

Predicting Used Car Prices

SPRINGBOARD CAPSTONE PROJECT

RYAN PAIK

The Problem

As a buyer, it is difficult to know what features of the vehicles and how they affect the price of a used vehicle. The focus is to predict the price of a used vehicle based on their features.



JUSTIN SULLIVAN/GETTY IMAGES

Who does this benefit?

Used Car Dealers



Individual Car buyers and
sellers

Car Features Used to Predict Price

- MSRP
- Listed Price
- Year
- Condition
- Cylinders
- Fuel
- Odometer
- Transmission
- Drive
- Size
- Type

Data acquired from Kaggle

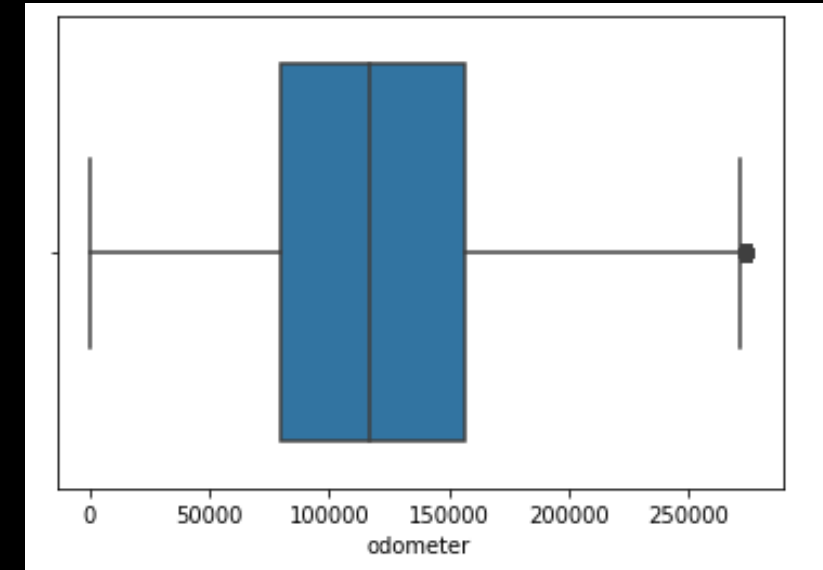
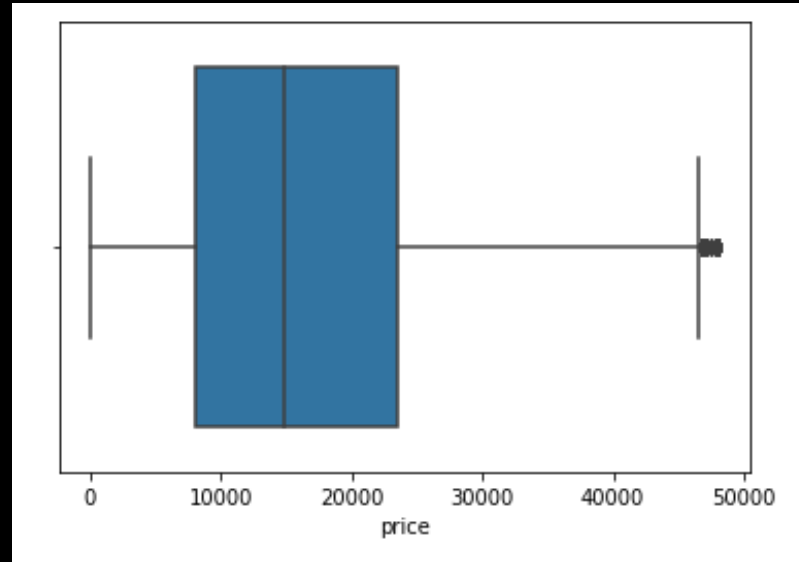
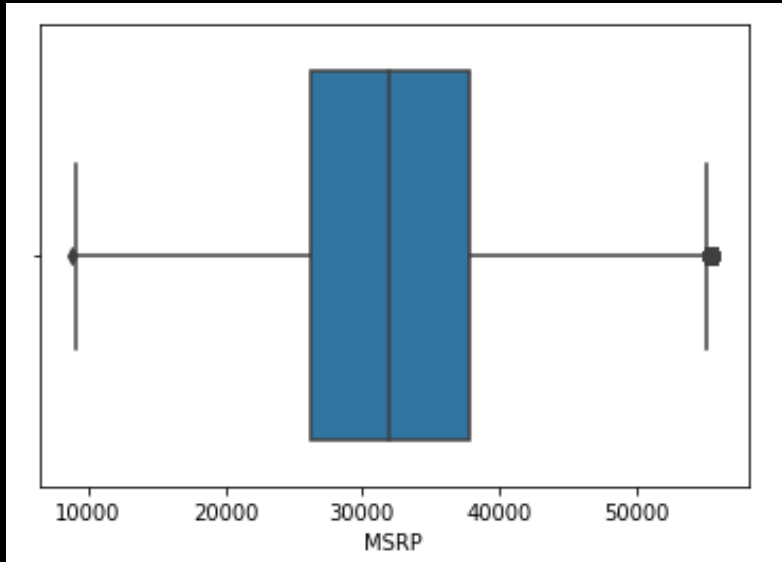
Exploratory Data Analysis

Merged Dataset

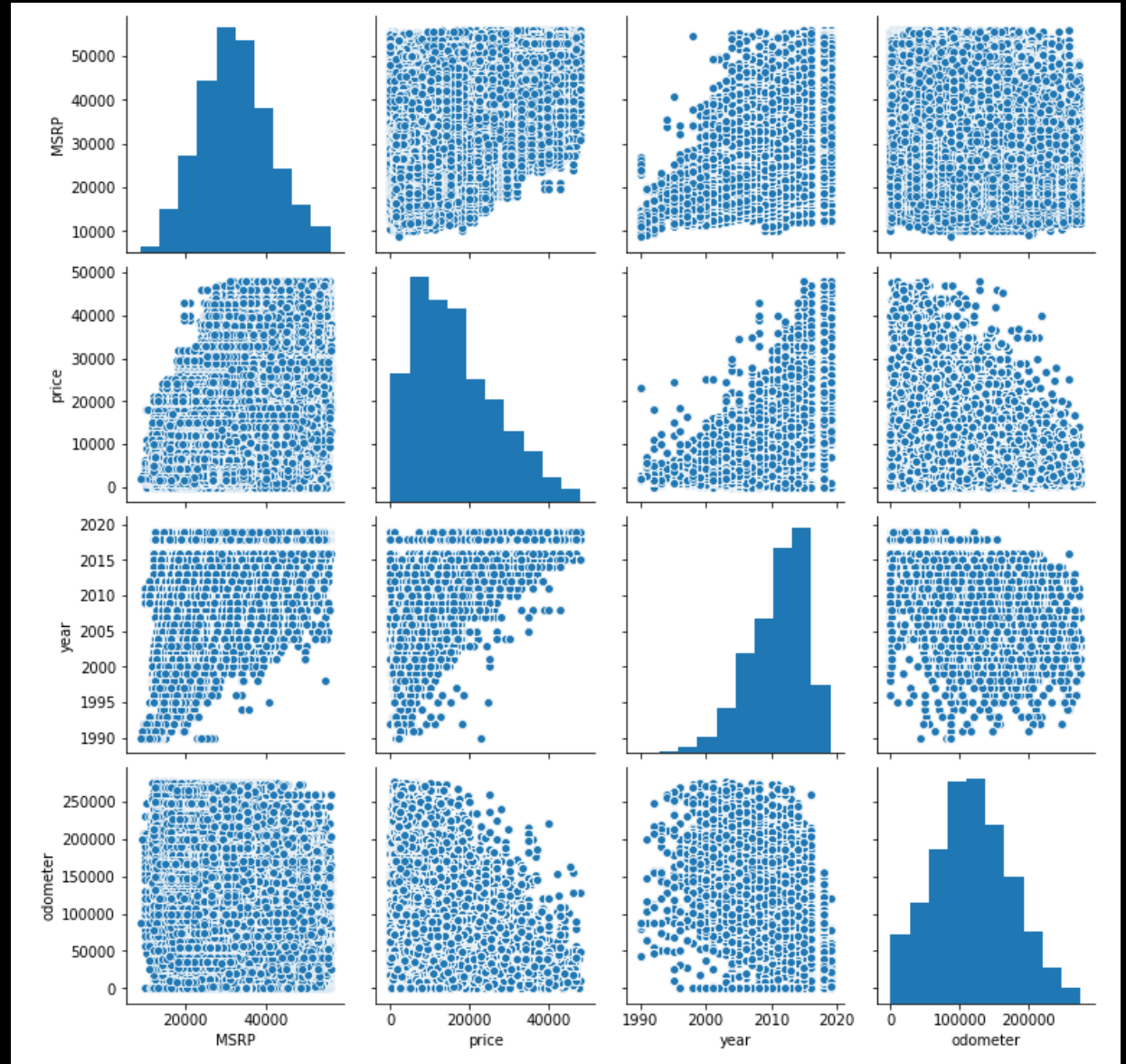
This dataset is the combination of 2 dataset: New car prices and used car prices

	Model	MSRP	price	year	condition	cylinders	fuel	odometer	transmission	drive	size	type
0	2019 acura rdx	40600	36250	2019	excellent	4 cylinders	gas	31250	automatic	4wd	mid-size	SUV
1	2019 acura rdx	45500	36250	2019	excellent	4 cylinders	gas	31250	automatic	4wd	mid-size	SUV
2	2019 acura rdx	43600	36250	2019	excellent	4 cylinders	gas	31250	automatic	4wd	mid-size	SUV
3	2019 acura rdx	37400	36250	2019	excellent	4 cylinders	gas	31250	automatic	4wd	mid-size	SUV
4	2019 acura rdx	42600	36250	2019	excellent	4 cylinders	gas	31250	automatic	4wd	mid-size	SUV
...
304065	2003 volvo xc90	35100	3500	2003	good	5 cylinders	gas	205000	automatic	4wd	full-size	SUV
304066	2003 volvo xc90	35100	2000	2003	fair	6 cylinders	gas	109000	automatic	4wd	mid-size	SUV
304067	1998 volvo s90	34300	2100	1998	good	6 cylinders	gas	144000	automatic	rwd	full-size	sedan
304068	1998 volvo s90	34300	2500	1998	good	6 cylinders	gas	160000	automatic	rwd	full-size	sedan
304069	1998 volvo s90	34300	3500	1998	good	6 cylinders	gas	170000	automatic	rwd	mid-size	sedan

Outliers Removed

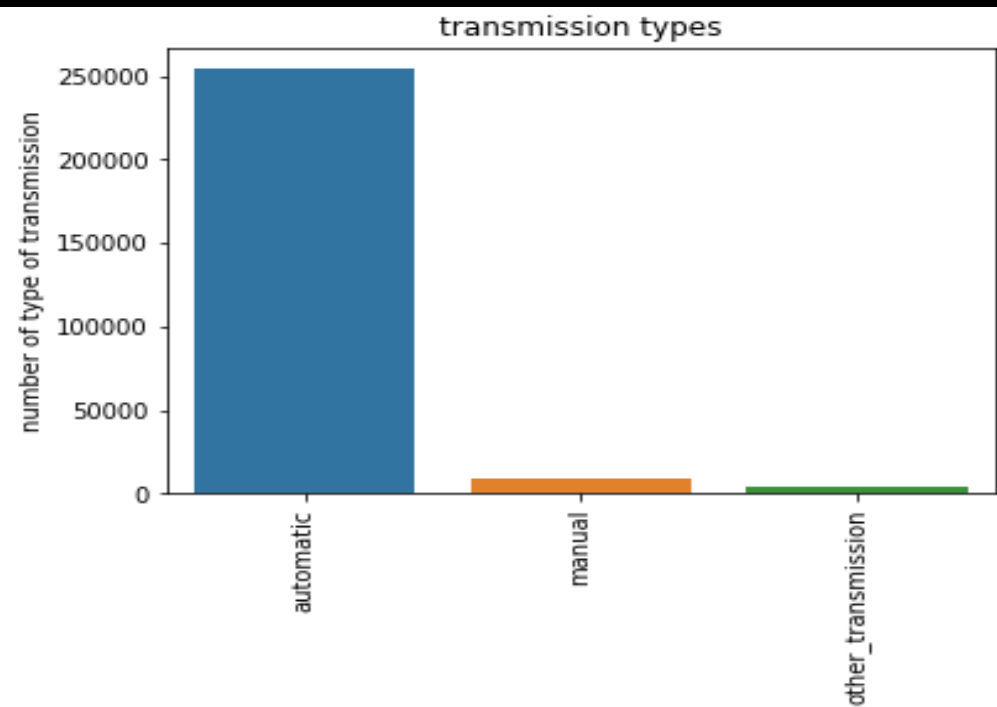
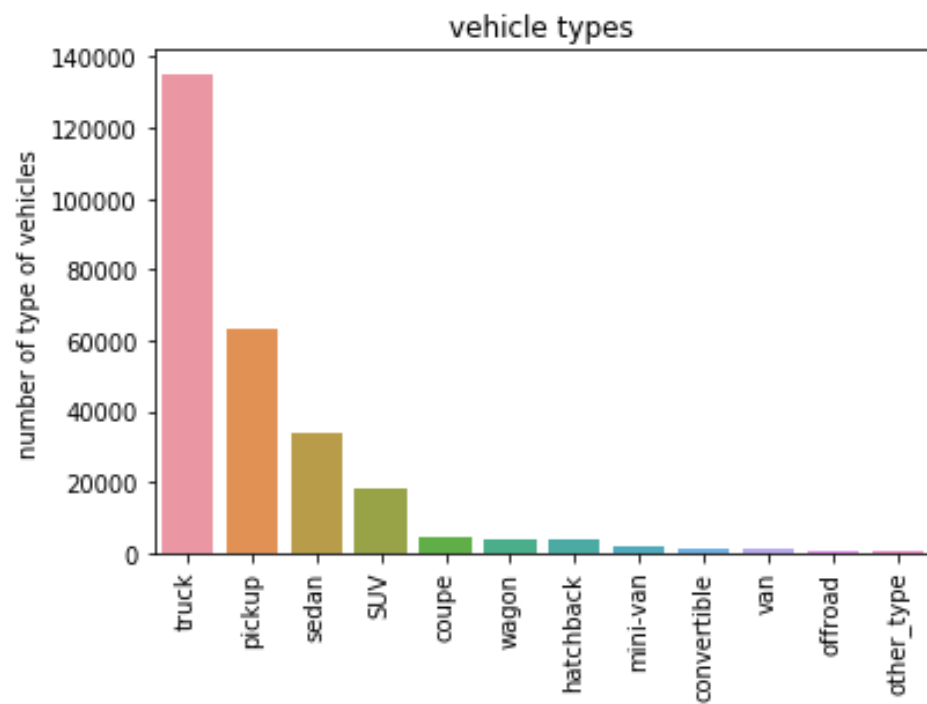


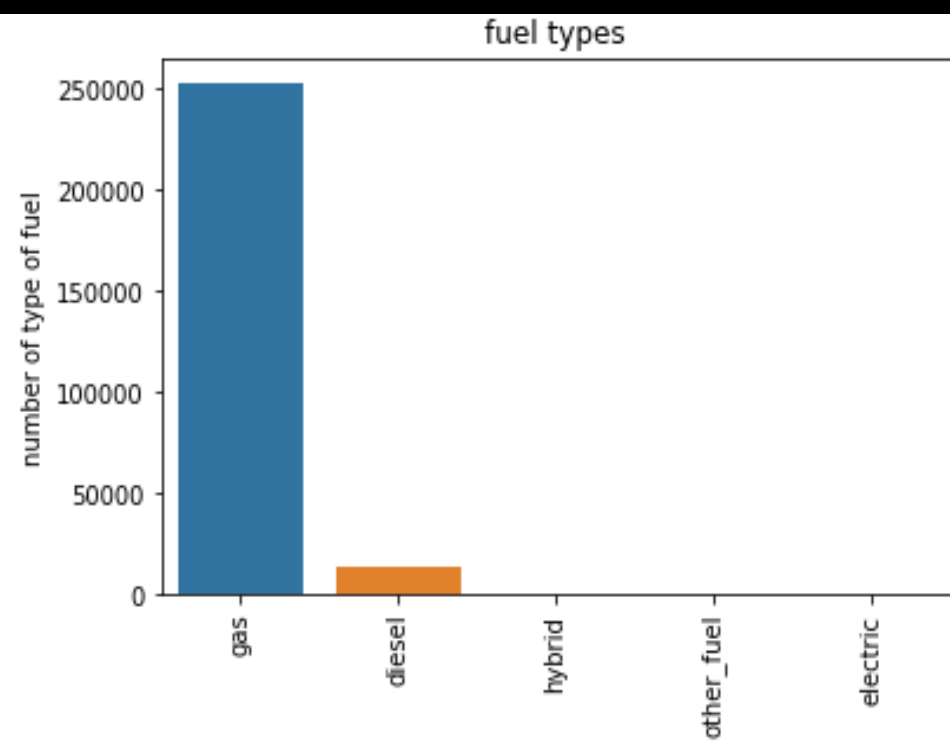
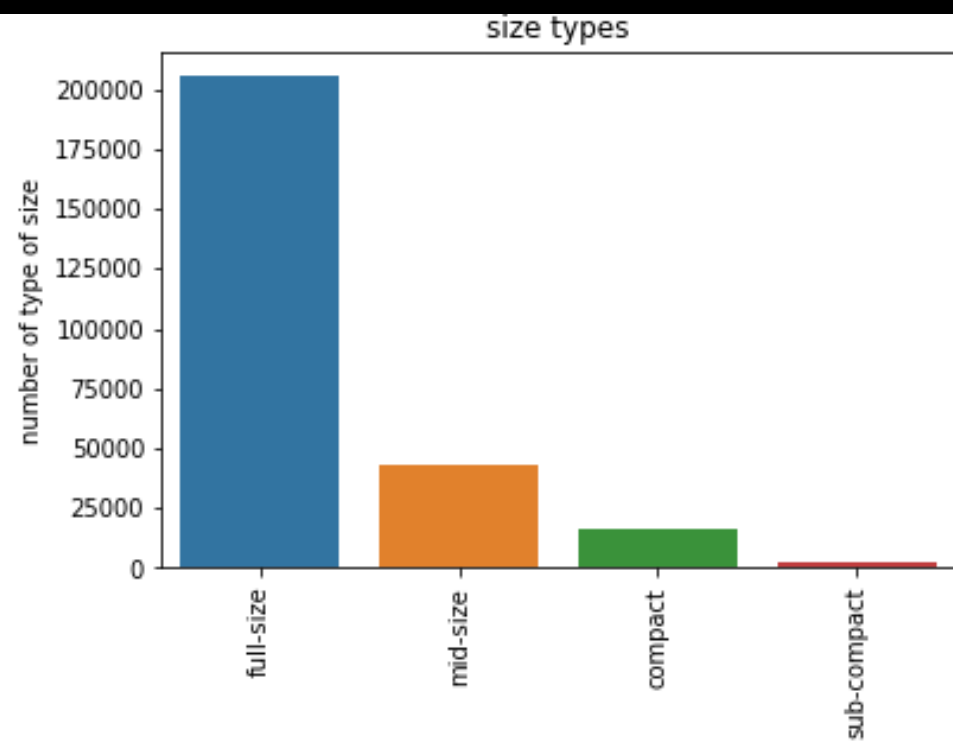
PairPlot

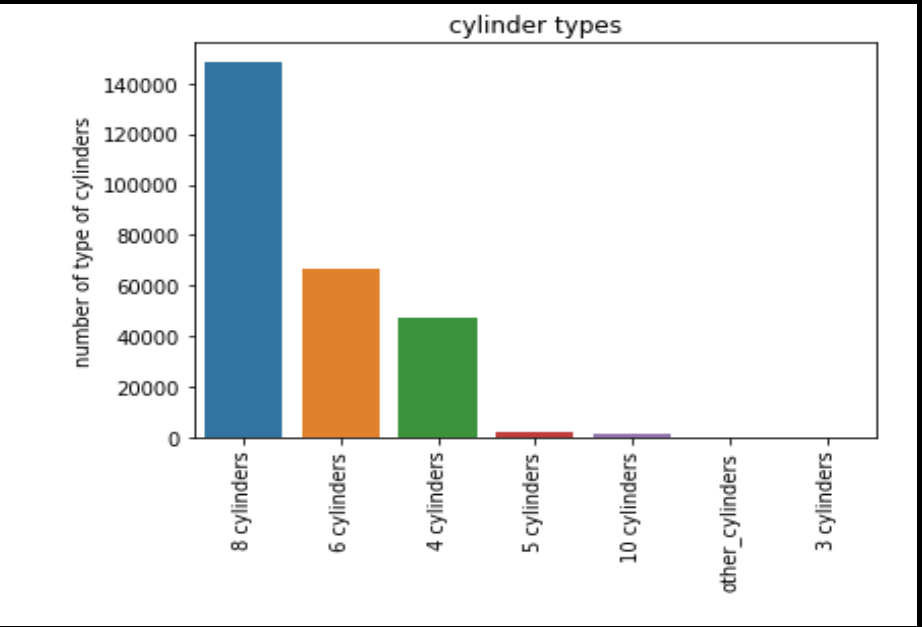
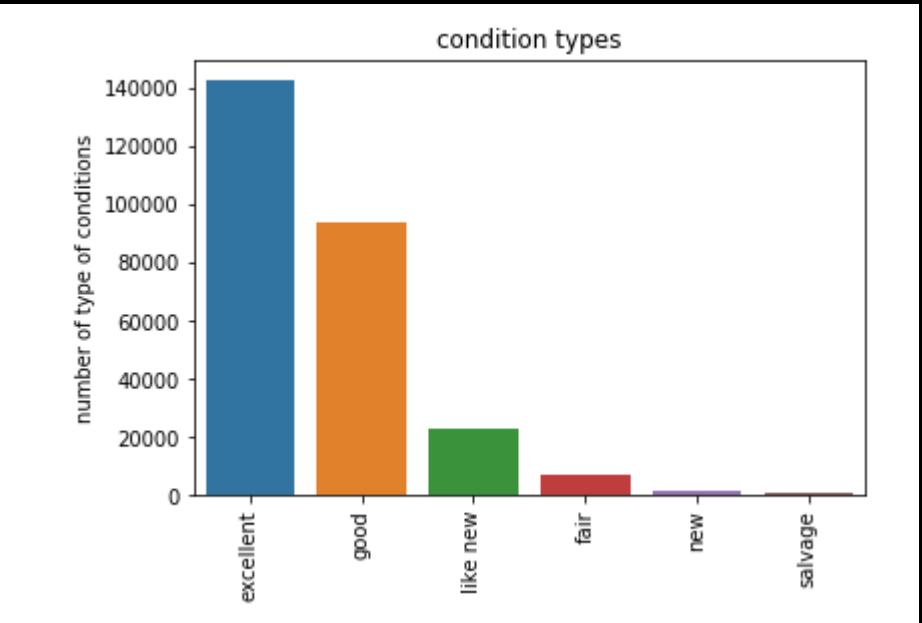
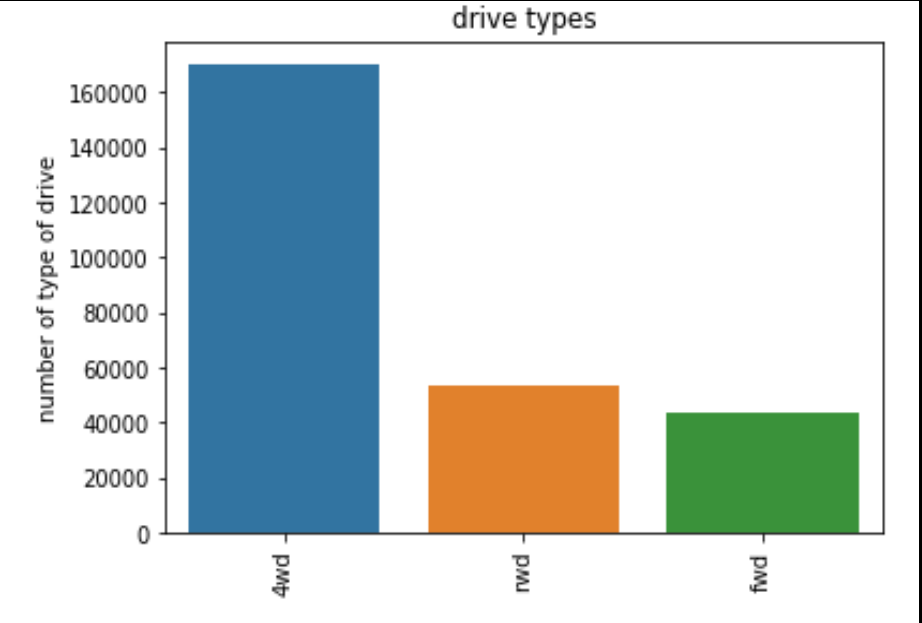


HeatMap









Machine Learning Models

Type: Supervised Learning – Regression

Linear Regression, Decision Tree,
Random Forest, Gradient Boosting

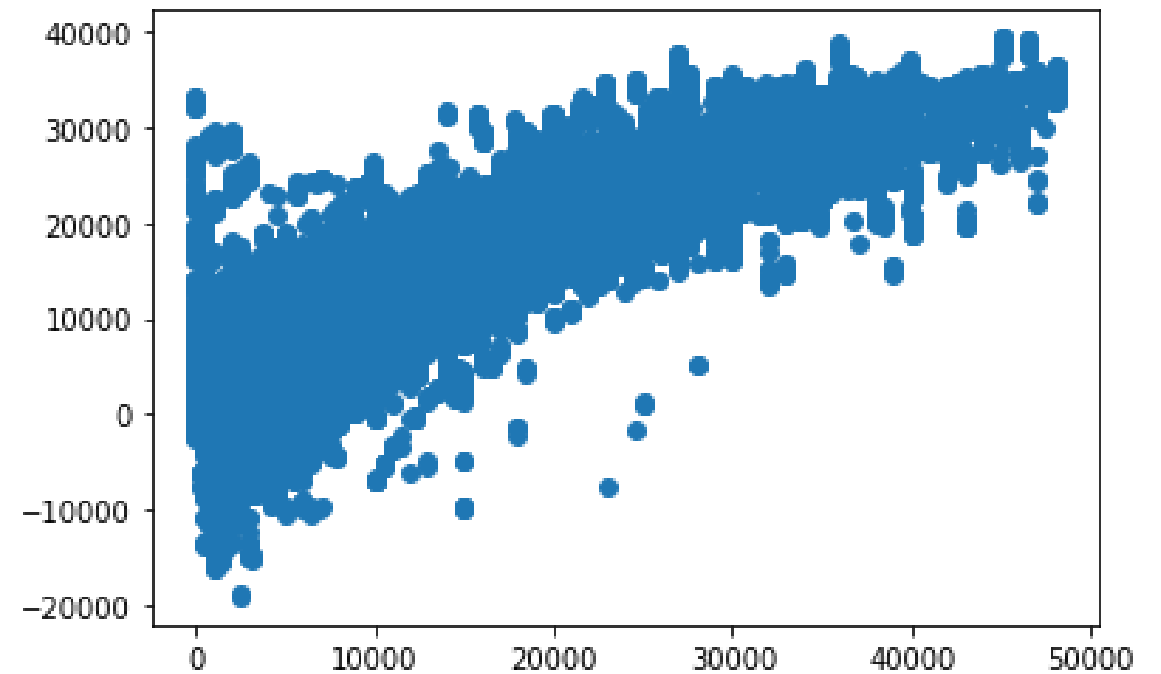
Standard Scaling

GridSearchCV

Linear Regression

Mean absolute error: 3910.62
Mean Squared Error: 27052164.51
Median Absolute Error: 3074.9
R2 score: 0.74

Average CV score: 0.74



Decision Tree

Mean absolute error: 89.03
Mean Squared Error: 740457.4
Median Absolute Error: 0.0
R2 score: 0.99

Average CV score: 0.99144

Hyperparameters:

Max depth: 10
Max features: 5

Train Score: 0.76
Test Score: 0.76

Average Validation Score: 0.728

Random Forest

Mean absolute error: 107.16

Mean Squared Error: 615529.52

Median Absolute Error: 0.0

R2 score: 0.99

Average CV score: 0.99297

Hyperparameters:

Max depth: 10

N estimators: 200

Train Score: 0.88

Test Score: 0.88

Average Validation Score: 0.884

Gradient Boosting

Mean absolute error: 3017.35

Mean Squared Error: 17363207.08

Median Absolute Error: 2272.67

R2 score: 0.84

Average CV score: 0.84

Hyperparameters:

Max depth: 10

N estimators: 300

Train Score: 0.99

Test Score: 0.99

Average Validation Score: 0.987

Conclusions

Model: Gradient Boosting

Future Improvements

- Improve vehicle feature grading system to be objective, rather than being subjective.
- Create a preferential input model for buyers