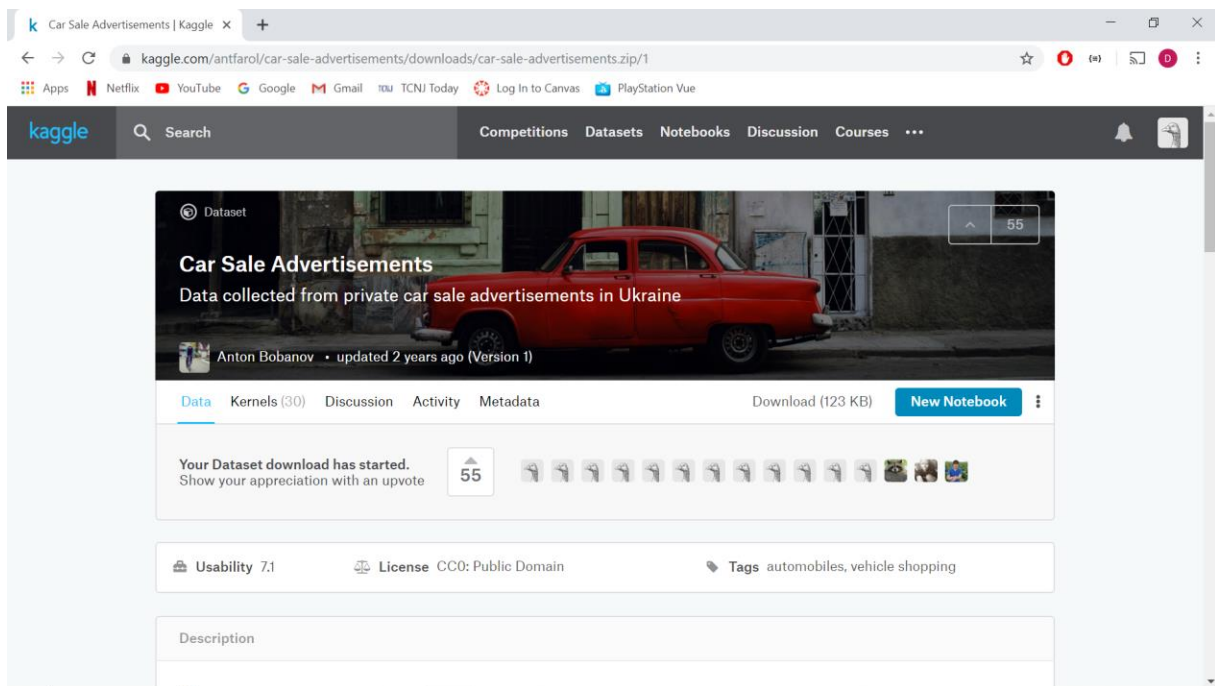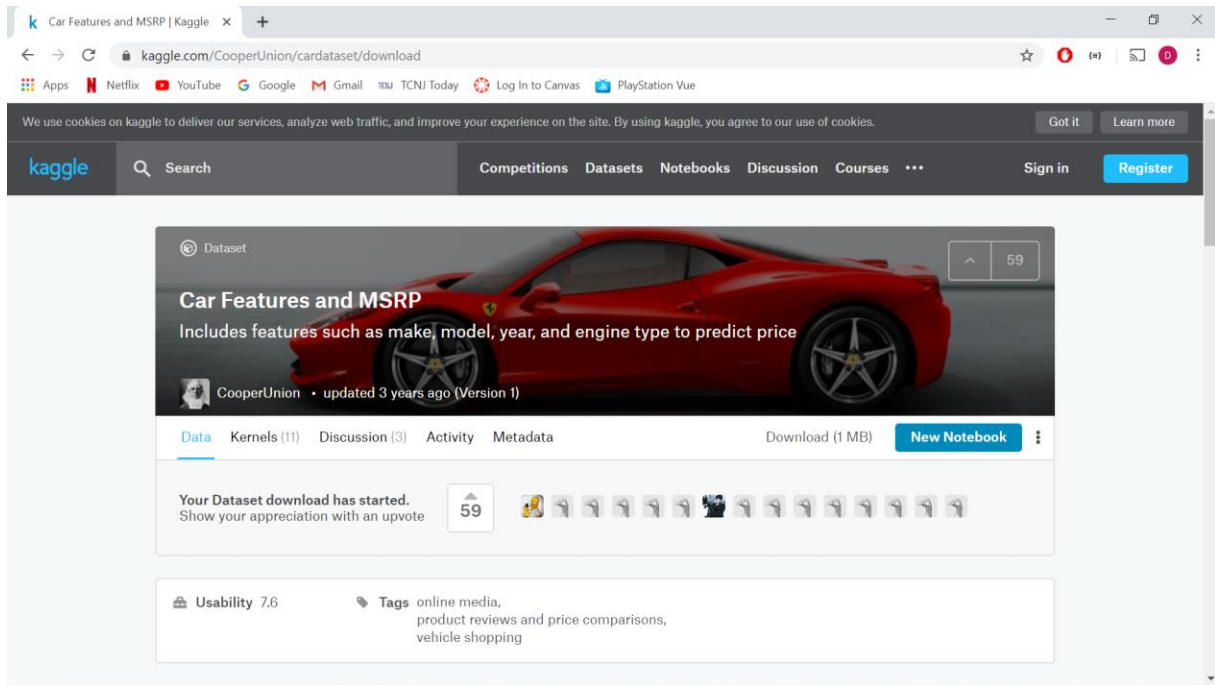**Danyelle Tolud**

**ETL Project Final Report**

**Extract**

The datasets were found on Kaggle:

Downloaded the CSV Files:

(car data and popularity)

data.csv :



(car advertising and sales data)

car_ad.csv :

converted  >>> CSV to JSON

car_ad.json :



## Transform

For the ETL process, the data was extracted from the original data.csv file and a converted csv to json car_ad.json file. Then the data was cleaned, transformed, and joined into one dataframe that included data from both datasets, joined by car make, model, and year to include popularity and price from the two different datasets. This was done in order to analyze car popularity and car price together as they are not both included in either dataset but are in the combined dataframe.

| File | Edit | View | Insert | Cell | Kernel | Widgets | Help | | Trusted | Python 3 ○ |

```
In [1]:   1  import pandas as pd
          2  from sqlalchemy import create_engine
```

### Store CSV into DataFrame

```
In [2]:   1  csv_file = "../data/data.csv"
          2  data_df = pd.read_csv(csv_file)
          3  data_df.head()
```

Out[2]:

| | Make | Model | Year | Engine Fuel Type | Engine HP | Engine Cylinders | Transmission Type | Driven_Wheels | Number of Doors | Market Category | Vehicle Size |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | BMW | 1 Series M | 2011 | premium unleaded (required) | 335.0 | 6.0 | MANUAL | rear wheel drive | 2.0 | Factory Tuner,Luxury,High-Performance | Compact |
| 1 | BMW | 1 Series | 2011 | premium unleaded (required) | 300.0 | 6.0 | MANUAL | rear wheel drive | 2.0 | Luxury,Performance | Compact |
| 2 | BMW | 1 Series | 2011 | premium unleaded (required) | 300.0 | 6.0 | MANUAL | rear wheel drive | 2.0 | Luxury,High-Performance | Compact |

| File | Edit | View | Insert | Cell | Kernel | Widgets | Help | | Trusted | Python 3 ○ |

### Create new data with select columns

```
In [3]:   1  new_data_df = data_df[['Make', 'Model', 'Year', 'Popularity' ]].copy()
          2  new_data_df.head()
```

Out[3]:

| | Make | Model | Year | Popularity |
|---|---|---|---|---|
| 0 | BMW | 1 Series M | 2011 | 3916 |
| 1 | BMW | 1 Series | 2011 | 3916 |
| 2 | BMW | 1 Series | 2011 | 3916 |
| 3 | BMW | 1 Series | 2011 | 3916 |
| 4 | BMW | 1 Series | 2011 | 3916 |

### Store JSON data into a DataFrame

```
In [4]:   1  json_file = "../data/car_ad.json"
          2  car_ad_df = pd.read_json(json_file)
          3  car_ad_df.head()
```

Out[4]:

| | body | car | drive | engType | engV | mileage | model | price | registration | year |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | crossover | Ford | full | Gas | 2.5 | 68 | Kuga | 15500.0 | yes | 2010 |
| 1 | sedan | Mercedes-Benz | rear | Gas | 1.8 | 173 | E-Class | 20500.0 | yes | 2011 |
| 2 | other | Mercedes-Benz | rear | Petrol | 5.5 | 135 | CL 550 | 35000.0 | yes | 2008 |
| 3 | van | Mercedes-Benz | front | Diesel | 1.8 | 162 | B 180 | 17800.0 | yes | 2012 |

```
In [5]:  1  new_car_ad_df = car_ad_df[["car", "model", "year", "price"]].copy()
         2  new_car_ad_df.head()
```

Out[5]:

|   | car | model | year | price |
|---|-----|-------|------|-------|
| 0 | Ford | Kuga | 2010 | 15500.0 |
| 1 | Mercedes-Benz | E-Class | 2011 | 20500.0 |
| 2 | Mercedes-Benz | CL 550 | 2008 | 35000.0 |
| 3 | Mercedes-Benz | B 180 | 2012 | 17800.0 |
| 4 | Mercedes-Benz | E-Class | 2013 | 33000.0 |

### Connect to local database

```
In [8]:  1  rds_connection_string = "postgres:password@localhost:5432/car_db"
         2  engine = create_engine(f'postgresql://{rds_connection_string}')
```

### Check for tables

```
In [9]:  1  engine.table_names()
```

### Use pandas to load csv converted DataFrame into database

```
In [10]:  1  new_car_ad_df.to_sql(name='car', con=engine, if_exists='append', index=False)
```

### Use pandas to load json converted DataFrame into database

```
In [11]:  1  new_car_ad_df.to_sql(name='model', con=engine, if_exists='append', index=False)
```

### Confirm data has been added by querying the customer_name table

- NOTE: can also check using pgAdmin

```
In [12]:  1  pd.read_sql_query('select * from car', con=engine).head()
```

Out[12]:

|   | car | model | year | price |
|---|-----|-------|------|-------|
| 0 | Ford | Kuga | 2010 | 15500.0 |
| 1 | Mercedes-Benz | E-Class | 2011 | 20500.0 |
| 2 | Mercedes-Benz | CL 550 | 2008 | 35000.0 |
| 3 | Mercedes-Benz | B 180 | 2012 | 17800.0 |
| 4 | Mercedes-Benz | E-Class | 2013 | 33000.0 |

### Confirm data has been added by querying the customer_location table ¶

```
In [13]:  1  pd.read_sql_query('select * from model', con=engine).head()
```
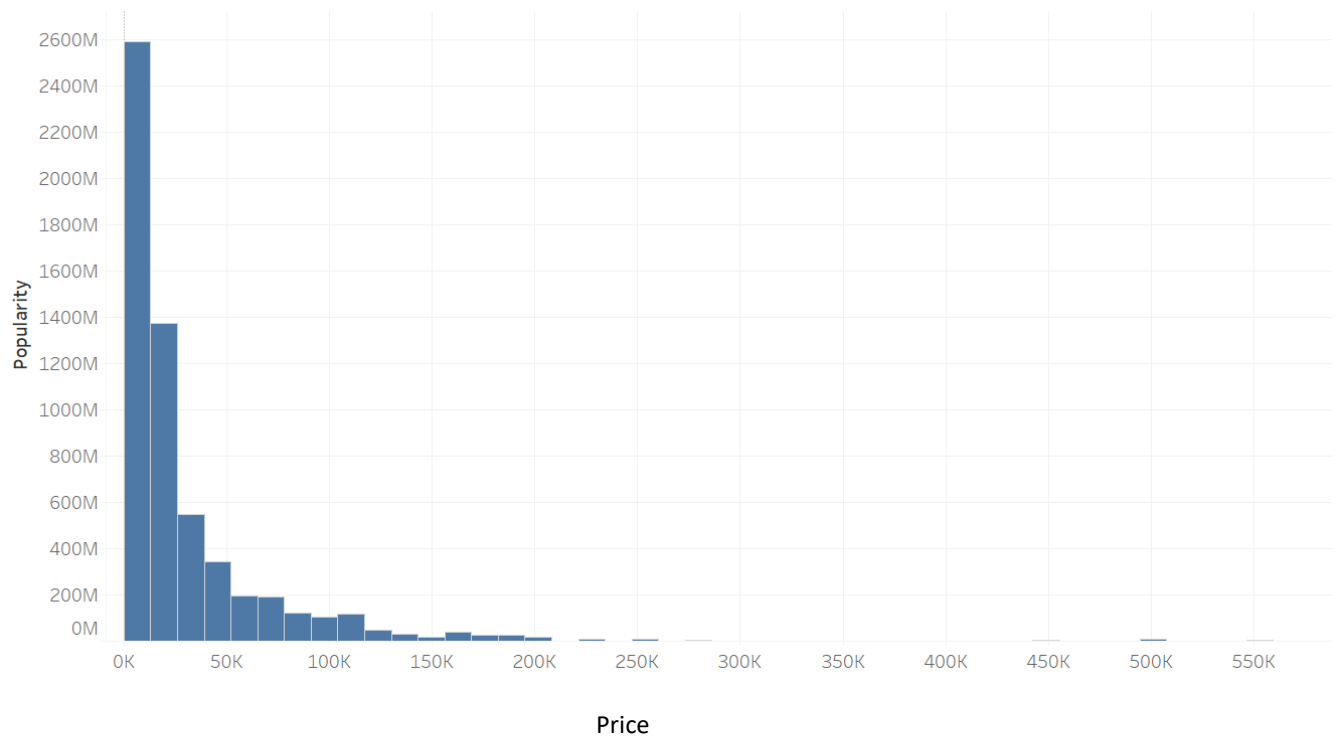
Out[13]:

|   | car | model | year | price |
|---|-----|-------|------|-------|
| 0 | Ford | Kuga | 2010 | 15500.0 |
| 1 | Mercedes-Benz | E-Class | 2011 | 20500.0 |

**Load**

Either dataset included some of the same data on cars, which made it possible to join the data together. From the cars.csv dataset we included car make, noted as simply as car, car model, year and popularity. From the car_ad.json dataset we included car make, model, and year as well, but also price. Because either set included popularity or price but not both, we cleaned the data and loaded them together to analyze both. We can now observe car popularity by price.



From the results, shown in the graph above, we can see a negative correlation—as price increases car popularity decreases. More expensive cars are less popular, and we can see that the most popular cars cost a lot less. This may not necessarily be because cheaper cars are more well-known, perform any better, or have any higher efficiency, etc. but maybe just because they are more affordable. A lot more people can afford cheaper cars, so this causes more expensive cars to be less "popular" and purchased less frequently. However, just because more expensive

cars are not bought as often does not mean they have a poor reputation or need increased advertising. People may be well-aware of an expensive car and want to purchase it, more than any other car, but still choose to purchase a cheaper car because it is within their budget. Advertisers should be aware of this when marketing ads, so instead of increasing advertising, they could make it more efficient. Given that price is strongly correlated to popularity, advertisers need focus on promoting affordability if they wish to increase popularity. Future analysis could include looking at cars sales by different demographics. Because cost is indicative of popularity, we could see how income or location ties into that. Then in terms of marketing, companies directly target the appropriate demographic to have more efficient means of advertising.