

Used Car Price Prediction and Recommender System

Trang Khong

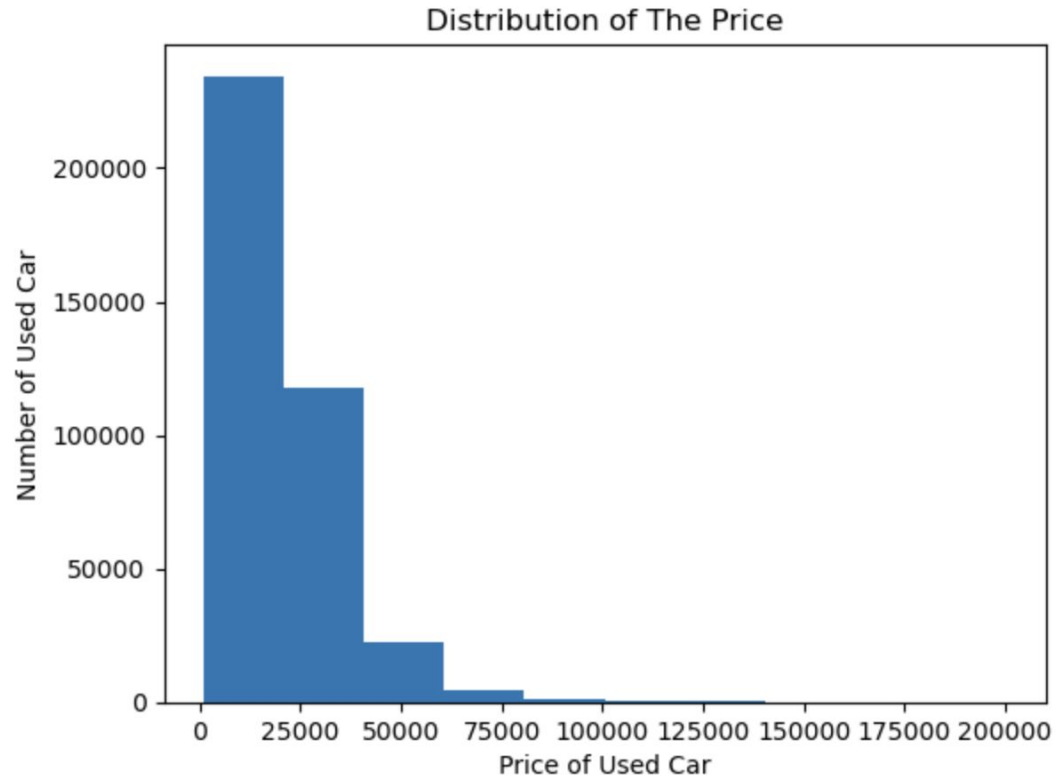
Problem Statement

- The price of used cars has increased after the coronavirus pandemic due to the issue of supply chain.
- Making a decision on buying and selling a car is challenging.
- The purpose of this project is to build a model to predict the price of used cars and a used car recommender system.
- Those models will help sellers when they are planning to sell their car and buyers to have an idea when they are looking for a used car.

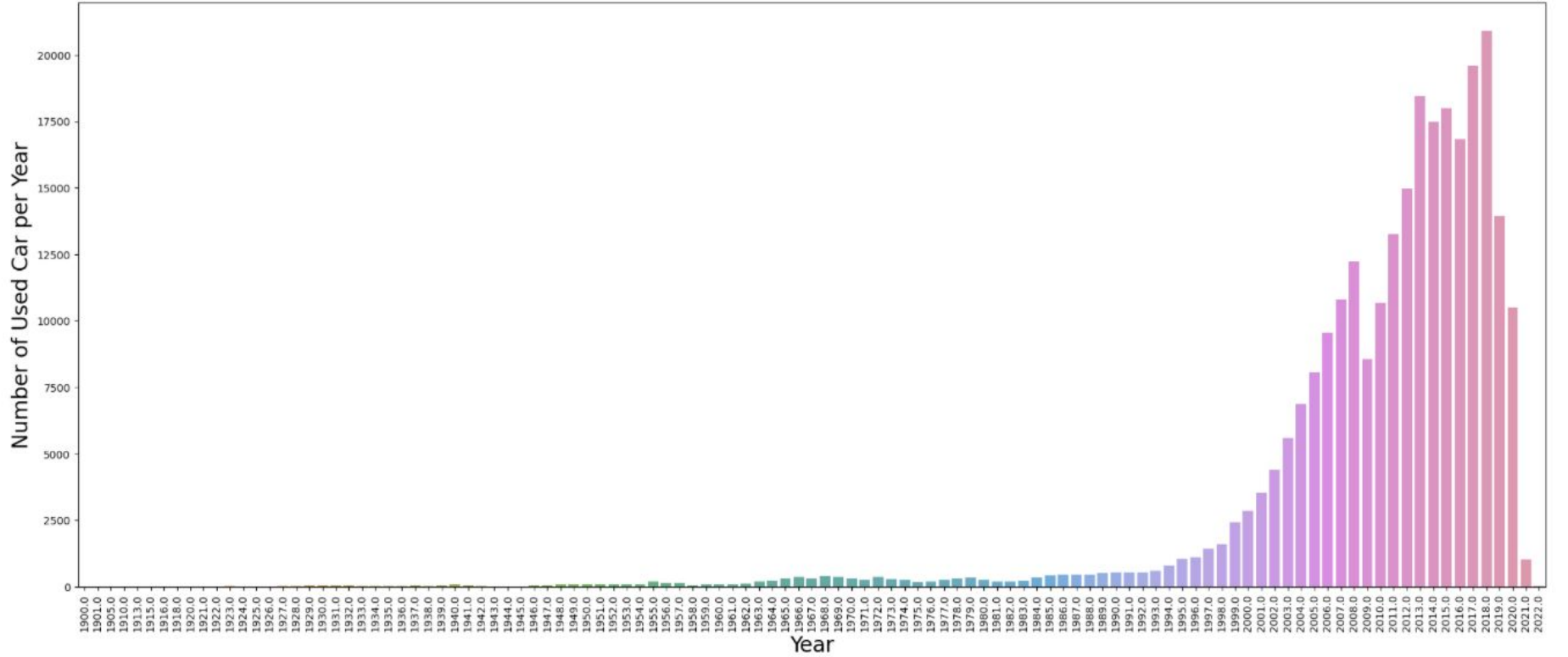
Analyzed Data

- Data was obtained from [Kaggle](#), collected by Austin Reese from Craigslist.
- The data have 26 columns.
- After cleaning, the data have 14 columns.
- The clean data will be use for imputing and modeling.
- List of features are applied to predict the price of used car are: year, manufacturers, condition, cylinders, fuel, odometer, title status, transmission, drive, size, type, color, state.

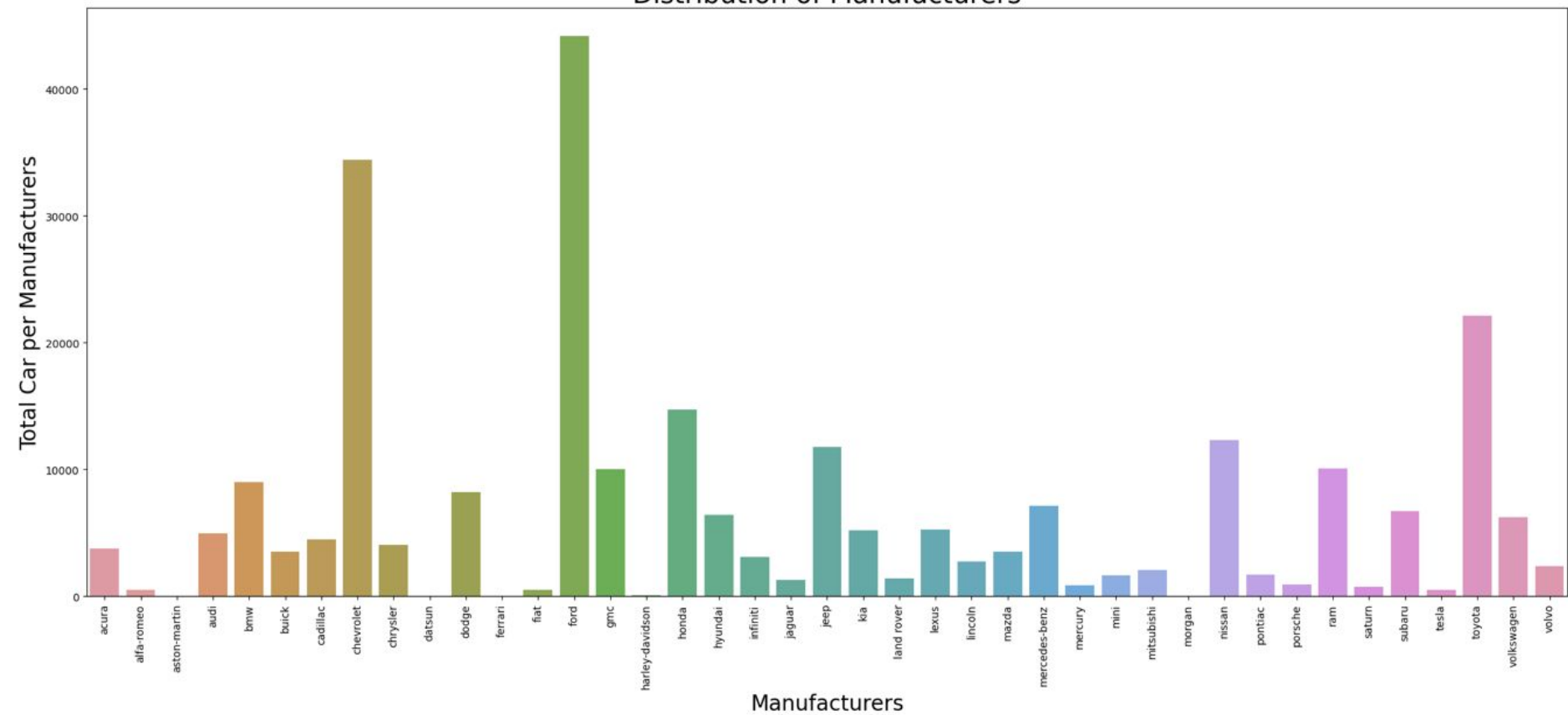
Distribution of Price



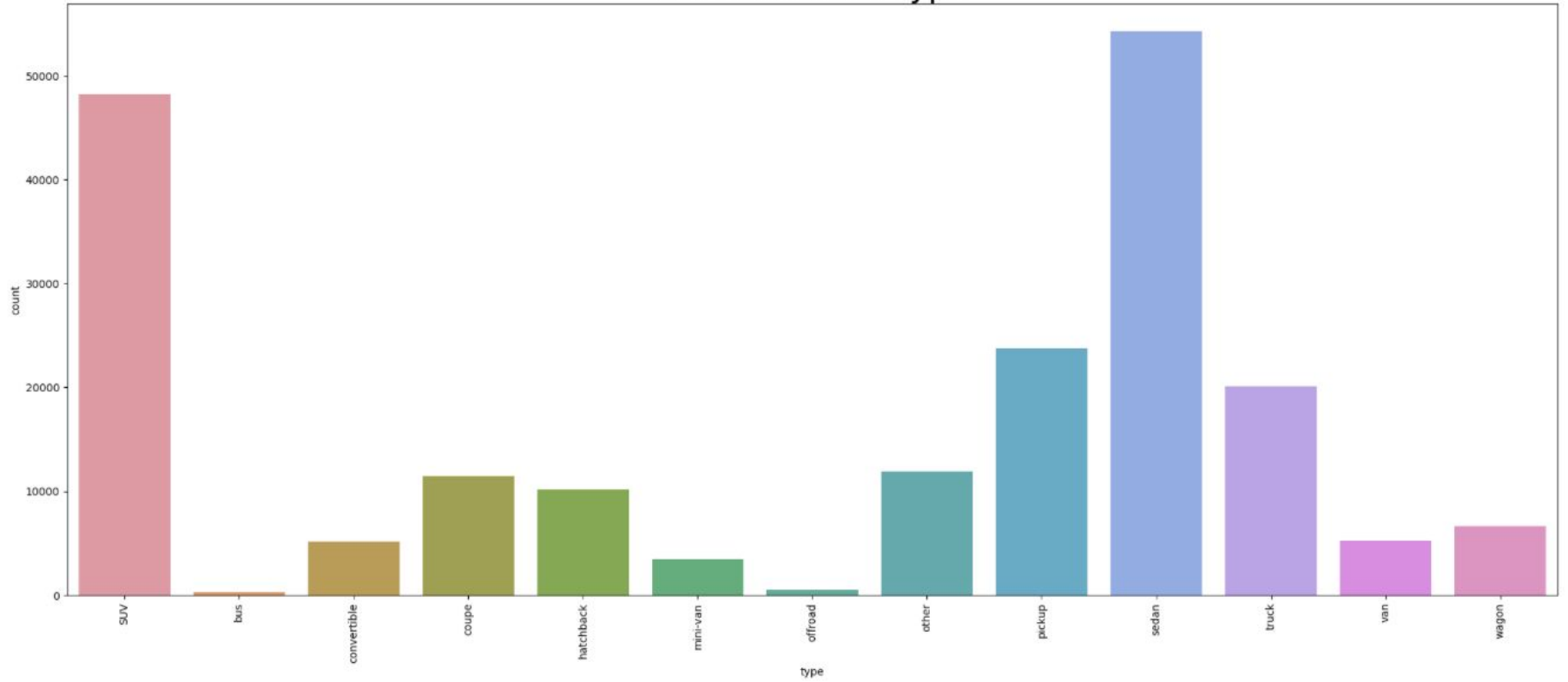
Distribution of Year



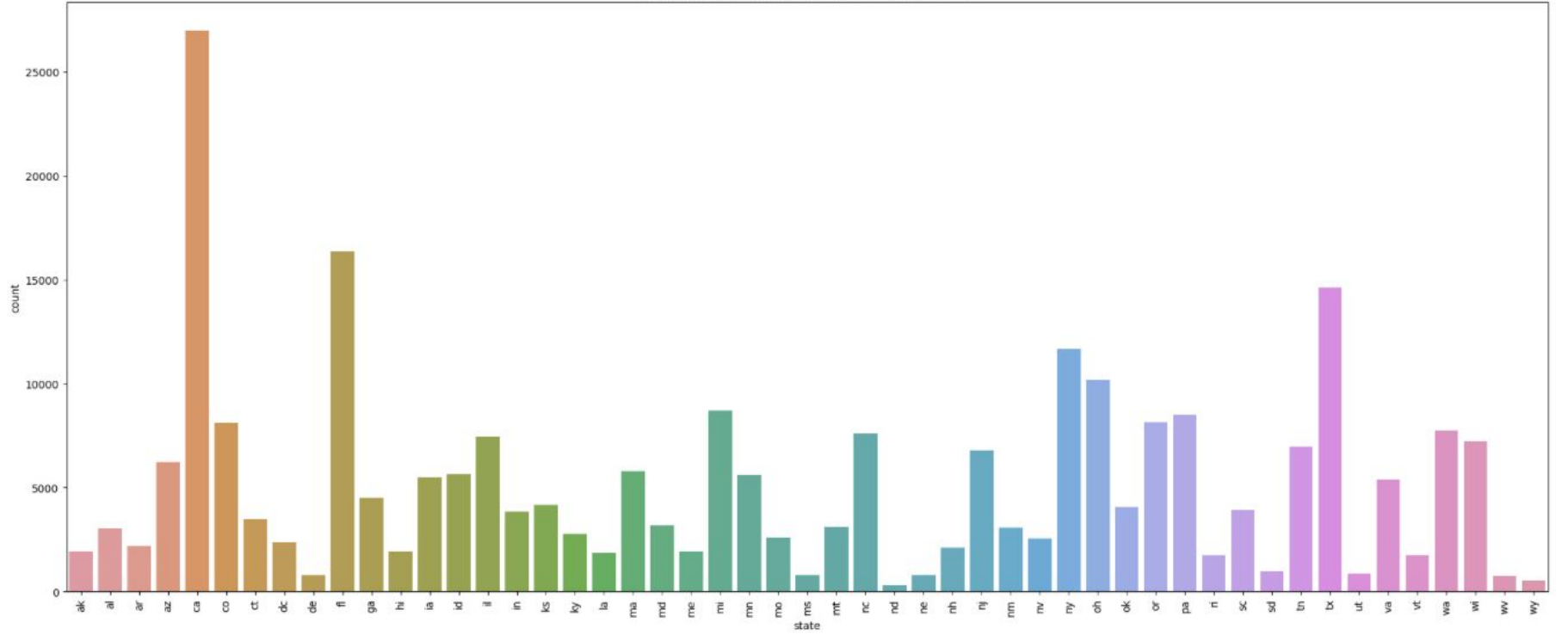
Distribution of Manufacturers



Distribution of Type



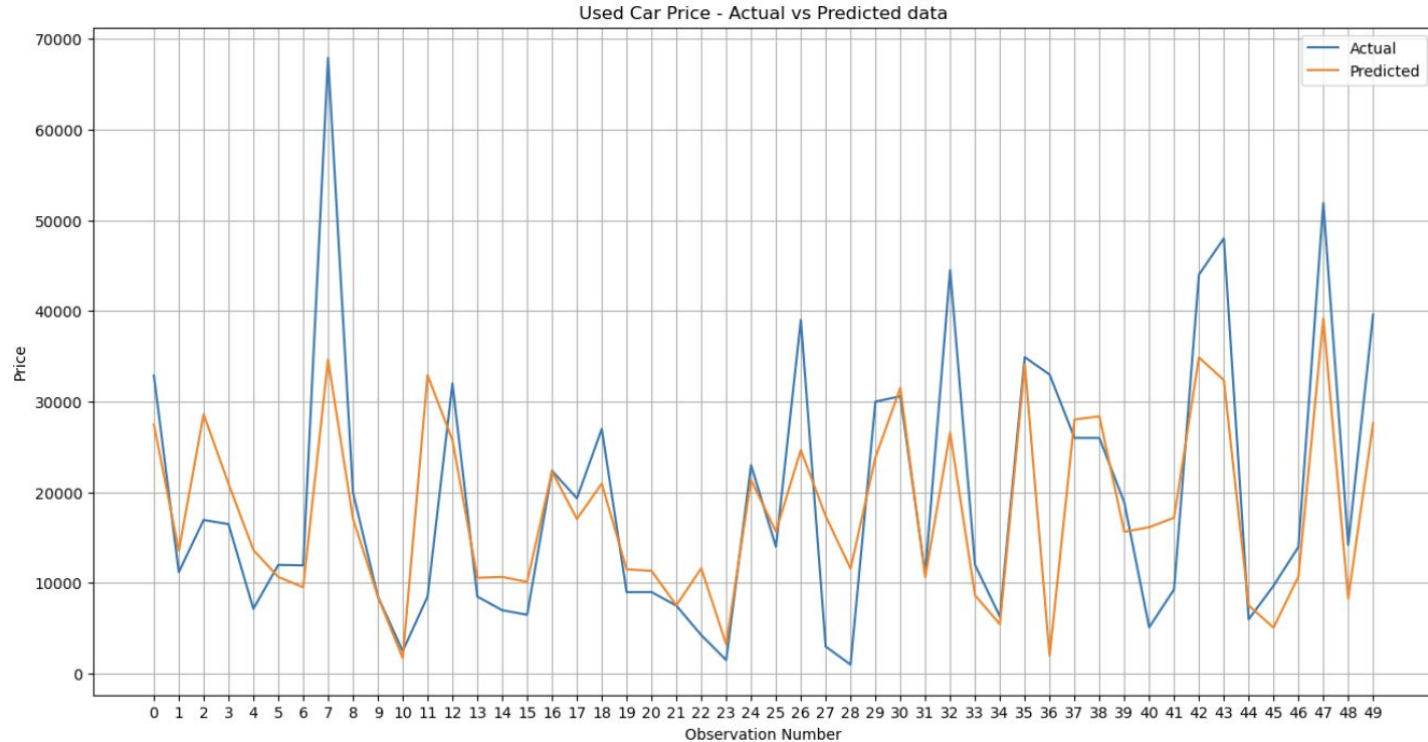
Distribution of States



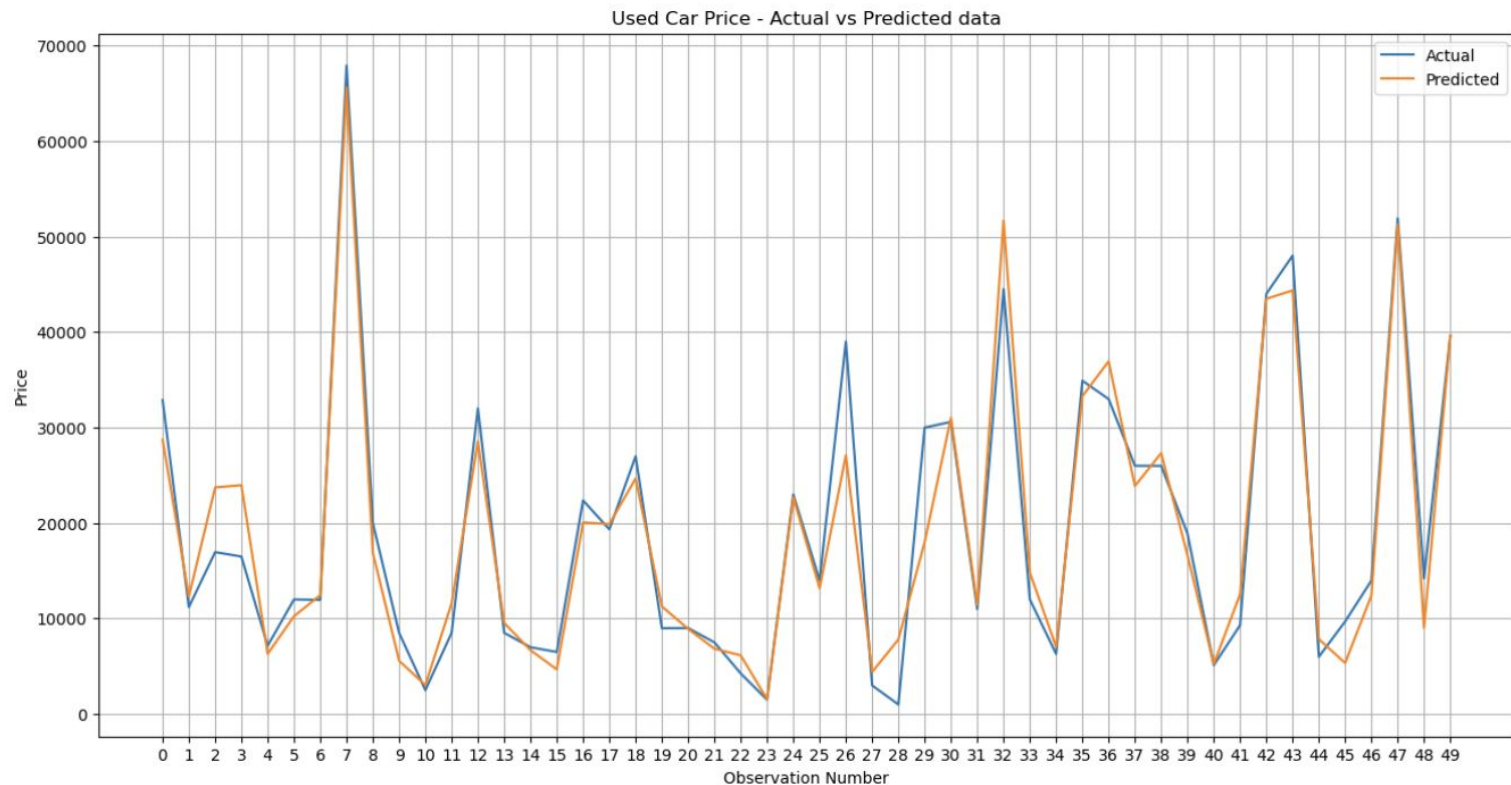
Imputing Method

- ❖ KNN imputer
- ❖ Mean and mode imputer base on different range of price
 - 1000 to 15000
 - 15000 to 50000
 - 50000 to 200000
- ❖ Different price of range give different mean and mode
- ❖ All models were compared to a baseline model.
- ❖ The used car price is evaluated with R Squared states how much of the variability in the data are explained by model, and Root Mean Squared Error (RMSE) represents (approximately) average difference between model predicted values and the actual values.

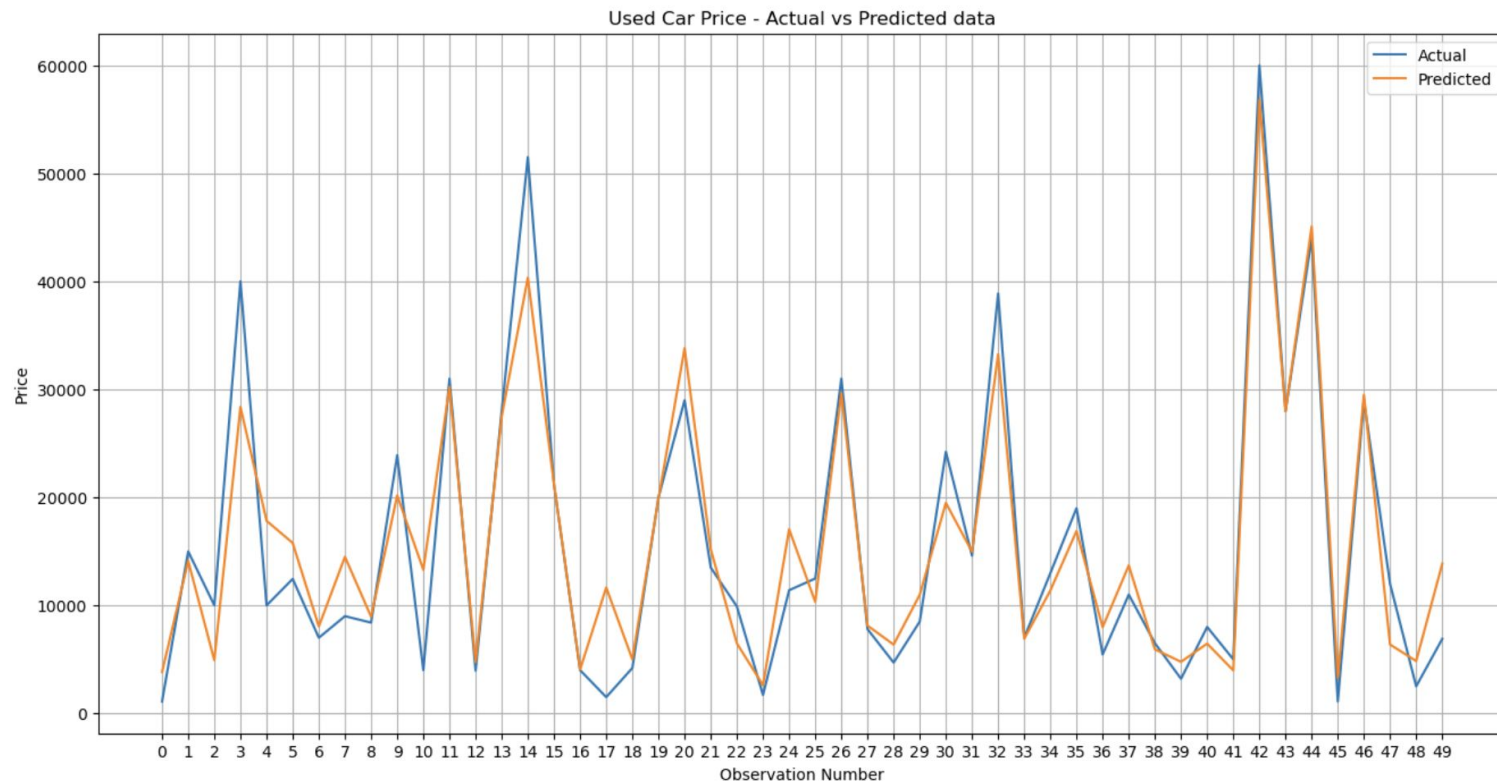
Linear Regression with Mean Imputing



Random Forest Regressor with Mean Imputer



Random Forest Regressor with KNN Imputer



Model Results with Two Different Imputer

Name of Model	Train score	Test score (R2)	RMSE
Baseline (Null model)	0.00	0.00	14504.270
Linear Regression (mean imputer)	0.514	0.521	10081.260
RandomForest Regressor (mean imputer)	0.982	0.874	5172.700
ExtraTrees Regressor (mean imputer)	0.952	0.874	5161.382
AdaBoost Regressor (mean imputer)	0.961	0.872	5217.493
GradientBoosting Regressor (mean imputer)	0.777	0.775	6913.230
RandomForest Regressor(pipeline)	0.982	0.874	5167.861
ExtraTrees Regressor(pipeline)	0.999	0.873	5179.484
RandomForest Regressor (KNN imputer)	0.925	0.836	5853.053
ExtraTrees Regressor (KNN imputer)	0.941	0.840	5781.103
AdaBoost Regressor (KNN imputer)	0.953	0.836	5867.411

Discussion

- The best performance model is Extra Trees Regressor with lowest RMSE and highest R2 score on both imputing methods.
- Linear Regression and Gradient Boosting have lowest train score and test score but the different between train and test score are small, those models are less overfit.
- The others models have train score higher than test score, those models are overfit but the test score are improved.
- Tuning hyperparameters can be apply to help improve the overfit.
- Keep the hyperparameter the same, train score, test score, and RMSE will be slightly different when apply the pipeline versus without applying the pipeline.
- Using the same model and hyperparameters, models with KNN imputer are less performance than mean imputer. However, those models with KNN imputer are less overfit than those models with mean imputer.

Price Recommender

```
price_recommender((10000, 20000), 7)
```

	price	manufacturer	type	year	condition	fuel	title_status	transmission	paint_color
0	12998.0	nissan	sedan	2018.0	excellent	gas	clean	automatic	white
1	16500.0	buick	SUV	2014.0	excellent	gas	clean	automatic	white
2	16900.0	jeep	SUV	2018.0	good	gas	clean	other	white
3	16990.0	honda	sedan	2015.0	good	other	clean	other	silver
4	17337.0	toyota	hatchback	2018.0	good	gas	clean	automatic	white
5	18500.0	audi	SUV	2015.0	excellent	diesel	clean	automatic	black
6	18995.0	land rover	SUV	2012.0	good	gas	clean	automatic	black

Brand Recommender

```
brand_recommender('toyota', (10000, 20000), 7)
```

	price	manufacturer	type	year	condition	fuel	title_status	transmission	paint_color
0	10995.0	toyota	SUV	2008.0	excellent	gas	clean	automatic	white
1	12987.0	toyota	sedan	2012.0	excellent	gas	clean	automatic	white
2	14998.0	toyota	SUV	2012.0	excellent	gas	clean	automatic	white
3	14999.0	toyota	SUV	2013.0	excellent	gas	clean	automatic	white
4	16999.0	toyota	truck	2010.0	good	gas	clean	manual	silver
5	17995.0	toyota	sedan	2014.0	like new	gas	clean	automatic	red
6	19500.0	toyota	SUV	2013.0	excellent	gas	clean	automatic	white

Used Car Recommender

```
all_recommender('acura',2007,'sedan',(5000, 10000),7)
```

	price	manufacturer	type	year	condition	fuel	title_status	transmission	paint_color
0	5500.0	acura	sedan	2007.0	excellent	gas	rebuilt	automatic	white
1	5950.0	acura	sedan	2007.0	excellent	gas	clean	automatic	grey
2	6900.0	acura	sedan	2007.0	excellent	gas	clean	automatic	blue
3	7500.0	acura	sedan	2007.0	excellent	gas	clean	automatic	black
4	7995.0	acura	sedan	2007.0	excellent	gas	clean	automatic	blue
5	7995.0	acura	sedan	2007.0	good	gas	clean	automatic	silver
6	8550.0	acura	sedan	2007.0	excellent	gas	clean	automatic	white

Cosine Similarity Recommender

	price	manufacturer	type	year	cylinders
10860	3000.0	nissan	SUV	2004.0	4
20084	3000.0	nissan	SUV	2004.0	4
10203	2995.0	nissan	SUV	2003.0	4
5895	3000.0	nissan	SUV	2003.0	4
32363	2800.0	nissan	SUV	2002.0	4
25463	2850.0	nissan	SUV	2007.0	4
24269	3290.0	nissan	SUV	2005.0	4
35024	3500.0	nissan	SUV	2003.0	4
19548	2195.0	nissan	SUV	2002.0	4
12262	2700.0	nissan	SUV	2010.0	4

Discussion

- Due to the limit of computer memory, the cosine similarity was not be able to run for this dataset. Limited the price of the data to limited the amount of data for cosine similarity.
- The cosine similarity works when limit the data to a small dataset.
- Brand recommender, Price recommender, and Used Car recommender can be use as a recommender system. Those recommender systems do not provide a perfect result but those recommender systems will give a list of used car base on some features that customers are looking for.

Conclusion

- The method use for data imputing has an important effect on model training and performance.
- The best performing model in this case was ExtratreesRegressor with mean imputed values.
- Cosine Similarity Recommender system will work with powerful computer or limited the data to certain size.

Bibliography

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

<https://www.cnbc.com/2023/09/08/used-vehicle-prices-may-have-bottomed-for-2023.html#:~:text=Used%20vehicle%20prices%20have%20been,high%20prices%20amid%20resilient%20demand.>

Thank you