# Linear Regression

Price of Used Cars Prediction

# Introduction

Linear Regression Model

Prediction Price of Used Cars Listing

# Methodology: Data

MSRP

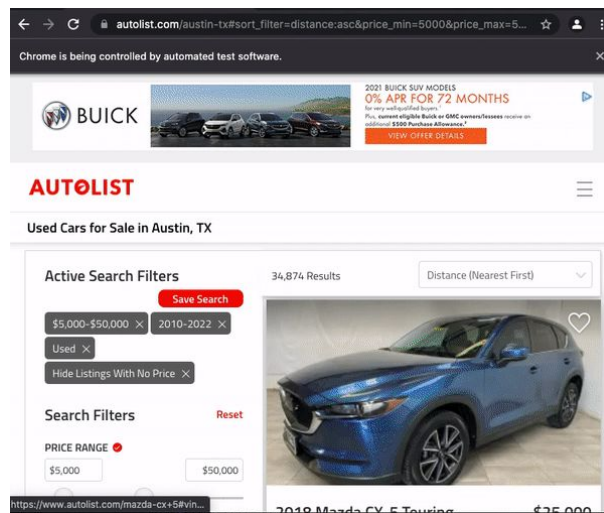https://www.cars.com/research/

MSRP for latest model of 453 cars

## Used Car Listings

https://www.autolist.com/
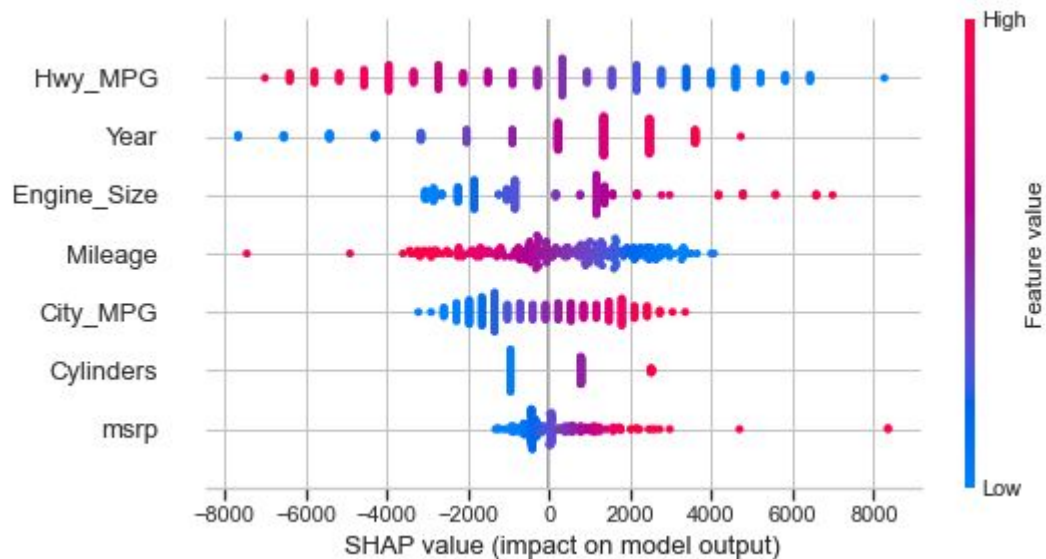
19 variables for 11809 used cars, including Price.

# Methodology: Libraries & Packages

- Selenium, requests & BeautifulSoup: Scrape and parse HTML
- sklearn.impute : Predictions to impute missing data
- Fuzzywuzzy - Search matches between different sources
- sklearn.linear_model : Create regression models
- yellowbrick.regressor: Drop influential datapoints
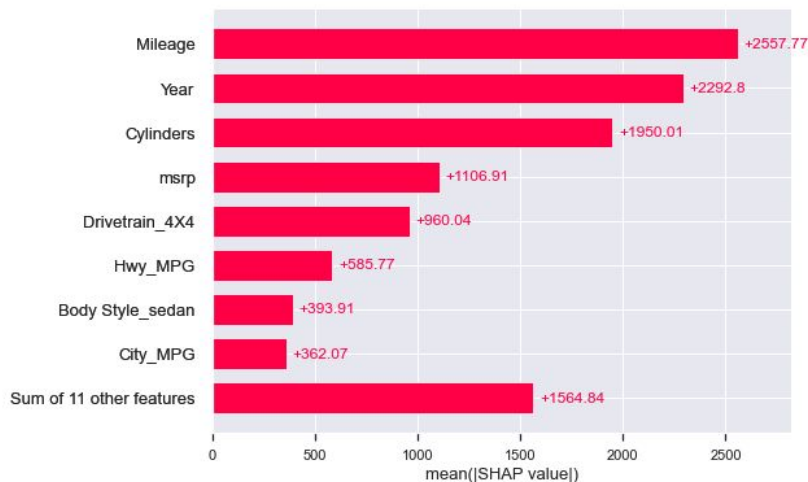- Shap: Visualization of influence of features in model.

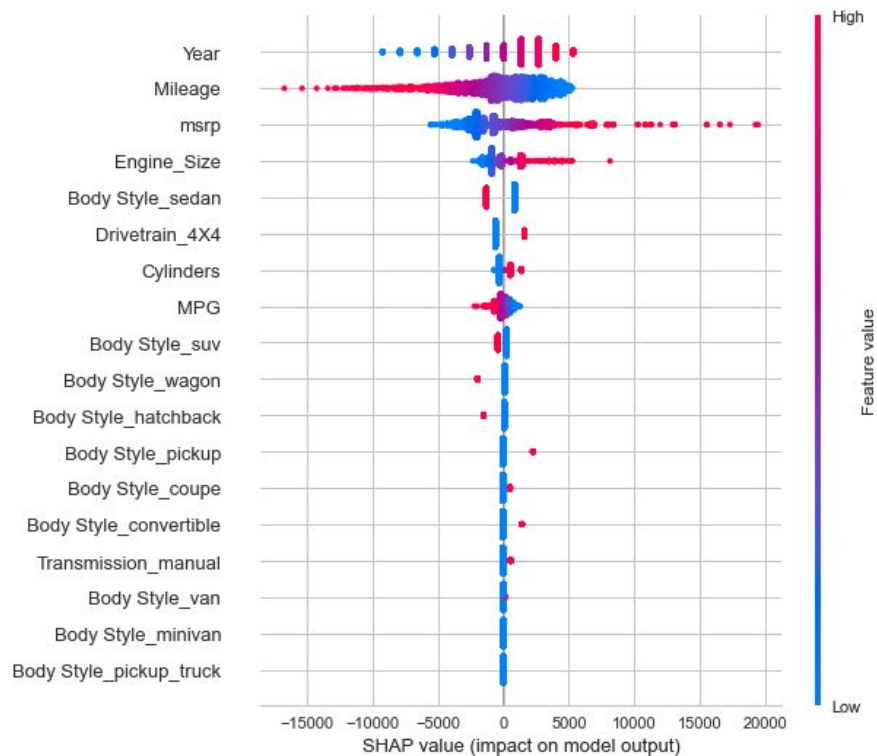# Numerical Features



Numerical Features' SHAP Values

# All Features



## Mean SHAP value of Features

SHAP Values of all Features

# Correlated Features & Regularization

Impute missing values

Imperfect Collinearity

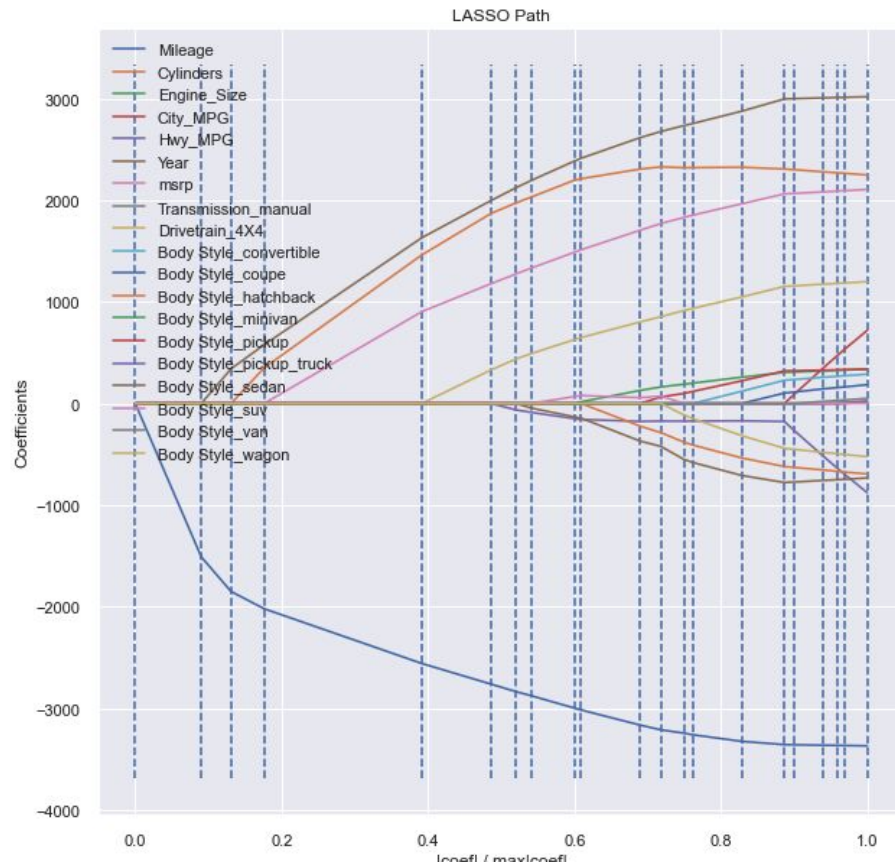Regularization

| | variables | vif |
|---|---|---|
| 0 | Mileage | 3.825799 |
| 1 | Cylinders | 40.772570 |
| 2 | Engine_Size | 7.444121 |
| 3 | City_MPG | 106.744936 |
| 4 | Hwy_MPG | 268.425848 |
| 5 | Year | 152.686841 |
| 6 | msrp | 6.901849 |

# Variables Dropped

- City MPG
- Manual Transmission
- Body Style - Minivan
- Body Style - Van
- Body Style - SUV
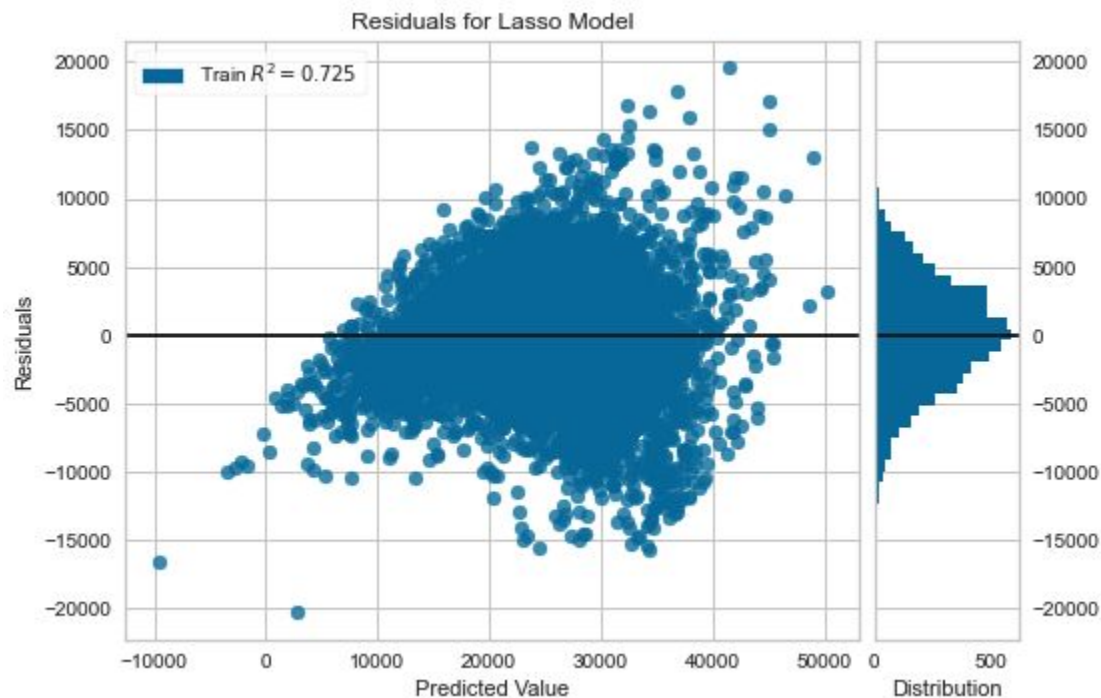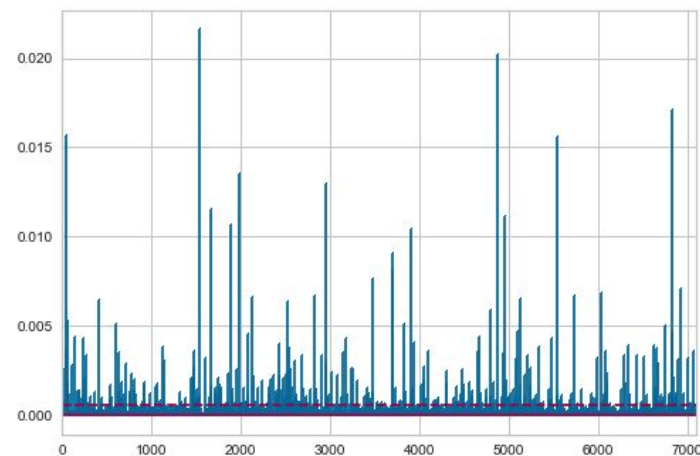- Body Style - Pickup Truck
- Body Style Hatchback



LASSO Path

# Regularization Results

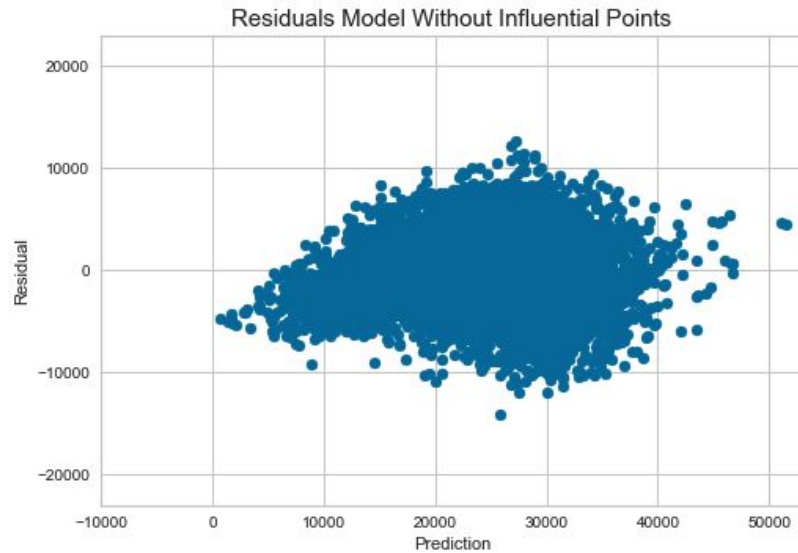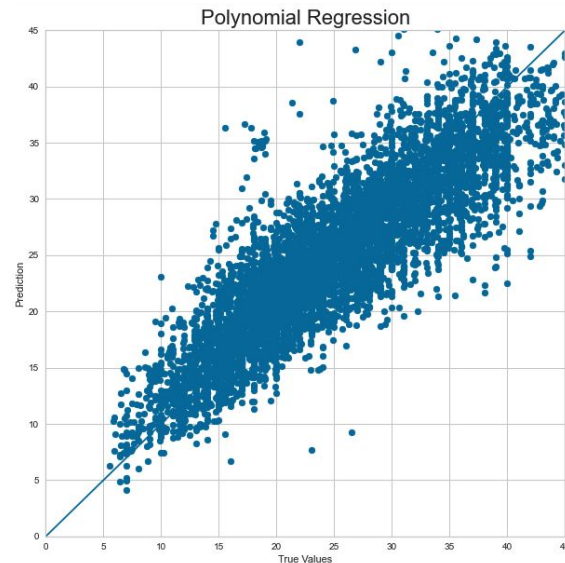| Model | Mean Squared Error | Mean Absolute Error | R^2 Score | Parameters |
|---|---|---|---|---|
| Ridge (L1) | 1.837455e^07 | 3335.628 | 0.725271 | Alpha: 12.9 |
| Lasso (L2) | 1.837435e^07 | 3335.626 | 0.725274 | Alpha: 2.0 |
| Elastic Net (L1 & L2) | 1.83735e^07 | 3335.624 | 0.725274 | Alpha:1.97 L1: 1 |

# Residual Analysis



Residuals for Lasso Model

Train $R^2 = 0.725$

399 Influential Points

# Regression without Influential Points



Penalized Linear Regression



Residuals Model Without Influential Points

| Linear Model | Mean Absolute Error | Root Mean Squared Error | R^2 Score | Parameters |
|---|---|---|---|---|
| With Influential Points | 1.837435e^07 | 3335.62 | 0.725 | Alpha: 2.0 (L1) |
| Without Influential Points | 1.357837e^07 | 2952.71 | 0.769 | Alpha: 2.0 (L1) |

# Polynomial Regression
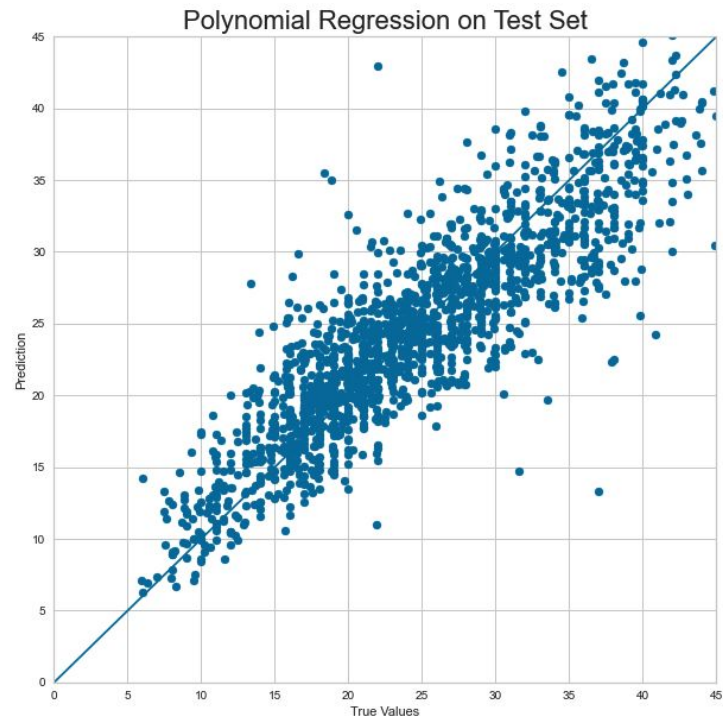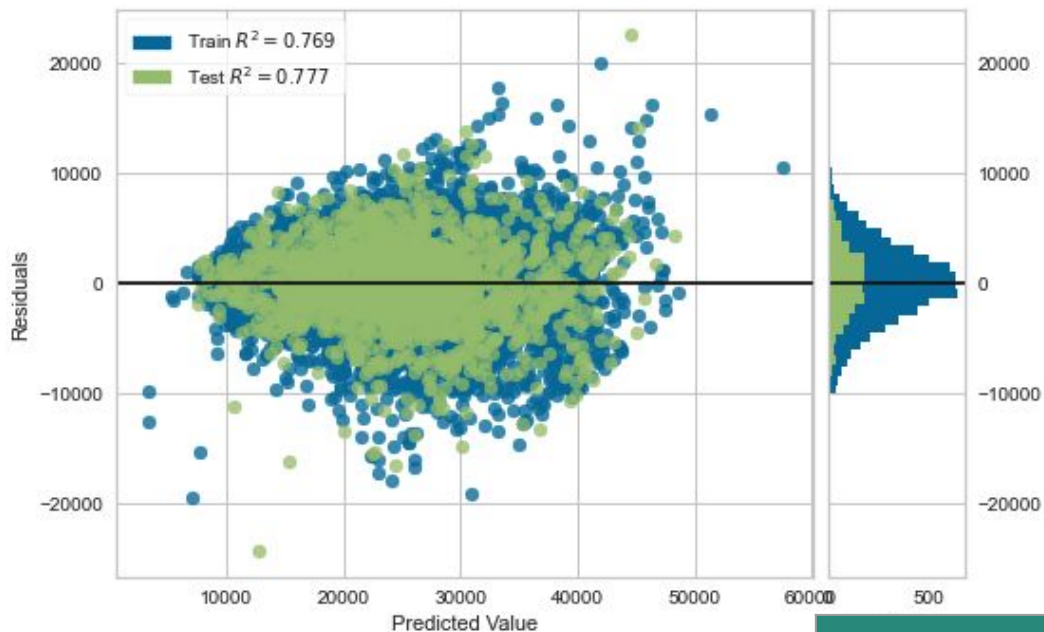


| Linear Model | Mean Absolute Error | Root Mean Squared Error | R^2 Score | Parameters |
|---|---|---|---|---|
| Penalized Linear Model | 1.357e^07 | 2953 | 0.769 | Lambda 2.0 (L1) |
| Polynomial Model | 1.201e^07 | 2714 | 0.796 | |

# Out of Sample Performance

| Linear Model | Mean Squared Error | Root Mean Squared Error | R^2 Score |
|---|---|---|---|
| Polynomial Model | 1.627e^07 | 3016 | 0.770 |

# Results: Unseen Data



| Linear Model | Mean Squared Error | Root Mean Squared Error | R^2 Score |
| --- | --- | --- | --- |
| Polynomial Model | 1.627e^07 | 3016 | 0.770 |

# Conclusions

- Low Bias


- High Variance


- RMSE: ~$3.000

# Future Work

Regularization in Polynomial Regression

Need features indicating level of car
(basic/premium)

MSRP for individual car.

Remove influential datapoints from SHAP
Summary Plot