

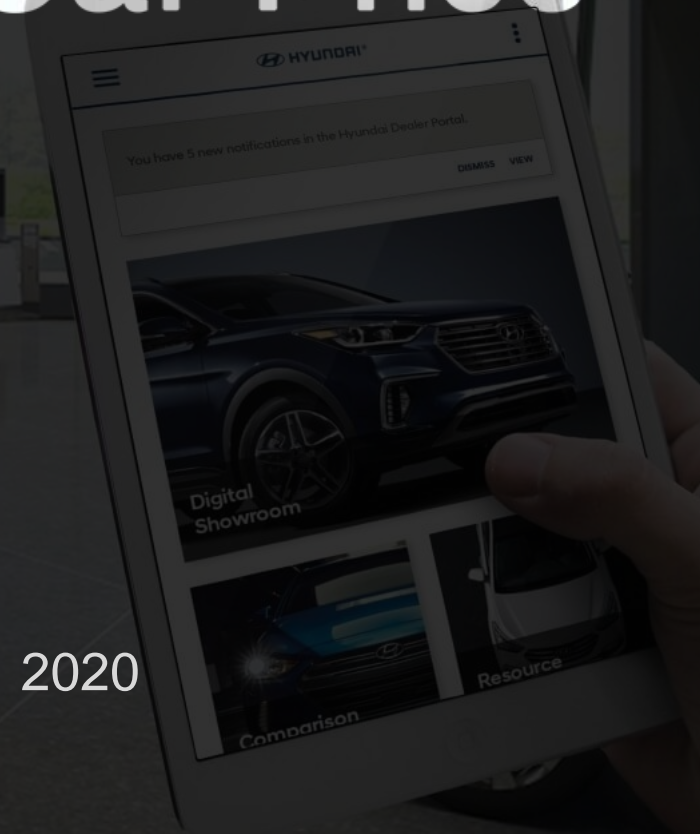
# Predicting the Used Car Price

From Craigslist

Yang Fei

Mentor: Kenneth Gil-Pasquel

Data Science Capstone Project 1, June 2020





What is the Target?

**Physical Dealer**



**Price Prediction**

predictions on the most popular Top 5 used car  
manufacturers on craigslist



**Online Shopping**



# Who might cares?



## Buyers

get a clear and concise price by observing an average standard.



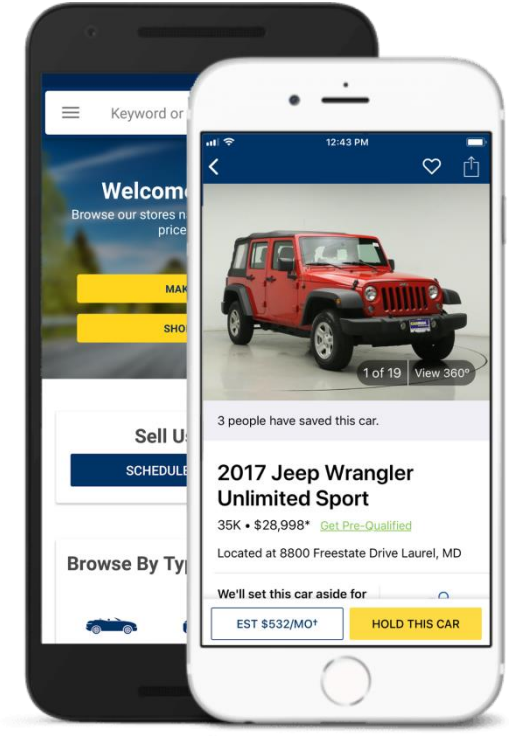
## Used Car Dealer

Make up their marketing plans by analyzing used car price trend and consumer psychology



## Car Manufacturers

Used car price performance may affect their future development and sales strategies



## Online APP Developer

Combine the prediction into their product to make it more competitive

# Where is the data from?

kaggle



Kaggle

<https://www.kaggle.com/austinreese/craigslist-carstrucks-data>

Craigslist

<https://www.craigslist.org/>  
the world's largest used car platform.

## Raw Data

509577 rows and 25 columns

### Year

1990-2020

Format  
csv

This data contains most all relevant information that Craigslist provides on car sales including columns like price, condition, manufacturer, latitude/longitude etc.

## DATA CLEANING & WRANGLING

- Deal with missing and duplicated data
- Drop the fake data
- Detecting & Filtering Outliers
- Add more columns

## Clean Data

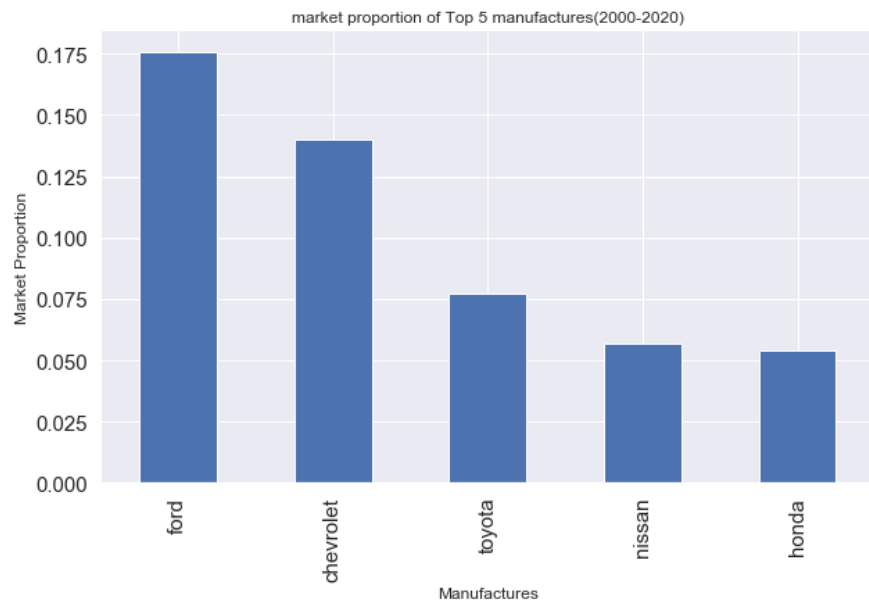
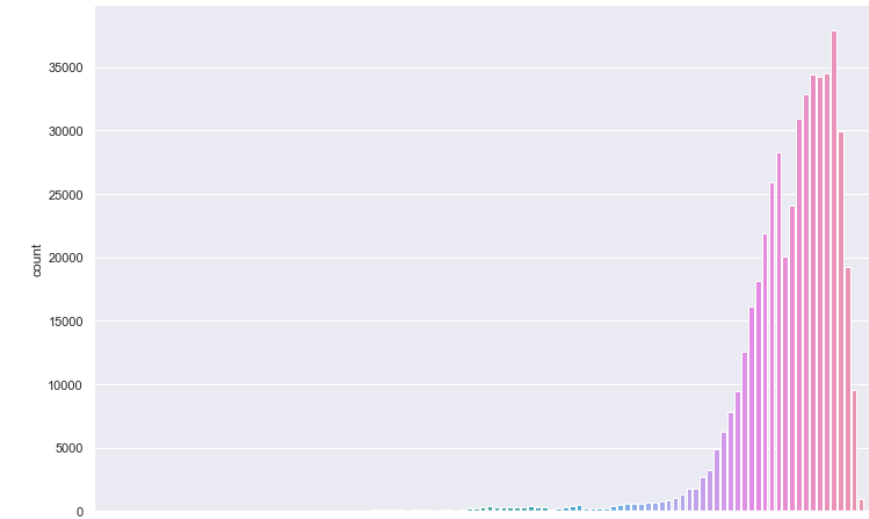
485862 rows and 15 columns

### Year

1990-2020

1. Add one additional column named 'age' which is calculated by 'year'
2. Add two columns name 'odometer\_class' and 'price\_class' for EDA Visualization.

# What is the research object?



## Research Period

Take the  
latest 20  
years

Take the top 5  
manufacturers  
which take the  
largest market  
proportion

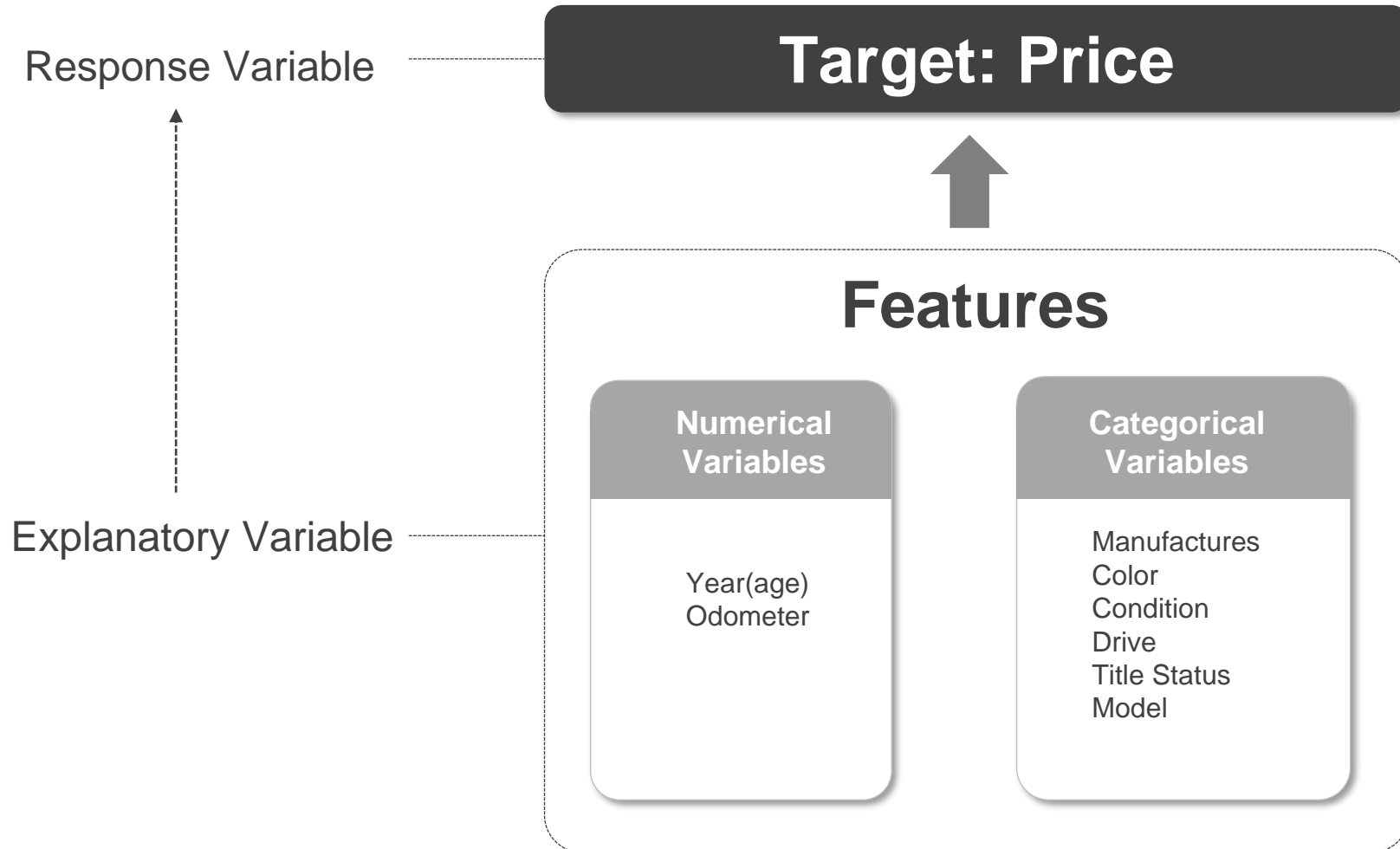
## Manufacturers

## Research Object

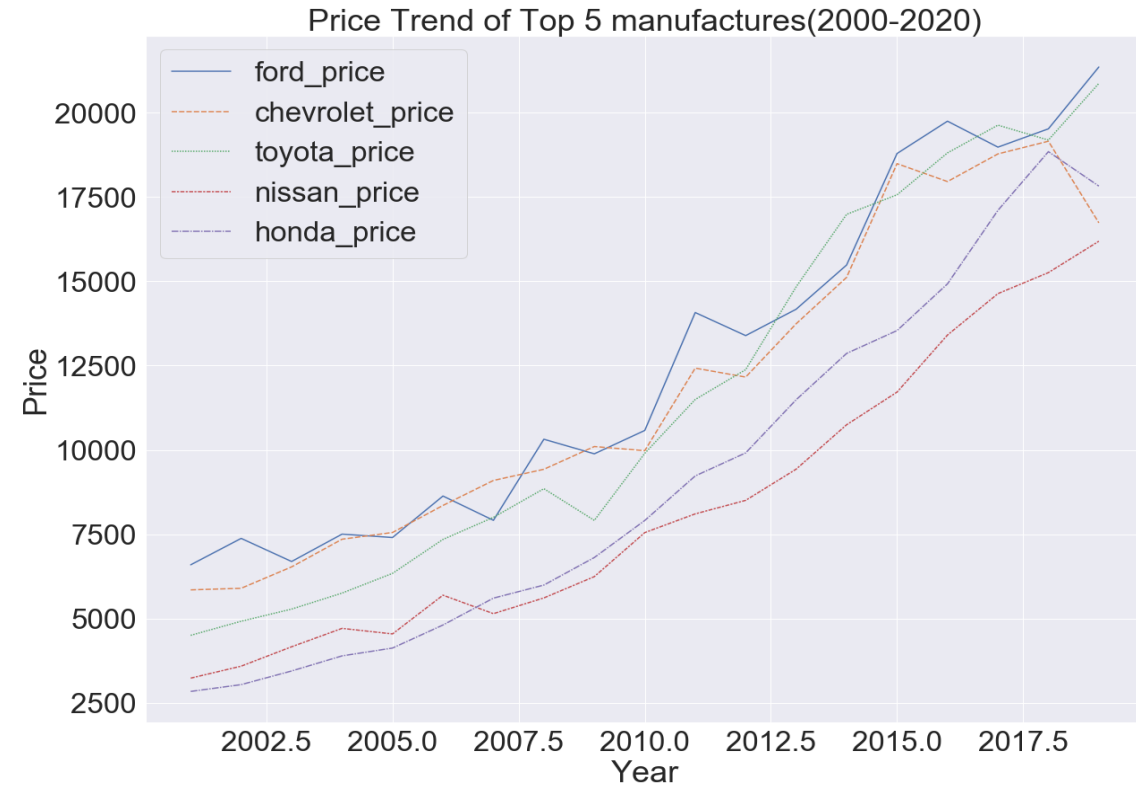
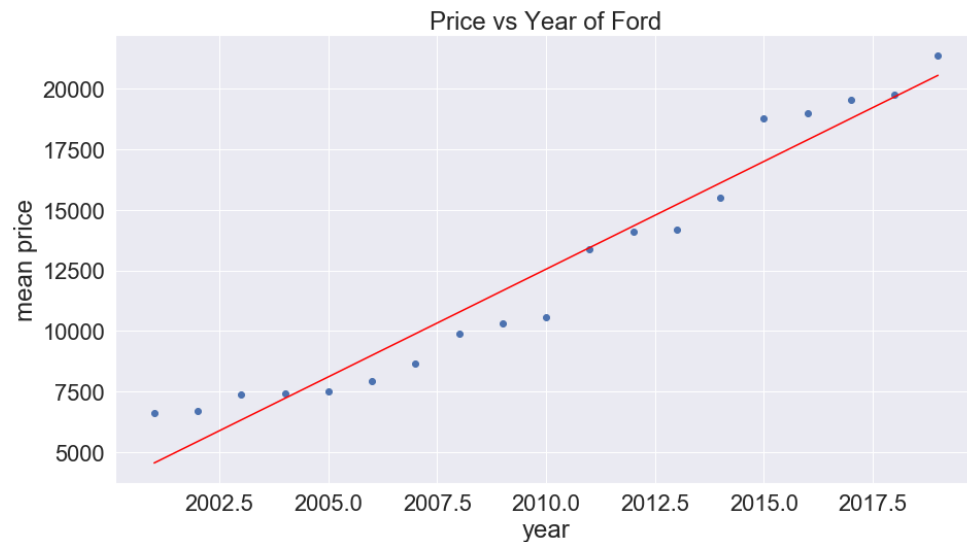
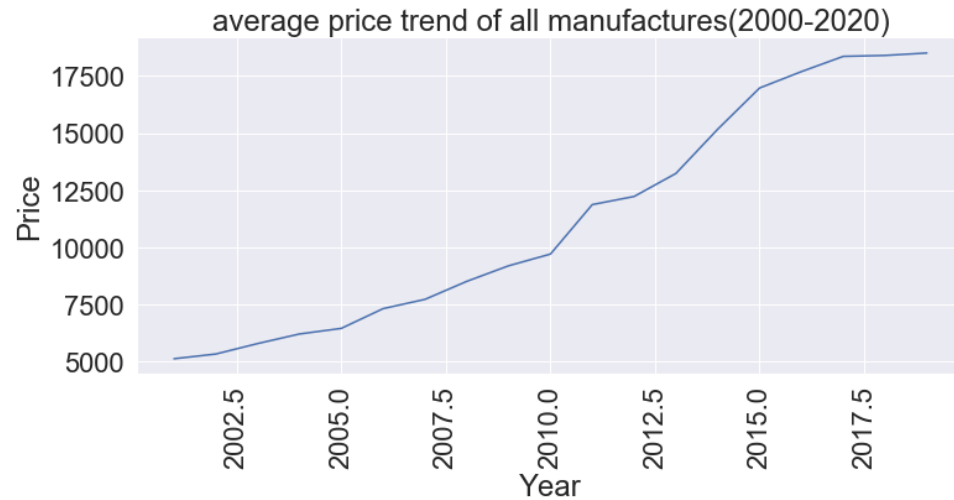
2000-2020

Ford  
Chevrolet Toyota  
Nissan  
Honda

# Data Exploration Analysis



# Year and Price

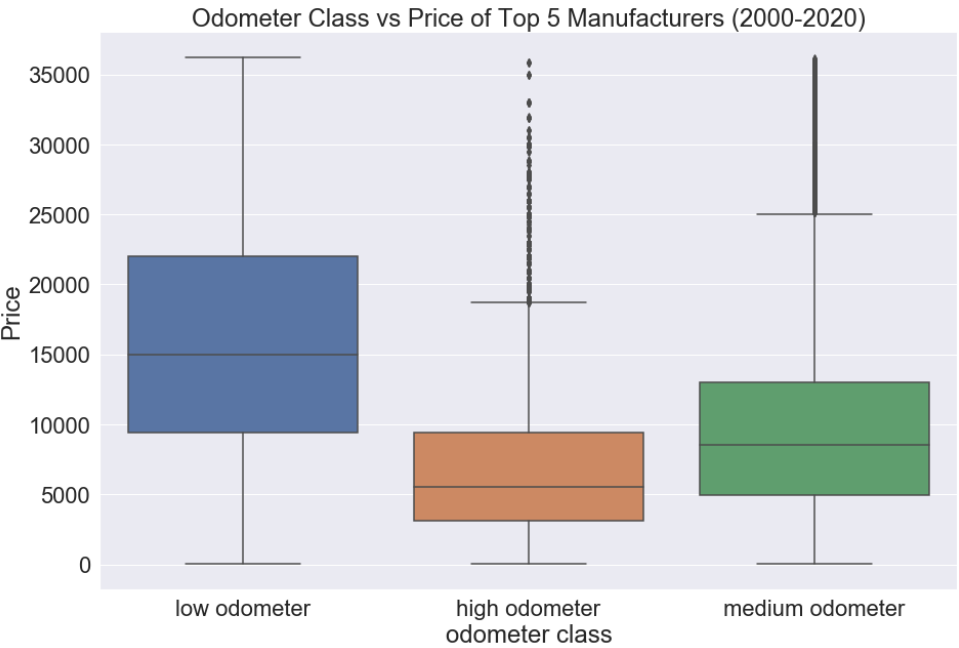


the average price of all used vehicles went up year by year. The increase speed is around 6% per year, which is a little higher than the inflation rate in America (3%). That's may be caused by online markets' developing. The used car price online is approaching to the physical dealer's gradually.

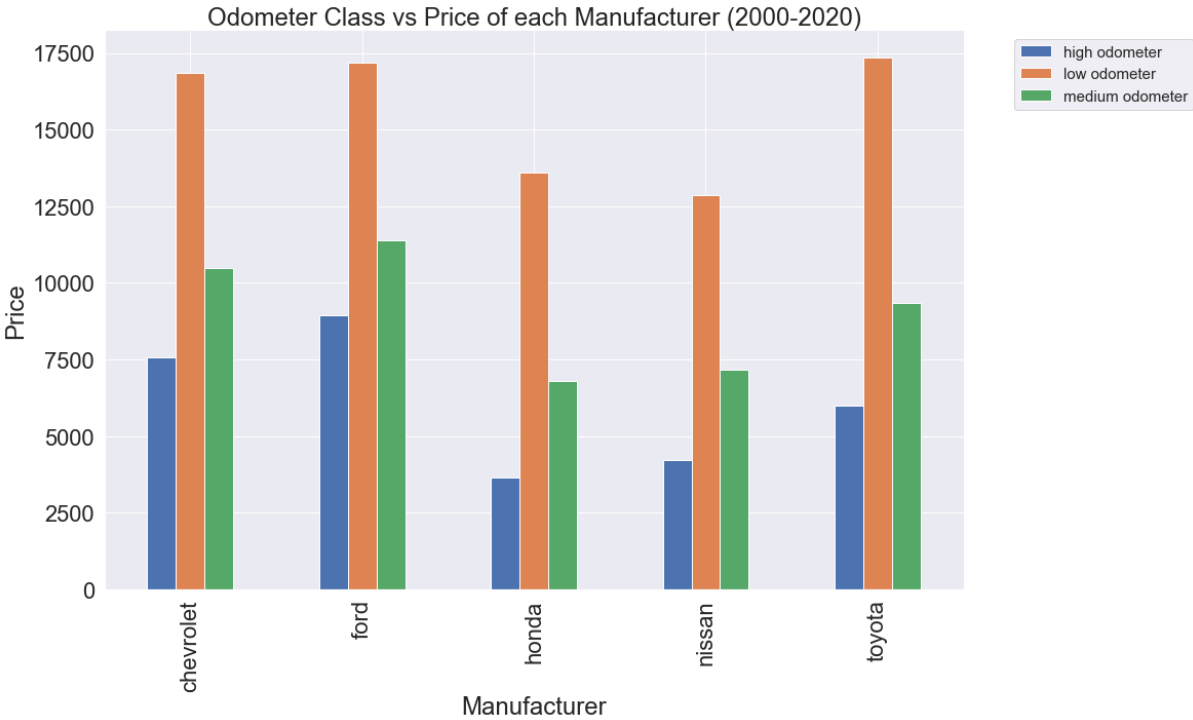
**Pearson value** is 0.5

There is a positive **linear** correlation between 'Year' and 'Price'.

# Odometer and Price



Vehicles with low odometers usually have a higher mean price as well as a wider price range.



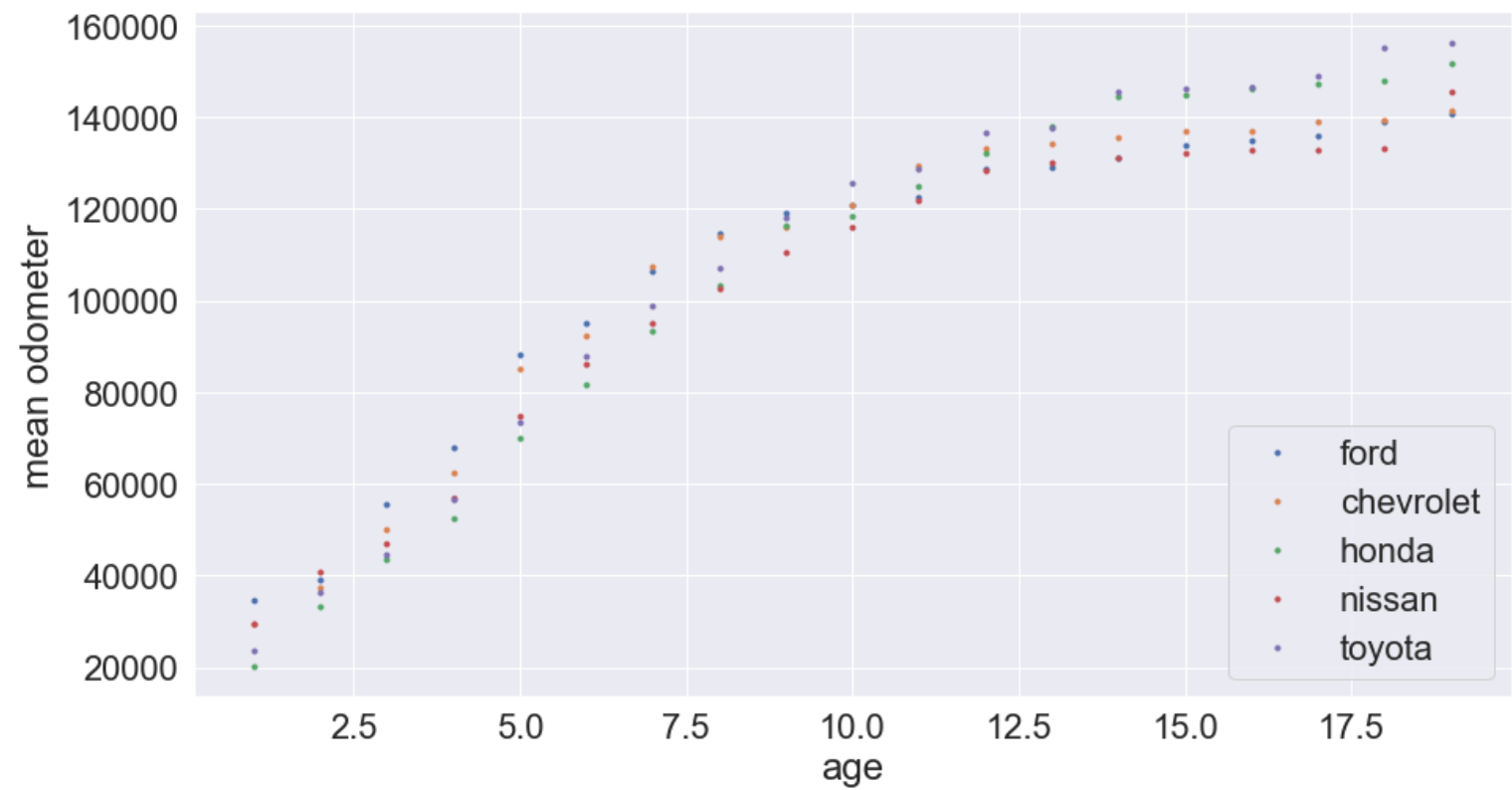
For different manufacturers, the relation between price and odometer is different. Toyota has the highest mean price of low kilometers(less than 10,000km). That's may be because of the statement that Japanese cars could held their value better than others in high odometer class.

**Ford Pearson value** is -0.18

It's not a linear correlation.



# Odometer and Year(Age)



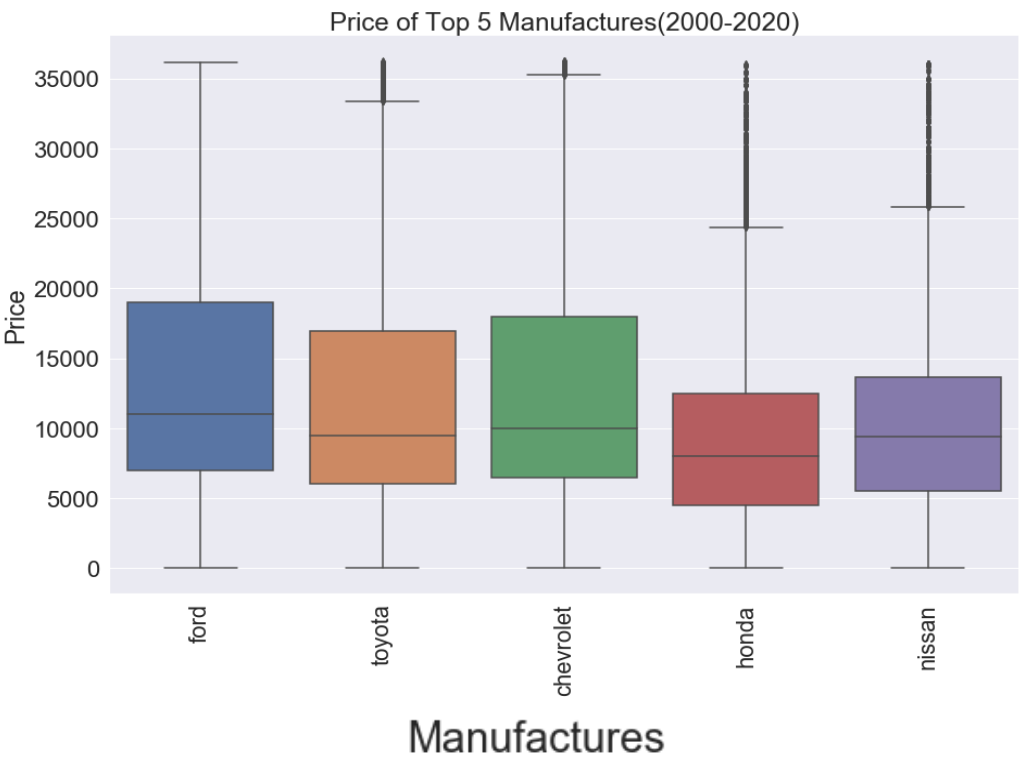
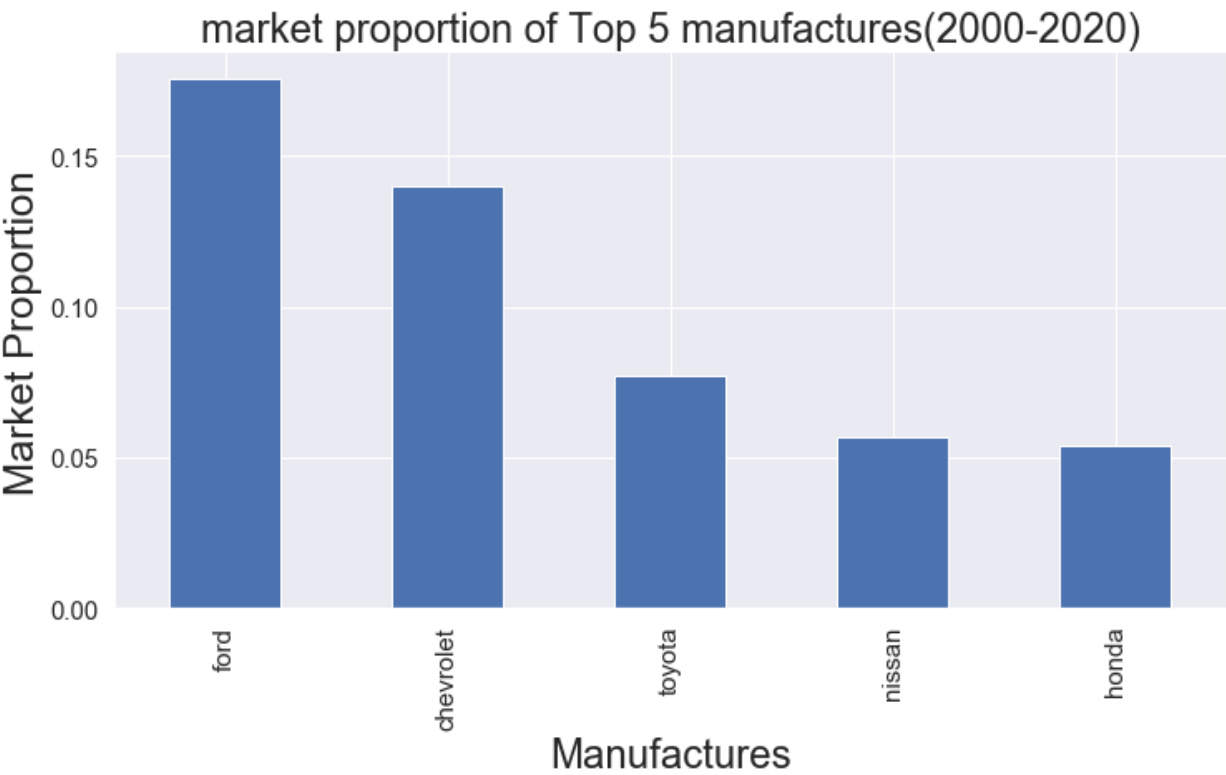
**Pearson Value** 0.5 indicate a positive linear relation between odometer and age.

What is remarkable is that the growth rate of odometers slowed down significantly among the vehicles which are over 12.5 years old.

# Manufacturers and Price

Combined with Market Proportion plot, it is obvious that their rankings are closely relevant.

**Ford** had the highest average price, followed by Chevrolet, Toyota, Nissan and Honda.

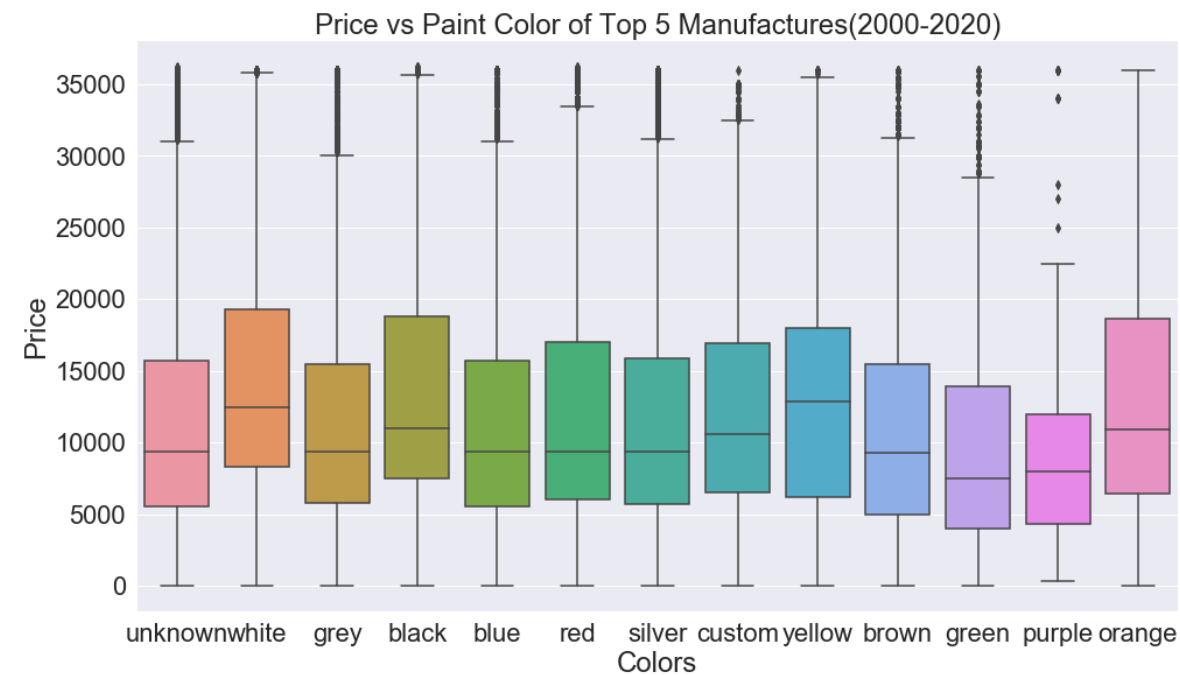
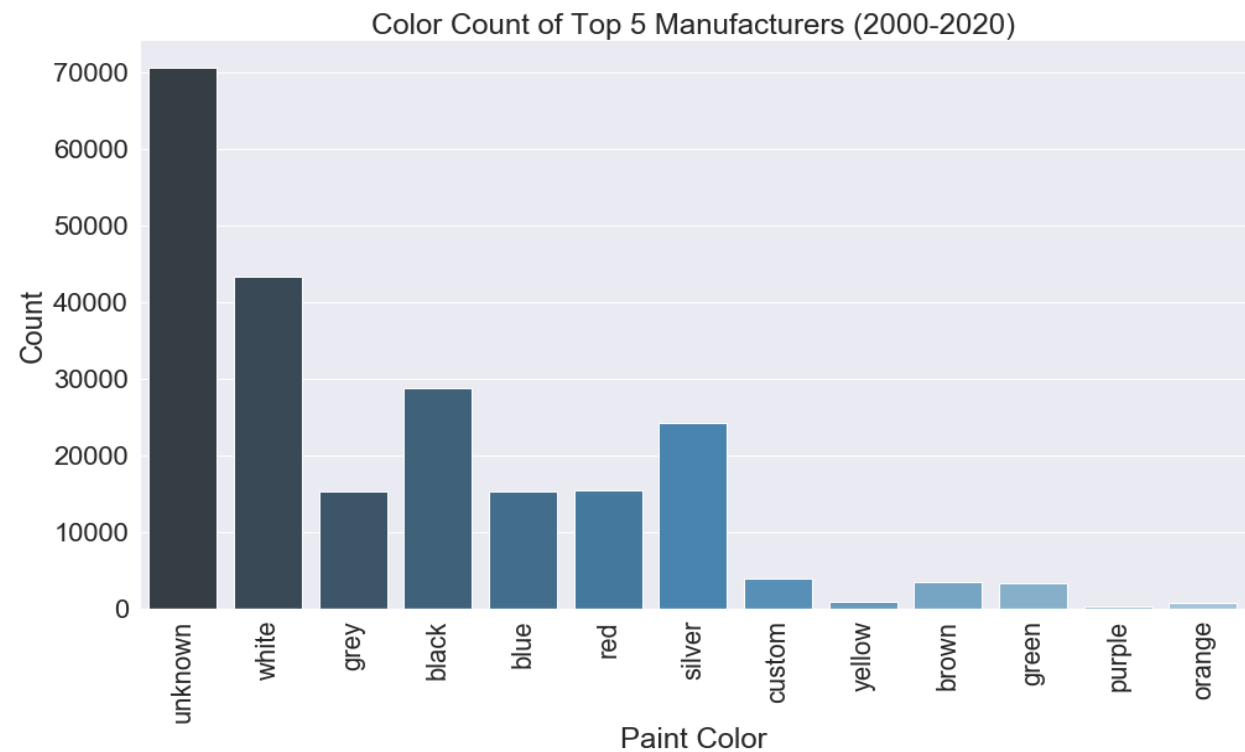


# Colors and Price

There are more than 13 species of paint color cars on sale during the past 20 years.

Top 3 colors which have the largest market proportion are white, black and silver. Probably because orange is a popular color for customers, however without enough supply in market.

**white, orange and black** usually could be sold at a higher price than other colors.



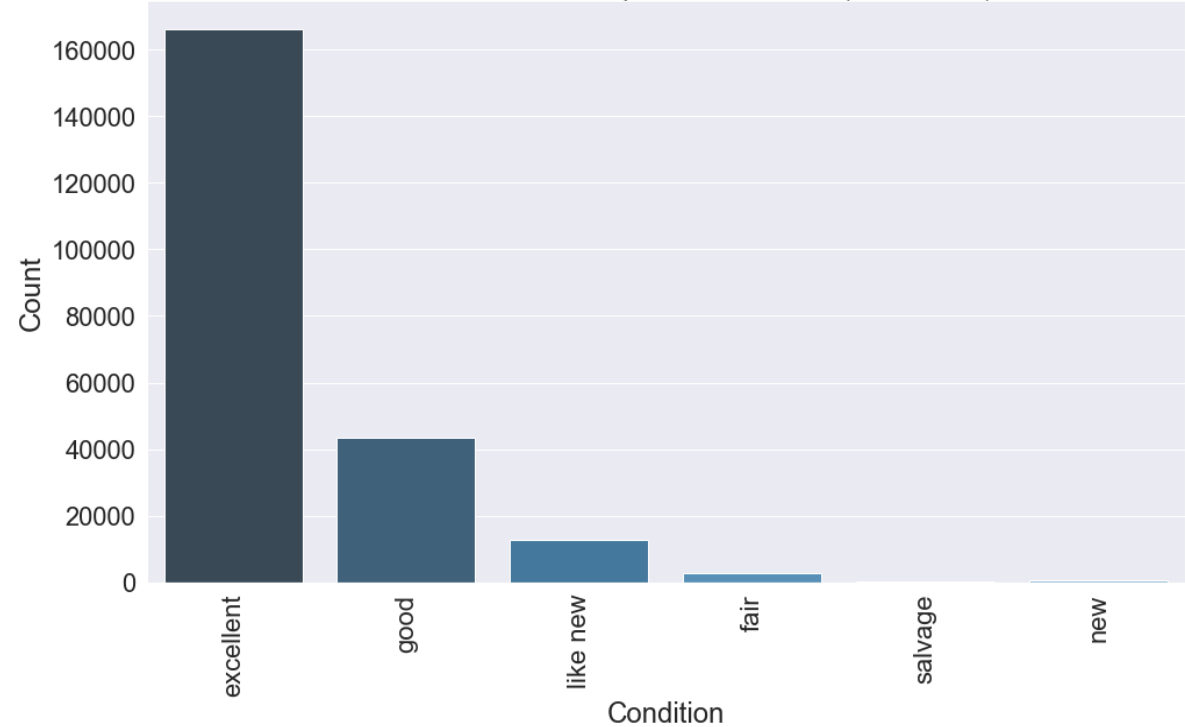
# Condition and Price

The vehicles labeled as ‘Like new’ and ‘excellent’ condition have the highest median prices.

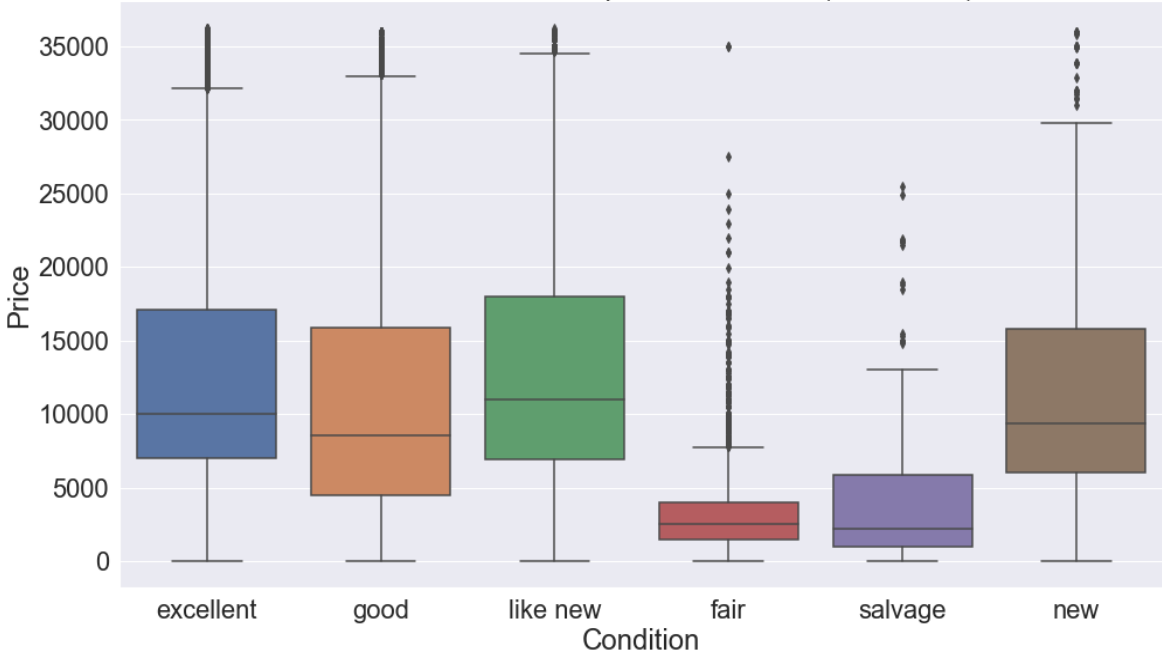
The ‘fair’ and ‘salvage’ ones have the smallest median prices as well as their price range.

Most vehicles on sale are described as **‘excellent condition’**.

Condition Counts of Top 5 Manufacturer (2000-2020)



Condition vs Price of Top 5 Manufactures(2000-2020)



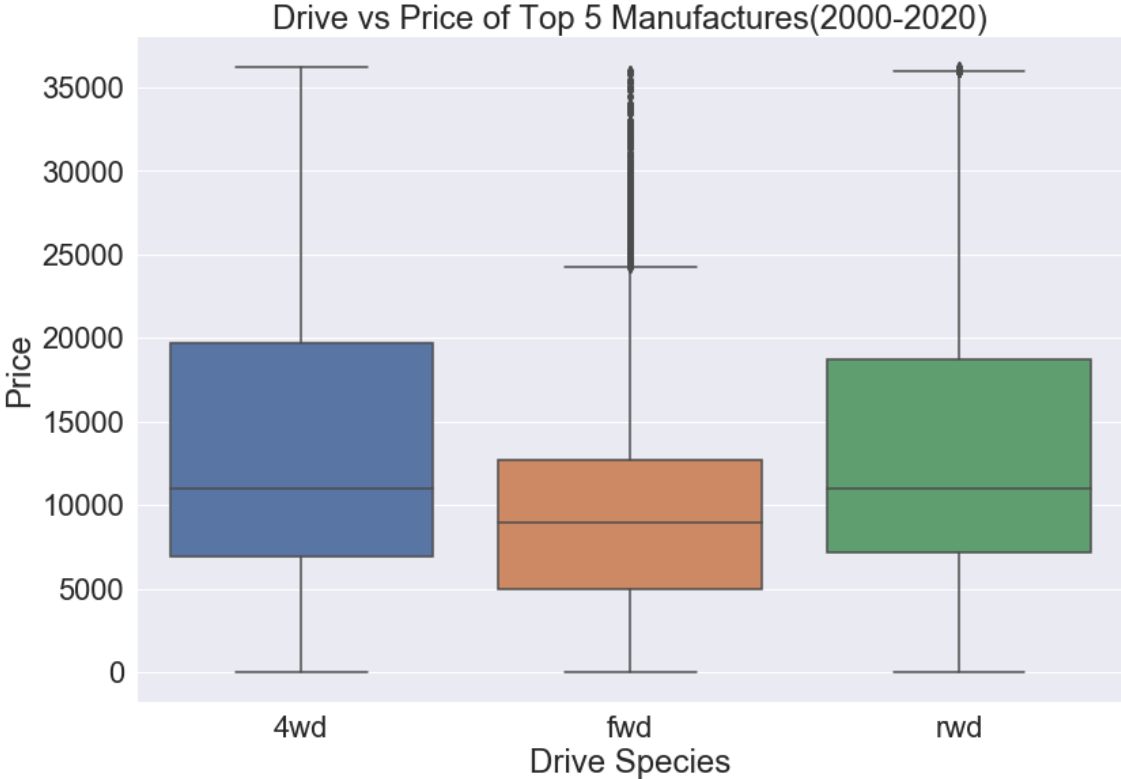
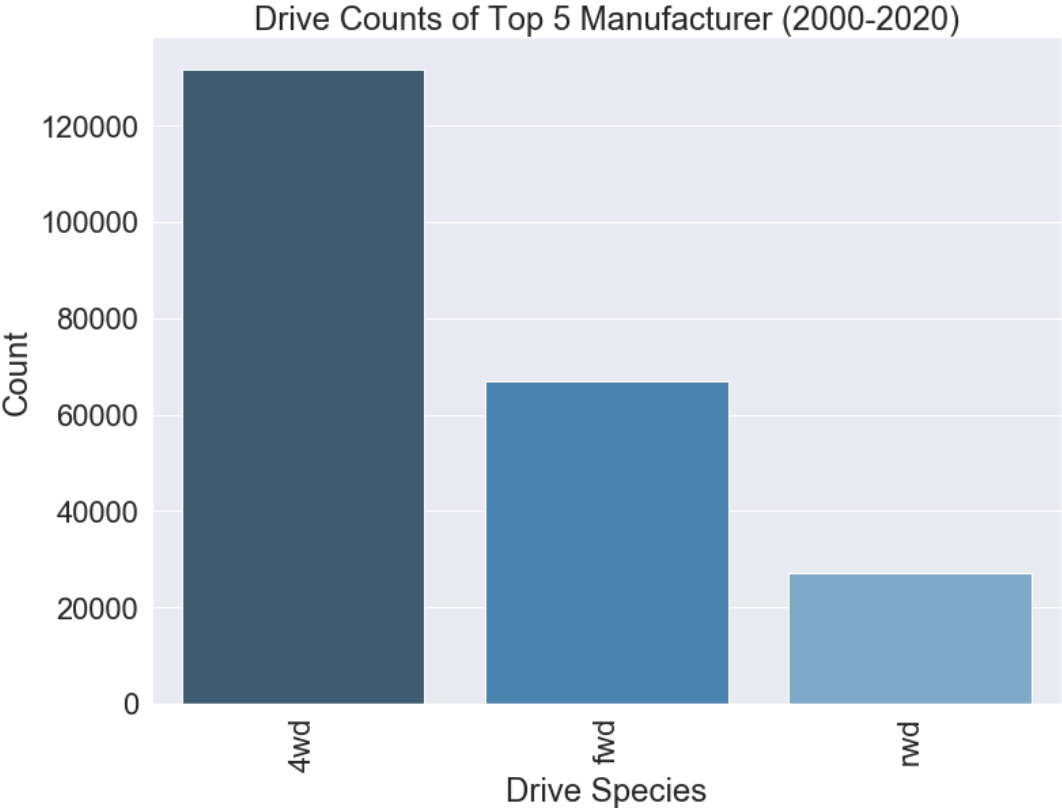


# Drive and Price

**4wd trucks** seem to have a bigger share of the market, which could explain why Ford and Chevrolet are so popular.

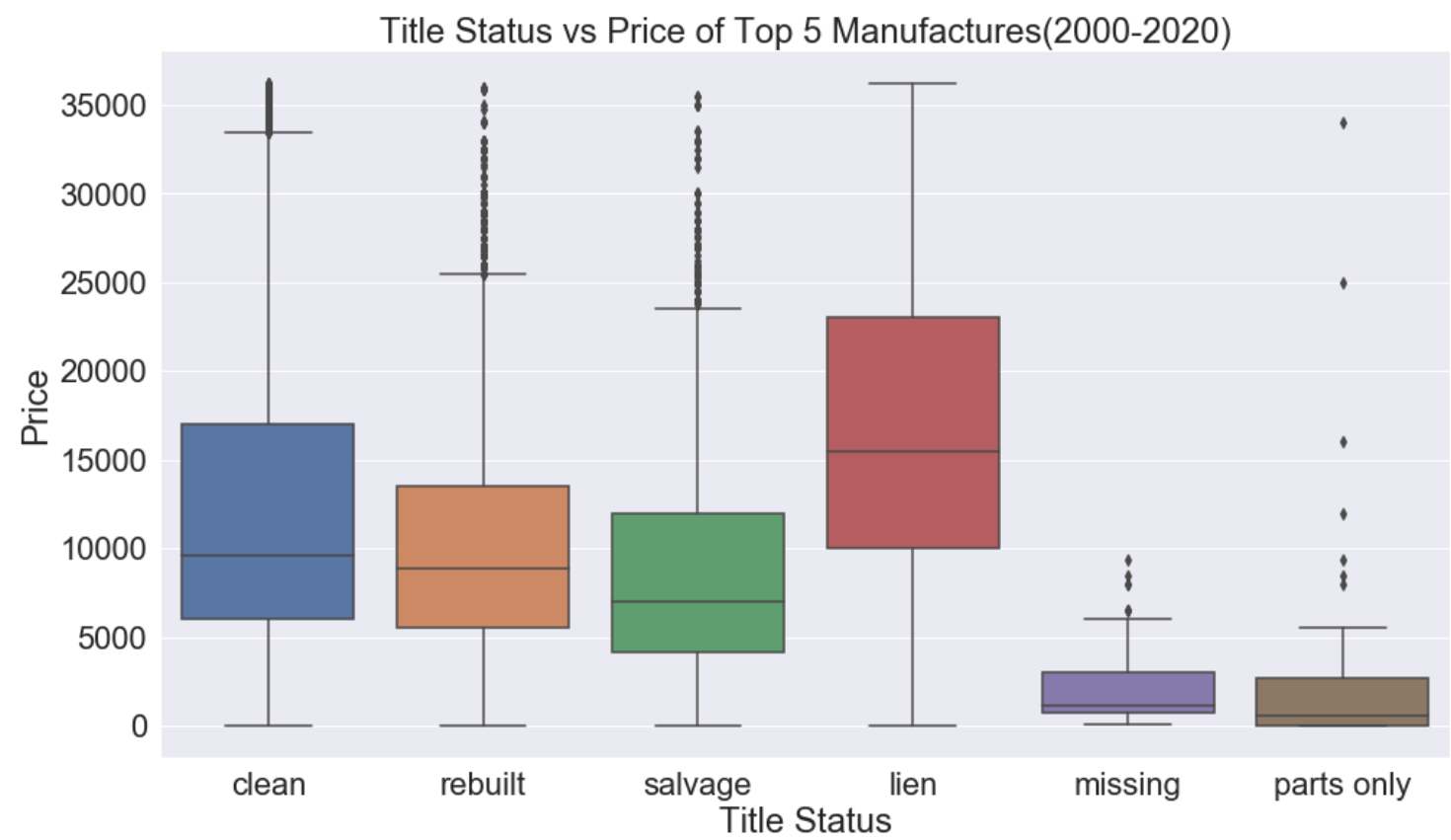
**4wd vehicles** always have a higher median price than other two categories.

**4wd** shares the largest proportion of the market.



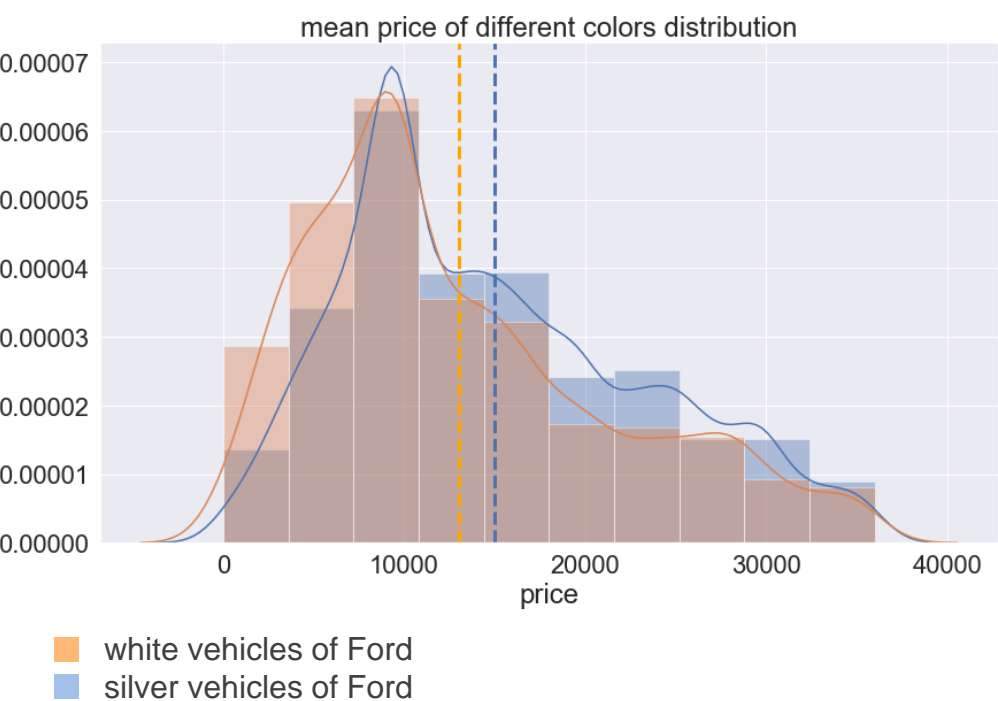
# Title Status and Price

**Lien vehicles** have the highest median price of all, which also have the largest price range..



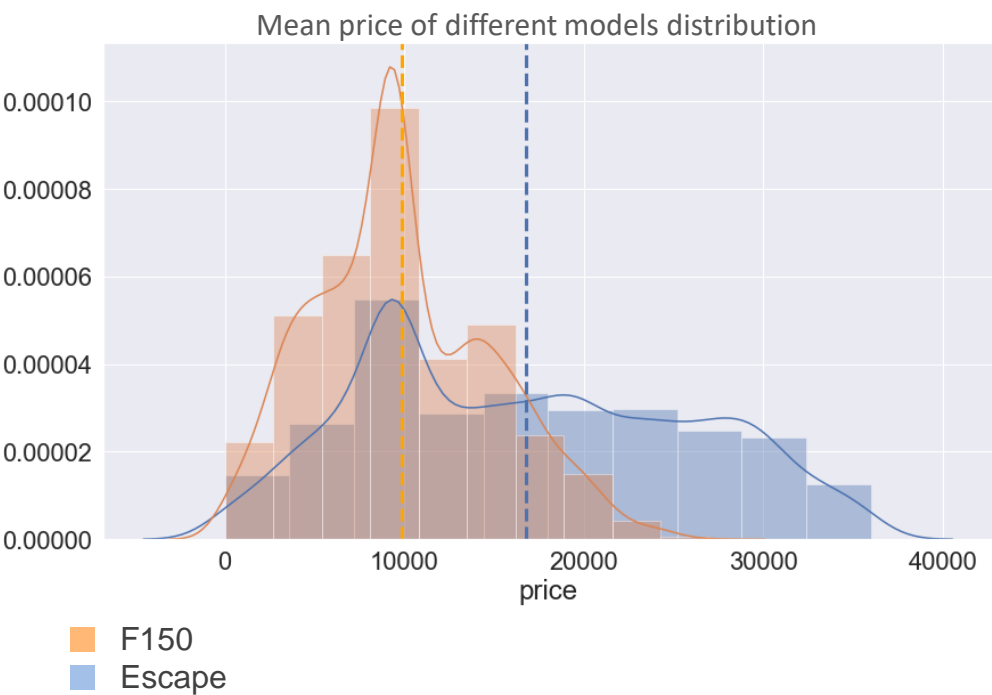
# Statistical Data Analysis

## Hypothesis Testing



H0: Samples from different colors share the same mean price  
p-value is 0.039

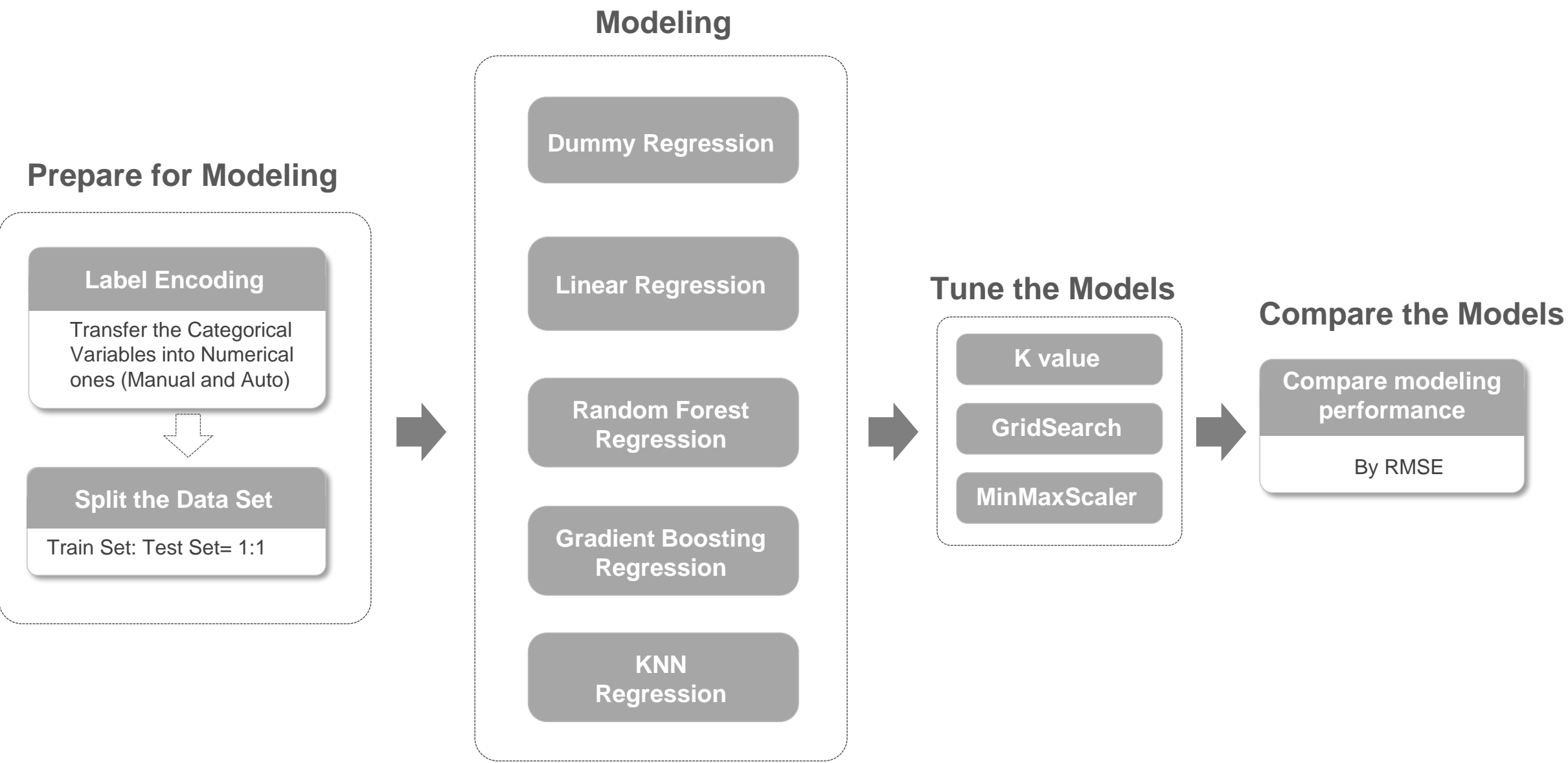
**Reject H0:** Samples from different colors share the different mean prices



H0: Samples from different models(F150 and Escape) share the same mean price  
p-value is 0.675

**Fail to reject H0:** There is no proof that samples from the two models share the same mean price.

# Modeling Ideas





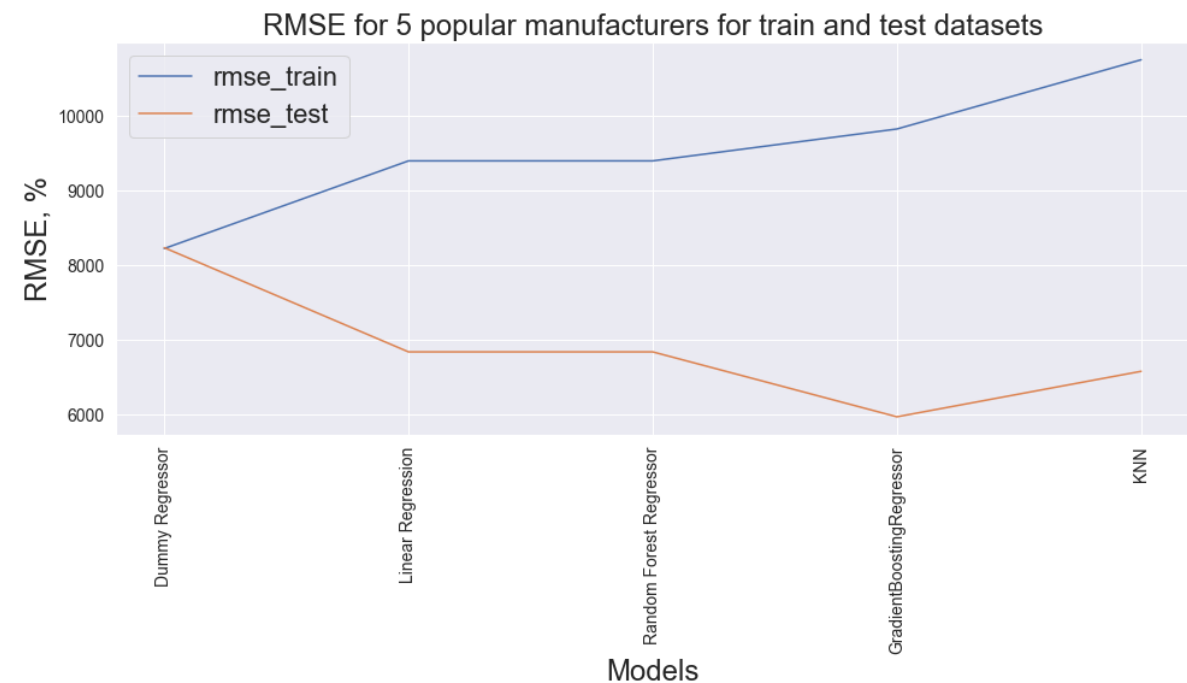
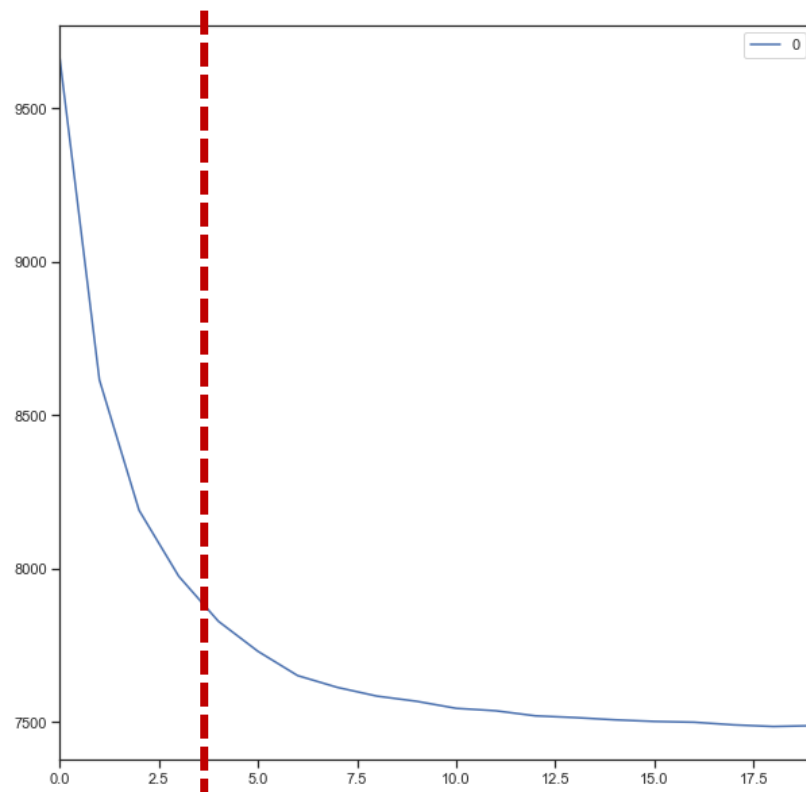
# Compare Models

Out of 5 models, the best models by the RMSE is Gradient Boosting Regression.

Determined by Random Forest modeling, price is more related with Manufacturers, Odometers, Condition, Paint Color.

Out of 24 features, we used only 7 features for the best model.

For KNN regression, the best K value should be 3.





# Thank You!

Yang Fei  
Email: [sophia.fei0302@gmail.com](mailto:sophia.fei0302@gmail.com)  
<https://www.linkedin.com/in/yang-fei-1a862194/>  
<https://github.com/fysophia0302>