# ENTITY

## Starting your First Data Science Portfolio Project!

Google Classroom Extra Activities for the modules from here on out, if done with passion and dedication, can lead to a complete portfolio project before you even reach DSO110! This means you can enter your data science interviews with TWO full real-world projects in your portfolio, the optimal numbers according to our data science mentors who are currently in the workforce. If you choose not to put this amount of effort into these activities then, at the very least, participating in these at least a bit will result in a substantial reusable code base that will give you a *huge* head start when you reach DSO110, the final project course module! This will give you more time to make your final ENTITY portfolio project the best it can be. Having multiple portfolio projects and/or higher quality projects will give you that competitive edge in the job market, which is what we want for ALL our ENTITY graduates!
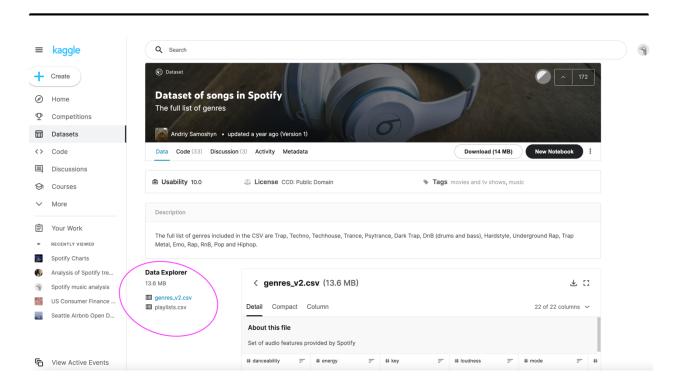
Here we go, let's get started!!

## Introduction

On Kaggle.com, search different topics that
   (a) interest you
   (b) relate to your academic or work history
   (c) seem like practical subjects for a data science project

Choose 5 dataset listings from different Kaggle URLs. Here is an example of one dataset listing:
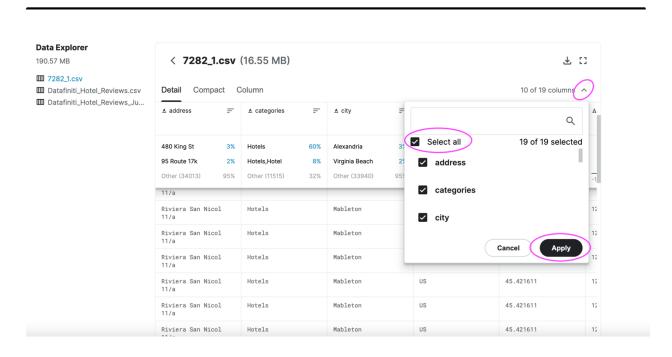
Although there are 2 CSV files, this is considered one dataset listing (see below):
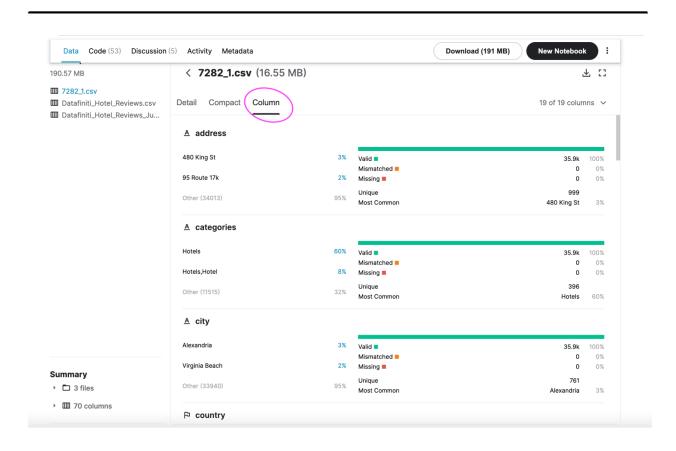
ENTITY



The 5 can be all one topic or different topics. Make sure each dataset or related datasets (like the Spotify example above) have categorical variables and numeric variables (floats/ints). You can check the datasets contents by *clicking dropdown arrow > checking "Select all" box > clicking "Apply"*.
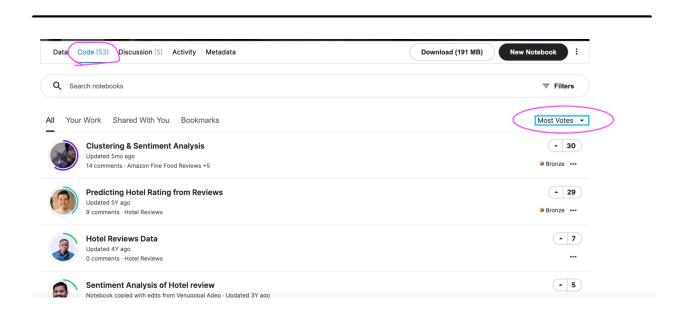
Also, make sure columns of interest do not have too many missing values. You can investigate by clicking "Column" tab.

This is not required, but a monumental plus. Check if the datasets have notebooks associated with them. This is the benefit of Kaggle.com and the ethos of open-source platforms! You can get great ideas by "standing on the shoulders of giants" and on many occasions, do not have to reinvent the wheel! You can find associated notebooks in the "Code" tab. Click "Most Votes" to find the highest rated notebooks. This is how data scientists operate in the real world—they help one another, use one another's code, and pay it forward on these platforms when they can!

**This is a time-consuming task, but a task you will inevitably have to do.** You will be glad you have done extensive dataset perusing months before you reach DSO110: Final Project.

After you have completed all of the above, you are ready to move on to the activities associated with each lesson module.

# Lesson 1

Out of the 5, choose a dataset that has columns that seem unnecessary and/or that could be renamed. Make a subset of that dataframe without those columns and/or rename the columns.
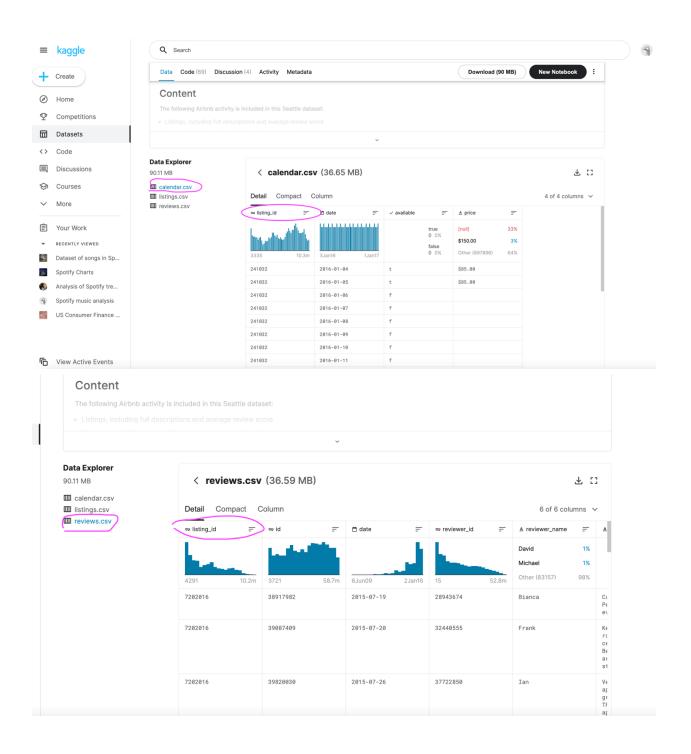
You may do this in Python and/or R.

# Lesson 2

*Option 1 (if dataset listing has multiple related datasets):*

Out of the 5, choose a dataset listing that has multiple datasets and merge them on the

appropriate columns. Example:

The above example with data on Airbnb reviews has two datasets that can be merged on *listing_id*.

*Option 2 (if dataset listings only have one dataset file or datasets don't have columns appropriate for merging):*

Search Kaggle for another dataset that has a related column and dataset subjects have value for comparison. For instance, a marketing dataset with zip codes can be merged with another dataset that has the zip code populations and average annual income. This is just one example of the many possible scenarios.

*Option 3 (if Option 1 and 2 are not viable or becoming a time suck):*

Choose a dataset and separate it into two subsets with one categorical column in common. Merge them back into one dataset.

You may do any of these options in Python and/or R.

## Lesson 3

*Option 1:*

Out of the 5 datasets, choose one and find a categorical column (string) where the groups have no inherent order (nominal) and dummy code that column.

*Option 2:*

Out of the 5 datasets, choose one and find a categorical column (string) where there are only 2 groups (binary), and change the groups to 0 and 1.

*Option 3:*

Out of the 5 datasets, choose one and find a column that is a continuous variable. Recode that column as a string providing a useful group name. For example, let's say there is an age column.

You could recode as the following:

age < 18 = "minor"
age > 18 = "adult"
age > 55 = "senior"

You may do any of these options in Python and/or R.

## Lesson 4

Out of the 5 datasets, choose one column of interest that is a continuous variable and create a histogram. Then, choose one column of interest that is a categorical variable and create a bar chart.

## Lesson 5

Out of the 5 datasets, choose one that has 3 or more continuous variables. Make a correlation heatmap.

You may do this in Python and/or R.

## Lesson 6

Out of the 5 datasets, choose one and make 2 visualizations with each representing 2 or more variables.

## Lesson 7

Peruse the internet and find an infographic that makes an impression on you.

## Lesson 8

# ENTITY

Watch this tutorial on interactive dashboard in Tableau:

https://www.youtube.com/watch?v=Nr31rv9tsJ8

## Lesson 9

Out of the 5 datasets, choose one. Identify columns of interest. Then decide whether they would serve better as independent or dependent variables in an analysis. Finally, use the mind maps provided in the curriculum to choose an appropriate statistical analysis or predictive model. Choose at least 2 statistical analyses or predictive models that can be performed on the dataset.

## Lesson 10

Watch this presentation on Storytelling with Data by Cole Knaflic. It does not have to be all in one sitting!

https://www.youtube.com/watch?v=8EMW7io4rSI