# Used Cars: The Data you never knew you needed



BY: REJANE BERINGER, SARA EDGAR AND MICHAEL GOOD

#### **Used Cars market**

 Vehicles are the second highest-priced asset consumers purchase, so it's very important for the economy.

• The used car market in the U.S. is already estimated at 41 million units annually.

#### **Used Cars market**

• Consumers can now find greater inventory of used vehicles online.

• More dealers are accepting the move to digital, which is a major behavioral shift for the market.

• A wider net can help you get a better price.

#### Covid-19 influence

Covid-19 has led to an increase in used car sales

> People avoid mass transportation

> People are more sensitive to auto cost in the recession.

## Why used cars?

- Large datasets
  - ☐ There is a lot of available information
- Many different parameters
  - □ What is important about the cars that are being sold?
- Location
  - ☐ Where is the best place to find a specific type of car
- Potentially confusing field

#### Sources:

Car Sales

https://www.kaggle.com/gagandeep16/car-sales

• US Used Cars Catalog

https://www.kaggle.com/lepchenkov/usedcarscatalog

Used Cars Dataset

https://www.kaggle.com/doaaalsenani/used-cars-dataets

#### ETL

#### **Extract**

• .csv from Kaggle.

#### **Transform**

- We will be using pandas to transform the dataset.
- We will select the important columns from each dataset, remove any irrelevant data or empty rows, and combine the datasets into the relevant tables.

#### Load

- PostgreSQL Relational database, expected tables to generate:
  - ➤ Important subset data from each dataset (x3)
  - ➤ Make & model Key table

#### Extraction

- Kaggle was key
  - List of links
  - Sifting through datasets not as simple as appears
- Data.world was slightly overwhelming

• CSVs reigned supreme

#### Transformation

- First Step:
  - Downloading the csvs (not as simple as it seems)
- Second Step:
  - Cleaning
    - De duplicating
    - Dropping/adding columns
    - Renaming Columns
    - Eliminating empty or N/A values
- Third Step:
  - Saving as new csv files

```
In [1]: import pandas as pd
         from sqlalchemy import create engine
         import os
        import csv
In [2]: #first file
        csv file = "Car sales.csv"
        car_sales1_df = pd.read_csv(csv_file, encoding="utf-8")
        car sales1 df.head()
Out[2]:
           Manufacturer | Model | Sales_in_thousands | _year_resale_value | Vehicle_type | Price_in_thousands | Engine_size | Horsepower | Wheelbase | Width | Length
         0 Acura
                        Integra 16.919
                                                  16.360
                                                                     Passenger
                                                                                  21.50
                                                                                                     1.8
                                                                                                                 140.0
                                                                                                                             101.2
                                                                                                                                       67.3
                                                                                                                                              172.4
         1 Acura
                        TL
                               39.384
                                                  19.875
                                                                     Passenger
                                                                                  28.40
                                                                                                    3.2
                                                                                                                 225.0
                                                                                                                                       70.3
                                                                                                                                              192.9
                                                                                                                             108.1
         2 Acura
                        CL
                               14.114
                                                  18.225
                                                                     Passenger
                                                                                  NaN
                                                                                                    3.2
                                                                                                                 225.0
                                                                                                                             106.9
                                                                                                                                       70.6
                                                                                                                                              192.0
         3 Acura
                        RL
                              8.588
                                                  29.725
                                                                                  42.00
                                                                                                    3.5
                                                                                                                210.0
                                                                                                                             114.6
                                                                                                                                       71.4
                                                                                                                                              196.6
                                                                     Passenger
         4 Audi
                        A4
                              20.397
                                                  22.255
                                                                                  23.99
                                                                                                    1.8
                                                                                                                 150.0
                                                                                                                             102.6
                                                                                                                                       68.2
                                                                                                                                              178.0
                                                                      Passenger
```

Out[3]:

	Manufacturer	Model	Vehicle_type	Engine_size	Horsepower	Wheelbase	Width	Length	Curb_weight	Fuel_capacity	Fuel_efficiency
0	Acura	Integra	Passenger	1.8	140.0	101.2	67.3	172.4	2.639	13.2	28.0
1	Acura	TL	Passenger	3.2	225.0	108.1	70.3	192.9	3.517	17.2	25.0
2	Acura	CL	Passenger	3.2	225.0	106.9	70.6	192.0	3.470	17.2	26.0
3	Acura	RL	Passenger	3.5	210.0	114.6	71.4	196.6	3.850	18.0	22.0
4	Audi	A4	Passenger	1.8	150.0	102.6	68.2	178.0	2.998	16.4	27.0
										***	***
152	Volvo	V40	Passenger	1.9	160.0	100.5	67.6	176.6	3.042	15.8	25.0
153	Volvo	S70	Passenger	2.4	168.0	104.9	69.3	185.9	3.208	17.9	25.0
154	Volvo	V70	Passenger	2.4	168.0	104.9	69.3	186.2	3.259	17.9	25.0
155	Volvo	C70	Passenger	2.3	236.0	104.9	71.5	185.7	3.601	18.5	23.0
156	Volvo	S80	Passenger	2.9	201.0	109.9	72.1	189.8	3.600	21.1	24.0

157 rows × 11 columns

4

In [5]: new\_car\_sales1\_df.insert(0,'id', range(1, 1+ len(new\_car\_sales1\_df)), True)

In [31]: new\_car\_sales1\_df

Out[31]:

	id	Manufacturer	Model	Vehicle_Type	Engine_Size	Horsepower	Wheelbase	Width	Length	Curb_Weight	Fuel_Capacity	Fuel_Efficiency
0	1	Acura	Integra	Passenger	1.8	140.0	101.2	67.3	172.4	2.639	13.2	28.0
1	2	Acura	TL	Passenger	3.2	225.0	108.1	70.3	192.9	3.517	17.2	25.0
2	3	Acura	CL	Passenger	3.2	225.0	106.9	70.6	192.0	3.470	17.2	26.0
3	4	Acura	RL	Passenger	3.5	210.0	114.6	71.4	196.6	3.850	18.0	22.0
4	5	Audi	A4	Passenger	1.8	150.0	102.6	68.2	178.0	2.998	16.4	27.0
										god.	***	
152	153	Volvo	V40	Passenger	1.9	160.0	100.5	67.6	176.6	3.042	15.8	25.0
153	154	Volvo	S70	Passenger	2.4	168.0	104.9	69.3	185.9	3.208	17.9	25.0
154	155	Volvo	V70	Passenger	2.4	168.0	104.9	69.3	186.2	3.259	17.9	25.0
155	156	Volvo	C70	Passenger	2.3	236.0	104.9	71.5	185.7	3.601	18.5	23.0
156	157	Volvo	S80	Passenger	2.9	201.0	109.9	72.1	189.8	3.600	21.1	24.0

157 rows × 12 columns

4

In [7]: new\_car\_sales1\_df.to\_csv("new\_car\_sales1")

```
In [8]: #second file
         cars df = pd.read csv('cars.csv')
In [9]: cars df.head()
Out[9]:
           manufacturer_name model_name transmission color odometer_value year_produced engine_fuel engine_has_gas engine_type engine_capacity ... |
                                                         silver 190000
         0 Subaru
                               Outback
                                            automatic
                                                                               2010
                                                                                              gasoline
                                                                                                          False
                                                                                                                          gasoline
                                                                                                                                      2.5
          1 Subaru
                                Outback
                                                         blue
                                                               290000
                                                                               2002
                                                                                                         False
                                                                                                                                      3.0
                                            automatic
                                                                                              gasoline
                                                                                                                          gasoline
          2 Subaru
                                                         red
                                                               402000
                                                                               2001
                                                                                                          False
                                                                                                                                      2.5
                                Forester
                                                                                                                          gasoline
                                            automatic
                                                                                              gasoline
                                                         blue
          3 Subaru
                                                               10000
                                                                               1999
                                                                                                          False
                                                                                                                                      3.0
                                Impreza
                                            mechanical
                                                                                              gasoline
                                                                                                                          gasoline
         4 Subaru
                                                         black 280000
                                                                               2001
                                                                                                          False
                                                                                                                                      2.5
                                Legacy
                                            automatic
                                                                                              gasoline
                                                                                                                          gasoline
         5 rows × 30 columns
```

```
In [10]: count = 0
          wanted columns = []
          for val in cars df.dtypes:
              if val != bool:
                  wanted columns.append(cars df.columns[count])
              count += 1
          wanted columns
Out[10]: ['manufacturer name',
           'model name',
           'transmission',
           'color',
           'odometer value',
           'year produced',
           'engine fuel',
           'engine type',
           'engine capacity',
           'body_type',
           'state',
           'drivetrain',
           'price usd',
           'location region',
           'number of photos',
           'up counter',
           'duration listed']
In [11]: cleaned cars = cars df[wanted columns]
In [12]: cleaned_cars = cleaned_cars.drop('number_of_photos',axis = 1)
          cleaned cars = cleaned cars.drop('engine fuel',axis = 1)
         cleaned_cars = cleaned_cars.drop('location_region',axis = 1)
```

```
In [13]:
    cleaned_cars = cleaned_cars.rename( columns = {
        'manufacturer_name': 'Manufacturer',
        'model_name': 'Model',
        'transmission'; 'Transmission',
        'color': 'Color',
        'odometer_value': 'Odometer',
        'year_produced': 'Year',
        'engine_type': 'Engine_Type',
        'engine_capacity': 'Engine_Capacity',
        'body_type': 'Body',
        'state': 'State',
        'drivetrain': 'Drivetrain',
        'price_usd': 'Price',
        'up_counter': 'Up_Counter',
        'duration_listed': 'List_Duration'
}

In [14]: cleaned_cars = cleaned_cars.drop('Up_Counter',axis = 1)

In [15]: cleaned_cars.insert(0,'id', range(1, 1+ len(cleaned_cars)), True)
```

#### In [16]: cleaned\_cars

Out[16]:

	id	Manufacturer	Model	Transmission	Color	Odometer	Year	Engine_Type	Engine_Capacity	Body	State	Drivetrain	Price	List_Dura
0	1	Subaru	Outback	automatic	silver	190000	2010	gasoline	2.5	universal	owned	all	10900.00	16
1	2	Subaru	Outback	automatic	blue	290000	2002	gasoline	3.0	universal	owned	all	5000.00	83
2	3	Subaru	Forester	automatic	red	402000	2001	gasoline	2.5	suv	owned	all	2800.00	151
3	4	Subaru	Impreza	mechanical	blue	10000	1999	gasoline	3.0	sedan	owned	all	9999.00	86
4	5	Subaru	Legacy	automatic	black	280000	2001	gasoline	2.5	universal	owned	all	2134.11	7
										•••			•••	
38526	38527	Chrysler	300	automatic	silver	290000	2000	gasoline	3.5	sedan	owned	front	2750.00	301
38527	38528	Chrysler	PT Cruiser	mechanical	blue	321000	2004	diesel	2.2	hatchback	owned	front	4800.00	317
38528	38529	Chrysler	300	automatic	blue	777957	2000	gasoline	3.5	sedan	owned	front	4300.00	369
38529	38530	Chrysler	PT Cruiser	mechanical	black	20000	2001	gasoline	2.0	minivan	owned	front	4000.00	490
38530	38531	Chrysler	Voyager	automatic	silver	297729	2000	gasoline	2.4	minivan	owned	front	3200.00	632

#### 38531 rows × 14 columns

In [17]: cleaned\_cars.to\_csv('cleaned\_cars.csv')

```
In [18]: #file 3
         csv file = "Resources/final cars datasets.csv"
         car sales2 df = pd.read csv(csv file, encoding="utf-8")
         car sales2 df.head()
Out[18]:
            Unnamed: 0 price mark
                                      model | year | mileage | engine_capacity | transmission | drive | hand_drive | fuel
                                              2003 80000
                                                                                              rhd
          0 0
                         80
                              nissan
                                       march
                                                            1240
                                                                                         2wd
                                                                            at
                                                                                                          gasoline
                                              2010 53000
                                                                                              rhd
                                                            1200
                                                                            at
                                                                                         2wd
                         110
                               nissan
                                       march
                                                                                                          gasoline
                         165
                                       lafesta
                                              2005 47690
                                                            2000
                                                                            at
                                                                                         2wd
                                                                                              rhd
                               nissan
                                                                                                          gasoline
                                      avensis 2008 130661
                                                            1990
          3 3
                         190
                              toyota
                                                                            at
                                                                                         2wd
                                                                                              rhd
                                                                                                          gasoline
                                                                                             rhd
                         190
                              daihatsu mira
                                              2006 66300
                                                            660
                                                                                         2wd
                                                                            at
                                                                                                          gasoline
In [19]: car_sales2_df= car_sales2_df.drop('Unnamed: 0', axis=1)
In [20]: car sales2 df.insert(0,'id', range(1, 1+ len(car sales2 df)), True)
In [21]: car_sales2_df
Out[21]:
                                                        engine_capacity transmission drive hand_drive fuel
                id
                     price mark
                                   model
                                          year mileage
                                          2003 80000
                                                                                     2wd rhd
          0
                     80
                                   march
                                                         1240
                                                                                                      gasoline
                           nissan
               2
                                          2010 53000
                                                         1200
                                                                                          rhd
                     110
                                                                                     2wd
                           nissan
                                   march
                                                                                                      gasoline
                                          2005 47690
                                                                                          rhd
                3
                     165
                           nissan
                                   lafesta
                                                        2000
                                                                        at
                                                                                     2wd
                                                                                                      gasoline
                                                                                          rhd
                     190
                           toyota
                                          2008 130661
                                                         1990
                                                                                     2wd
                4
                                   avensis
                                                                                                      gasoline
                5
                     190
                          daihatsu
                                   mira
                                           2006 66300
                                                                                     2wd rhd
                                                                                                      gasoline
                                                        996
          2313 2314 1400
                                           2009 121000
                          toyota
                                   vitz
                                                                                     2wd
                                                                                          rhd
                                                                                                      gasoline
          2314 2315 1400
                                          2003 101000
                                                        3000
                                                                                          rhd
                          toyota
                                   estima
                                                                        at
                                                                                     2wd
                                                                                                      gasoline
          2315 2316 1400 subaru
                                   r2
                                           2005 101000
                                                        660
                                                                        cvt
                                                                                          rhd
                                                                                     2wd
                                                                                                      gasoline
          2316 2317 1400 honda
                                           2000 170000
                                                                                         rhd
                                                        660
                                                                        at
                                                                                     4wd
                                                                                                      gasoline
         2317 2318 1400 toyota
                                   estima t 2005 72320
                                                        3000
                                                                                          rhd
                                                                        at
                                                                                     2wd
                                                                                                      gasoline
         2318 rows × 11 columns
```

In [22]: car sales2 df.to csv('car sales2 df.csv')

In [26]: pd.read\_sql\_query('select \* from cars', con=engine).head()

Out[26]:

	id	Manufacturer	Model	Transmission	Color	Odometer	Year	Engine_Type	Engine_Capacity	Body	State	Drivetrain	Price	List_Duration
(	1	Subaru	Outback	automatic	silver	190000	2010	gasoline	2.5	universal	owned	all	10900.00	16
•	2	Subaru	Outback	automatic	blue	290000	2002	gasoline	3.0	universal	owned	all	5000.00	83
2	2 3	Subaru	Forester	automatic	red	402000	2001	gasoline	2.5	suv	owned	all	2800.00	151
;	4	Subaru	Impreza	mechanical	blue	10000	1999	gasoline	3.0	sedan	owned	all	9999.00	86
4	1 5	Subaru	Legacy	automatic	black	280000	2001	gasoline	2.5	universal	owned	all	2134.11	7

In [27]: new\_car\_sales1\_df.to\_sql(name='vehicles', con=engine, if\_exists='append', index=False)

In [35]: pd.read\_sql\_query('select \* from vehicles', con=engine).head()

Out[35]:

	id	Manufacturer	Model	Vehicle_Type	Engine_Size	Horsepower	Wheelbase	Width	Length	Curb_Weight	Fuel_Capacity	Fuel_Efficiency
0	1	Acura	Integra	Passenger	1.8	140.0	101.2	67.3	172.4	2.639	13.2	28.0
1	2	Acura	TL	Passenger	3.2	225.0	108.1	70.3	192.9	3.517	17.2	25.0
2	3	Acura	CL	Passenger	3.2	225.0	106.9	70.6	192.0	3.470	17.2	26.0
3	4	Acura	RL	Passenger	3.5	210.0	114.6	71.4	196.6	3.850	18.0	22.0
4	5	Audi	A4	Passenger	1.8	150.0	102.6	68.2	178.0	2.998	16.4	27.0

.

```
In [29]: car_sales2_df.to_sql(name='sara_cars', con=engine, if_exists='append', index=False)
In [30]: pd.read_sql_query('select * from sara_cars', con=engine).head()
Out[30]:
            id price mark
                            model year mileage engine_capacity transmission drive hand_drive fuel
         0 1 80
                                    2003 80000
                                                                             2wd
                                                                                 rhd
                    nissan
                            march
                                                 1240
                                                                at
                                                                                             gasoline
                                    2010 53000
         1 2 110
                                                 1200
                                                                             2wd rhd
                    nissan
                            march
                                                                at
                                                                                             gasoline
         2 3
               165
                    nissan
                            lafesta
                                    2005 47690
                                                 2000
                                                                at
                                                                             2wd rhd
                                                                                             gasoline
         3 4
               190
                            avensis 2008 130661
                                                                                 rhd
                                                 1990
                                                                at
                                                                            2wd
                    toyota
                                                                                             gasoline
         4 5
                    daihatsu mira
                                    2006 66300
                                                660
               190
                                                                             2wd rhd
                                                                                             gasoline
                                                                at
```

#### Load

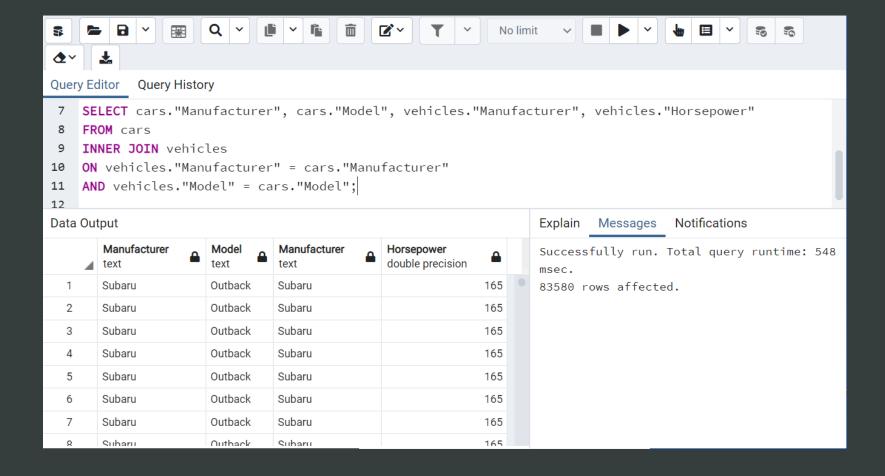
• Creating cars database

• Writing schema for the tables

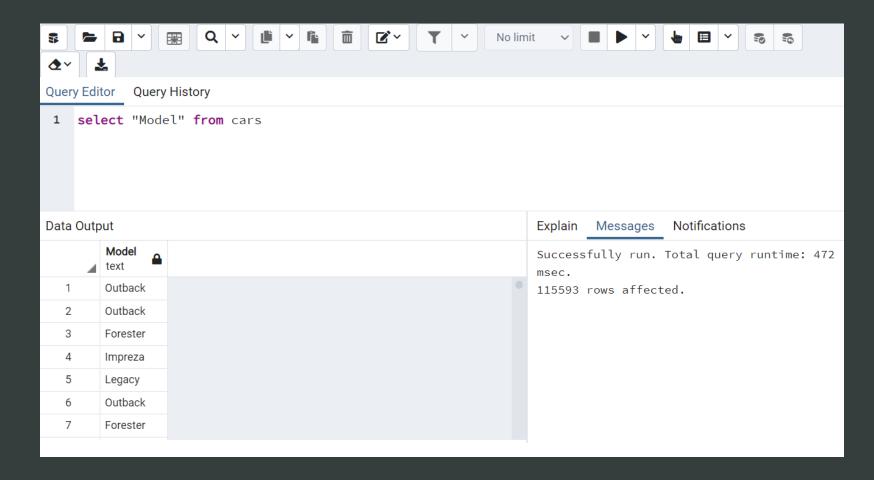
Connecting pandas dataframes to SQL

- Creating queries for new tables
  - Using inner join to create new tables

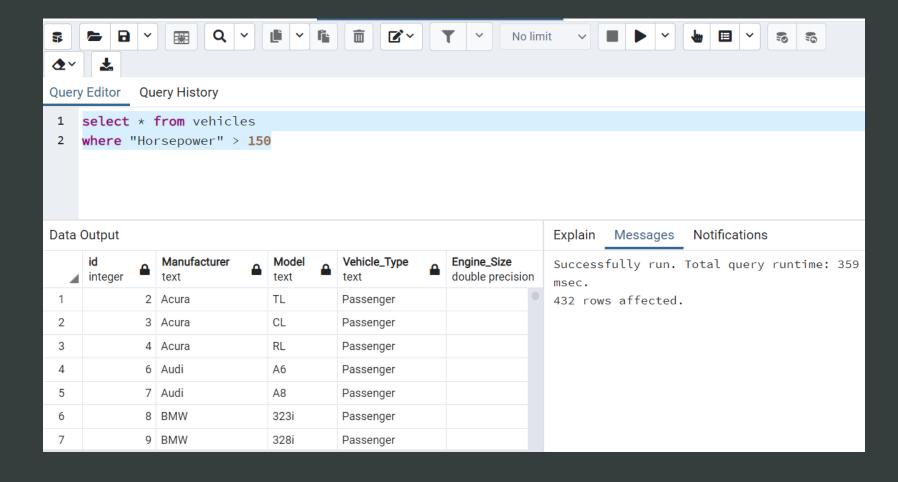
#### SQL



### SQL



#### SQL



#### **Issues Encountered**

Connecting Data Frames to SQL

 Discovering that inserting the "tables" through pandas and creating the schema from scratch in SQL would cause MAJOR issues

Combining Tables

 Needing to use "" to emphasize uppercase letters when identifying the columns in the tables within SQL

# Questions

