



# PREDICTING CAR PRICE

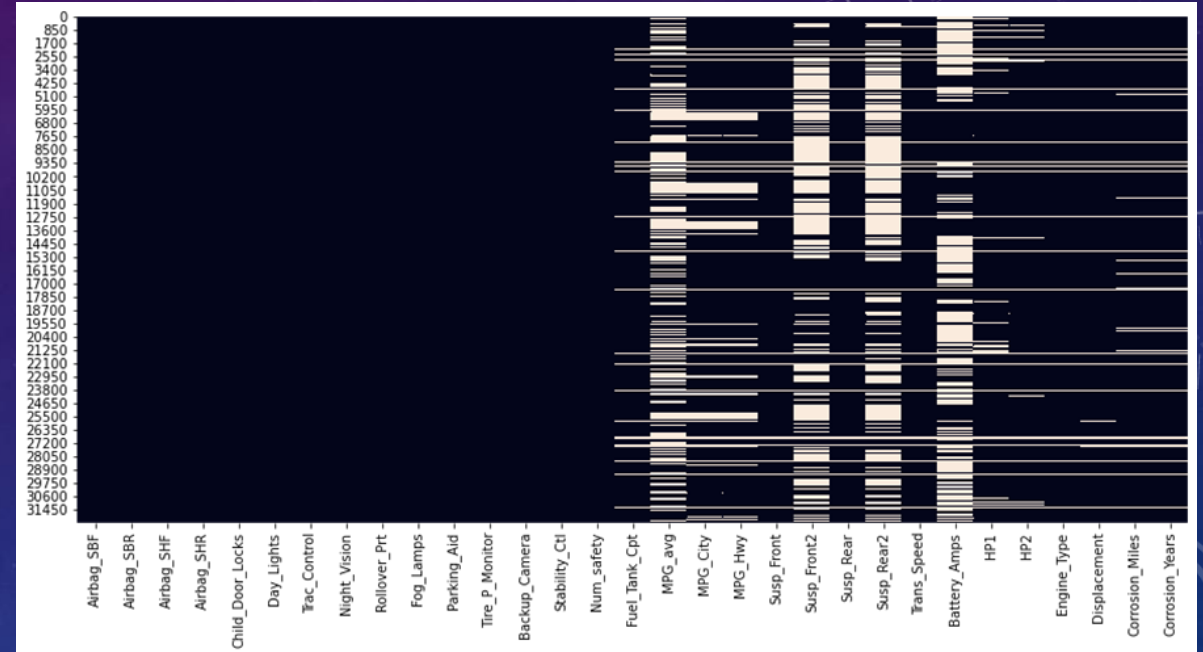
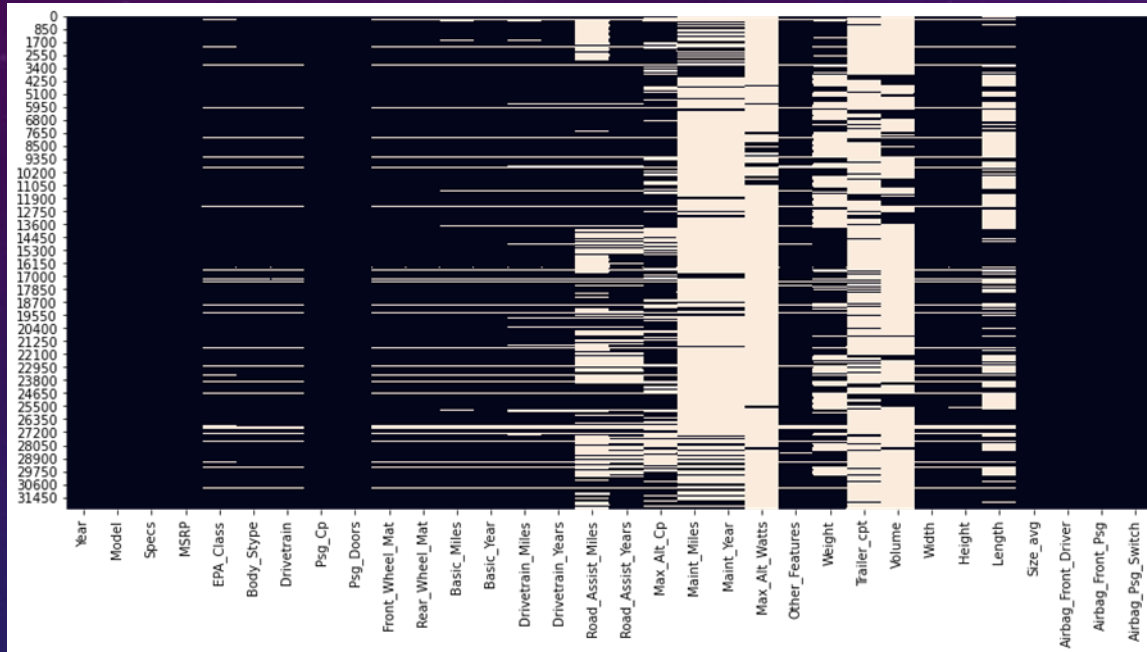
AUTO MACHINE LEARNING

Dongtao Jiang, Data Science Career Track, Springboard  
8/2/2020

# DATASET

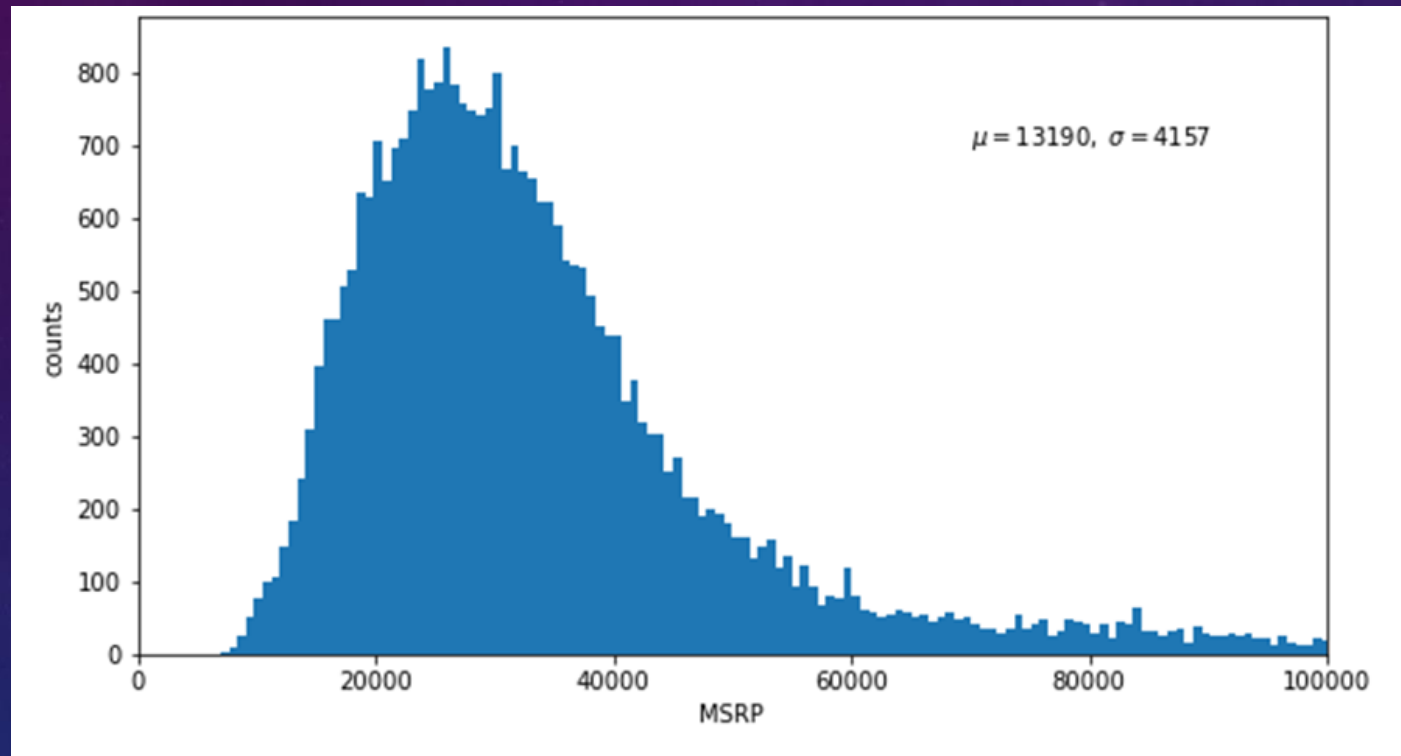
- Data source: [thearconnection.com](http://thearconnection.com)
- Features
  - Dimensions
  - Fuel economy
  - Performance specs
  - Safety features
  - Warranty

# MISSING VALUES



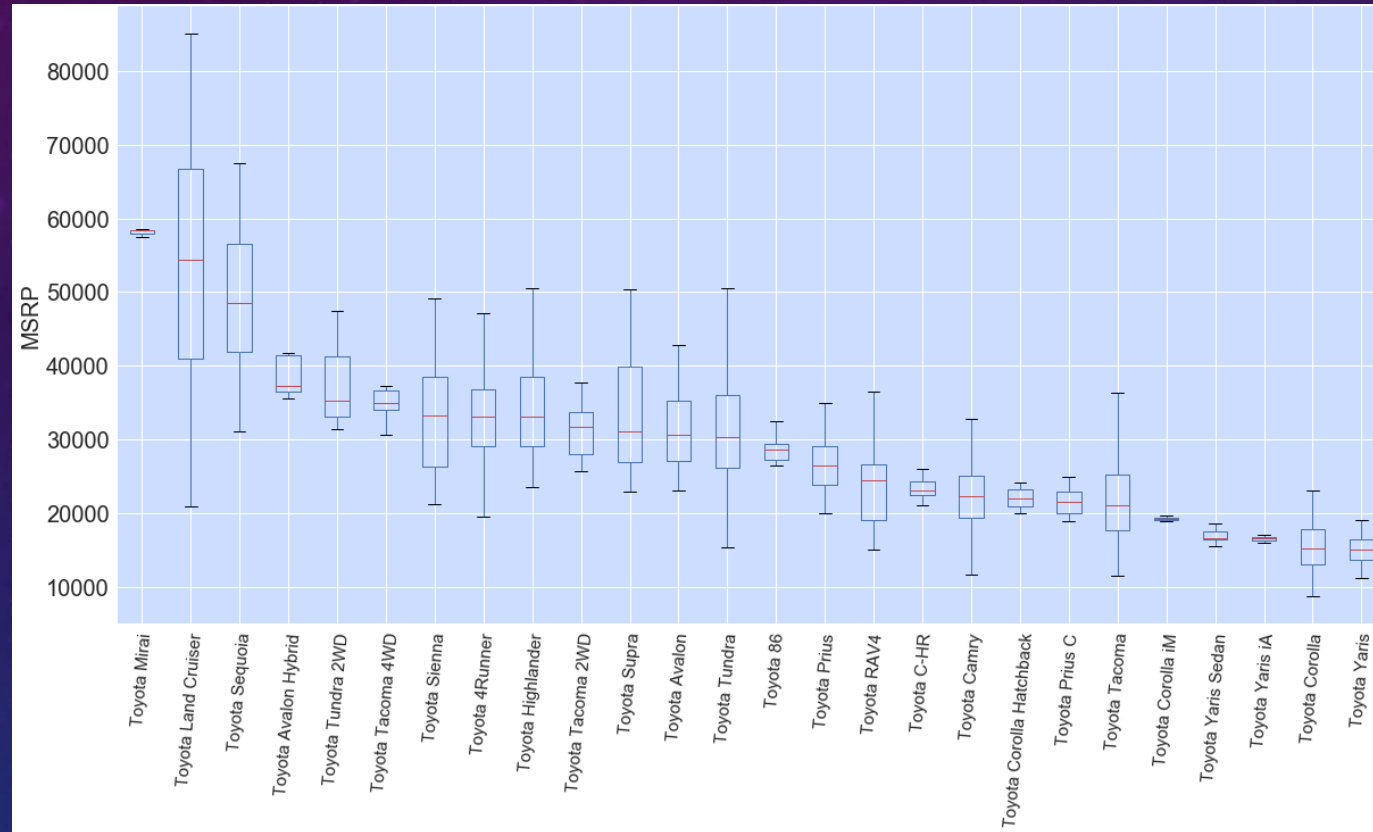
- 15.7% of total dataset are missing.
- Top missing features: Max\_Alt\_Watts Maint\_Miles, Maint\_Year, Volume

# PRICE DISTRIBUTION



- Average price is \$13200.
- Standard deviation is 4160.

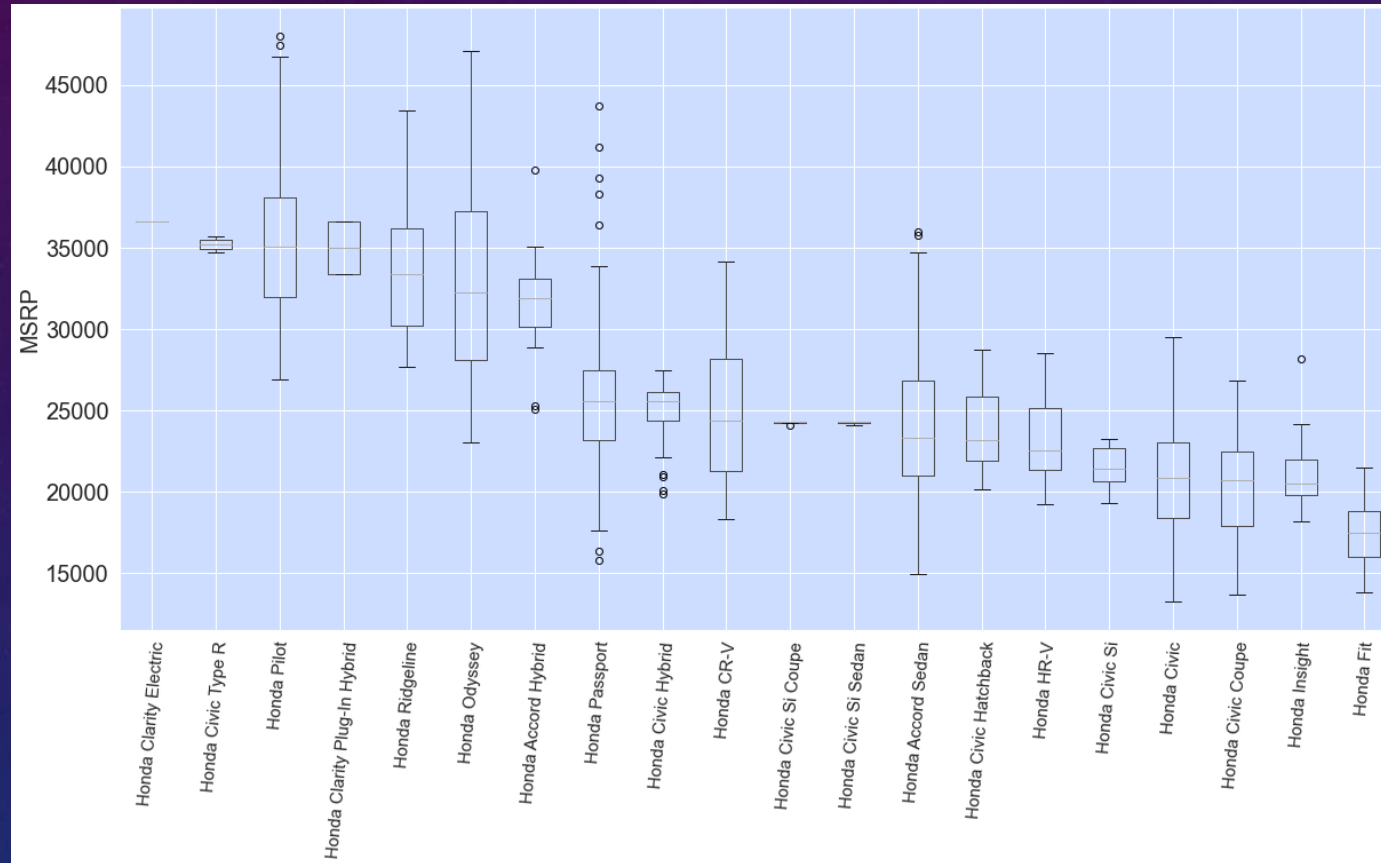
# TOYOTA PRICES SORTING – BOX PLOT



- Luxury model to basic model in sequence.

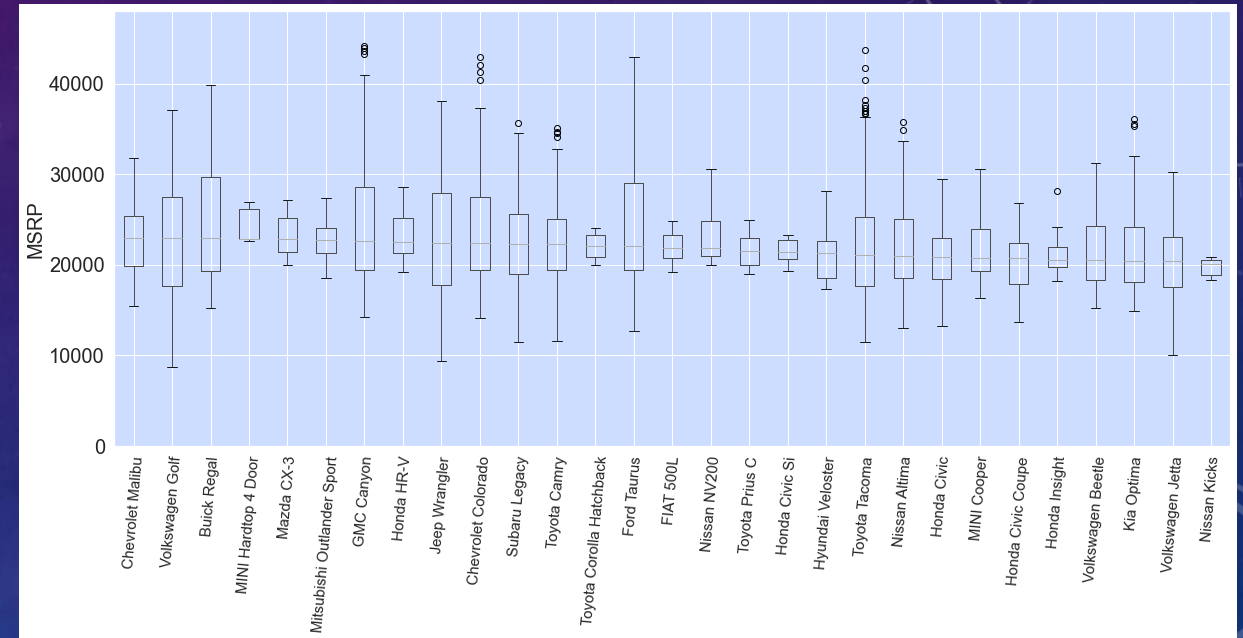
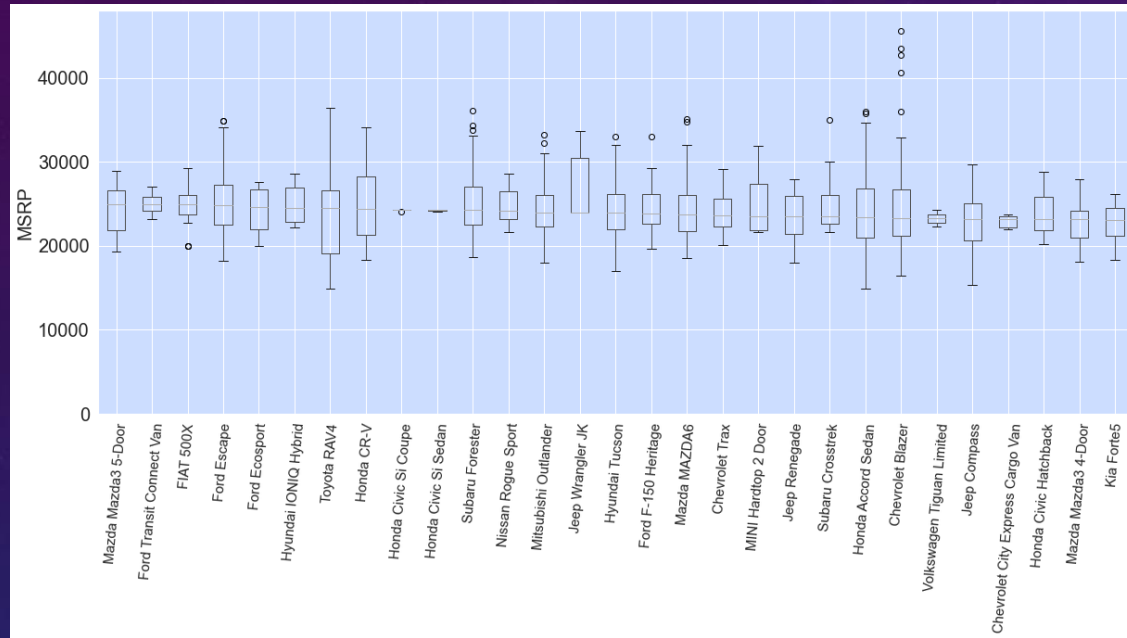


# HONDA PRICES SORTING – BOX PLOT



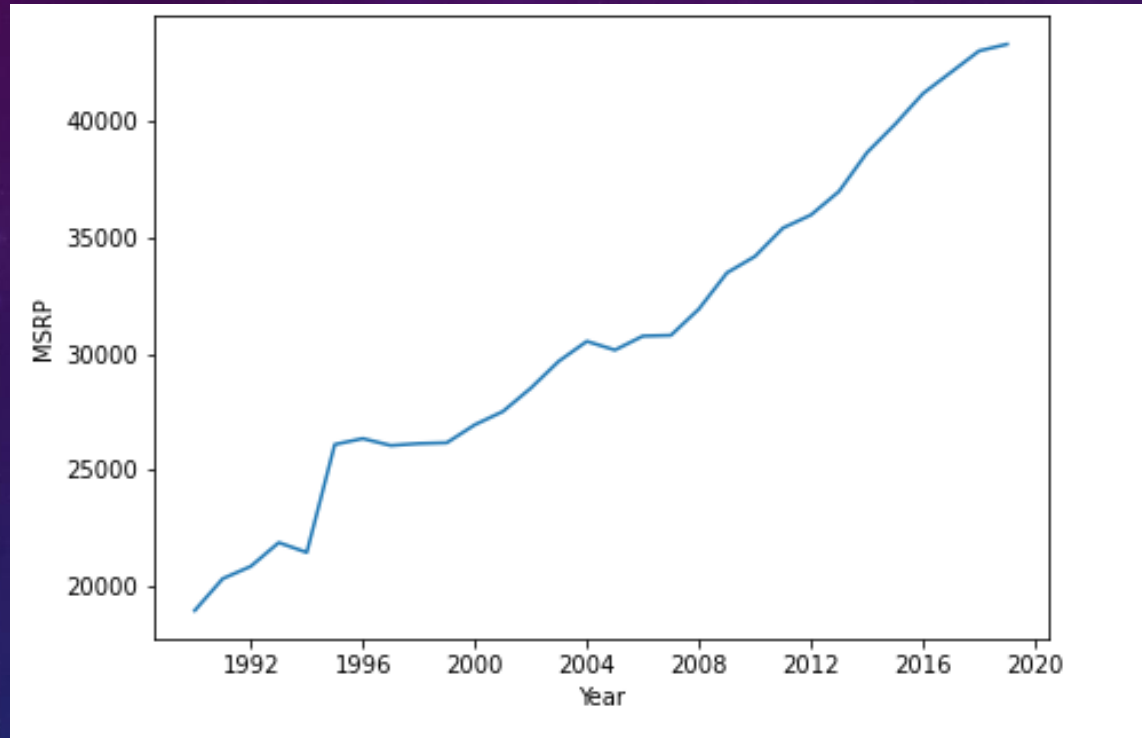
- Luxury model to basic model in sequence.

# BUYING GUIDE FOR BUDGET-TIGHT CUSTOMERS



- Selection of low-end models priced about \$22,000.

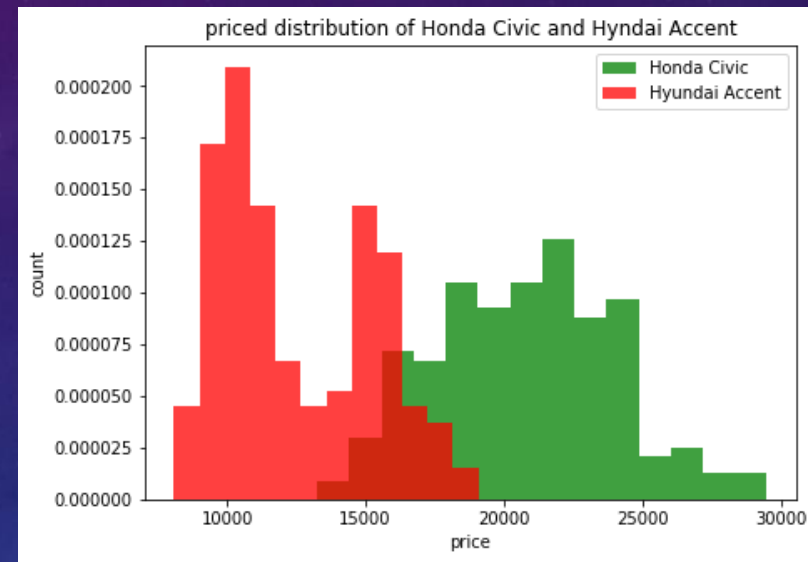
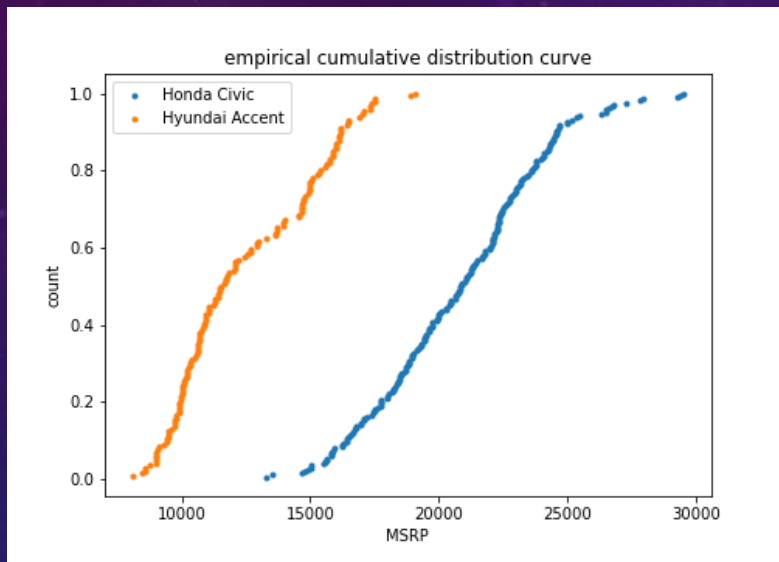
# EVOLUTION OF AVERAGE PRICE OVER ALL MODELS



- Car prices have been constantly increasing generally linearly over the last two decades.
- This indicates year is an important feature to predict prices.



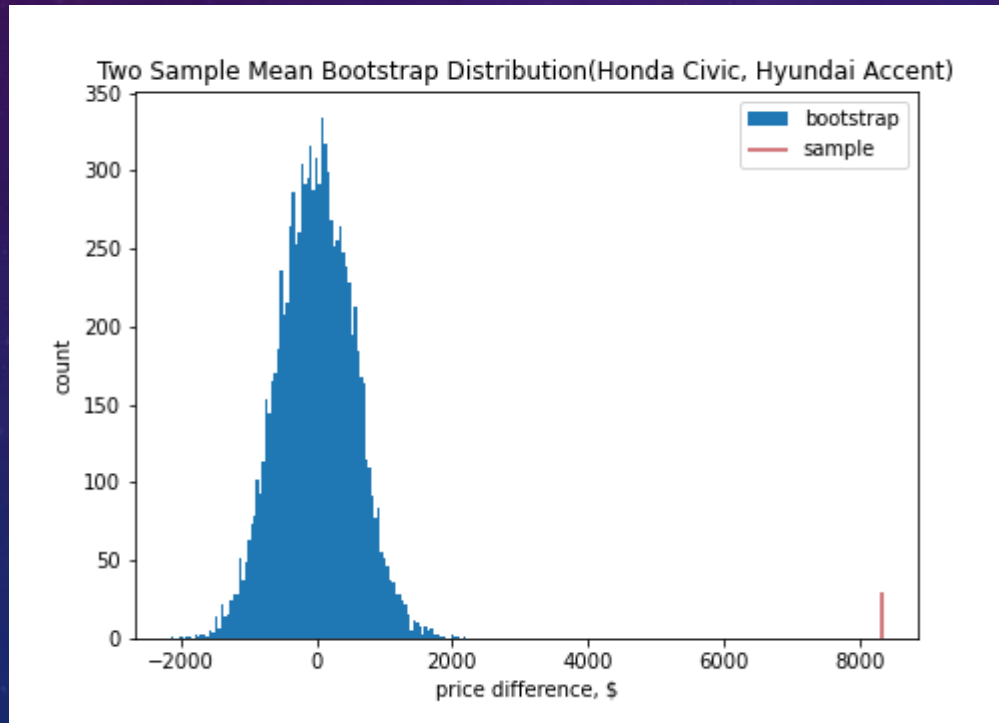
# PRICE DISTRIBUTION FOR TWO CAR MODELS



- Difference in prices of the two low-end car models is significant.
- A small fraction of price overlap exists.

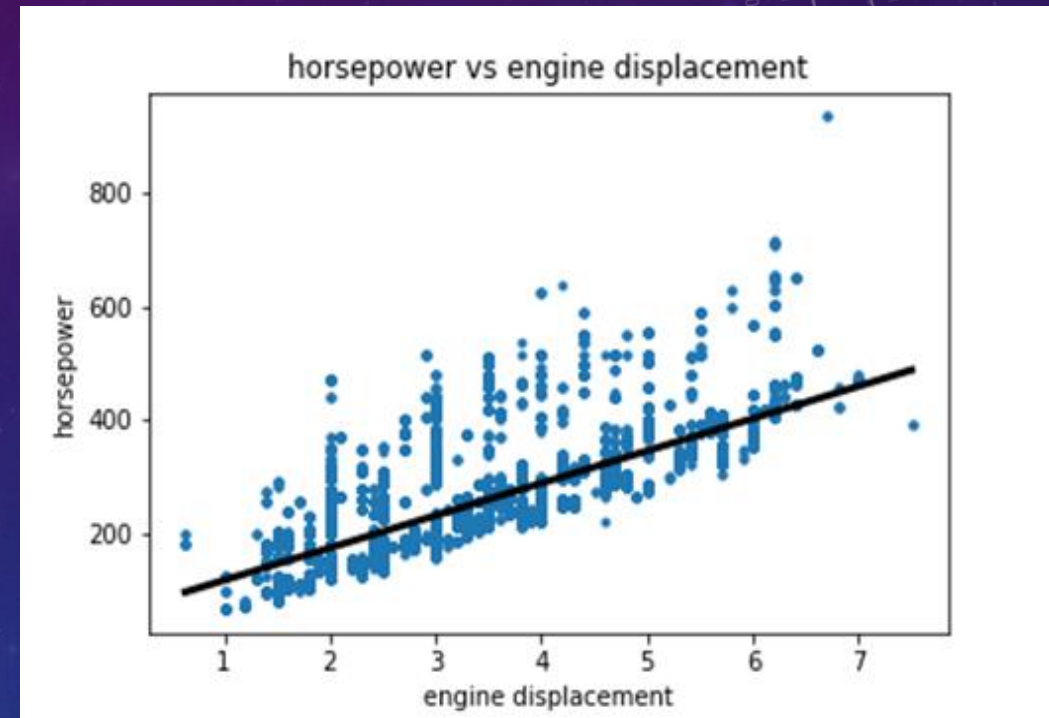
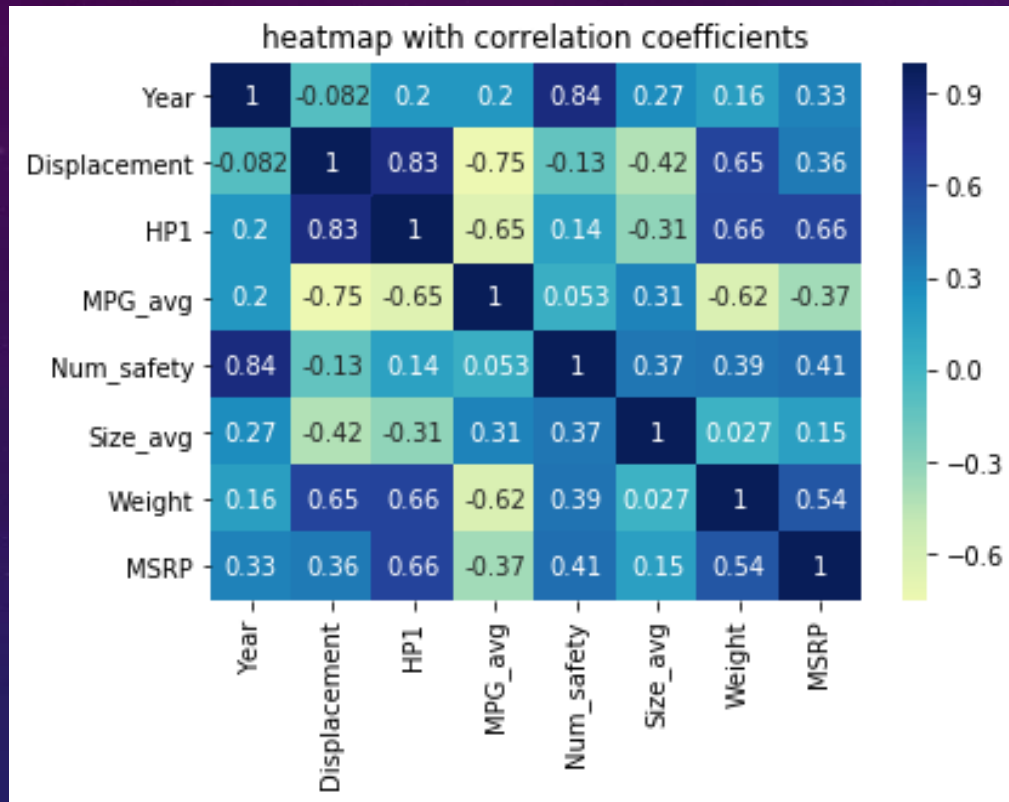
# TWO-SAMPLE HYPOTHESIS TEST – BOOTSTRAP APPROACH

- $H_0$ : There is no difference in the mean price between Hyundai Accent and Honda Civic.
- $H_a$ : There is obvious difference in the mean price between Hyundai Accent and Honda Civic.
- $\alpha$ : 5%



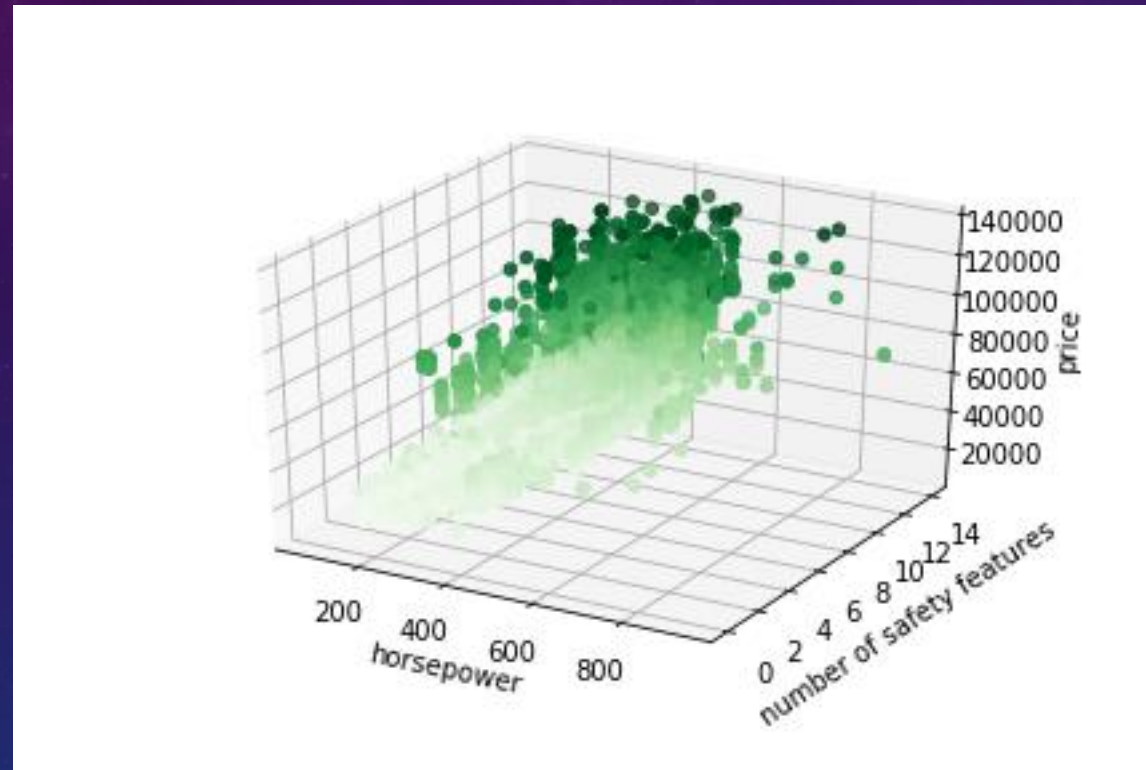
- Results:
  - p-value: 0
  - Null hypothesis rejected.
  - Approve the sharp difference in mean.

# CORRELATION INVESTIGATIONS



- Strong correlation between horsepower and engine displacement agrees well with physics and engineering principles.

# 3-D VISUALIZATION



- 3-D plot is utilized to visualize the effect of two features combined on prices.



# MACHINE LEARNING MODELS

- **Linear Regression**

One of the most well-known and well understood algorithms in statistics and machine learning.

- **Ridge Regression**

Useful to mitigate the problem of multicollinearity in linear regression.

- **Lasso Regression**

Ideal for producing simpler models.

- **Decision Tree**

Easy to interpret.

- **Random Forest**

Merges multiple decision trees to get a more accurate and stable prediction.



# MACHINE LEARNING

Models	r2_score	Negative mean_absolute_error
Random Forest	0.986	-1443
Linear regression	0.968	-1002483
Lasso Regression	0.967	-3262
Ridge Regression	0.964	-3525
Decision Tree	0.908	-5470

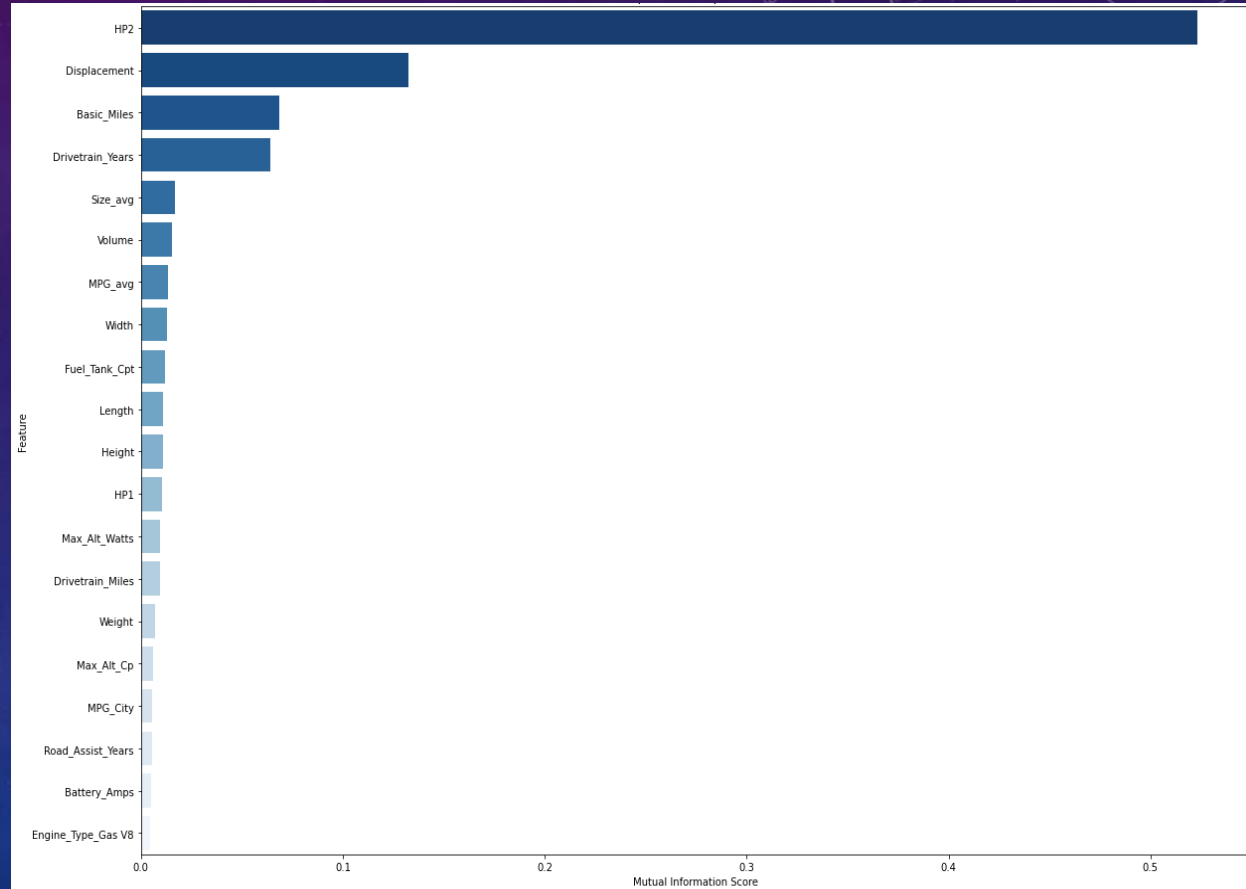
- Most of the models used give excellent goodness of fit that their r2 scores are mostly above 0.96. This indicate these models are of good choice.
- Random Forest gives the best performance with r2 being very close to 0.99. Decision Tree has the poorest performance with a r2\_score of 0.908.

# DIFFERENT METRICS

- Negative mean absolute error compared with  $r^2$ .
- Decision Trees still tops the perforce.
- Scoring for Linear Regression is, however, exceedingly low.
- Performance of algorithms varies with metrics.

# FEATURE IMPORTANCE

Rank	Feature	Importance
0	HP2	0.523067
1	Displacement	0.132579
2	Basic_Miles	0.068280
3	Drivetrain_Years	0.064147
4	Size_avg	0.016712
5	Volume	0.015565
6	MPG_avg	0.013185
7	Width	0.013044
8	Fuel_Tank_Cpt	0.011811
9	Length	0.011152
10	Height	0.011051



# FEATURE IMPORTANCE

- The HP<sub>2</sub> (torque spec) has the highest importance, which is at least 4 times higher than the rest.
- The 2<sup>nd</sup> highest is Displacement.
- This result agrees well with engineering know-how.
- Engine is the most dominating part for the major performance of car for example the car's lifetime, speed, driving smoothness, horsepower, fuel efficiency, etc.



# CONCLUDING REMARKS

- Hypothesis testing quantified significant difference in the mean price between two low-end popular car models: Hyundai Accent and Honda Civic.
- Exploratory data analysis was conducted to visualize missing values over all dataset, provide buying guide for low-income customers by extracting all lowly-priced car models and sorting in order.
- The pair plot and heatmap indicate a positive correlation between horsepower and engine displacement, which agrees well with physics and engineering principles.
- Machine learning algorithm Missforest was successfully used to impute missing values, which is one of the reasons we end up with very high predicting accuracy.
- Five models were experimented including linear regression, ridge regression, lasso regression, decision trees and random forest. Except lower performance of decision trees, all the other models deliver very good  $r^2$  scores higher than 96%. The best model is random forest that scored close to 99%.
- Feature importance analysis revealed the torque spec (HP1) and displacement are the most important factor determining the car prices. This result indicates the power of random forest because these two features are related to the heart of car: engine.



# FUTURE WORK

- Acquire more domain knowledge for the purpose of feature engineering
  - remove unnecessary features
  - create new features
- Tune models hyperparameters to marginally improve performance
- Obtain data for newer car models to test the model