask expert mentor, questions about your project or any specific topic. You can also search for topics that have already been asked by other students in your Nanodegree. On the [Knowledge Platform](#), our expert mentor staff will give you with ongoing assistance.

Common Student Questions

Question: I still haven't completely figured out how to clean the data completely and how much data to drop or keep for analysis. How to decide this?

Answer: One of the next projects you will be doing would be [Wrangle and Analyze Data](#) which specifically focuses on data cleaning and using a structure to clean a dataset. Just keep going and keep honing your skills in Data Analysis.

Question: How can I improve my code readability and code structure

Answer: You can write the code following the [PEP8 style guidelines](#). Also, instead of using single line comments, you can use multi-line comments in python to increase the readability. [Read more about it here](#)

Question: Are duplicate rows in these datasets significant/ need to be dropped?

Answer: You can read about the relevant answer on this [Knowledge Hub answer by Abdulaziz](#)

Code Functionality



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

Things done well

The submitted code works well as it doesn't produce errors when run. Also, it's sufficient to reproduce the results described.

The coding structure and logical flow of your coding practices are impressive. Keep up the good work.

Additional Readings

[Jupyter magic to handle notebook exceptions](#)

[Jupyter Notebook Tips and Tricks](#)

[What is the right way to debug in iPython notebook?](#)



- 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.

You have shown a good understanding of built-in functions and the correct application thereof

You have shown a good understanding of built-in functions and the correct application thereof.

Learning Notes

Here are a couple of links with useful pandas and numpys built-in functions and methods:

[Pandas Cheat Sheet](#)

[Numpy Cheat Sheet](#)



- 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.

Things done well

- You have given a brief introduction about the dataset you will be working on. That's a really good practice for a reader (or yourself when you will revisit the project in future) to get an idea about the dataset and its related fields.
- You have also posed the questions you will be exploring at the beginning of the project. This is good practice as it lays the ground work for the rest of the project and give you a direction to think and analyze the dataset even before delving deep into it.
- The structure of your notebook is clean and has a logical flow. Different sections are clearly shown for each one of the steps of the data wrangling process.
- Your code is properly commented and contain good variable names which is making your code easy to read. You have incorporated many markdowns and in-code comments. Markdowns are very important as this allows the reader to follow along with the intentions of the author. Good job!

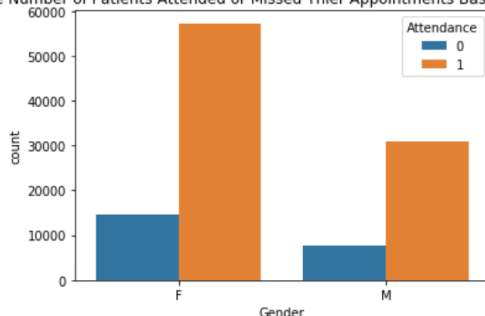
Things that Require Changes

- You should create any custom defined function (as it is one of the rubric requirement - The code makes use of functions to avoid repetitive code) making your job easier and making your code look cleaner. For example - for the multiple visualizations you plotted of the same kind, you can create a function and invoke that function whenever you want to plot that kind of visualization. Please [refer this link](#) to know more about how to create functions to avoid repetitive code.

```
1 temp = sns.countplot(x=appoints_df['Gender'], hue=appoints_df['Attendance'], data=appoints_df)
2 temp.set_title("The Number of Patients Attended or Missed Thier Appointments Based on Gender")
```

```
Text(0.5,1,'The Number of Patients Attended or Missed Thier Appointments Based on Gender')
```

The Number of Patients Attended or Missed Thier Appointments Based on Gender



For example - This code snippet can be enclosed inside a user defined function

Additional Readings

[Six steps to more professional data science code](#)
[Reviewing Data Science Projects](#)

Quality of Analysis



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

Things done well

You have done a great job with the questions you have posed! Also, each of these questions are explored thoroughly.

In EDA, great questions help you focus on relevant parts of your data and direct your analysis towards meaningful insights. Questions should be measurable, clear and concise. They should be designed to either qualify or disqualify potential solutions to your specific problem or opportunity. You have gone above and beyond in this area :)

Additional Readings

[Your Data Won't Speak Unless You Ask It The Right Data Analysis Questions](#)
[How Do Data Scientists Ask the Right Questions?](#)

Data Wrangling Phase



The project documents any changes that were made to clean the data, such as merging multiple files, handling missing values, etc.

The structure of your notebook is clean and has a logical flow. Comments are used to clearly identify each one of the steps of the data wrangling process. You have correctly identified issues regarding the datasets.

Comments:

Data wrangling is the process of cleaning and unifying messy and complex data sets for easy access and analysis. It is about targeting a field, row or column in a data set, and applying an action such as joining, parsing, cleansing, consolidating, or filtering to create the desired output, which will then be used down the road. Data wrangling involves activities like:

- Remove unused columns.
- Remove duplicate rows.
- Change data formats (date columns)
- Discard missing values.

Benefits:

- it makes your data useful
- it can be organized into a standardized and repeatable process that moves and transforms data sources into a common format, which can be reused multiple times.

Additional Readings

[The Growing Importance Of Data Cleaning](#)

[Data Cleaning Using Python Pandas](#)

[Data Cleaning with Python and Pandas: Detecting Missing Values](#)

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.

The analysis makes good use of both single (1D) and multiple (2D) variable explorations to investigate different features and the relations between these features in the dataset. Questions are investigated from a single variable perspective as well as a multiple-variable perspective. Good Job!

Learning note:

Below is a table that allows us to see the distinction between 1D and 2D explorations:

Summary: Differences between univariate and bivariate data.

Univariate Data	Bivariate Data
<ul style="list-style-type: none"> • involving a single variable 	<ul style="list-style-type: none"> • involving two variables
<ul style="list-style-type: none"> • does not deal with causes or relationships 	<ul style="list-style-type: none"> • deals with causes or relationships
<ul style="list-style-type: none"> • the major purpose of univariate analysis is to describe 	<ul style="list-style-type: none"> • the major purpose of bivariate analysis is to explain
<ul style="list-style-type: none"> • central tendency - mean, mode, median • dispersion - range, variance, max, min, quartiles, standard deviation. • frequency distributions • bar graph, histogram, pie chart, line graph, box-and-whisker plot 	<ul style="list-style-type: none"> • analysis of two variables simultaneously • correlations • comparisons, relationships, causes, explanations • tables where one variable is contingent on the values of the other variable. • independent and dependent variables
Sample question: How many of the students in the freshman class are female?	Sample question: Is there a relationship between the number of females in Computer Programming and their scores in Mathematics?

Additional Readings

[1D and 2D Variable Explorations](#)

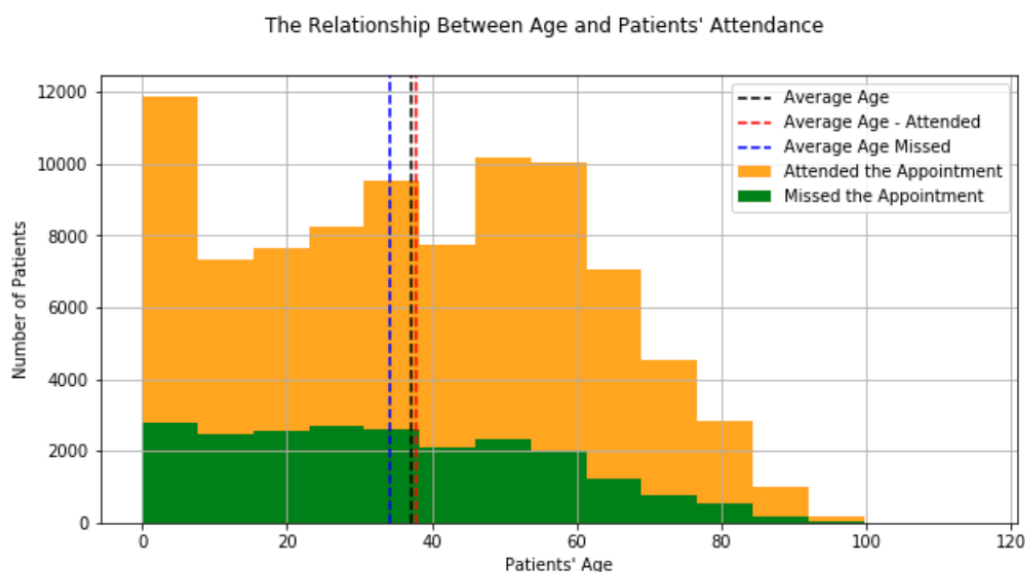
[A Comprehensive Guide to Data Exploration](#)



- 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.

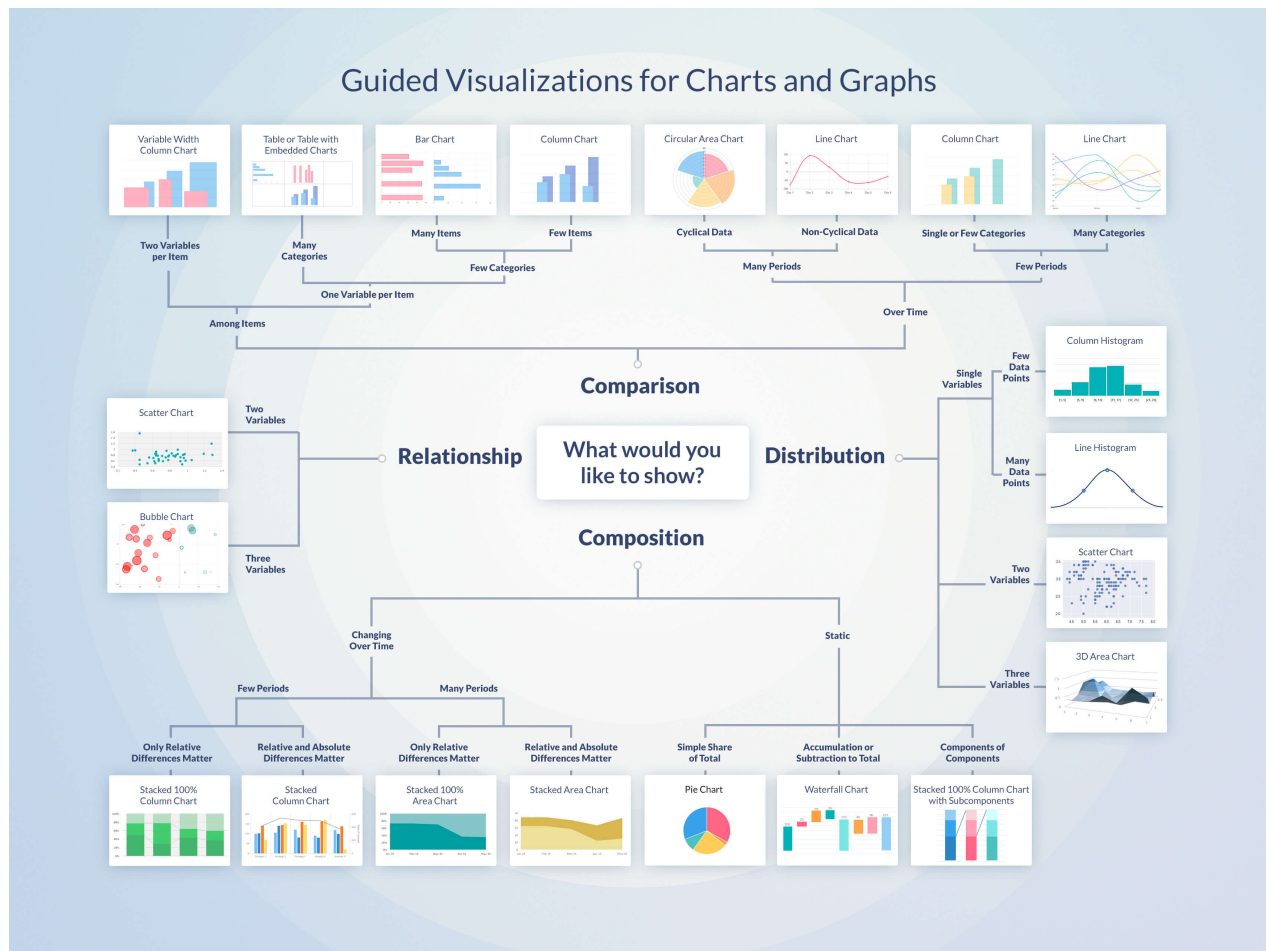
The visualizations appropriately support the investigation of the data set.

Visualizing data requires a lot of patience and determination because it's not easy selecting the best visualization to match with a given data type. Well enough, the project rightly builds descriptive visualizations using a variety of plots. I really found the below visualization interesting especially Average dotted lines plotted.



Learning note:

Below is the screenshot I have provided in relation to plotting data. I have found this very useful when deciding what types of plots to use. I hope this will be useful for you in future too:



Additional Readings

[A step-by-step guide to Data Visualizations in Python](#)

[Data visualization with different Charts in Python](#)

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.

Good Job adding Conclusions and Limitations to your project summarising your findings and figuring out what issues this dataset lacked or what you dropped affecting the overall analysis

Additional Readings

[How to Plan and Organize a Data Science/Analytics Project?](#)

[Write Discussion and Conclusion of a project](#)

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.

Well done!

It is very important to communicate the results adequately; however, it is also very important to describe each activity, analysis, or graph. This will allow your audience to understand what you are doing and how you are doing it. Moreover, the reasoning makes your work organized, formal, and sophisticated.



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.

Visualization presented clearly depict the data and represent the questions posed. The plots are well structured and easy to interpret. The plots have clearly represented titles and labels. Good job.

Comments:

One of the most important steps in creating an impactful visualization is making sure all of its elements are labeled appropriately. The text components of a graph give your reader visual clues that help your data tell a story and should allow your graph to stand alone, outside of any supporting narrative.

 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

Rate this review

START

