

---

# Used Car Price Prediction

Data science and Artificial Intelligence Society

Jiho Lee, Sungjun Choi, Earl Chae, Haejeong Cho



# Predicting the Price of Used Cars

1. Data Extraction
2. Feature Engineering
3. Model Performance Analysis

<https://github.com/jayhoneylee527/Used-Car-Price-Estimator>

---

# Data Extraction

- Data scraped from **Cargurus.com**
- Used “uszipcode” package to get zip codes all over Illinois
- Extracted used car data from Cargus at different region in Illinois by randomly selecting zip codes.

```
"""
We want to scrape used car listings in IL, so scrape zipcodes all over IL by population in descending order
"""

zipcodes = []
state_list = ['IL'] |

# Gather zipcodes from each state
search = SearchEngine()
for st in state_list:
    res = search.by_state(state=st, returns = 0) # Return all zipcodes in each state
    for r in res:
        zipcodes.append(r.zipcode)

# Shuffle zipcode list - For randomization
random.shuffle(zipcodes)
```

# Data

- We successfully extracted 10,942 used car data
- There are 30 columns

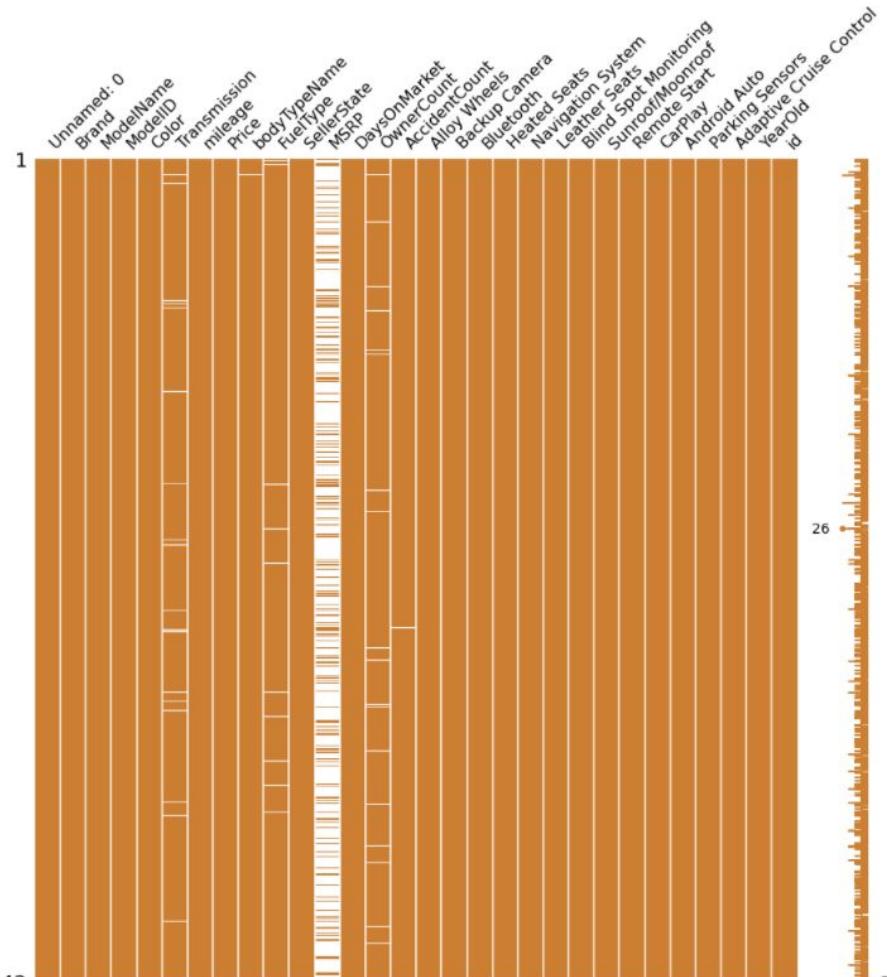
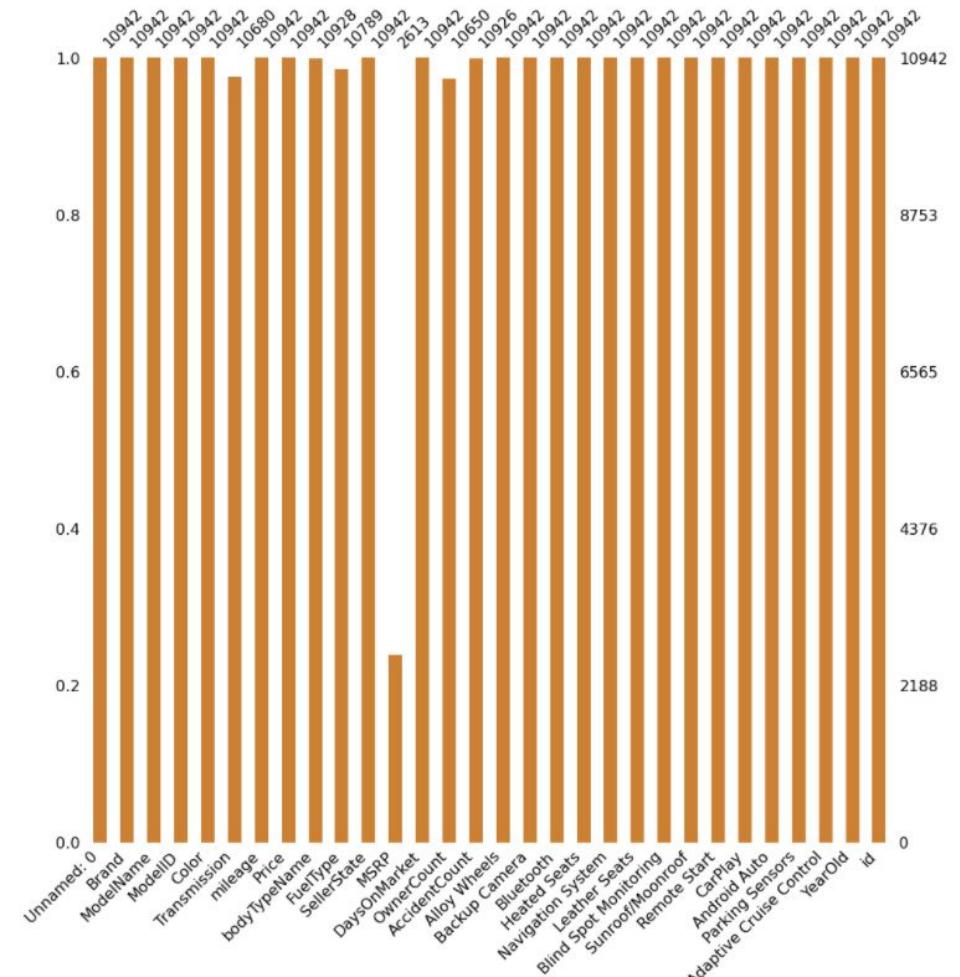
```
spec_tbl_df [10,942 x 30] (s3: spec_tbl_df/tbl_df/tbl/data.frame)
$ ...1 : num [1:10942] 0 1 2 3 4 5 6 7 8 9 ...
$ Brand : chr [1:10942] "Ford" "Chrysler" "Toyota" "RAM" ...
$ ModelName : chr [1:10942] "F-150" "Pacifica" "4Runner" "2500" ...
$ ModelID : chr [1:10942] "d337" "d177" "d290" "d2102" ...
$ Color : chr [1:10942] "GRAY" "GRAY" "BLUE" "BLACK" ...
$ Transmission : chr [1:10942] "Automatic" "Automatic" "Automatic" "Automatic" ...
$ mileage : num [1:10942] 25111 94774 179200 26438 93724 ...
$ Price : num [1:10942] 36990 20990 21990 53580 25995 ...
$ bodyTypeName : chr [1:10942] "Pickup Truck" "Minivan" "SUV / Crossover" "Pickup Truck" ...
$ FuelType : chr [1:10942] "Flex Fuel Vehicle" "Gasoline" "Gasoline" "Gasoline" ...
$ SellerState : chr [1:10942] "IL" "IL" "IL" "IL" ...
$ MSRP : num [1:10942] 37990 20990 NA NA NA ...
$ DaysOnMarket : num [1:10942] 211 11 203 37 9 232 68 23 69 23 ...
$ OwnerCount : num [1:10942] 1 1 2 2 1 1 1 1 1 1 ...
$ AccidentCount : num [1:10942] 0 1 0 0 0 0 0 0 0 0 ...
$ Alloy Wheels : num [1:10942] 1 1 0 0 1 1 1 0 1 ...
$ Backup Camera : num [1:10942] 1 1 1 0 1 1 1 1 1 ...
$ Bluetooth : num [1:10942] 1 1 1 1 1 1 1 0 1 1 ...
$ Heated Seats : num [1:10942] 1 0 1 1 0 1 1 0 0 1 ...
$ Navigation System : num [1:10942] 0 0 1 0 0 1 1 0 1 1 ...
$ Leather Seats : num [1:10942] 0 0 1 0 0 1 1 0 0 1 ...
$ Blind Spot Monitoring : num [1:10942] 0 0 0 1 0 1 0 0 1 1 ...
$ Sunroof/Moonroof : num [1:10942] 0 0 1 0 0 0 0 0 1 0 ...
$ Remote Start : num [1:10942] 1 1 0 1 0 1 1 1 1 0 ...
$ CarPlay : num [1:10942] 1 0 0 0 0 0 0 0 1 1 ...
$ Android Auto : num [1:10942] 0 0 0 0 0 0 0 0 1 1 ...
$ Parking Sensors : num [1:10942] 0 0 0 1 0 1 1 0 0 0 ...
$ Adaptive Cruise Control: num [1:10942] 0 0 0 0 0 1 0 0 0 0 ...
$ Yearold : num [1:10942] 4 5 10 3 5 4 5 2 3 5 ...
$ id : num [1:10942] 0 1 2 3 4 5 6 7 8 9 ...
```

---

# Checking Null Value

- There are 6 columns that has Null Value
- MSRP(manufacturer's suggested retail price) has 76.12% Null value

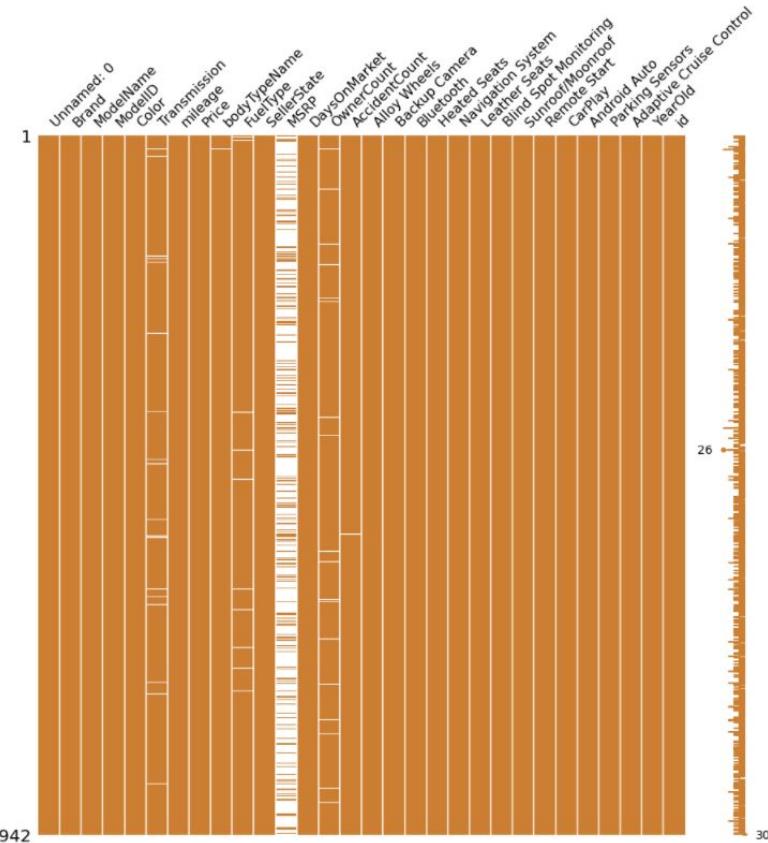
column: Unnamed: 0	Percent of NaN value: 0.00%
column: Brand	Percent of NaN value: 0.00%
column: ModelName	Percent of NaN value: 0.00%
column: ModelID	Percent of NaN value: 0.00%
column: Color	Percent of NaN value: 0.00%
column: Transmission	Percent of NaN value: 2.39%
column: mileage	Percent of NaN value: 0.00%
column: Price	Percent of NaN value: 0.00%
column: bodyTypeName	Percent of NaN value: 0.13%
column: FuelType	Percent of NaN value: 1.40%
column: SellerState	Percent of NaN value: 0.00%
column: MSRP	Percent of NaN value: 76.12%
column: DaysOnMarket	Percent of NaN value: 0.00%
column: OwnerCount	Percent of NaN value: 2.67%
column: AccidentCount	Percent of NaN value: 0.15%
column: Alloy Wheels	Percent of NaN value: 0.00%
column: Backup Camera	Percent of NaN value: 0.00%
column: Bluetooth	Percent of NaN value: 0.00%
column: Heated Seats	Percent of NaN value: 0.00%
column: Navigation System	Percent of NaN value: 0.00%
column: Leather Seats	Percent of NaN value: 0.00%
column: Blind Spot Monitoring	Percent of NaN value: 0.00%
column: Sunroof/Moonroof	Percent of NaN value: 0.00%
column: Remote Start	Percent of NaN value: 0.00%
column: CarPlay	Percent of NaN value: 0.00%
column: Android Auto	Percent of NaN value: 0.00%
column: Parking Sensors	Percent of NaN value: 0.00%
column: Adaptive Cruise Control	Percent of NaN value: 0.00%
column: YearOld	Percent of NaN value: 0.00%
column: id	Percent of NaN value: 0.00%



---

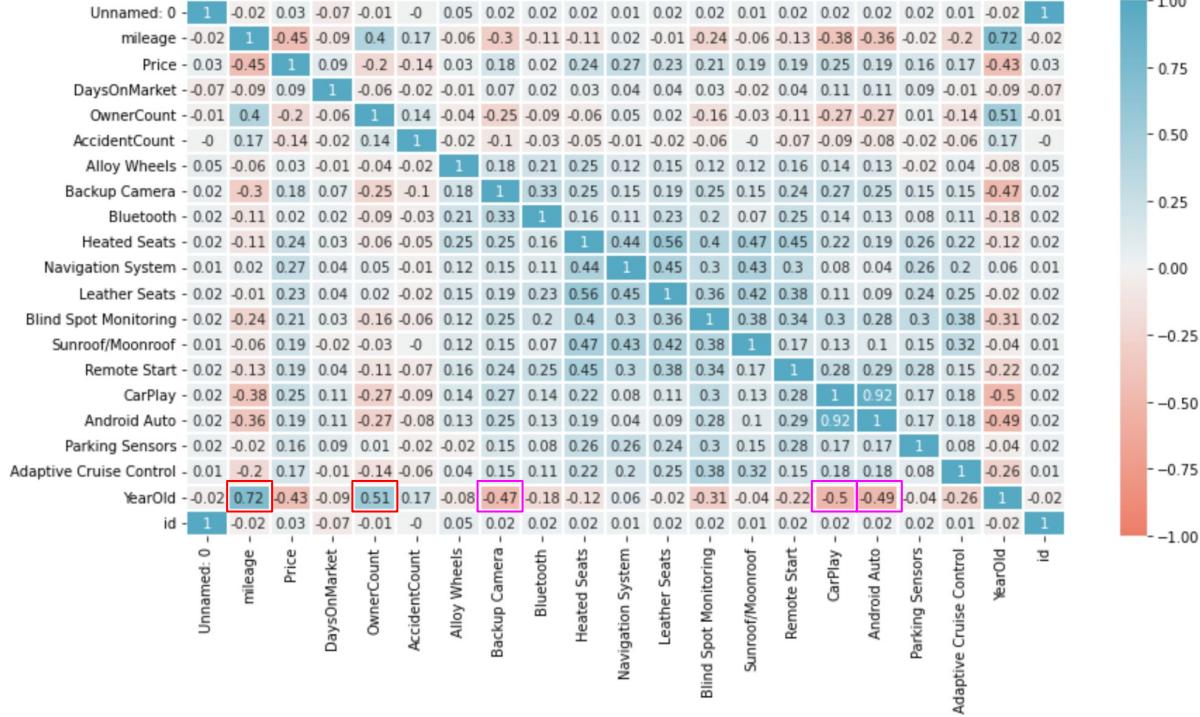
# Recap from Week 1

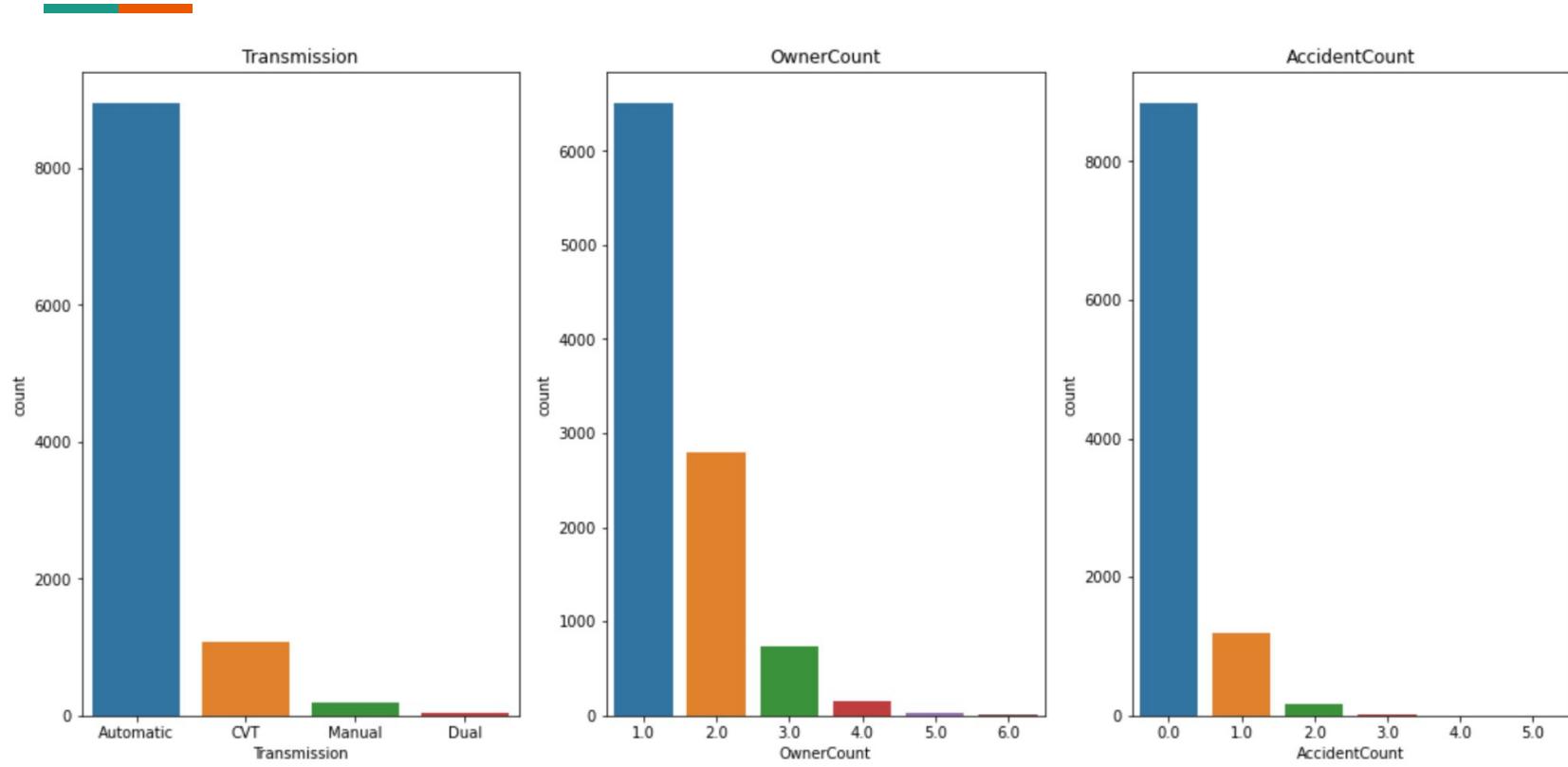
1. Data extracted from Cargurus.com using web scraping method
2. Checked Null Data
  - a. 76.12% of null data for MSRP  
(manufacturer's suggested retail price)

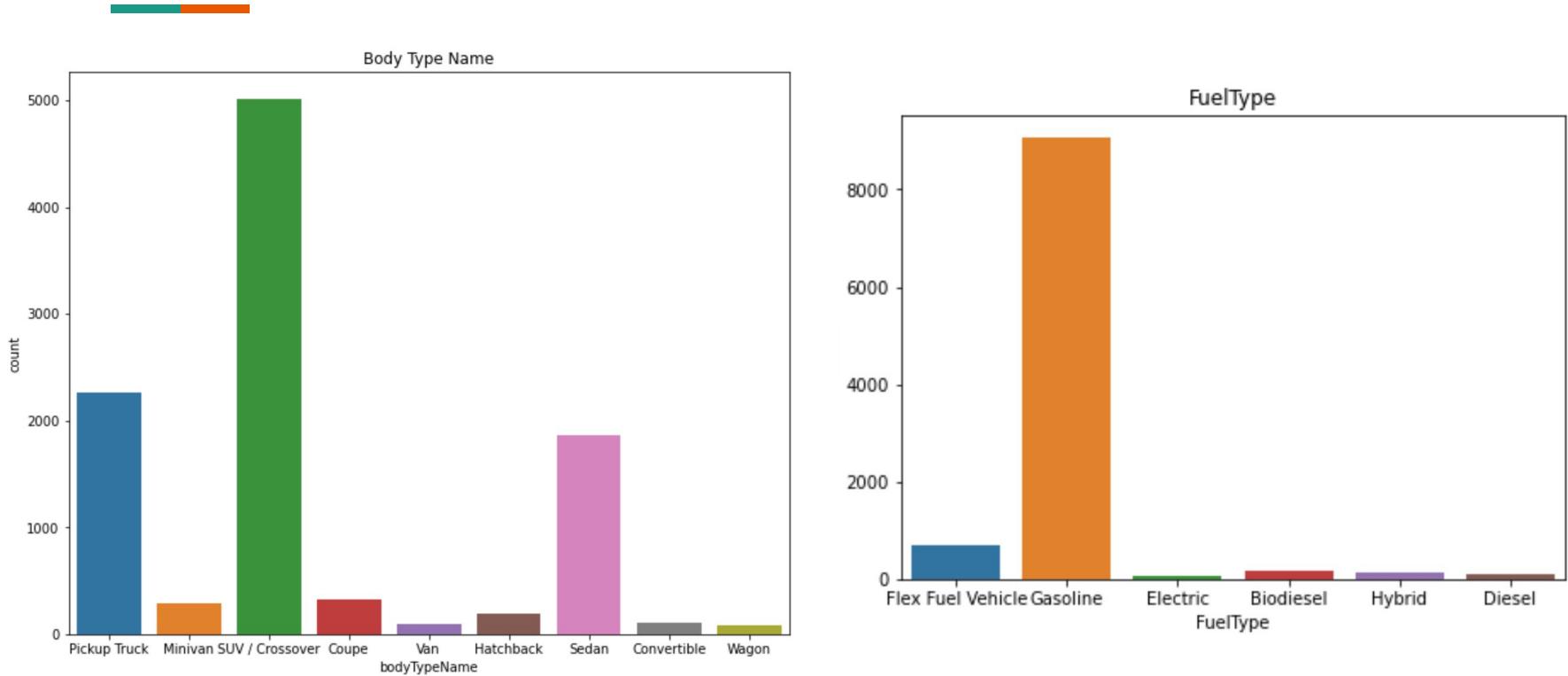


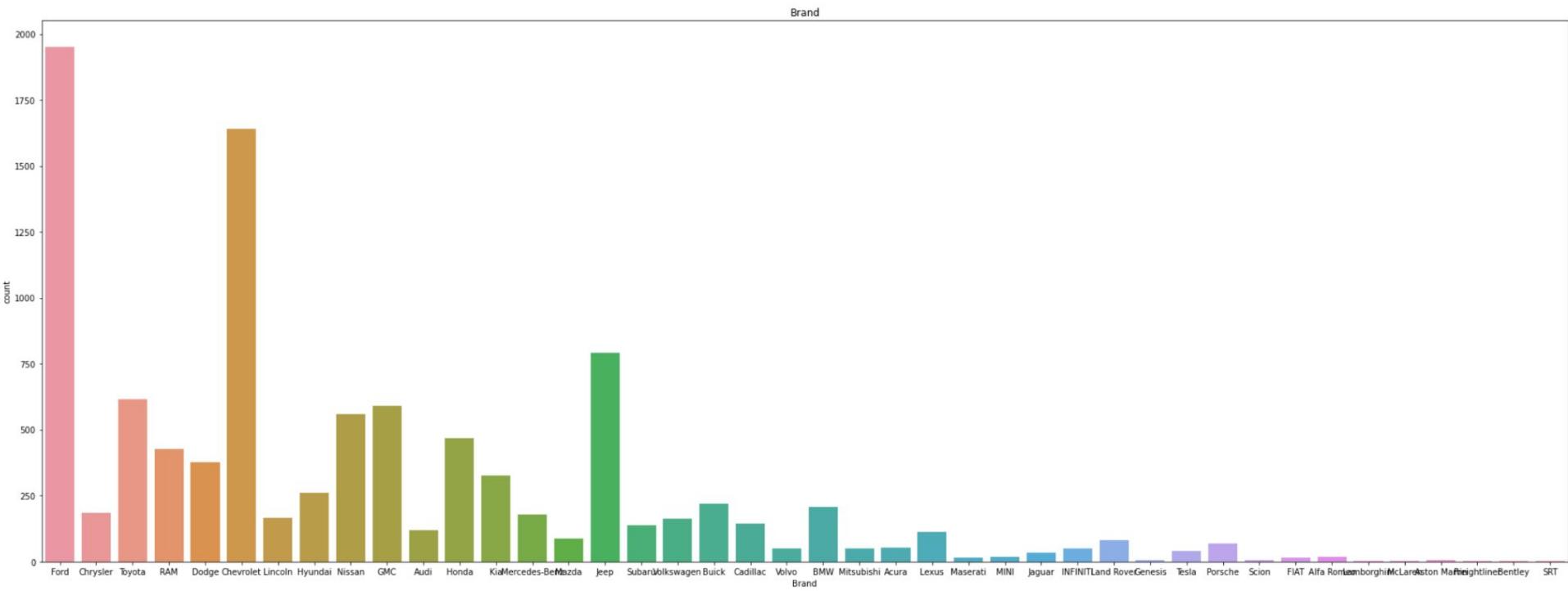
# EDA

Correlation Matrix









# Data Cleaning

- Dropped MSRP and all the rows that have null value.
  - Dropped 703 entry from 10,942 u: car data
- 0 % null value for all the columns after data cleaning

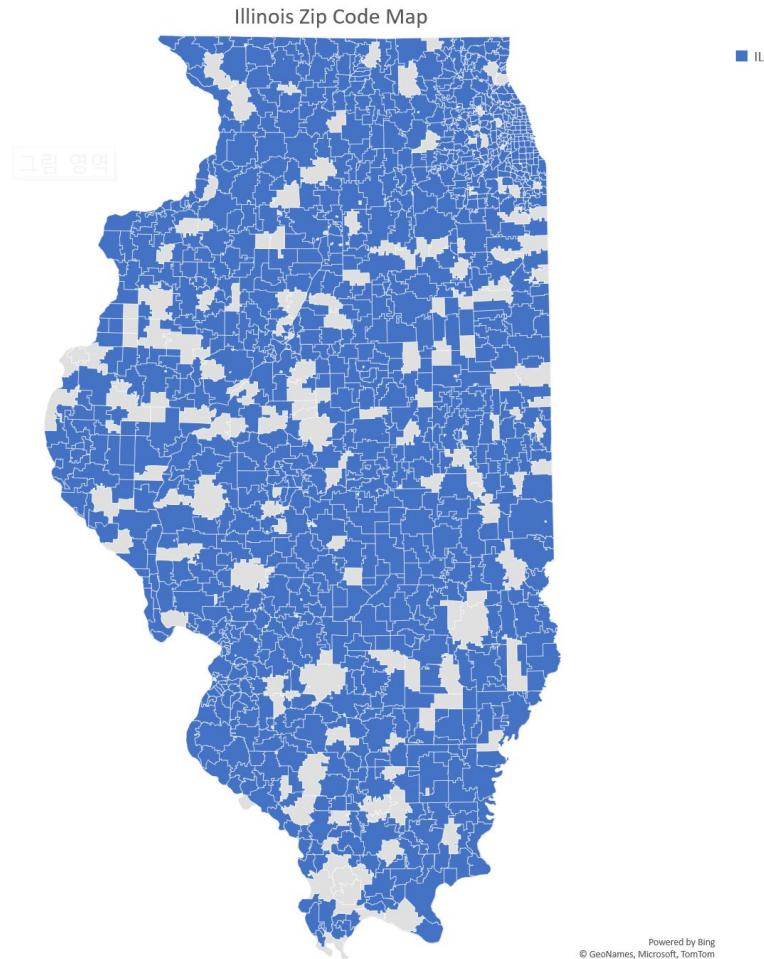
```
column: Unnamed: 0 Percent of NaN value: 0.00%
column: Brand Percent of NaN value: 0.00%
column: ModelName Percent of NaN value: 0.00%
column: ModelID Percent of NaN value: 0.00%
column: Color Percent of NaN value: 0.00%
column: Transmission Percent of NaN value: 2.39%
column: mileage Percent of NaN value: 0.00%
column: Price Percent of NaN value: 0.00%
column: bodyTypeName Percent of NaN value: 0.13%
column: FuelType Percent of NaN value: 1.40%
column: SellerState Percent of NaN value: 0.00%
column: MSRP Percent of NaN value: 76.12%
column: DaysOnMarket Percent of NaN value: 0.00%
column: OwnerCount Percent of NaN value: 2.67%
column: AccidentCount Percent of NaN value: 0.15%
column: Alloy Wheels Percent of NaN value: 0.00%
column: Backup Camera Percent of NaN value: 0.00%
column: Bluetooth Percent of NaN value: 0.00%
column: Heated Seats Percent of NaN value: 0.00%
column: Navigation System Percent of NaN value: 0.00%
column: Leather Seats Percent of NaN value: 0.00%
column: Blind Spot Monitoring Percent of NaN value: 0.00%
column: Sunroof/Moonroof Percent of NaN value: 0.00%
column: Remote Start Percent of NaN value: 0.00%
column: CarPlay Percent of NaN value: 0.00%
column: Android Auto Percent of NaN value: 0.00%
column: Parking Sensors Percent of NaN value: 0.00%
column: Adaptive Cruise Control Percent of NaN value: 0.00%
column: YearOld Percent of NaN value: 0.00%
column: id Percent of NaN value: 0.00%
```

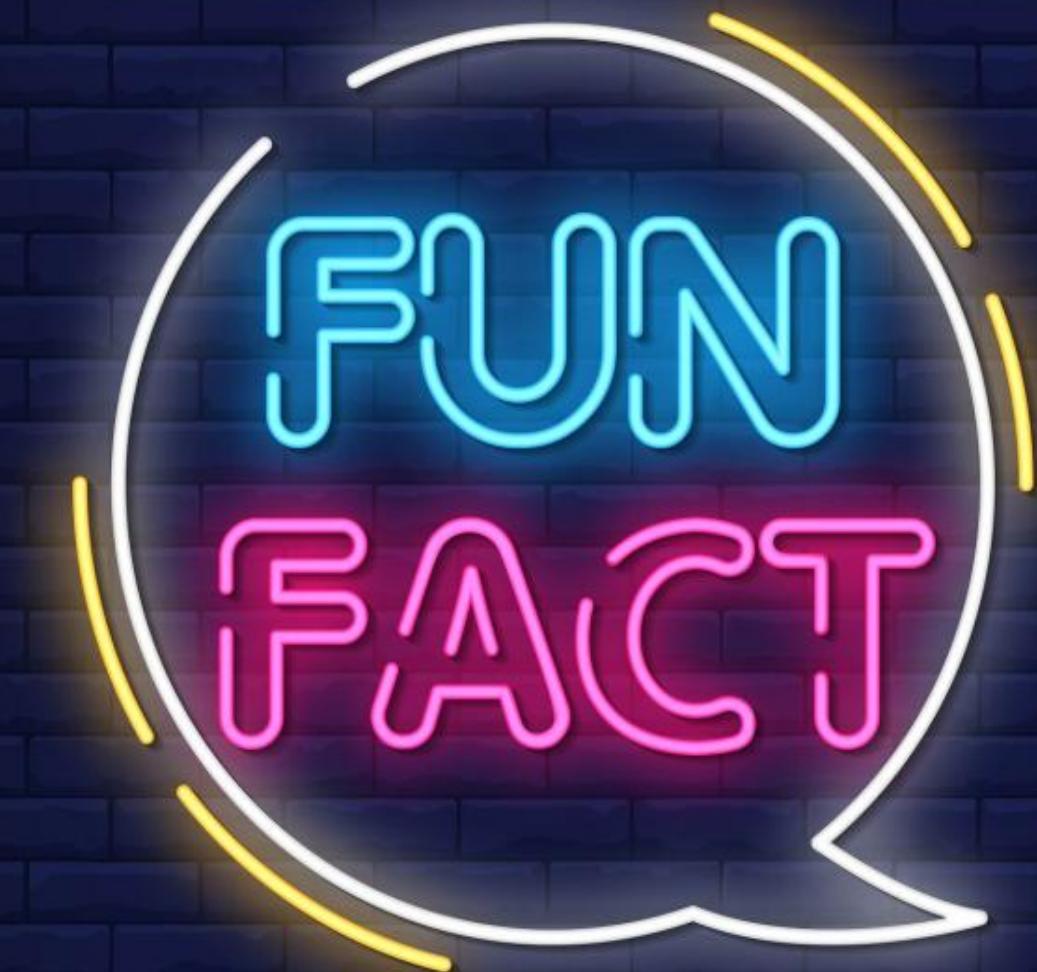


```
column: Unnamed: 0 Percent of NaN value: 0.00%
column: Brand Percent of NaN value: 0.00%
column: ModelName Percent of NaN value: 0.00%
column: ModelID Percent of NaN value: 0.00%
column: Color Percent of NaN value: 0.00%
column: Transmission Percent of NaN value: 0.00%
column: mileage Percent of NaN value: 0.00%
column: Price Percent of NaN value: 0.00%
column: bodyTypeName Percent of NaN value: 0.00%
column: FuelType Percent of NaN value: 0.00%
column: SellerState Percent of NaN value: 0.00%
column: DaysOnMarket Percent of NaN value: 0.00%
column: OwnerCount Percent of NaN value: 0.00%
column: AccidentCount Percent of NaN value: 0.00%
column: Alloy Wheels Percent of NaN value: 0.00%
column: Backup Camera Percent of NaN value: 0.00%
column: Bluetooth Percent of NaN value: 0.00%
column: Heated Seats Percent of NaN value: 0.00%
column: Navigation System Percent of NaN value: 0.00%
column: Leather Seats Percent of NaN value: 0.00%
column: Blind Spot Monitoring Percent of NaN value: 0.00%
column: Sunroof/Moonroof Percent of NaN value: 0.00%
column: Remote Start Percent of NaN value: 0.00%
column: CarPlay Percent of NaN value: 0.00%
column: Android Auto Percent of NaN value: 0.00%
column: Parking Sensors Percent of NaN value: 0.00%
column: Adaptive Cruise Control Percent of NaN value: 0.00%
column: YearOld Percent of NaN value: 0.00%
column: id Percent of NaN value: 0.00%
```

# Zip Codes Mapping

- This is the map of Illinois zip codes that we randomly selected
- Used car datas are from these zip codes



A neon sign is mounted on a dark blue brick wall. The sign features the words "FUN" in blue and "FACT" in pink, both enclosed within a white-outlined speech bubble shape. The neon tubes are glowing brightly against the dark background.

**FUN**  
**FACT**

Brand	ModelName	ModelID	Color	Transmission	mileage	Price	bodyTypeName	FuelType	SellerState	MSRP	DaysOnMarket	OwnerCount
-------	-----------	---------	-------	--------------	---------	-------	--------------	----------	-------------	------	--------------	------------

Ford	Focus	d346	GRAY		NaN	168252.0	3990.0	Sedan	Gasoline	IN	NaN	17	2.0
AccidentCount	Alloy Wheels	Backup Camera	Bluetooth	Heated Seats	Navigation System	Leather Seats	Blind Spot Monitoring	Sunroof/Moonroof	Remote Start	CarPlay	Android Auto	Parking Sensors	Adaptive Cruise Control
0.0	1	0	1	0	1	0	0	0	0	0	0	0	0



YearOld

11

Brand	ModelName	ModelID	Color	Transmission	mileage	Price	bodyTypeName	FuelType	SellerState	MSRP	DaysOnMarket	OwnerCount	YearOld	
Lamborghini	Aventador	d2125	ORANGE	Automatic	13077.0	450000.0	Coupe	Gasoline	IL	NaN	121	2.0	6	Adaptive Cruise Control
AccidentCount	Alloy Wheels	Backup Camera	Bluetooth	Heated Seats	Navigation System	Leather Seats	Blind Spot Monitoring	Sunroof/Moonroof	Remote Start	CarPlay	Android Auto	Parking Sensors	0.0	0



JIDD MOTORS.COM  
1313 RAND RD, DES PLAINES, IL 60016



LEATHER  
INTERIOR

AWD



NAVIGATION



POWER  
WINDOWS



POWER  
DOOR  
LOCKS

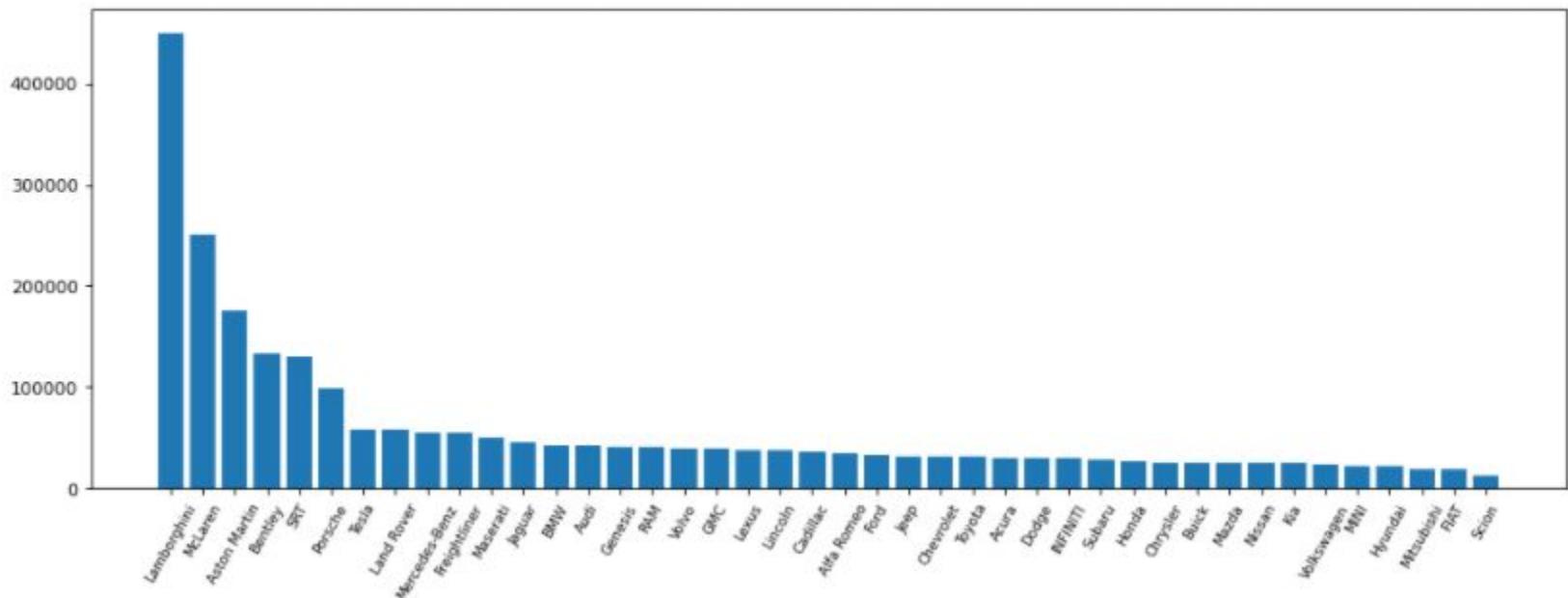


A/C

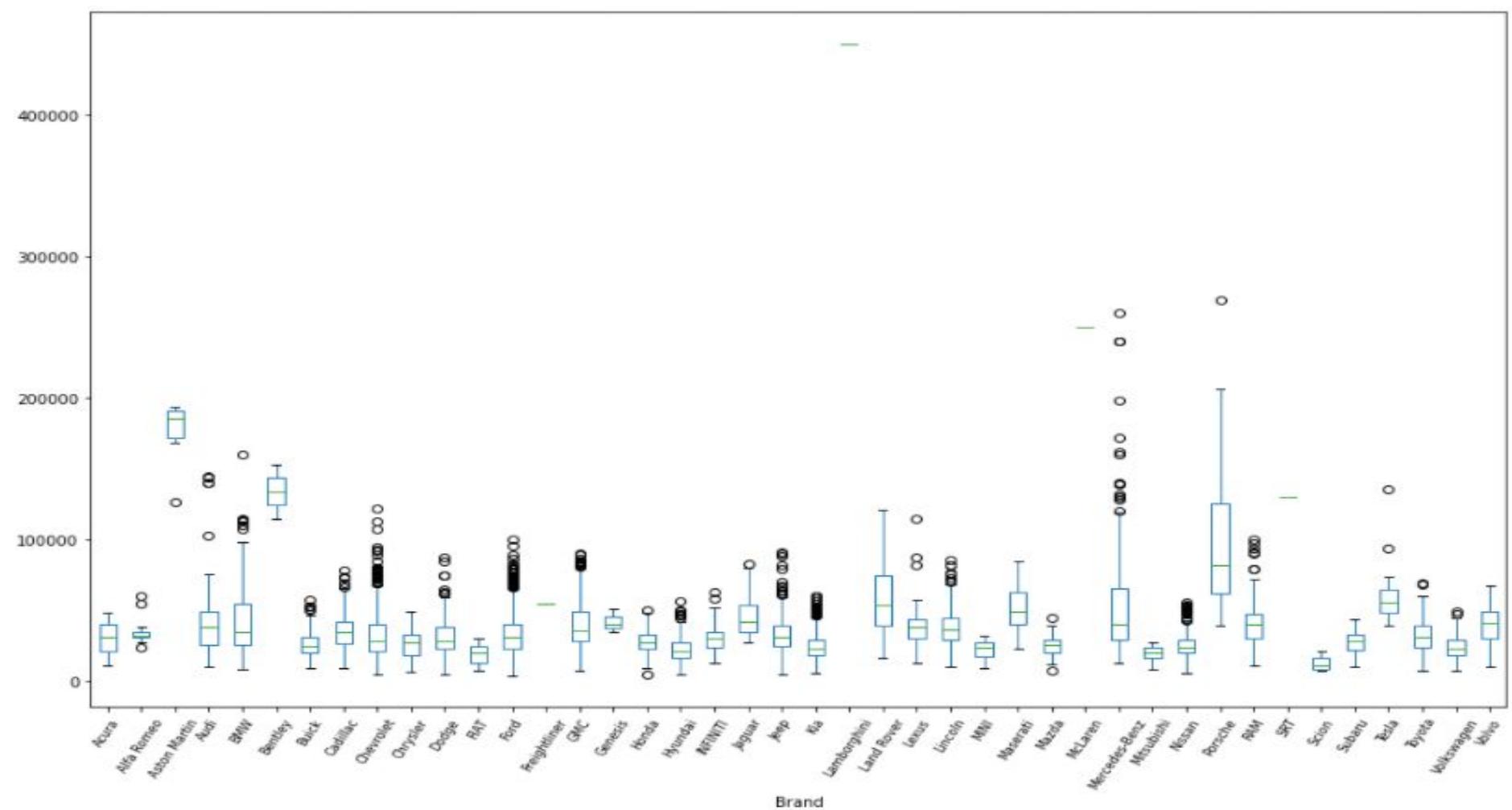
---

# Check Target(Price) Based Distribution

Brands vs Average Price



Brand vs Price



---

## Dropping Unnecessary Columns and Split the Data

- Drop Price column
- Drop Seller State because there are other States besides IL because zipcodes we used could be used at another state or another possibility is that region near the borderline of the States dealer posted their vehicle at Illinois.
- Split the dataset with the ratio of 70:30
- Check if model ID and model name match to see if we can drop model name
- Result came out true, so drop model ID

```
# Train test split - 30% of whole dataset into test data
X = df.drop(['Price', 'SellerState'], axis=1)
y = df.Price
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=123456789)
```

```
#Check if model id and model name match
df.loc[df.ModelID == 'd622', 'ModelName']

# Drop modelName because modelName can be found using modelId
X_train.drop(['id', 'modelName'], axis=1, inplace=True)
X_test.drop(['id', 'modelName'], axis=1, inplace=True)
```

---

# Target Encoding and One hot encoding

- Replace columns that has categorical value with the mean of target value : Brand, ModelID, Color, bodyTypeName
- Apply one hot encoding with Transmission and FuelType

```
# Target Encoding
target_enc_cols = ['Brand', 'ModelID', 'Color', 'bodyTypeName']
enc = ce.target_encoder.TargetEncoder(return_df=True, handle_unknown='value')
X_train[target_enc_cols] = enc.fit_transform(X_train[target_enc_cols], y_train)
X_test[target_enc_cols] = enc.transform(X_test[target_enc_cols])
```

```
# One Hot Encoding
one_hot_cols = ['Transmission', 'FuelType']
enc = OneHotEncoder(handle_unknown='ignore')
encoded = enc.fit_transform(X_train[one_hot_cols]).toarray()
X_train[enc.get_feature_names_out()] = encoded
X_train.drop(one_hot_cols, axis=1, inplace=True)

encoded = enc.transform(X_test[one_hot_cols]).toarray()
X_test[enc.get_feature_names_out()] = encoded
X_test.drop(one_hot_cols, axis=1, inplace=True)
```



# Model Fitting

---

## 1. Linear Regression

```
#modeling Linear Regression

model = LinearRegression(fit_intercept=True, n_jobs=4)
reg = model.fit(X_train, y_train)
pred = model.predict(X_test)

#R^2 score
print("Linear Regression R^2 Score for test set: ", reg.score(X_test, y_test))
print("Linear Regression MSE: ", mean_squared_error(pred, y_test))
print("Linear Regression MAE: ", mean_absolute_error(pred, y_test))

Linear Regression R^2 Score for test set:  0.7936733576176747
Linear Regression MSE:  72671238.66087796
Linear Regression MAE:  4553.894298360479
```

```
model.intercept_
model.coef_.shape
#unable to visualize because it is 32 dimension
(32,)
```

- Unable to visualize because we have 32 parameters.

## 2. Polynomial Regression Degree 2

```
# Polynomial Regression - degree 2
poly = PolynomialFeatures(degree= 2)
X_train_ = poly.fit_transform(X_train)
X_test_ = poly.fit_transform(X_test)

model = LinearRegression(fit_intercept=True, n_jobs=4)
reg = model.fit(X_train_, y_train)
pred = model.predict(X_test_)

print("Polynomial Regression R^2 Score for test set: ", reg.score(X_test_, y_test))
print("Polynomial Regression MSE: ", mean_squared_error(pred, y_test))
print("Polynomial Regression MAE: ", mean_absolute_error(pred, y_test))
```

Polynomial Regression R^2 Score for test set: 0.8232675511760047

Polynomial Regression MSE: 62247734.07503561

Polynomial Regression MAE: 4477.578155412296

Polynomial  
Linear  
Regression

$$y = b_0 + b_1x_1 + b_2x_1^2 + \dots + b_nx_1^n$$

### 3. Random Forest

```
# With the best parameter
rf = RandomForestRegressor(n_estimators= 150, min_samples_split=2, max_depth=None, random_state=42)
rf.fit(X_train, y_train)
pred = rf.predict(X_test)
print("Test MSE: ", mean_squared_error(pred, y_test))
print("Test MAE: ", mean_absolute_error(pred, y_test))

pred = rf.predict(X_train)
print("Train MSE: ", mean_squared_error(pred, y_train))
print("Train MAE: ", mean_absolute_error(pred, y_train))

# Show feature importances
feature_dict = dict(zip(rf.feature_names_in_, rf.feature_importances_.round(4)))
display(sorted(feature_dict.items(), key=lambda x:x[1], reverse=True))

#accuracy
from sklearn.metrics import accuracy_score
print('정확도 : {:.4f}'.format(model.score(X_test, y_test)))

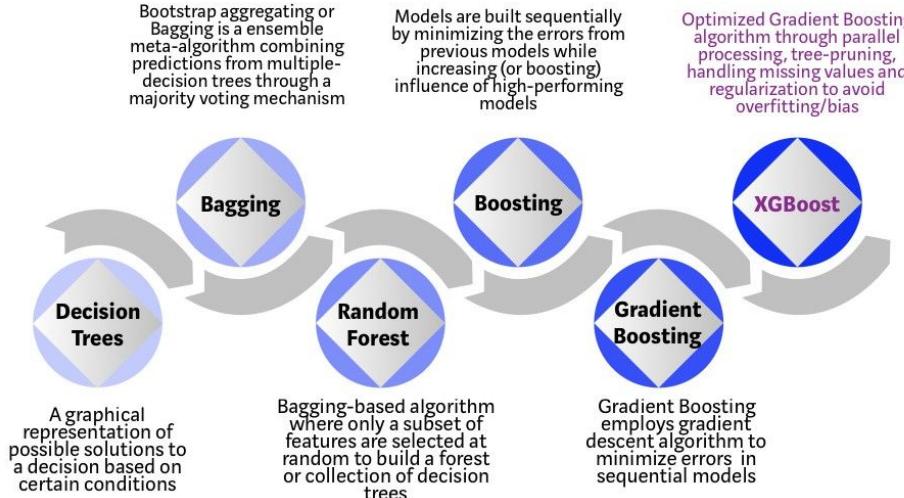
정확도 : 0.7937
```

```
Test MSE:  47375925.13666701
Test MAE:  3504.609691049727
Train MSE:  7823178.279621231
Train MAE:  1254.8186176855559
[('ModelID', 0.5816),
 ('mileage', 0.1488),
 ('YearOld', 0.0876),
 ('Color', 0.04),
 ('DaysOnMarket', 0.0361),
 ('Brand', 0.0176),
 ('OwnerCount', 0.0128),
 ('bodyTypeName', 0.0104),
 ('Leather Seats', 0.0098),
 ('Backup Camera', 0.0068),
 ('Bluetooth', 0.0068),
 ('Heated Seats', 0.0052),
 ('Adaptive Cruise Control', 0.0051),
 ('Navigation System', 0.0043),
 ('Sunroof/Moonroof', 0.0039),
 ('Remote Start', 0.0031),
 ('Blind Spot Monitoring', 0.0028),
 ('CarPlay', 0.0028),
 ('Parking Sensors', 0.0028),
 ('Alloy Wheels', 0.0025),
 ('Android Auto', 0.0019),
 ('FuelType_Biodiesel', 0.0017),
 ('AccidentCount', 0.0012),
 ('FuelType_Gasoline', 0.0012),
 ('FuelType_Diesel', 0.001),
 ('FuelType_Flex Fuel Vehicle', 0.0007),
 ('FuelType_Electric', 0.0005),
 ('Transmission_Automatic', 0.0004),
 ('Transmission_Manual', 0.0003),
 ('Transmission_CVT', 0.0002),
 ('FuelType_Hybrid', 0.0001),
 ('Transmission_Dual', 0.0)]
```

# Future Plans

Xgboost - decision-tree-based ensemble Machine Learning algorithm

*dmlc*  
**XGBoost**





## XGBoost(Extreme Gradient Boost)

XGBoost is an optimized distributed gradient boosting library designed to be highly efficient, flexible and portable. It implements machine learning algorithms under the Gradient Boosting framework. XGBoost provides a parallel tree boosting (also known as GBDT, GBM) that solve many data science problems in a fast and accurate way.

---

# XGBoost(Extreme Gradient Boost)

Pros :

1. GBM Improved speed compare to GBM
  - a. use parallel tree boosting
2. Efficient, Flexible, and Portable
3. Prevent overfitting
4. Optimal model can be made by tuning hyperparameter

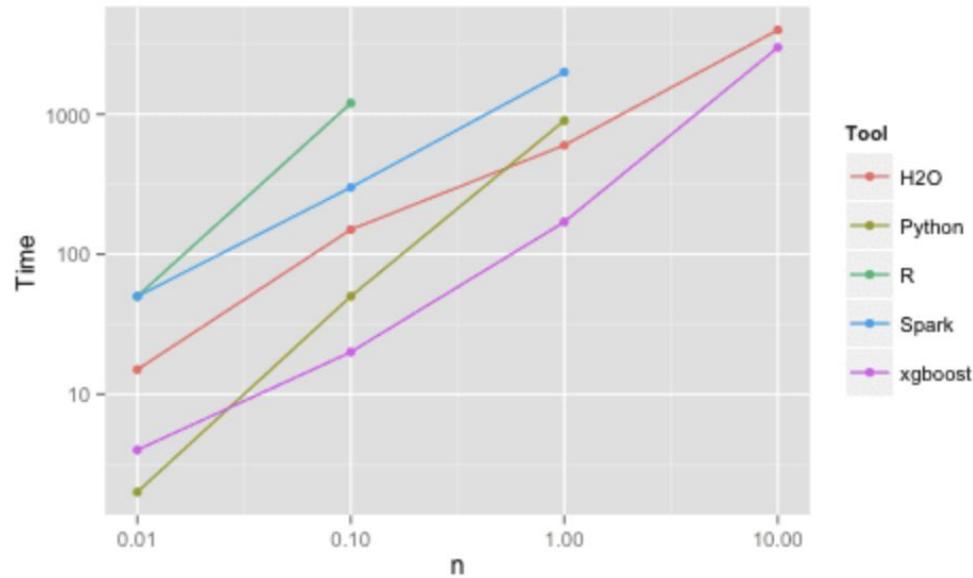
Cons:

1. XGBoost doesn't work well on dispersed data
2. Sensitive to outliers
  - a. XGBoost classifier is forced to fix every error
3. Complicated hyperparameter

---

# Execution Speed

- benchmarks comparison Random Forest Implementation



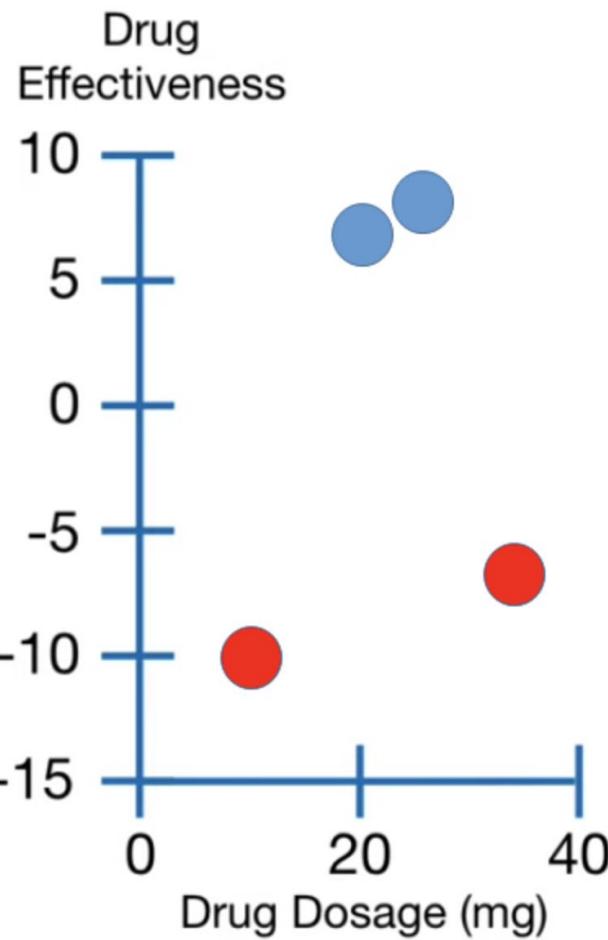
---

## Example

Simple data with 4 observations

X : Drug Dose  
Y = Drug Effectiveness

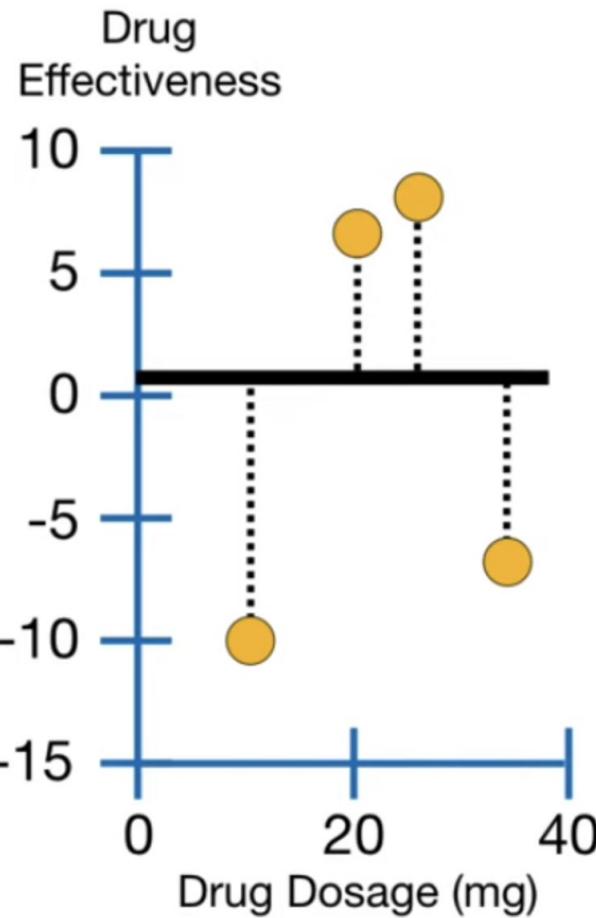
Using this data, we are going to make  
drug dosage effectiveness model





Initial Prediction Value (Black line): 0.5

We think this is because computer predicts either 0 or 1 so they set it to middle point which is 0.5 as base score





We want to maximize our **Gain**

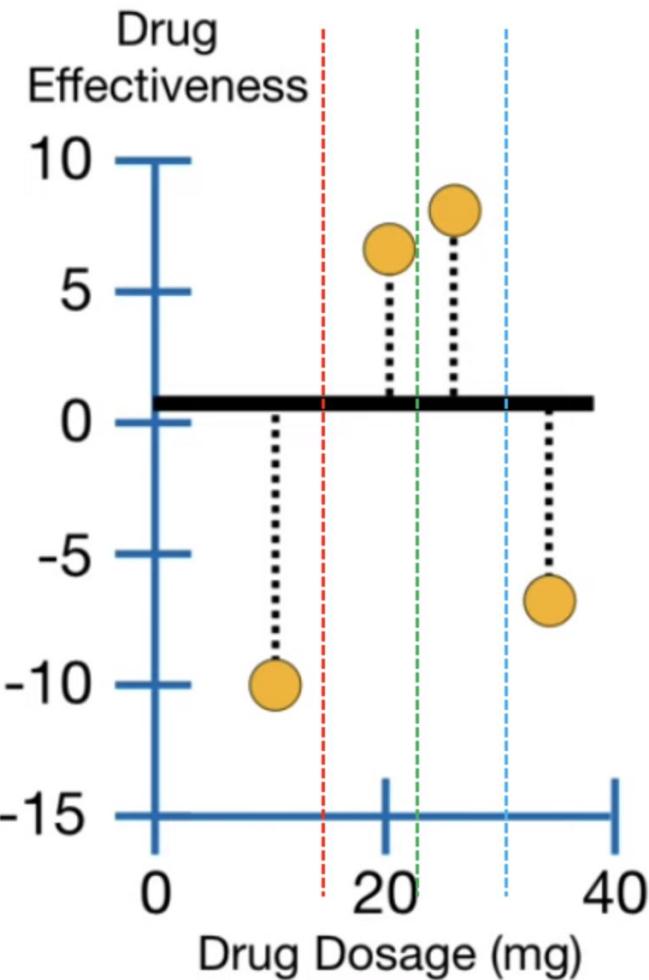
- Gain = (New Similarity Score) - (Old Similarity Score)

Similarity Score = (Sum of Residuals)<sup>2</sup> / (Number of Residuals + lambda)

We can consider lambda as 0 here.  
Lambda is just a parameter to prevent overfitting.

So in this case, our old similarity score is

$$(-10.5 + 6.5 + 7.5 - 7.5)^2 / (4 + 0) = 4$$



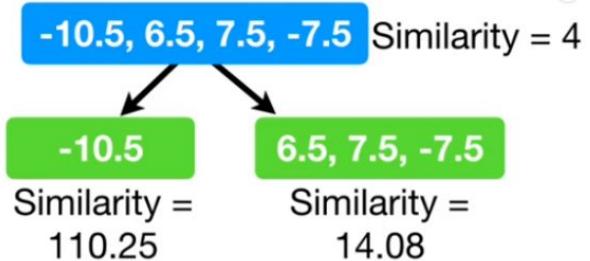
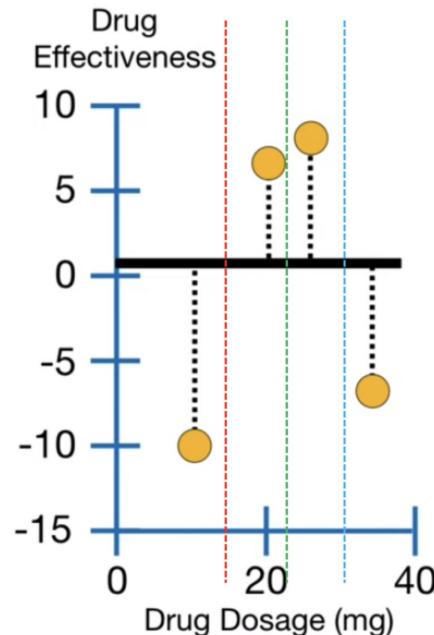
- Gain = (New Similarity Score) - (Old Similarity Score)

Dividing it by the red line

- $(-10.5)^2 / 1 = 110.25$  (New similarity Score)
- $(6.5 + 7.5 - 7.5)^2 / 3 = 14.08$  (New similarity score)

Using this information, we calculate the Gain

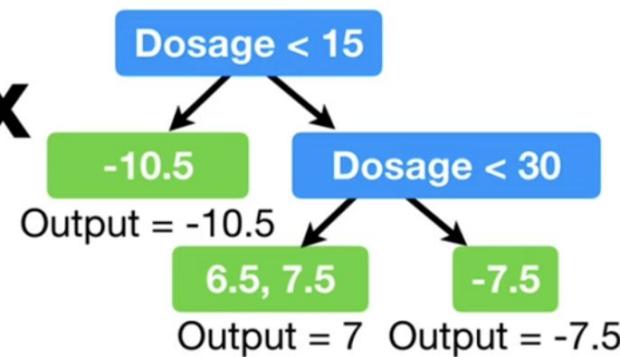
- $(110.25 + 14.08) - 4 = 120.33$



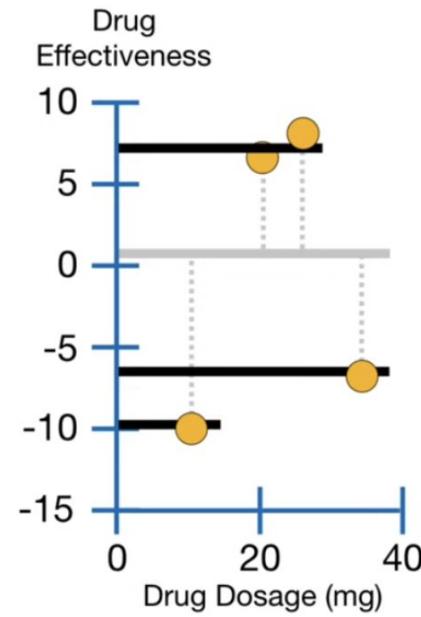
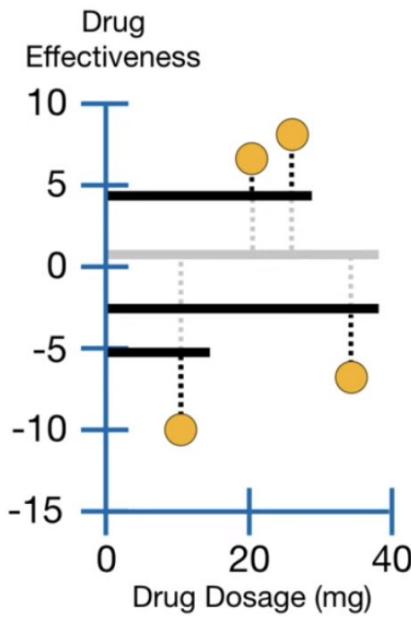
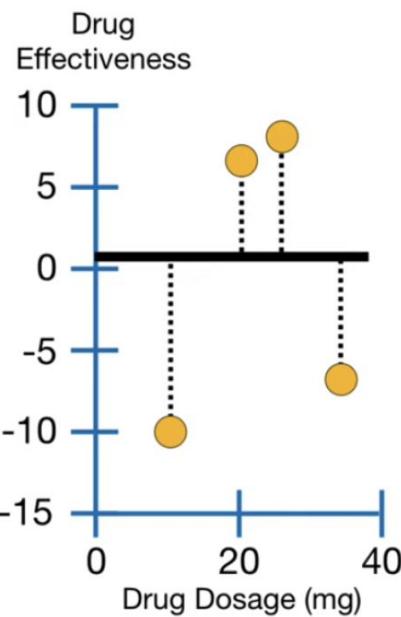
# Decision Tree

Predicted Drug  
Effectiveness  
**0.5**

+ Learning Rate **X**



—



# Used Car Price Prediction Using XGBoost Regressor

In [119...]

```
"""
XGB Regressor - with Randomized Grid Search
"""

xgb = XGBRegressor(objective='reg:squarederror', eval_metric='rmse')
params = {'n_estimators':[700,800,900], 'max_depth':[5,6], 'learning_rate':[0.1]}
model = RandomizedSearchCV(xgb, params, n_iter=5, cv=5, scoring='neg_mean_squared_error', n_jobs=-1, verbose=0)
model.fit(X_train, y_train)

print(model.best_params_)
pred = model.predict(X_test)
print("MSE: ", mean_squared_error(pred, y_test))
print("MAE: ", mean_absolute_error(pred, y_test))

{'n_estimators': 800, 'max_depth': 5, 'learning_rate': 0.1}
MSE:  39510885.08189203
MAE:  3184.4337905248008
```

정확도 : 86.8579

In [139...]

```
# With the best parameters searched
xgb = XGBRegressor(objective='reg:squarederror', eval_metric='rmse', n_estimators=1000
                   , max_depth=6, learning_rate=0.1)
xgb.fit(X_train, y_train)
xgb_pred = xgb.predict(X_test)
print("XGB MSE: ", mean_squared_error(xgb_pred, y_test))
print("XGB MAE: ", mean_absolute_error(xgb_pred, y_test))
```

XGB MSE: 39493181.22963319  
XGB MAE: 3191.6426135698953

# 1. Linear Regression

```
#modeling Linear Regression

model = LinearRegression(fit_intercept=True, n_jobs=4)
reg = model.fit(X_train, y_train)
pred = model.predict(X_test)

#R^2 score
print("Linear Regression R^2 Score for test set: ", reg.score(X_test, y_test))
print("Linear Regression MSE: ", mean_squared_error(pred, y_test))
print("Linear Regression MAE: ", mean_absolute_error(pred, y_test))
```

```
Linear Regression R^2 Score for test set: 0.7936733576176747
Linear Regression MSE: 72671238.66087796
Linear Regression MAE: 4553.894298360479
```

## 2. Polynomial Regression Degree 2

```
# Polynomial Regression - degree 2
poly = PolynomialFeatures(degree= 2)
X_train_ = poly.fit_transform(X_train)
X_test_ = poly.fit_transform(X_test)

model = LinearRegression(fit_intercept=True, n_jobs=4)
reg = model.fit(X_train_, y_train)
pred = model.predict(X_test_)

print("Polynomial Regression R^2 Score for test set: ", reg.score(X_test_, y_test))
print("Polynomial Regression MSE: ", mean_squared_error(pred, y_test))
print("Polynomial Regression MAE: ", mean_absolute_error(pred, y_test))
```

Polynomial Regression R<sup>2</sup> Score for test set: 0.8232675511760047

Polynomial Regression MSE: 62247734.07503561

Polynomial Regression MAE: 4477.578155412296

Polynomial  
Linear  
Regression

$$y = b_0 + b_1 x_1 + b_2 x_1^2 + \dots + b_n x_1^n$$

### 3. Random Forest

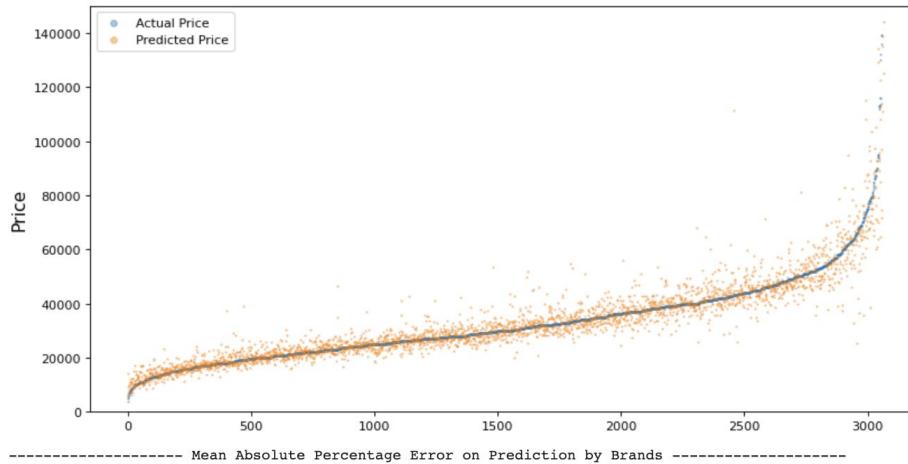
```
# With the best parameter
rf = RandomForestRegressor(n_estimators= 150, min_samples_split=2, max_depth=None, random_state=42)
rf.fit(X_train, y_train)
pred = rf.predict(X_test)
print("Test MSE: ", mean_squared_error(pred, y_test))
print("Test MAE: ", mean_absolute_error(pred, y_test))

pred = rf.predict(X_train)
print("Train MSE: ", mean_squared_error(pred, y_train))
print("Train MAE: ", mean_absolute_error(pred, y_train))

# Show feature importances
feature_dict = dict(zip(rf.feature_names_in_, rf.feature_importances_.round(4)))
display(sorted(feature_dict.items(), key=lambda x:x[1], reverse=True))
```

정확도 : 82.9708

### Actual Listed Price vs Predicted Price



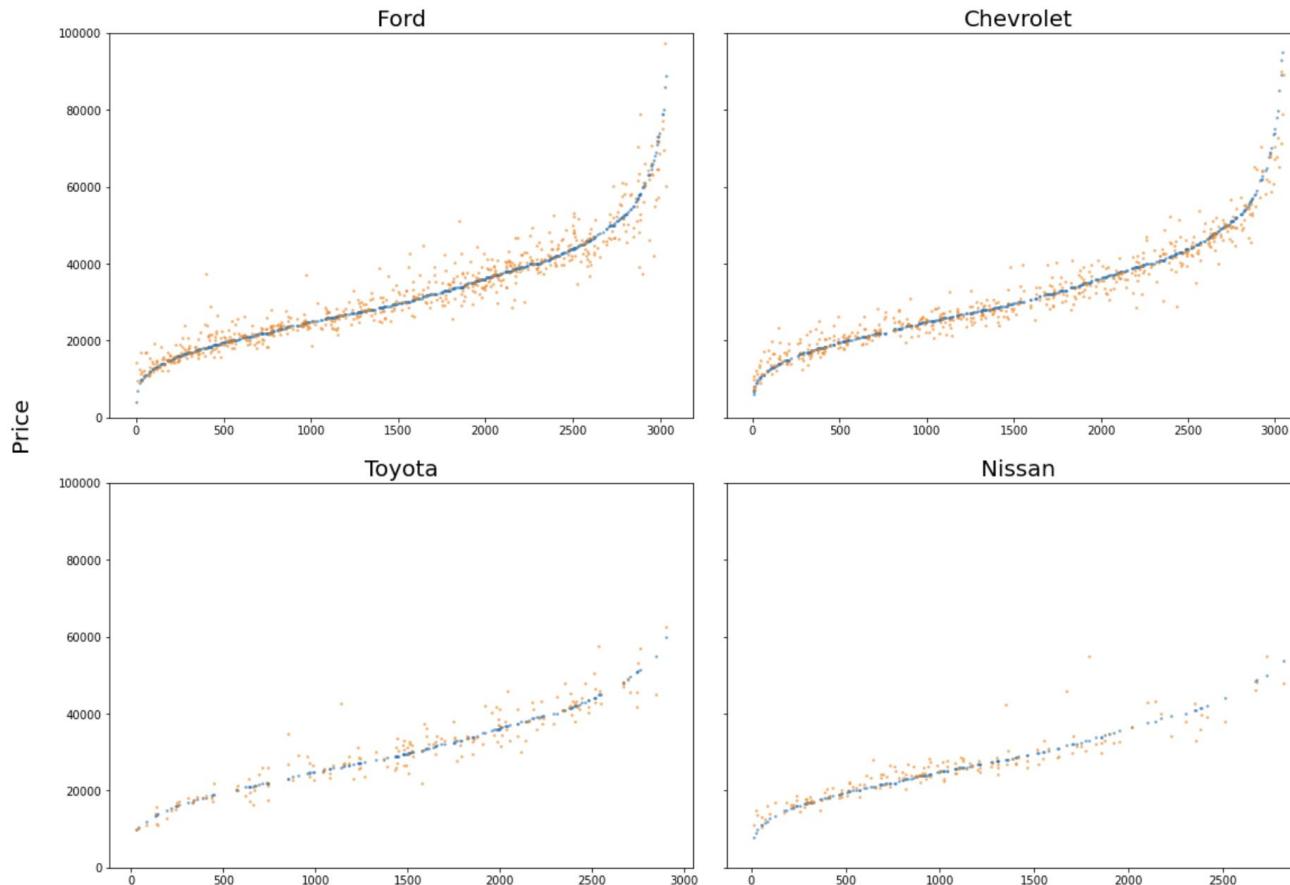
----- Mean Absolute Percentage Error on Prediction by Brands -----

Ford : 9.27%  
Chevrolet : 8.55%  
Toyota : 8.40%  
Nissan : 8.70%

----- Mean Absolute Percentage Error on Prediction by Body Types -----

Pickup Truck : 9.02%  
Sedan : 10.38%  
SUV / Crossover : 8.25%  
Wagon : 22.56%  
Hatchback : 13.46%  
Van : 14.26%  
Minivan : 6.82%  
Coupe : 16.39%  
Convertible : 15.44%

## Major Brands Actual Price vs Predicted Price



---

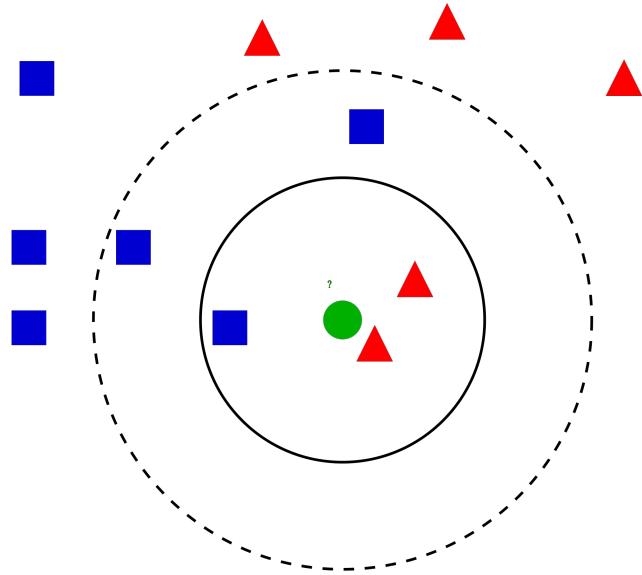
# Improving the model

1. Enlarged the data set!
  - a. From 1,000 zip codes from IL  
-> 3,000 zip codes from 5 states (IL, IA, WI, MO, IN)
2. Tried KNN with different scaling methods
3. Dropped Android Auto
  - a. Due to multicollinearity determined by VIF (variance inflation factor)
  - b. High market share of Apple CarPlay
    - i. Either merge the two or drop Android Auto

variables	VIF
0 CarPlay	12.465760
1 Android Auto	12.215435
2 YearOld	6.697340
3 mileage	6.702101



variables	VIF
0 CarPlay	1.184677
1 YearOld	6.690461
2 mileage	6.698865



---

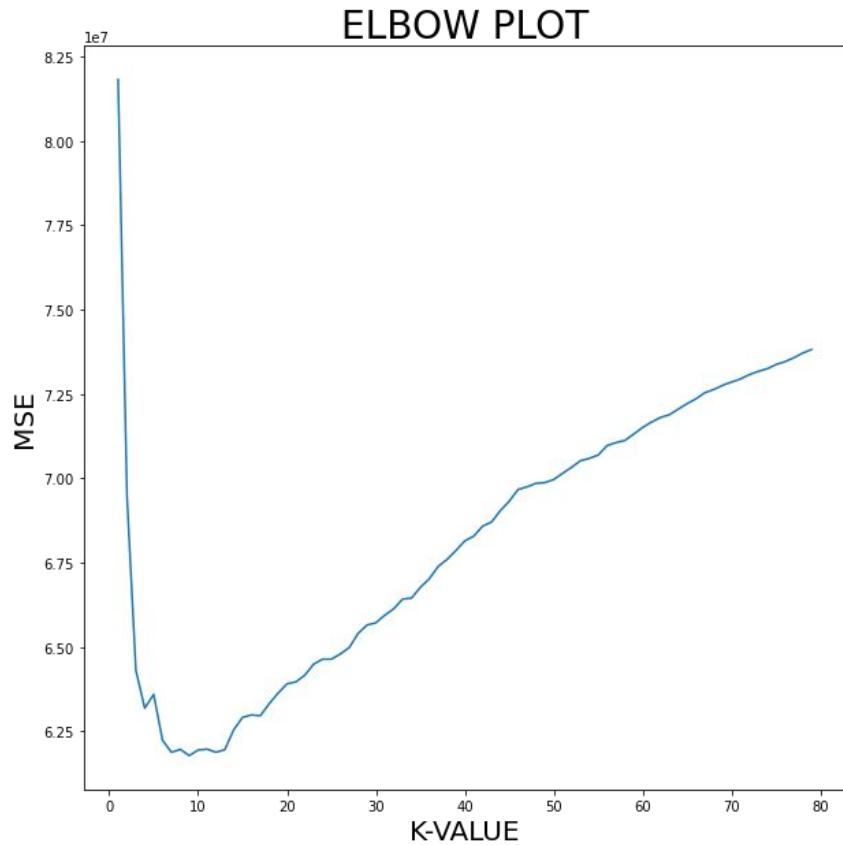
# K Nearest Neighbors

**MinMax Scaling** : choose minimum and maximum values of a feature to rescale values

**Standard Scaling** : scaling features to be approximately standard normally distributed.

**Robust Scaling** : overcome outliers within the data. This is acquired by using the first and third quartile values in the scaling process.

Best K selected : 8



---

## K Nearest Neighbors - MinMax Scaling

With our best K selected (8), MinMax scale range of (0,1)

KNN Score for test set after scailing: 0.6356390551835813  
KNN MSE: 100904088.59185754  
KNN MAE: 5711.1753761418595



## K Nearest Neighbors - Standard Scaling

With our best K selected (8)

KNN Score for test set after scaling: 0.7614678844856192

KNN MSE: 66057754.15362708

KNN MAE: 4676.748361096184



## K Nearest Neighbors - Robust Scaling

With our best K selected (8)

KNN Score for test set after scaling: 0.8138833082892858

KNN MSE: 51542118.92348536

KNN MAE: 3843.2880171950565

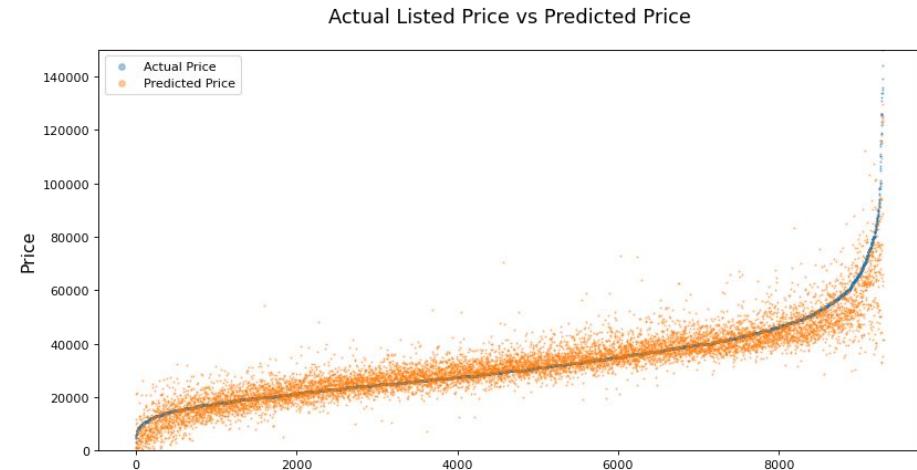
Results are lower when using Standard and MinMax scaling. We think this is because of weight problem.

# Model Performance

---

## Linear Regression

	Before	After
R^2	0.7937	0.7851
MSE	72,671,238	59,288,734
MAE	4554	4154.19



# Model Performance

---

## KNN Regression

	Before Scaling	After Scaling
R^2	0.7762	0.8138
MSE	61,969,064	51,542,118
MAE	4643.04	3843.28

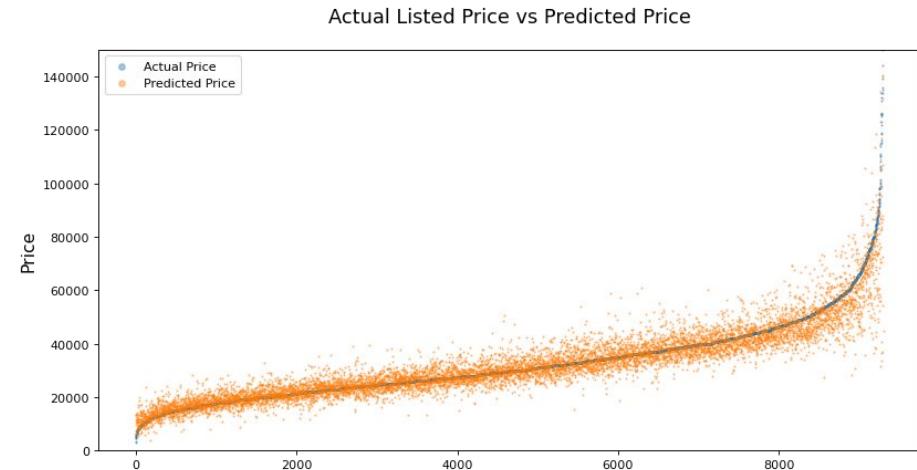


# Model Performance

---

## Polynomial Regression

	Before	After
R^2	0.8232	0.8331
MSE	62,247,734	46,202,202
MAE	4477.58	3443.15



# Model Performance

---

## Random Forest

	Before	After
R^2	0.83	0.8924
MSE	39,510,885	29,780,137
MAE	3184	2724.70

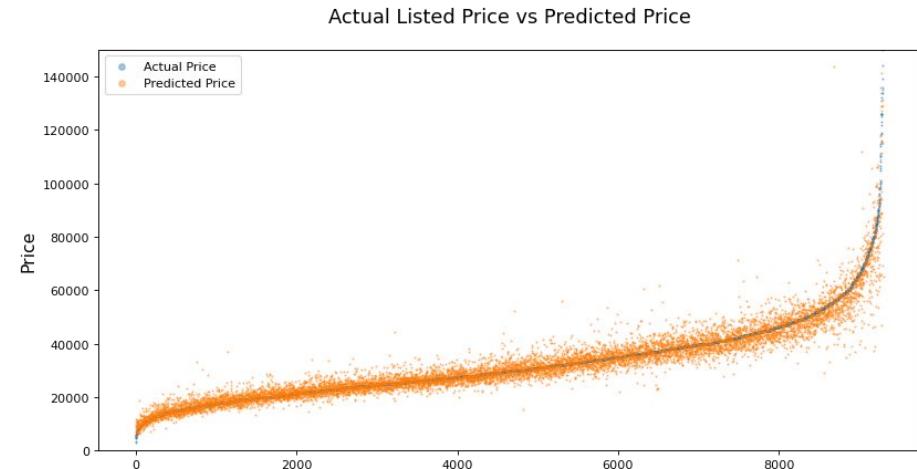


# Model Performance

---

## XGBoost Regression

	Before	After
R^2	0.8686	0.8930
MSE	39,493,181	29,623,373
MAE	3191.64	2584.23



# Feature Importance

