

[Return to Classroom](#)

Investigate a Dataset

REVIEW

HISTORY

Meets Specifications

Very impressive submission! I can see your hard work reflected in your project 🏆 Congratulations on achieving this and good luck on your way to master data analysis 😊

Code Functionality

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

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.

Good work

Suggestion

Here are a few Pandas built-in methods that are very useful for exploring variables in this project:

- [Boolean-Indexing](#)
- [Group-by](#)
- [Value-Counts](#)
- [Series.map](#)
- [Working-with-text-data](#)

The code makes use of functions to avoid repetitive code. The code contains good comments and variable names, making it easy to read.

Excellent job! solid code and well documented 😊

Quality of Analysis

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

Awesome job! As always we want you to go above and beyond, here are some suggestions of interesting questions for this dataset:

1. How is popularity trending over time?
2. How are revenues trending over time?
3. How is runtime trending over time?

4. Do top ratings movies always generate big revenue?
5. Do higher budget movies always generate big revenue?
6. Is there any impact of vote count on revenue?
7. Can we provide a list of the most popular directors based on ratings?
8. Can we provide a list of directors that generates big revenue?
9. What are typical runtimes for directors? Is there a duration preferred by directors?
10. Is there a relation between popularity and revenue for directors? etc.

Project: Investigating the Profitability of Films

Table of Contents

- [Introduction](#)
- [Data Wrangling](#)
- [Exploratory Data Analysis](#)
- [Conclusions](#)

Movie Data Analysis

Introduction

This dataset was taken from The Movie Database (TMDb), "a community built movie and TV database." (www.themoviedb.org/about) Each row corresponds to a movie and includes a range of data about each film. Relevant data to be used in the following analysis include the following variables:

- original_title
- genres
- release_date
- release_year
- budget_adj (budget in terms of 2010 dollars)
- revenue_adj (revenue in terms of 2010 dollars)

In this report, I explore the following questions:

1. How has the profitability of making films changed over time?
2. How does profitability vary for films released during different months?
3. How does a film's budget relate to its profitability?
4. How does a film's genre relate to its profitability?

Throughout my analysis film profitability (as calculated by subtracting each film's adjusted budget from its adjusted revenue) will be the dependent variable, while release year, release month, budget, and genre will be the independent variables.

Data Wrangling Phase

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

Good work in implementing a Data Wrangling Phase

Suggestion

The most important aspect of Data Wrangling is to clean or transform the data preparing it for analysis.

One main issue is having missing data while conducting analysis, which can provide skew/bias results. Luckily there are a few methods that Pandas provide to deal with these issues:

- The first thing to do is to always [Identify the missing values](#) within the dataset. The few steps after this explain how to deal with the missing data
- If there are columns with a few rows of missing data the [Dropna method](#) could be used to drop the missing rows.
- If there are rows with missing data the [Fillna-method](#) can be used instead of dropping them completely (This method can vary with the data and the project)
- The final option is if there are way too many missing values within a column it is best to drop the column completely using the [Drop-column-method](#)

Data Wrangling does not only involve Identifying and dealing with missing values but also involves in transforming the data to a more effective state to target the analysis. Here are other wrangling methods:

- **Binning or Cutting** Groups continuous or numerical values into smaller groups or 'bins'
- **Pandas-Dummies** Transforms categorical data into dummy/indicator variables

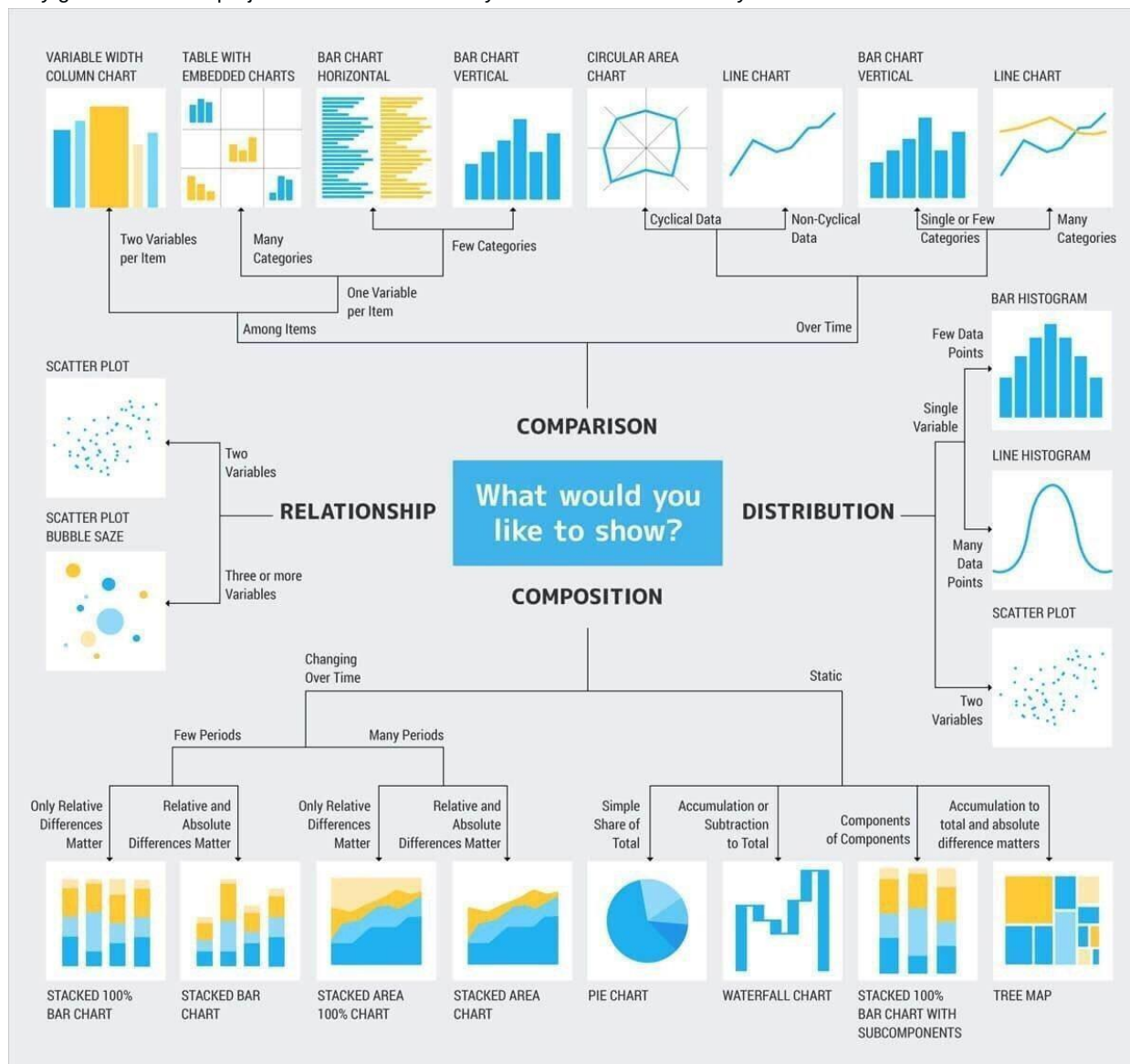
Exploration Phase

The project investigates the stated question(s) from multiple angles. At least three variables are investigated using both single-variable (1d) and multiple-variable (2d) explorations.

The project's visualizations are varied and show multiple comparisons and trends. Relevant statistics are computed throughout the analysis when an inference is made about the data.

At least two kinds of plots should be created as part of the explorations.

Very good ! for future projects let me recommend you [these](#) tools to choose your visualizations



Conclusions Phase

The results of the analysis are presented such that any limitations are clear. The analysis does not state or imply that one change causes another based solely on a correlation.

Congratulations, your project is super impressive 😊

Communication

Reasoning is provided for each analysis decision, plot, and statistical summary.

Fantastic 🙌

Visualizations made in the project depict the data in an appropriate manner that allows plots to be readily interpreted.



 [DOWNLOAD PROJECT](#)