# Car MSRP Prediction with Multiple Linear Regression

CIS 4321

Beshoy Sadek, Sherleen Lee, Xiao Tong Gan, and Somar Dakak

May 14, 2021

# Introduction

The problem for the manufacturers is to find a balance between what consumers want and pricing the particular vehicle
- Important to Know
  - which features and specs make the Manufacturer Suggested Retail Price
  - which features largely determine the price of the car
  - produce and sell cars that meet market demands for its target consumer group

The purpose of our group project is to figure out what features have an influence on a vehicle's Manufacturer Suggested Retail Price (MSRP) in the Car Features and MSRP dataset
- will be able to
  - know what features are most important to consumers by predicting an appropriate MSRP given a specific set of features
  - understand how those features affect the MSRP the most

# Data Collection and Variable Description

Collecting Data:
- found our dataset on Kaggle
  - could not find the real source of the set
  - all the records were collected from Twitter and Edmunds
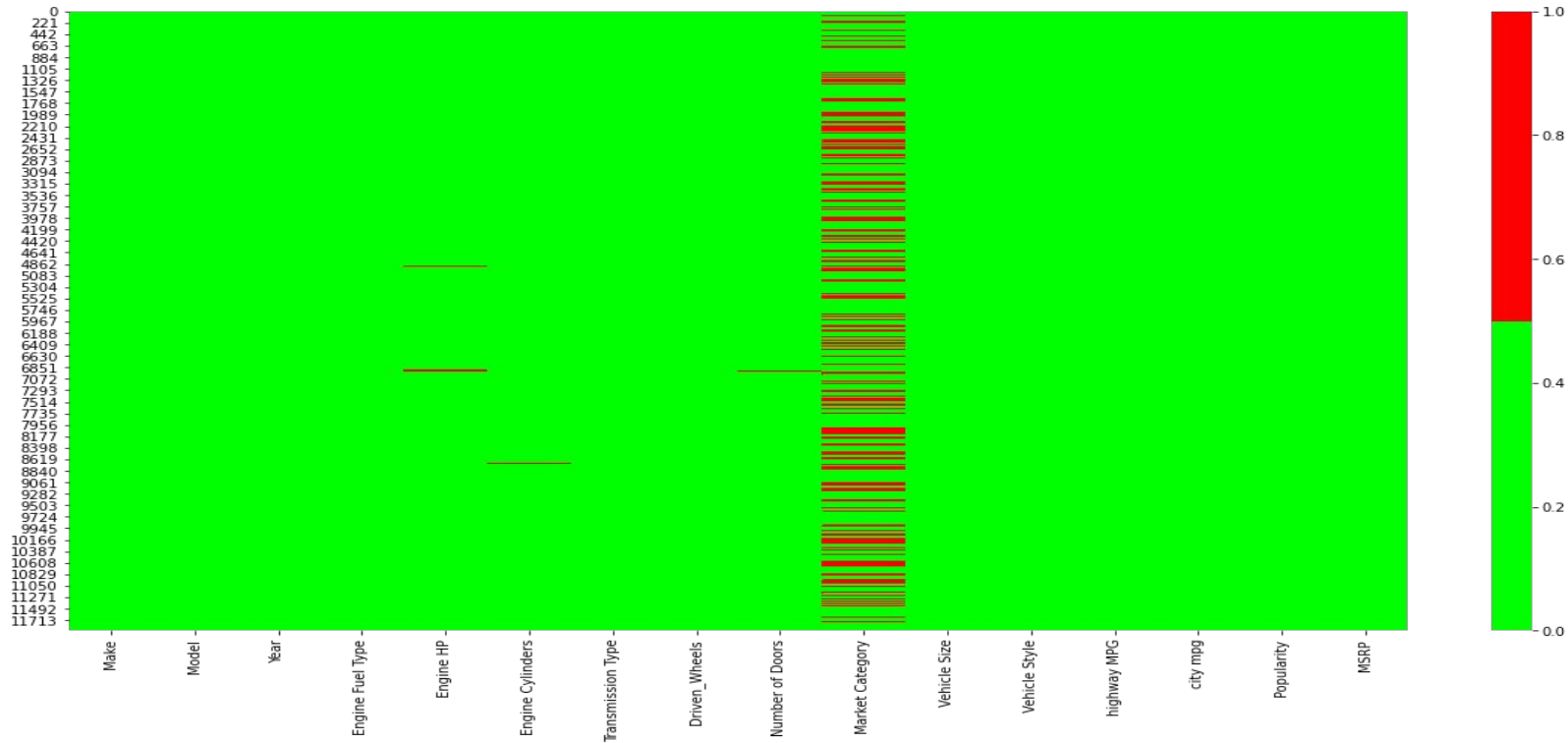
Dataset has 16 variables which are divided to 8 numerical and 8 categorical variables
- numerical variables: Year, Engine HP, Engine Cylinders, Number of Doors, highway mpg, city mpg, Popularity, and MSRP
- categorical variables: Make, Model, Engine Fuel Type, Transmission type, Driven Wheels, Market Category, Vehicle Size, and Vehicle Style
  - categorical variables such as Transmission type, Vehicle Size, and Vehicle Style, are transformed to dummy variables
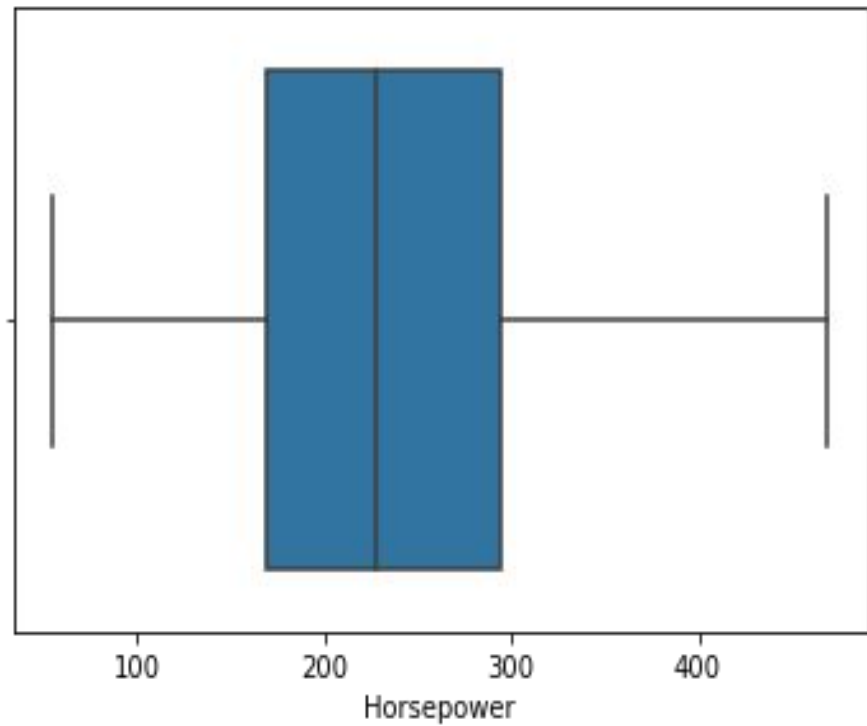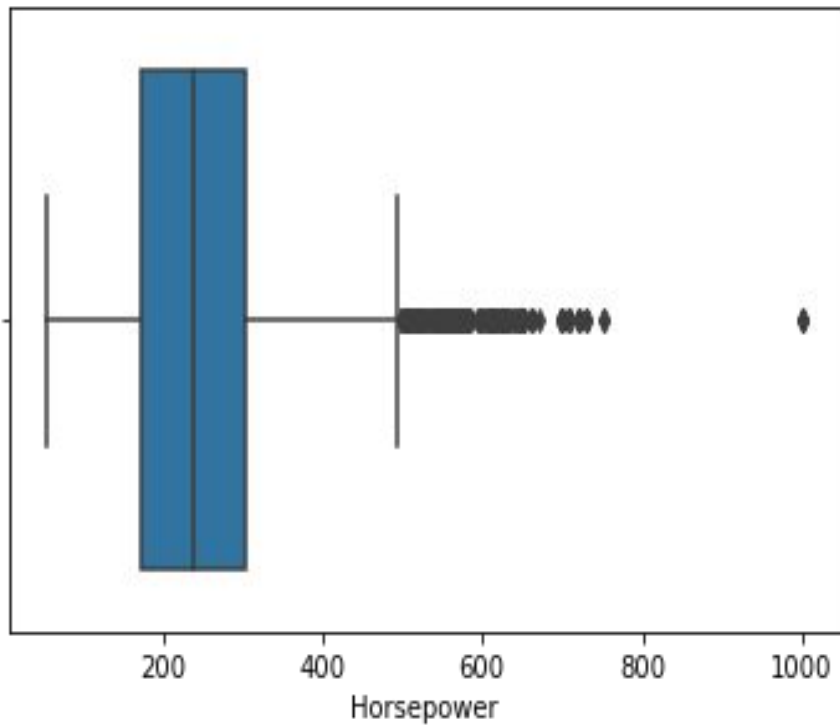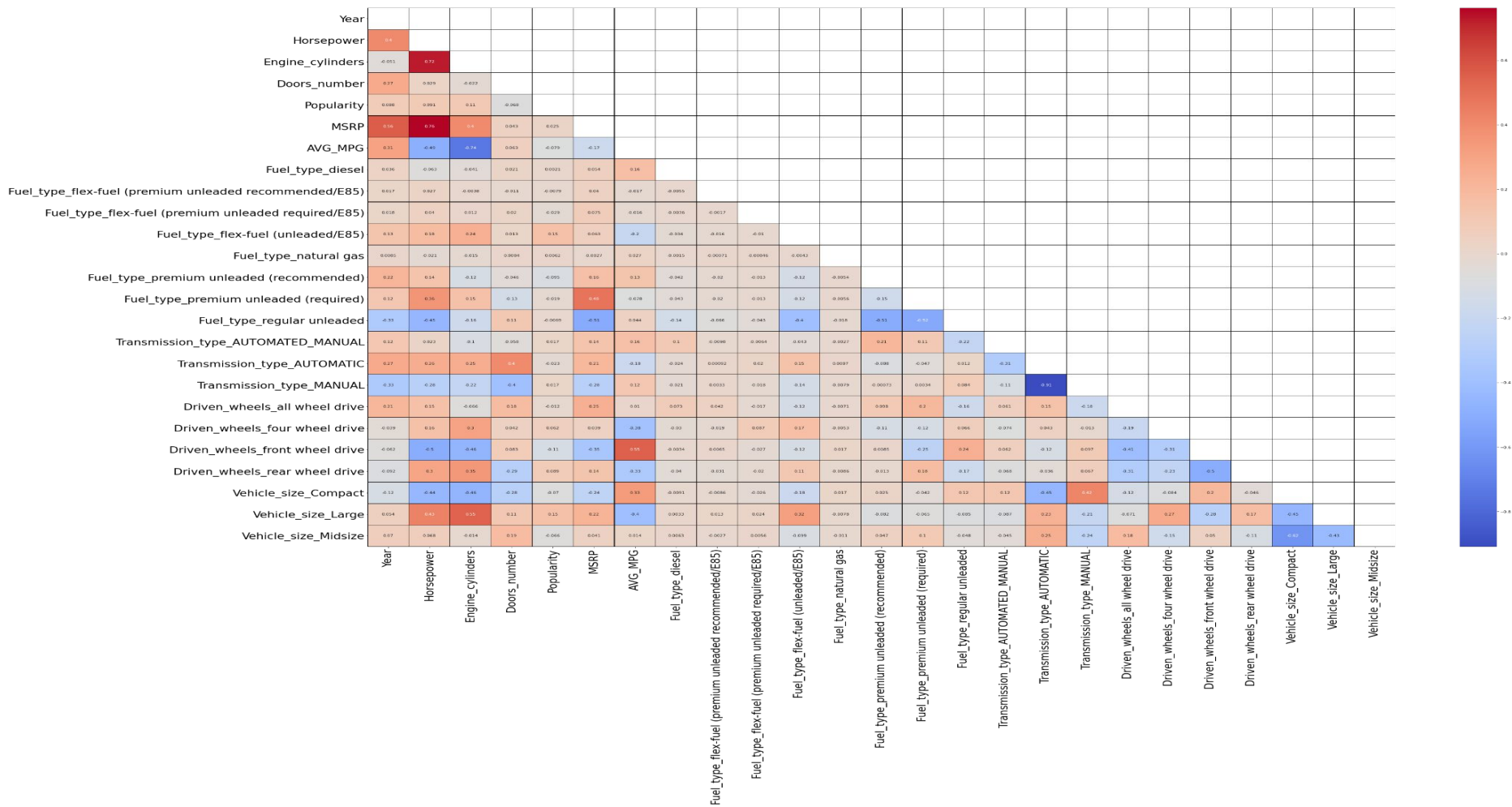
# Data Collection and Variable Description

- Make : the brand of the vehicle
- Model : the model of the vehicle
- Year : the year of manufacture of the vehicle
- Engine fuel type : the type of you that the vehicle operates on
- Engine HP : the engine horsepower
- Engine cylinders :the number of cylinders in the combustion engine
- Transmission type: what kind of transmission the vehicle has
- Driven_Wheels : the amount of wheels that receive power directly from the engine
- Number of doors : how many doors the vehicle has

- Market category : the market category to which the vehicle belongs
- Vehicle size : the size of the vehicle based on its volume capacity
- Vehicle style : the style of the vehicle based on its physical shape
- Highway MPG : the number of miles the vehicle drives per gallon at highway speed
- City MPG : the number of miles the vehicle drives per gallon at a city speed
- Popularity : how popular the car is among people
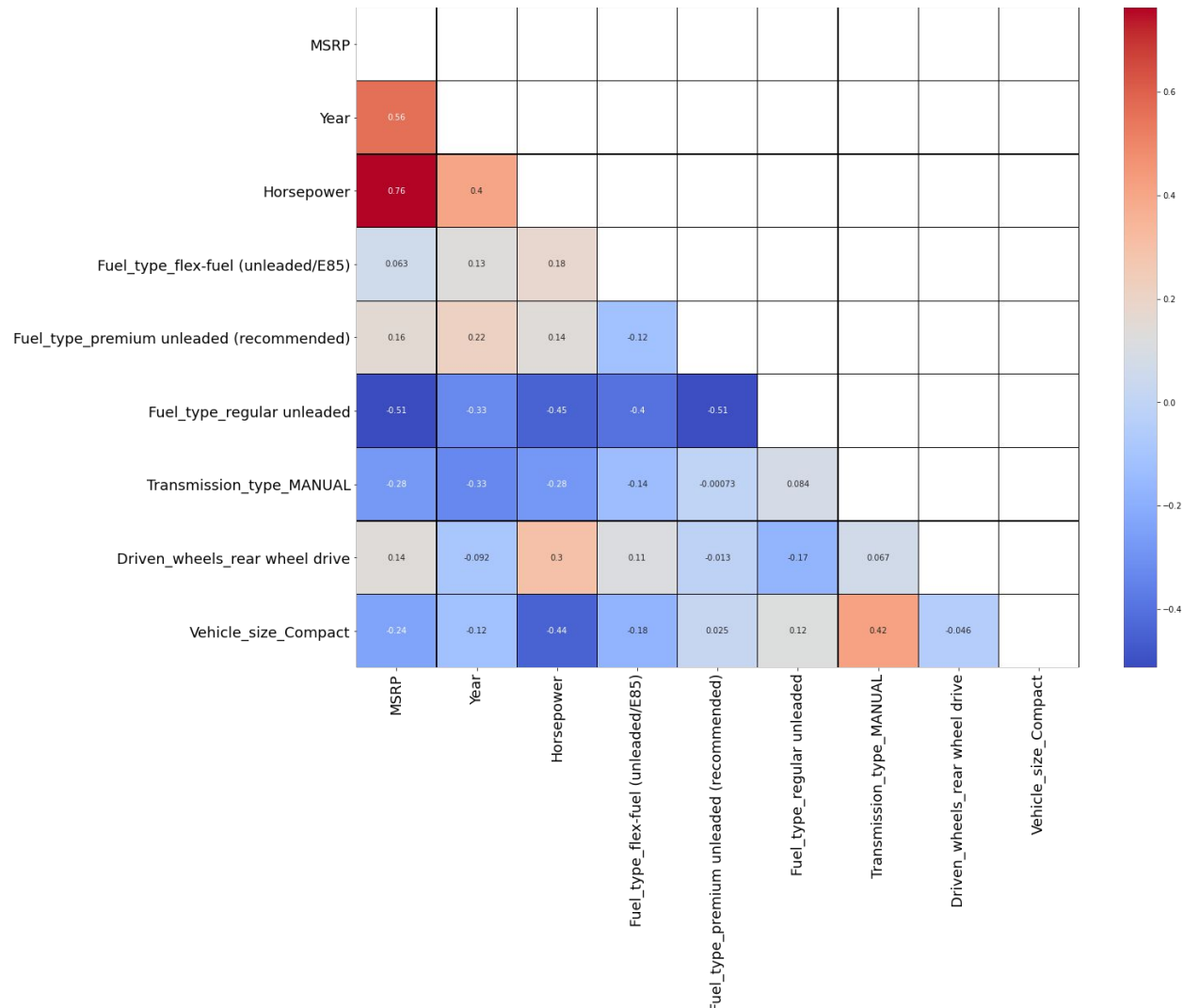- MSRP :the manufacturer's suggested retail price of a vehicle

# Descriptive Analysis

**Descriptive Analysis**

Correlation heatmap (lower-triangular matrix). Columns are numbered to match the row labels below.

| Variable | Year | Horsepower | Engine_cylinders | Doors_number | Popularity | MSRP | AVG_MPG | Fuel_type_diesel | Fuel_type_flex-fuel (premium unleaded recommended/E85) | Fuel_type_flex-fuel (premium unleaded required/E85) | Fuel_type_flex-fuel (unleaded/E85) | Fuel_type_natural gas | Fuel_type_premium unleaded (recommended) | Fuel_type_premium unleaded (required) | Fuel_type_regular unleaded | Transmission_type_AUTOMATED_MANUAL | Transmission_type_AUTOMATIC | Transmission_type_MANUAL | Driven_wheels_all wheel drive | Driven_wheels_four wheel drive | Driven_wheels_front wheel drive | Driven_wheels_rear wheel drive | Vehicle_size_Compact | Vehicle_size_Large | Vehicle_size_Midsize |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Year | | | | | | | | | | | | | | | | | | | | | | | | | |
| Horsepower | 0.4 | | | | | | | | | | | | | | | | | | | | | | | | |
| Engine_cylinders | -0.051 | 0.72 | | | | | | | | | | | | | | | | | | | | | | | |
| Doors_number | 0.27 | 0.029 | -0.022 | | | | | | | | | | | | | | | | | | | | | | |
| Popularity | 0.088 | 0.091 | 0.11 | -0.068 | | | | | | | | | | | | | | | | | | | | | |
| MSRP | 0.55 | 0.76 | 0.6 | 0.043 | 0.025 | | | | | | | | | | | | | | | | | | | | |
| AVG_MPG | 0.31 | -0.45 | -0.74 | 0.063 | -0.079 | -0.17 | | | | | | | | | | | | | | | | | | | |
| Fuel_type_diesel | 0.036 | -0.063 | -0.041 | 0.021 | 0.0021 | 0.054 | 0.16 | | | | | | | | | | | | | | | | | | |
| Fuel_type_flex-fuel (premium unleaded recommended/E85) | 0.017 | 0.027 | -0.0038 | -0.013 | -0.0079 | 0.04 | -0.017 | -0.0055 | | | | | | | | | | | | | | | | | |
| Fuel_type_flex-fuel (premium unleaded required/E85) | 0.018 | 0.04 | 0.012 | 0.02 | -0.029 | 0.075 | -0.016 | 0.0036 | -0.0017 | | | | | | | | | | | | | | | | |
| Fuel_type_flex-fuel (unleaded/E85) | 0.14 | 0.18 | 0.24 | 0.015 | 0.15 | 0.063 | -0.2 | -0.036 | -0.016 | -0.01 | | | | | | | | | | | | | | | |
| Fuel_type_natural gas | 0.0081 | -0.021 | -0.015 | 0.0094 | 0.0062 | -0.0027 | 0.027 | -0.0015 | -0.0071 | -0.00046 | -0.0042 | | | | | | | | | | | | | | |
| Fuel_type_premium unleaded (recommended) | 0.22 | 0.14 | -0.12 | -0.046 | -0.095 | 0.16 | 0.13 | -0.042 | -0.02 | -0.013 | -0.12 | -0.0054 | | | | | | | | | | | | | |
| Fuel_type_premium unleaded (required) | 0.12 | 0.36 | 0.15 | 0.13 | 0.019 | 0.49 | -0.078 | -0.043 | 0.02 | -0.013 | -0.12 | -0.0056 | 0.15 | | | | | | | | | | | | |
| Fuel_type_regular unleaded | -0.33 | -0.45 | -0.16 | 0.13 | -0.0065 | -0.51 | 0.044 | -0.14 | -0.066 | -0.045 | -0.4 | -0.018 | -0.51 | -0.52 | | | | | | | | | | | |
| Transmission_type_AUTOMATED_MANUAL | 0.12 | 0.022 | -0.1 | 0.058 | 0.017 | 0.14 | 0.16 | 0.1 | 0.0098 | 0.0054 | 0.043 | 0.0027 | 0.21 | 0.11 | 0.22 | | | | | | | | | | |
| Transmission_type_AUTOMATIC | 0.17 | 0.26 | 0.25 | 0.4 | -0.023 | 0.21 | -0.19 | -0.024 | 0.00002 | 0.02 | 0.15 | 0.0097 | -0.086 | -0.047 | 0.012 | -0.31 | | | | | | | | | |
| Transmission_type_MANUAL | -0.33 | -0.28 | -0.22 | -0.4 | 0.017 | -0.28 | 0.12 | 0.021 | 0.0033 | -0.018 | -0.14 | -0.0079 | -0.00073 | 0.0034 | 0.084 | -0.11 | -0.91 | | | | | | | | |
| Driven_wheels_all wheel drive | 0.21 | 0.15 | -0.066 | 0.18 | -0.013 | 0.25 | 0.01 | 0.073 | 0.042 | -0.017 | -0.13 | -0.0071 | 0.093 | 0.3 | 0.061 | 0.15 | -0.18 | -0.19 | | | | | | | |
| Driven_wheels_four wheel drive | -0.039 | 0.16 | 0.3 | 0.042 | 0.062 | 0.039 | -0.28 | -0.03 | -0.019 | 0.17 | -0.0053 | -0.11 | -0.12 | 0.079 | 0.093 | -0.013 | -0.19 | | -0.19 | | | | | | |
| Driven_wheels_front wheel drive | -0.062 | -0.5 | -0.46 | 0.085 | -0.11 | -0.35 | 0.55 | -0.0054 | 0.0065 | -0.027 | -0.12 | 0.017 | 0.0085 | -0.25 | 0.24 | 0.097 | -0.12 | -0.41 | -0.31 | | | | | | |
| Driven_wheels_rear wheel drive | 0.092 | 0.3 | 0.35 | -0.29 | 0.096 | 0.14 | -0.33 | 0.04 | 0.031 | -0.02 | 0.11 | 0.0086 | 0.013 | 0.18 | 0.17 | 0.066 | -0.31 | -0.23 | 0.5 | | | | | | |
| Vehicle_size_Compact | -0.13 | -0.44 | -0.46 | -0.28 | -0.07 | -0.24 | 0.33 | -0.0051 | -0.0086 | -0.026 | -0.18 | 0.017 | 0.025 | -0.042 | 0.12 | 0.12 | -0.45 | 0.43 | -0.13 | -0.084 | 0.2 | -0.046 | | | |
| Vehicle_size_Large | 0.054 | 0.42 | 0.55 | 0.11 | 0.15 | 0.22 | -0.4 | 0.0033 | 0.013 | 0.024 | 0.32 | -0.0078 | -0.092 | -0.063 | -0.085 | -0.067 | 0.22 | 0.21 | -0.071 | 0.27 | 0.28 | 0.17 | -0.45 | | |
| Vehicle_size_Midsize | 0.07 | 0.068 | -0.014 | 0.19 | -0.066 | 0.041 | 0.014 | 0.0063 | -0.0027 | 0.0056 | -0.096 | 0.011 | 0.047 | 0.1 | 0.048 | -0.045 | -0.25 | -0.24 | 0.18 | -0.15 | 0.05 | -0.11 | -0.62 | -0.43 | |

# Data Analysis

- Use multiple linear regression

    - to predict one variable based on the known information about other variables, and

    - examine how multiple independent variables are related to one dependent variable

- Used recursive model elimination (RFE) to compare models with different sizes and selected Model 8* with highest adjusted $R^2$ of 0.75

- Compared the performance of the model using the training and test sets

Model 8* (Year, Horsepower, Fuel_type_flex-fuel (unleaded/E85), Fuel_type_premium unleaded (recommended), Fuel_type_regular unleaded, Transmission_type_MANUAL, Driven_wheels_rear wheel drive, Vehicle_size_Compact)

# Model 8

| Final Features For The Best Model(Model 8) | Coefficient |
|---|---|
| Year: | 745.196675 |
| Horsepower: | 150.668859 |
| Fuel_type_flex_fuel(unleaded/E85): | -17650.075979 |
| Fuel_type_premium unleaded (recommended): | -12529.747319 |
| Fuel_type_regular unleaded: | -16154.122128 |
| Transmission_type_MANUAL: | -2799.322581 |
| Driven_wheel_rear wheel drive: | -2139.430713 |
| Vehicle_size_Compact: | 3808.463290 |

# Summary of Findings

- Result of the analysis: predict the MSRP of a car with a 75% accuracy
- We were able to predict a car's MSRP with a given set of features
- We also concluded that the RFE is not the best model to predict a car's MSRP based on the 75% accuracy that the RFE method yielded.
  - We were able to accurately predict a car's msrp after plugging in values for the features listed in Model 8

| Car Features | Inputs |
|---|---|
| Year: | 2011 |
| Horsepower: | 335 |
| Fuel type: | Premium |
| Transmission type: | Manual |
| Driven Wheels: | Rear |
| Vehicle Size: | Compact |
| Predicted MSRP: | $46,549.31 |

# Implications

- At 75% accuracy in predicting MSRP, manufacturers can:
    - Adjust features and specs according to the preferred MSRP,
    - Budget early on,
    - Decide on how to market the car, and
    - Have a better understanding and control of the factors contributing to the MSRP

# Limitations

- Lack of numerical variables
- Unknown values which needed to be dropped individually
- Irrelevant variables such as market category, vehicle style, and popularity.

# Future steps

- Improve our results with a dataset that has more complete data and more relevant variables
- Use a different model for our prediction



**Dataset**

**Car Features and MSRP**

Includes features such as make, model, year, and engine type to predict price

CooperUnion • updated 4 years ago (Version 1)

339

# Classification

The problem: Filling a missing Engine Cylinder value on a car shopping website.

The approach: Using 3 different Binary Classification methods to solve this problem.

Used Methods: Decision Tree, Random Forest, and Naive Bayes

Results: All Methods yielded impressive Classification abilities of 97% and above accuracy scores.

# Analysis and Results



Full Decision Tree

Pruned Decision Tree

## Cross Validation

### Decision Tree

| | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| **4 Cylinders** | 99% | 99% | 99% | 4156 |
| **8 Cylinders** | 97% | 97% | 97% | 1662 |
| **Accuracy** | | | 98% | 5818 |
| **Macro Avg** | 98% | 98% | 98% | 5818 |
| **Weighted Avg** | 98% | 98% | 98% | 5818 |

### Random Forest

| | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| **4 Cylinders** | 98% | 98% | 98% | 4156 |
| **8 Cylinders** | 96% | 94% | 95% | 1662 |
| **Accuracy** | | | 97% | 5818 |
| **Macro Avg** | 97% | 96% | 97% | 5818 |
| **Weighted Avg** | 97% | 97% | 97% | 5818 |

### Pruned Decision Tree

| | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| **4 Cylinders** | 100% | 98% | 99% | 4156 |
| **8 Cylinders** | 96% | 99% | 97% | 1662 |
| **Accuracy** | | | 98% | 5818 |
| **Macro Avg** | 98% | 99% | 99% | 5818 |
| **Weighted Avg** | 98% | 98% | 98% | 5818 |

### Naïve Bayes

| | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| **4 Cylinders** | 99% | 99% | 99% | 4156 |
| **8 Cylinders** | 96% | 98% | 97% | 1662 |
| **Accuracy** | | | 98% | 5818 |
| **Macro Avg** | 98% | 98% | 98% | 5818 |
| **Weighted Avg** | 98% | 98% | 98% | 5818 |

ROC Curve