



CAR PRICE PREDICTION PROJECT

Submitted by:
Asheem Siwach

ACKNOWLEDGMENT

I would like to thank Flip Robo Technologies for providing me with the opportunity to work on this project from which I have learned a lot. I am also grateful to Mr. Shubham Yadav for his constant guidance and support.

Reference sources are :-

- Google
- Stackoverflow.com
- Analytics vidhya
- Notes and repository from DataTrained

INTRODUCTION

- **Business Problem Framing**

With the covid 19 impact in the market, we have seen lot of changes in the car market. Now some cars are in demand hence making them costly and some are not in demand hence cheaper. One of our clients works with small traders, who sell used cars. With the change in market due to covid 19 impact, our client is facing problems with their previous car price valuation machine learning models. So, they are looking for new machine learning models from new data. We have to make car price valuation model.

- **Review of Literature**

The goal of this project is to build a model which can be used to predict the car price using various features available. This project is more about data exploration, finding the better insights from data and using the skills and techniques to build an efficient model which can predict the prices after the recession faced during covid-19 crisis. Since we scrape a good amount of data that are related to the price of car, we can do better data exploration and derive some interesting features using the available columns.

Analytical Problem Framing

- Mathematical/ Analytical Modeling of the Problem

In this project we have created a model to predict the "SalePrice" of cars which are for resale in Used Cars Market. So "SalePrice" is our target feature and it is continuous in nature. Therefore we are dealing with Regression type of problem.

This project is based on 2 main phases –

- 1). Data Collection
- 2). Model Building

Data Collection Phase –

You have to scrape at least 5000 used cars data. You can scrape more data as well, it's up to you.

More the data better the model

In this section You need to scrape the data of used cars from websites (Olx, cardekho, Cars24 etc.)

You need web scraping for this. You have to fetch data for different locations. The number of

columns for data doesn't have limit, it's up to you and your creativity. Generally, these columns are Brand, model, variant, manufacturing year, driven kilometers, fuel, number of owners, location and at last target variable Price of the car. This data is to give you a hint about important variables in used car model. You can make changes to it, you can add or you can remove some columns, it completely depends on the website from which you are fetching the data.

Try to include all types of cars in your data for example- SUV, Sedans, Coupe, minivan, Hatchback.

Model Building :

After collecting the data, you need to build a machine learning model. Before model building do all data pre-processing steps. Try different models with different hyper parameters and select the best model.

Follow the complete life cycle of data science. Include all the steps like

1. Data Cleaning
2. Exploratory Data Analysis
3. Data Pre-processing
4. Model Building
5. Model Evaluation
6. Selecting the best model

- **Data Sources and their formats**

Data sources for our model are online Car Re-sale websites like Car24. We scraped our data using the web scrapping tool Selenium. We scrapped almost 5000 rows of data and then saved all data into different excel files so that it can be used for model building.

After saving into excel files we load it into our Jupyter notebook and now let's have a look on data formats –

```
In [4]: # Now combining all the dataframes
data = pd.concat([dataframe1,dataframe2,dataframe3,dataframe4,dataframe5,dataframe6],ignore_index=True)
data
```

Out[4]:

	Model	History	Owner	Kms_Driven	Fuel_Type	Transmission	Registration	Insurance	Price
0	2021 Hyundai NEW SANTRO 1.1 550KM/AUTOMATIC	Non-Accident	1st	3084 K	Petrol	AUTOMATIC	DL-60 a-xxxx	Zero Depreciation	₹9,04,750
1	2020 Maruti Suzuki DZIX 1.2 K12 388KM/	Non-Accident	1st	3920 K	Petrol	MANUAL	HR-23 a-xxxx	Third party	₹5,00,500
2	2017 Ford Eco 1.2 V AT AUTOMATIC	Non-Accident	1st	3927 K	Petrol	AUTOMATIC	TL-30 a-xxxx	Comp	₹6,67,150
3	2020 Suzuki Swift 1.0 699KM/	Non-Accident	1st	3942 K	Petrol	MANUAL	CH-15 a-xxxx	Comp	₹3,08,300
4	2019 Hyundai Elite i20 Ngrin Executive 1.2 388KM/	Non-Accident	1st	3984 K	Petrol	MANUAL	DL-60 a-xxxx	Zero Dep	₹6,25,000
...
5264	2015 Hyundai Elite i20 SPORTZ (i) 1.4 388KM/	Non-Accident	1st	5511 K	Diesel	MANUAL	UP-32 a-xxxx	Comp	₹4,33,000
5265	2017 Tata NEXON 1.2+ 1.5 388KM/	Non-Accident	1st	5710 K	Diesel	MANUAL	UP-32 a-xxxx	Zero Dep	₹6,20,500
5266	2019 Hyundai Venue 1.4 CRDi EX MT	Non-Accident	1st	5754 K	Diesel	Auto	UP-32 a-xxxx	Full	₹7,02,500
5267	2018 Maruti Suzuki Swift AUTOMATIC	Non-Accident	1st	6015 K	Diesel	AUTOMATIC	UP-32 a-xxxx	Comp	₹5,35,000
5268	2011 Toyota Innova 2.5 V 7510 699KM/	Non-Accident	1st	20120 K	Diesel	MANUAL	UP-13 a-xxxx	Comp	₹4,75,000

After loading the data we checked data formats -

```
In [13]: 1 # checking datatypes and other information
2 dataset.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4790 entries, 0 to 4789
Data columns (total 9 columns):
 #   Column              Non-Null Count  Dtype  
---  --   ---
 0   Model               4790 non-null   object  
 1   History             4789 non-null   object  
 2   Owner               4790 non-null   object  
 3   Kms_driven          4790 non-null   float64  
 4   Fuel_Type           4789 non-null   object  
 5   Transmission        4549 non-null   object  
 6   Registration        4789 non-null   object  
 7   Insurance           3626 non-null   object  
 8   Price              4790 non-null   object  
dtypes: float64(1), object(8)
memory usage: 356.9+ KB
```

- There is 2 type of data present - 1.) Float 2.) Object in nature
- Insurance and Transmission features contains null values & 1 row is present in dataset with null values in all features.

Only Kms_driven is a continuous variable and rest of all are categorical variables.

• Data Preprocessing Done

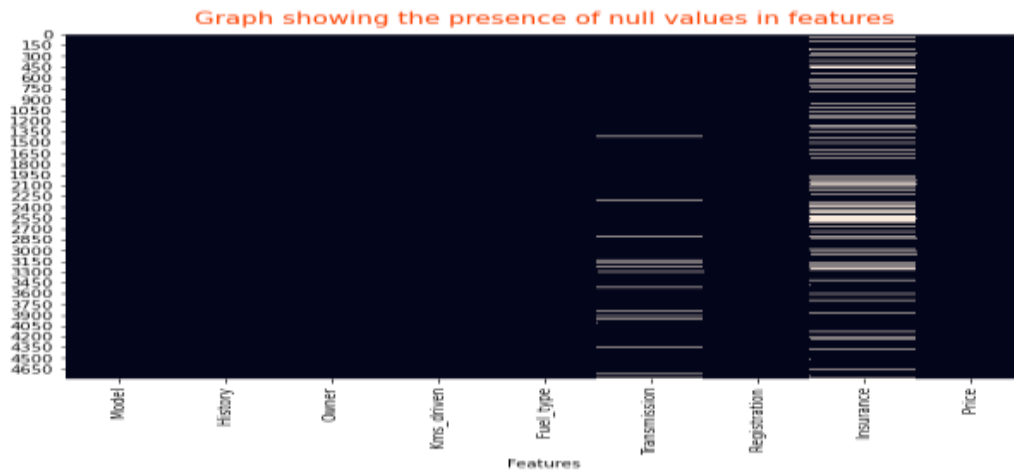
1. Checking Null Values:

Now after loading the data and checking the formats we moved to data preprocessing step. In this step we have to change data into proper formats and then do some data cleaning by treating null values.

We have used heatmap visualization technique to see the presence of null values in our dataset along with some method to find number of null values.

```
In [34]: 1 # Checking null values using heatmap
2 for i in dataset:
3     if dataset[i].isnull().any():
4         print("{} -> Having null values = {}".format(i, dataset[i].isnull().sum()))
5
6 # Heatmap
7 plt.figure(figsize=(10,8))
8 sea.heatmap(dataset.isnull(),cbar=False)
9 plt.title("Graph showing the presence of null values in features",color="orange",pad=10,fontsize=15,loc="center")
10 plt.xlabel("Features")
11 plt.show()

History -> Having null values = 1
Owner -> Having null values = 1
Kms_driven -> Having null values = 1
Fuel_Type -> Having null values = 1
Transmission -> Having null values = 241
Registration -> Having null values = 1
Insurance -> Having null values = 1264
```



Using the graph we can see that null values present in Transmission and Insurance variables.

2. Data Cleaning :

For data cleaning ,first we drop the duplicated values present in dataset.

```
1 1 # After combining data let's check duplicated data
2 dataset.duplicated().sum()
3 419
```

Model Variable -After removing the duplicated values we splited the Model variable into many useful variable which gives exact information about one character.

```
In [19]: 1 # Now using split function finding out the followings
2
3 dataset['Make_yr']=dataset['Model'].str.split().str[0] # make year of vehicle
4 dataset['Brand_name']=dataset['Model'].str.split().str[1] # brand to which vehicle belongs
5 dataset['model']=dataset['Model'].str.split().str[2] # model of vehicle
6 dataset['Variant']=dataset['Model'].str.split().str[3] # Variant of vehicle

In [20]: 1 # Let's convert make_yr into interger
2 dataset['Make_yr']=dataset['Make_yr'].astype(int)
3
4 # We can find one more feature showing the age of vehicle
5 dataset['Years_old']=2021-dataset['Make_yr']
```

Transmission Variable -We find out the most frequent value in Transmission variable and fill null values with that.

manual transmission is most preferred.

```
In [38]: 1 # There are only 2 types of transmissions either automatic or manual.
2 # So filling null values with most repeated type of transmission
3 dataset['Transmission'].fillna("MANUAL",inplace=True)
```

Registration Variable – In this variable some useless data was present ,so using the split method we extract only first 2 digit data which represent to which State car belongs.

```
In [32]: 1 # Based on the RTO we will extract the first 2 alphabets as these have some effect on the price
2 dataset['RTO']=dataset['Registration'].str.split("-").str[0]
3
4 dataset['RTO'].unique()

Out[32]: array(['DL', 'HR', 'CH', 'UP', 'PB', 'AS', 'GJ', 'WB', 'AP', 'MH', 'UK',
               'HP', 'KL', 'WA', 'HJ', 'KA'], dtype=object)
```

Insurance variable -Now we move towards Insurance variable and find out that there are some values with same meaning but different representation. Firstly we treated these values and again based on the most frequent value fill the null values.

```
In [34]: 1 # Check unique values in insurance_type feature
2 dataset['Insurance'].unique()

Out[34]: array(['Zero Depreciation', 'Third_party', 'Comp', 'Zero Dep', nan,
               'Third Party', 'Comprehensive'], dtype=object)

In [35]: 1 # Create dictionary to replace all values with these 3 values
2 val_to_replace={'Zero Depreciation':'Zero Depreciation','Zero Dep':'Zero Depreciation',
3               'Comprehensive':'Comprehensive','Comp':'Comprehensive',
4               'Third Party':'Third Party','Third_party':'Third Party','3rd Party':'Third Party'}
5
6 # Mapping values to be replaced in list
7 dataset['Insurance']=dataset['Insurance'].map(val_to_replace)

In [37]: 1 # Let's see null values in Insurance feature
2 dataset['Insurance'].isnull().sum()

Out[37]: 1163

In [38]: 1 # Fill null values using the most repeated character
2 dataset['Insurance'].fillna("Third Party",inplace=True)
```

SalePrice Variable – In this variable due to some sign like “,” ,”₹” it was not giving us continuous values. So we removed these signs and convert into continuous variable.

```
In [39]: 1 # As we have Re-sale price mentioned in our dataset and it is object in nature.
2 # Convert it into number(int,float)
3 price=[]
4 for i in dataset['Price']:
5     price.append(i.replace("₹","").replace(",",""))
6
7 dataset['SalePrice']=price
8 dataset['SalePrice']=dataset['SalePrice'].astype(float)
```

So we have completed the process of data cleaning. Now will drop some unnecessary columns.

```
In [55]: 1 # Copying the dataset
2 model_dataset=dataset.copy()
3
4 # removing unnecessary columns
5 model_dataset.drop(columns=['Make_yr','Model','Price','Registration','History'],inplace=True)
```

3. Checking for Outliers :

```
In [68]: 1 # Checking outliers
2 model_dataset.plot.box(figsize=(12,6))
3 plt.title("BoxPlot for Outliers",color='indianred',fontsize=18,pad=11)
4 plt.show()
```



Kms_driven and SalePrice variables contains outliers. Let's treat them using Zscore method. But before treating them convert categorical data into numerical data.

We are using Ordinal Encoding method to convert categorical data into numerical data –

```
In [61]: 1 from sklearn.preprocessing import OrdinalEncoder
2
3 # Using ordinal encoder to encode categorical data into numerical
4 ol=OrdinalEncoder()
5 for i in cols:
6     model_dataset[[i]]=ol.fit_transform(model_dataset[[i]])

In [62]: 1 # checking outliers using zscore method
2 from scipy.stats import zscore
3 z=np.abs(zscore(model_dataset))
4 z

Out[62]: array([[0.5682677, 1.42411539, 0.73644487, ..., 2.16054217, 1.22487635,
0.0674221 ],
[0.5682677, 1.41100608, 0.73644487, ..., 1.82642836, 0.35973705,
0.47848498],
[0.5682677, 1.41083162, 0.73644487, ..., 0.82408693, 1.22487635,
0.50203136],
...,
[0.5682677, 0.07330777, 1.35787489, ..., 1.49231455, 1.94730109,
0.86357595],
[0.5682677, 0.13908305, 1.35787489, ..., 1.15820074, 1.94730109,
0.37689848],
[0.5682677, 5.16132098, 1.35787489, ..., 1.10059592, 1.94730109,
0.09383772]])

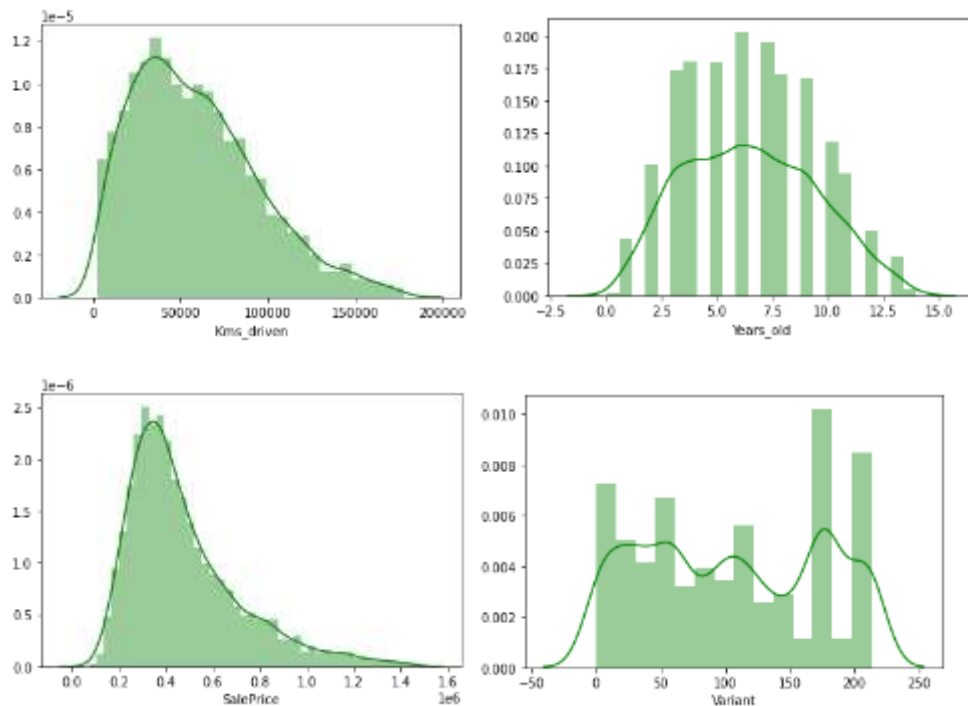
In [65]: 1 # Checking percentage loss of data
2 print("Percentage loss of data in treating outliers: {}%".format(round((model_dataset.shape[0]-card
4
Percentage loss of data in treating outliers: 3.38%
```

Removed 3.3% of outliers from dataset.

4. skewness in data –

```
In [66]: 1 # Checking skewness
2 cardata_new.skew()

Out[66]: Owner          1.189748
Kms_driven             0.715876
Fuel_type             -0.685095
Transmission          -2.831751
Insurance              0.040706
Brand_name             0.146087
model                 0.002558
Variant               0.059977
Years_old              0.181510
RTO                   0.212382
SalePrice              1.298426
dtype: float64
```

Treated skewness using the PowerTransformer() technique-

```
In [68]: 1 # Removing skewness
2 from sklearn.preprocessing import PowerTransformer
3
4 pt = PowerTransformer()
5
6 cardata_new['Transmission'] = pt.fit_transform(cardata_new[['Transmission']])
7 cardata_new['Owner'] = pt.fit_transform(cardata_new[['Owner']])
8 cardata_new['Brand_name'] = pt.fit_transform(cardata_new[['Brand_name']])
```

5. Scaling Data :

Now in the last step of data preprocessing we scale our dataset into equal range because sometimes our model can be biased towards higher values present in dataset to eliminate this risk we are using MinMaxScaler to scale our data.

```
In [71]: 1 # now rescaling our data so that model will not be biased
2 from sklearn.preprocessing import MinMaxScaler
3 min_scale = MinMaxScaler()
4
5 for i in cardata_new.columns:
6     cardata_new[i] = min_scale.fit_transform(cardata_new[i])

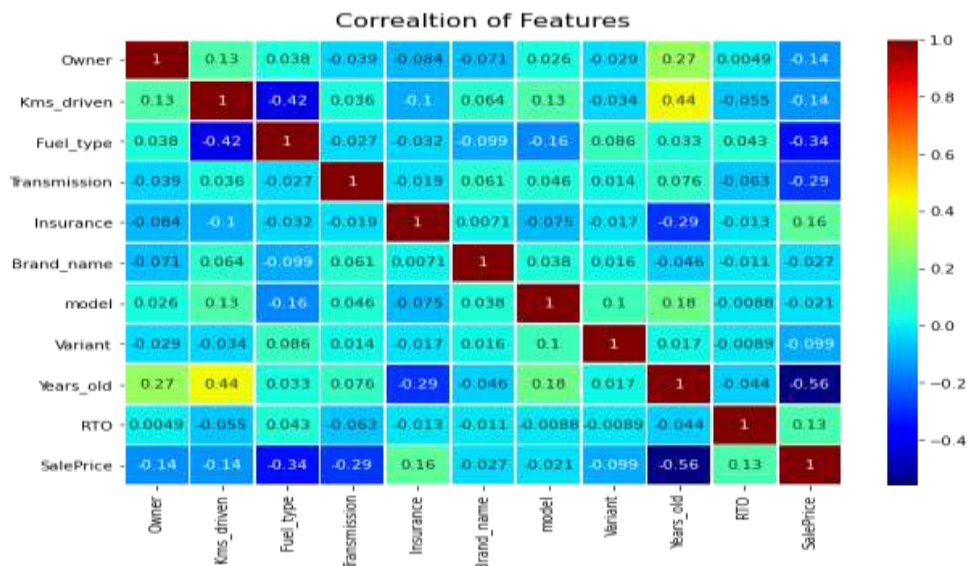
In [72]: 1 cardata_new.head(5)

Out[72]:
```

	Owner	kms_driven	Fuel_type	Transmission	Insurance	Brand_name	model	Variant	Years_old	RTD	SalePrice
0	0.0	0.000445	1.0	0.0	1.0	0.373034	0.541303	0.071362	0.000000	0.153040	0.298101
1	0.0	0.009432	1.0	1.0	0.5	0.686604	0.112782	0.262011	0.071429	0.384815	0.423083
2	0.0	0.009472	1.0	0.0	0.0	0.327950	0.458647	0.000390	0.280714	0.153040	0.429560

- Data Inputs- Logic- Output Relationships

```
In [69]: 1 # Plotting correlational plot
2 plt.figure(figsize=(10,7))
3 sns.heatmap(cardata_new.corr(), annot=True, fmt='.2g', linewidths=0.5,
4             linecolor='white', cmap='jet')
5 plt.title("Correlation of Features", fontsize=16, pad=11)
6 plt.show()
```



Observations:

SalePrice having 50% negative correlation with Years_old feature which indicated that more the old less will be the SalePrice of car and vice-versa.

Variant is least correlated with SalePrice which indicated less fluctuations in SalePrice based on Variant of car.

Years_old and Kms_driven are sharing positive correlation of 40% which is quite realistic as more the years more will be the distance travelled.

- Hardware and Software Requirements and Tools Used

There is no such requirement for hardware ,but I have used intel i5 8th generation processor.

Software: Jupyter Notebook (Anaconda 3)

Language : Python 3.9

Libraries used in project:

- a. Pandas
- b. Numpy
- c. Matplotlib
- d. Seaborn
- e. Sklearn
- f. scipy
- g. selenium
- h. time

Model/s Development and Evaluation

- Splitting Dataset

First of all we are splitting data into dependent and independent variable.

```
In [73]: 1 X=cardata_new.iloc[:, :-1]
          2 y=cardata_new['SalePrice']
          3 print("Splitted dependent & Independent data successfully.")
          Splitted dependent & Independent data successfully.
```

Now will split data into training and testing data using train_test_split method –

```
In [75]: 1 # now lets provide the best random state to the model
          2 x_train,x_test,y_train,y_test=train_test_split(X,y,test_size=0.25,random_state=final_random_state)
          3
          4
          5 # Success
          6 print("Training and testing split was successful.")
          Training and testing split was successful.
```

- Identification of possible problem-solving approaches (methods)

As we already mentioned above we have continuous target variable so based on the type of variable we are using Regression approach for our model and will use some algorithms like

- a) Linear Regression
- b) Ridge Regression
- c) Support Vector Regressor
- d) DecisionTree Regressor
- e) KNeighbors Regressor
- f) RandomForest Regressor
- g) AdaBoost Regressor
- h) GradientBoosting Regressor

- Testing of Identified Approaches (Algorithms)

After initiating the instances for our model we will create a function which will give us all the metrics score we are using in our model building by passing the model into it.

```
1 # Importing the required libraries
2
3 # Model = {} # algorithm used
4 # R2Score = {} # R2 score of algorithm
5 # MAE = {} # Mean absolute error
6 # MSE = {} # Mean squared error
7 # RMSE = {} # Root mean squared error
8 # CVScore = {} # Mean of cross val score
9 # Std = {} # Standard deviation in cross val
10
11
12 def r2score(model):
13     print(f"Model: {model}")
14     Model.append(str(model))
15     # Training score
16     model.fit(x_train, y_train)
17     print(f"Training score: {model.score(x_train, y_train)*100}")
18     y_pred = model.predict(x_test)
19     # R2 score value
20     r2 = r2_score(y_test, y_pred)
21     print(f"R2 score = {round(r2*100, 2)}%")
22     R2Score.append(round(r2*100, 2))
23     # Mean absolute error
24     mae = mean_absolute_error(y_test, y_pred)
25     print(f"Mean absolute error = {mae}")
26     MAE.append(mae)
27     # Mean squared error
28     mse = mean_squared_error(y_test, y_pred)
29     print(f"Mean squared error = {mse}")
30     MSE.append(mse)
31     # Root mean squared error
32     rmse = np.sqrt(mean_squared_error(y_test, y_pred))
33     RMSE.append(np.sqrt(mean_squared_error(y_test, y_pred)))
34
35     # Cross validation
36     cv_score = cross_val_score(model, x, y, cv=5, scoring='r2')
37     CVScore.append(round(cv_score.mean()*100, 2))
38     # Standard deviation
39     print(f"Standard deviation = {2*std(cv_score)}")
40     std.append(cv_score.std())
41     print(f"\n", 30*"=")
```

- Run and Evaluate selected models

Evaluations of Algorithms

```
In [86]: 1 # Table view of result of each metrics from above algorithms
2 evaluations = pd.DataFrame({"Model":Model,"R2_Score":R2Score,"MAE":MAE,"MSE":MSE,
3                             "RMSE":RMSE,"CV_Score":CVScore,"Std_dev":Std})
4 evaluations

Out[86]:
```

	Model	R2 Score	MAE	MSE	RMSE	CV Score	Std_dev
0	LinearRegression()	54.051	0.087556	0.013420	0.115848	46.353	0.057337
1	Ridge()	54.049	0.087512	0.013421	0.115849	46.367	0.057066
2	DecisionTreeRegressor()	60.807	0.048438	0.005606	0.074871	72.461	0.061656
3	SVR()	68.391	0.069881	0.009232	0.096063	61.930	0.040841
4	RandomForestRegressor()	60.060	0.035693	0.002903	0.053682	64.536	0.025321
5	KNeighborsRegressor()	66.042	0.066268	0.009918	0.099500	55.130	0.046838
6	AdaBoostRegressor()	50.929	0.096763	0.014332	0.119717	41.175	0.131246
7	GradientBoostingRegressor()	82.658	0.049714	0.005065	0.071169	78.761	0.015046

After sending the model into our created function, it created list for us and using these list we created a dataframe of all algorithms and metrics which we used for training.

Among these algorithms ,RandomForest Regressor and GradientBoosting Regressor are giving better result then others. So we will tune them using RandomizedSearchCV HyperParameter tunnig technique and find best result out of them.

• Hyperparameter tuning –

1). RandomForest Regressor:

```
In [88]: 1 rf_rdcv = RandomizedSearchCV(rf_reg, param_distributions,
2                                   scoring='r2', n_jobs=-1, verbose=2, cv=5)
3 rf_rdcv.fit(x_train,y_train)

Fitting 5 folds for each of 10 candidates, totalling 50 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done 25 tasks | elapsed: 3.4min
[Parallel(n_jobs=-1)]: Done 50 out of 50 | elapsed: 4.2min finished

Out[88]: RandomizedSearchCV(cv=5, estimator=RandomForestRegressor(), n_jobs=-1,
                             param_distributions={'bootstrap': [True],
                                                  'criterion': ['mse', 'max'],
                                                  'max_depth': [80, 90, 100, None],
                                                  'max_features': ['auto', 'sqrt'],
                                                  'min_samples_leaf': [1, 2, 4],
                                                  'min_samples_split': [8, 16, 32],
                                                  'n_estimators': [100, 200, 500, 1000]},
                             scoring='r2', verbose=2)

In [89]: 1 print(rf_rdcv.best_params_)
2 print(rf_rdcv.best_score_)

{'n_estimators': 1000, 'min_samples_split': 8, 'min_samples_leaf': 4, 'max_features': 'auto', 'max_depth': None, 'criterion':
'mse', 'bootstrap': True}
0.6174880153313118

In [90]: 1 # let's now try best parameter
2 rfr = RandomForestRegressor(n_estimators=100, min_samples_split=8, min_samples_leaf=4,
3                             max_features='auto', max_depth=90, criterion='mse',
4                             bootstrap=True, random_state=final_random_state)
5 rfr.fit(x_train,y_train)
6 print("Training score = %.2f"%rfr.score(x_train,y_train))
7 rfr_pred = rfr.predict(x_test)
8 print("Coefficient of determination = %.2f"%r2_score(y_test,rfr_pred))
9 print("\nMean Absolute error = %.2f"%mean_absolute_error(y_test,rfr_pred))
10 print("\nMean Squared error = %.2f"%mean_squared_error(y_test,rfr_pred))
11 cvscore_val_score(rfr, X,y, cv=5,scoring='r2')
12 print("\nCv valid. score = %.2f"%cv.mean())
13 print("Standard deviation = %.2f"%cv.std())

Training score = 0.94
Coefficient of determination = 0.89

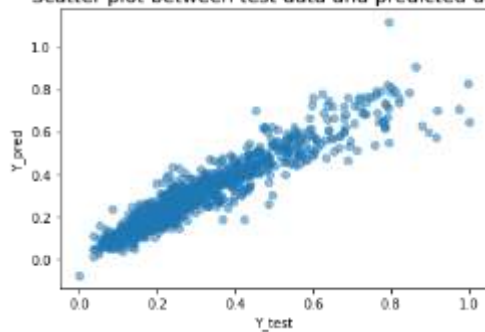
Mean Absolute error = 0.038
Mean Squared error = 0.003

Cv valid. score = 0.82
Standard deviation = 0.03
```

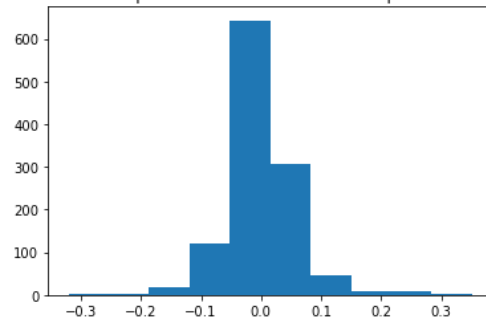
2). GradientBoostingRegressor:

```
In [94]: 1 gbr_cv = RandomizedSearchCV(gbr_reg, param_distributions=gradient_parameters,
2   scoring='r2', n_jobs=-1, verbose=0, cv=5)
3
4 gbr_cv.fit(x_train,y_train)
5
6 Fitting 5 folds for each of 10 candidates, totalling 50 fits
7
8 [Parallel(n_jobs=-1)]: Using backend joblib backend with 8 concurrent workers.
9 [Parallel(n_jobs=-1)]: from 35 out of 50 | elapsed: 2.0min remaining: 50.2s
10 [Parallel(n_jobs=-1)]: Done 50 out of 50 | elapsed: 11.5min finished
11
12 Out[94]: RandomizedSearchCV(cv=5, estimator=GradientBoostingRegressor(), n_jobs=-1,
13   param_distributions={'criterion': ['mse', 'mae'],
14   'learning_rate': [1, 0.5, 0.2, 0.1,
15   0.05, 0.02, 0.01,
16   0.001],
17   'max_depth': range(5, 15, 2),
18   'max_features': ['auto', 'sqrt'],
19   'min_samples_split': range(200, 1000, 200),
20   'n_estimators': [100, 500, 1000, 2000,
21   10000]},
22   scoring='r2', verbose=2)
23
24 In [95]: 1 print(gbr_cv.best_params_)
25 2 print(gbr_cv.best_score_)
26
27 ('n_estimators': 10000, 'min_samples_split': 400, 'max_features': 'auto', 'max_depth': 7, 'learning_rate': 0.5, 'criterion': 'mse')
28 0.8752076586866666
29
30 In [96]: 1 # let's now try best parameter
31 2 gbr = GradientBoostingRegressor(n_estimators=1000, min_samples_split=1000, learning_rate=0.5,
32   max_features='sqrt', max_depth=7, criterion='mse')
33 3 gbr.fit(x_train,y_train)
34 4 print("Training score = %.2f"%gbr.score(x_train,y_train))
35 5 gbr_pred = gbr.predict(x_test)
36 6 print("Coef. of determination = %.2f"%r2_score(y_test,gbr_pred))
37 7 print("Mean Absolute error = %.3f"%mean_absolute_error(y_test,gbr_pred))
38 8 print("Mean Squared error = %.3f"%mean_squared_error(y_test,gbr_pred))
39 9 cvscore_val_score(gbr, X,y, cv=5, scoring='r2')
40 10 print("Cross vald. score = %.2f"%cv.mean())
41 11 print("Standard deviation = %.2f"%cv.std())
42
43 Training score = 0.93
44 Coef. of determination = 0.89
45
46 Mean Absolute error = 0.039
47 Mean Squared error = 0.003
48
49 Cross vald. score = 0.84
50 Standard deviation = 0.03
51
52 In [97]: 1 # let's plot the plot of predicted and actual data
53 2 plt.scatter(x_test,gbr_pred, alpha=0.5)
54 3 plt.xlabel('y_test')
55 4 plt.ylabel('y_pred')
56 5 plt.title('Scatter plot between test data and predicted data',fontsize=15)
57 6 plt.show()
```

Scatter plot between test data and predicted data



Distribution plot of difference in actual and predicted data



After testing various algorithms and hypertunning the best two algorithms we find that both RandomForestRegressor and GradientBoostingRegressor are giving better results but we are finalising the GradientBoostingRegressor algorithm for our model as there is minimum difference in coefficient of determination and cross validation score and also less variance in result in the higher SalePrice.

The selected model gives coefficient of determination value =0.89 with minimum mean absolute error of 0.039 and cross validation score 0.84 with a least standard deviation of 0.03.

- **Finalising and Saving the Model**

```
In [99]: 1 import joblib
         2 joblib.dump(gbr, 'Carprice_prediction.obj') # model and filename

Out[99]: ['Carprice_prediction.obj']

In [100]: 1 gbr_model = joblib.load('Carprice_prediction.obj')
          2 gbr_model

Out[100]: GradientBoostingRegressor(criterion='mse', learning_rate=0.5, max_depth=11,
                                     max_features='sqrt', min_samples_split=1000)

In [101]: 1 # Predicting the values
          2 gbr_predictions=gbr_model.predict(X)
          3 print('Predictions of Random Forest Regressor: ',gbr_predictions)

Predictions of Random Forest Regressor: [0.3479031 0.41334629 0.44043025 ... 0.4482615 0.39320814 0.45234612]
```

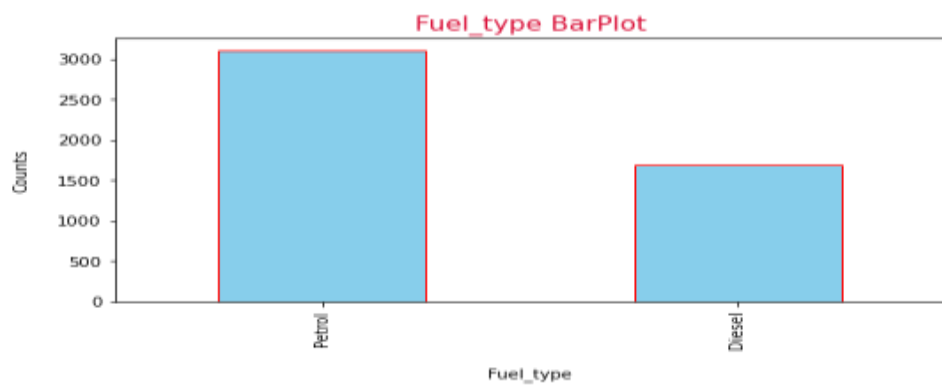
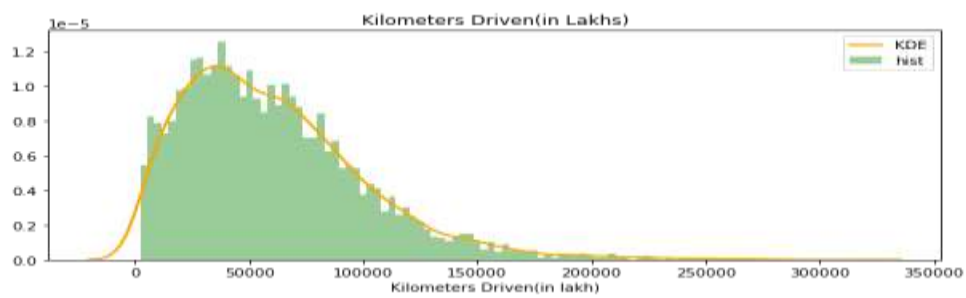
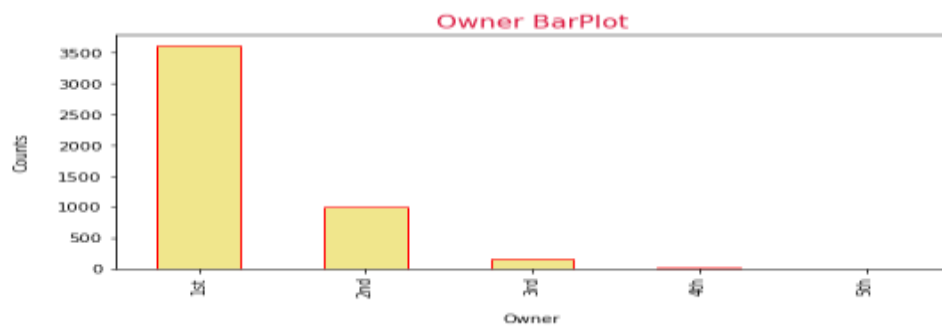
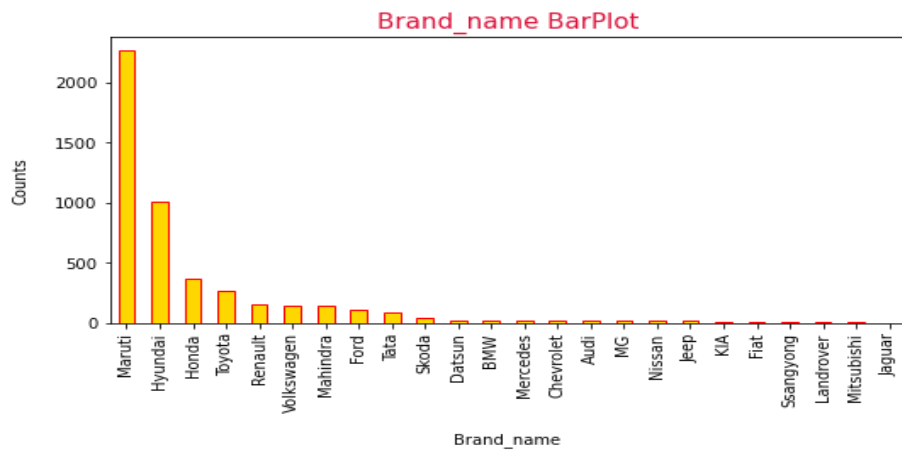
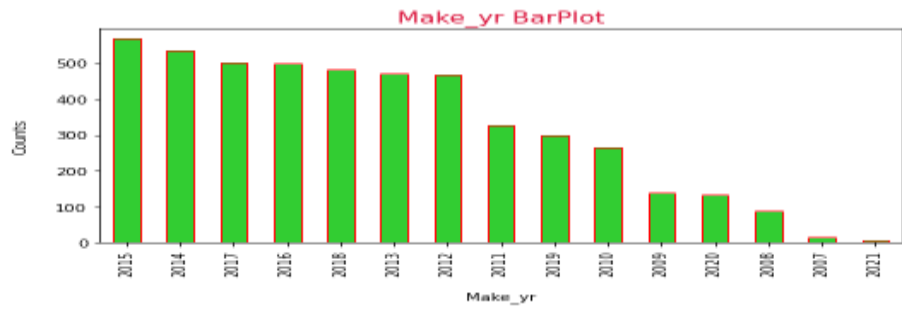
- **Key Metrics for success in solving problem under consideration**

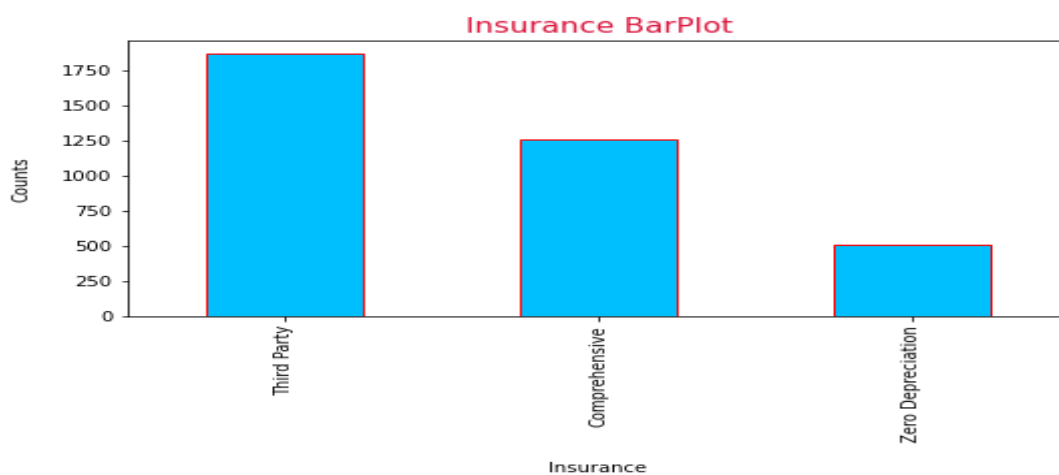
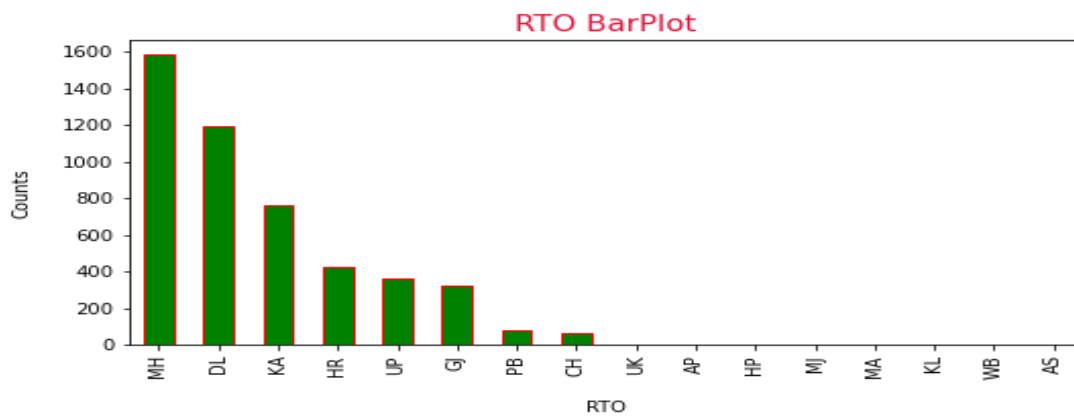
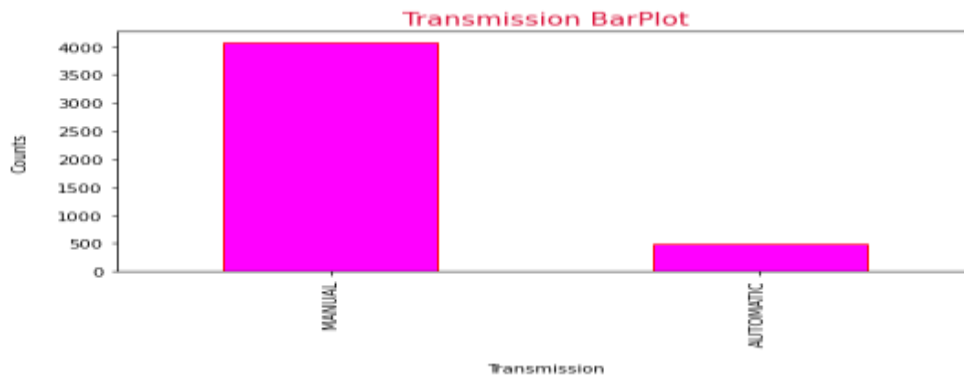
In this model we have used -

R2score ,
mean absolute error,
mean squared error,
rootmean squared error ,
cross validation score ,
Standard deviation to find the best result out of any algorithm.

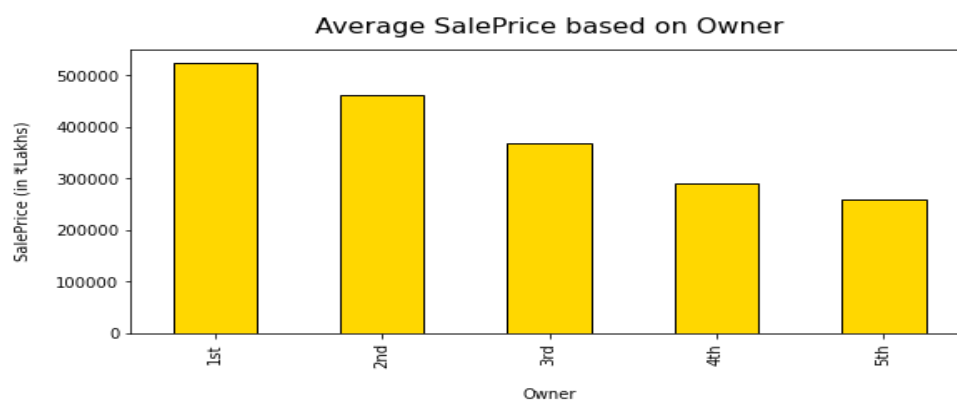
- **Visualizations**

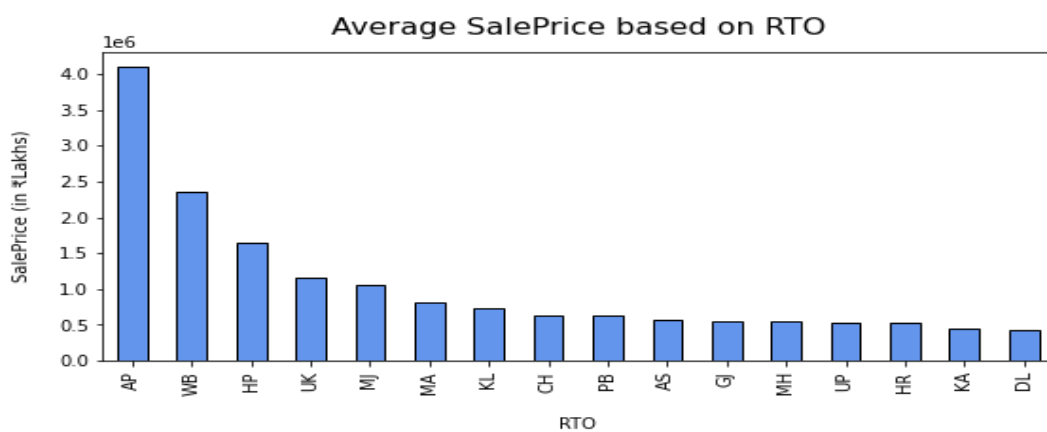
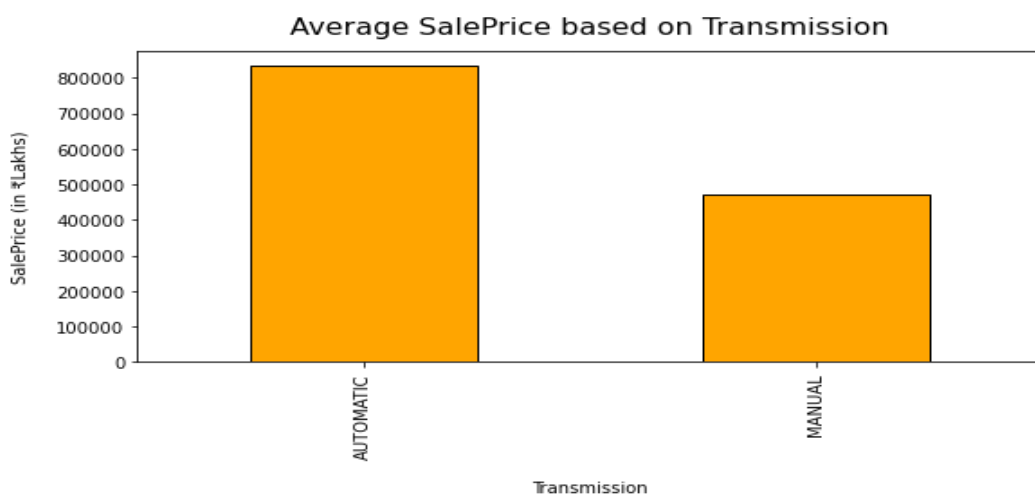
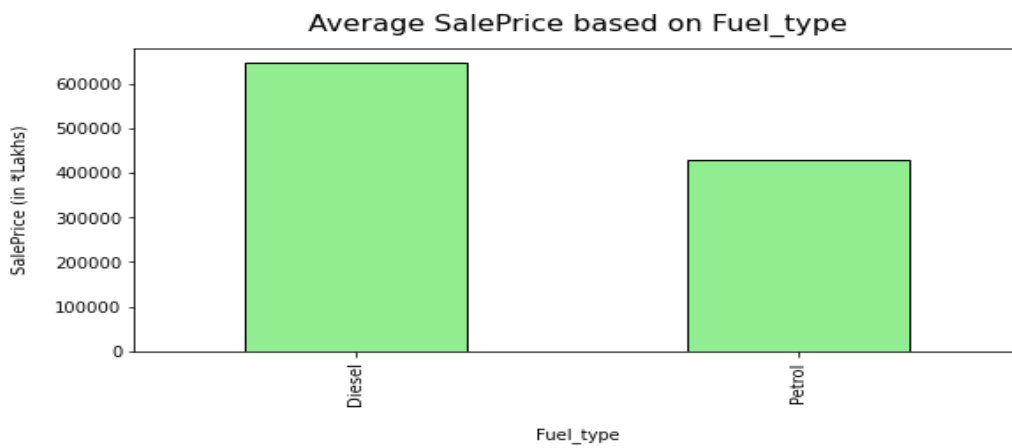
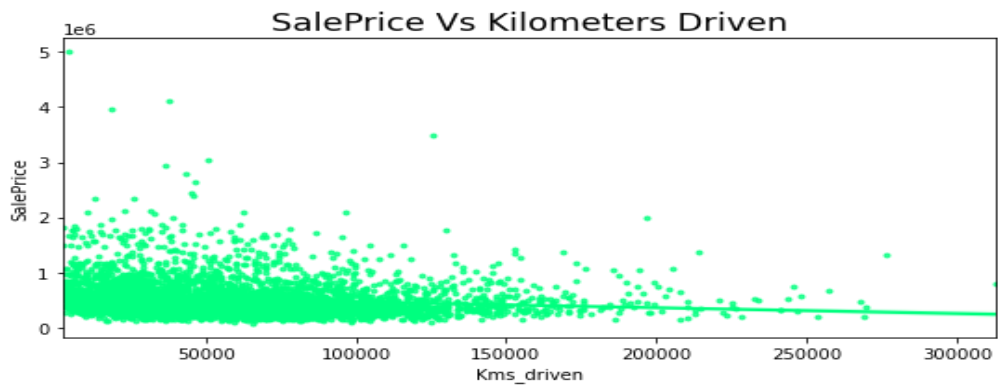
Checking data distribution in each variable –

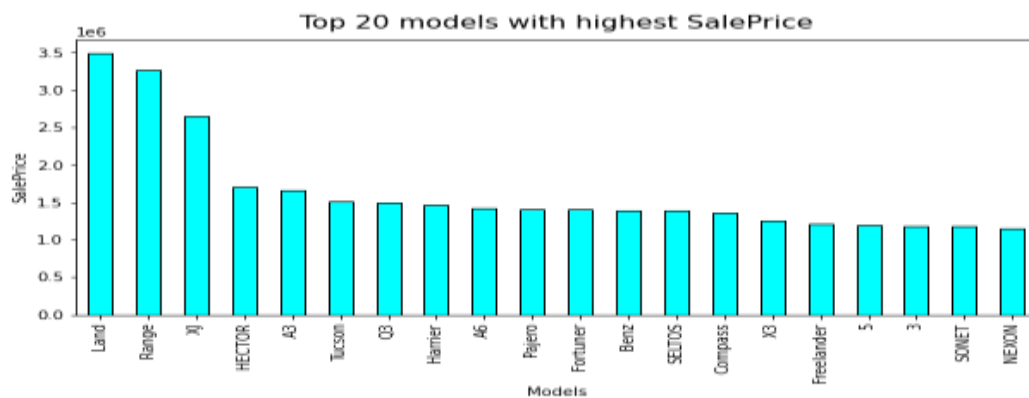
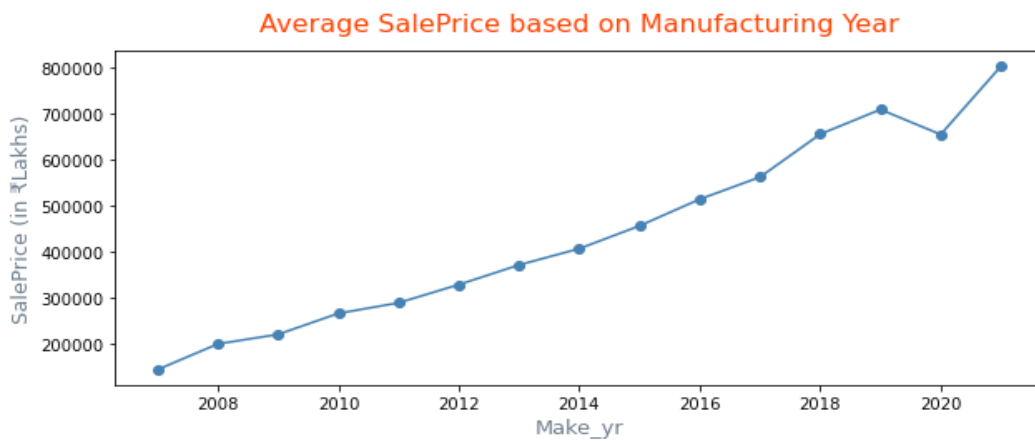
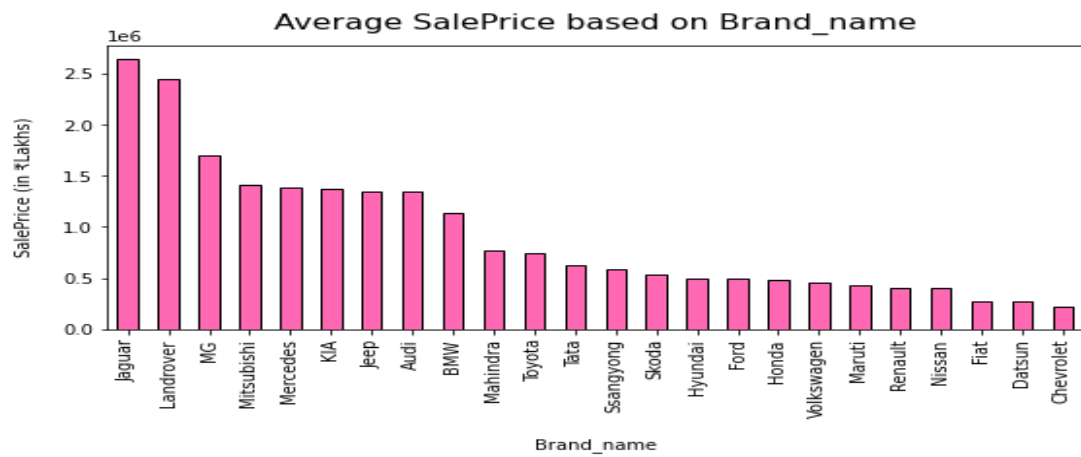
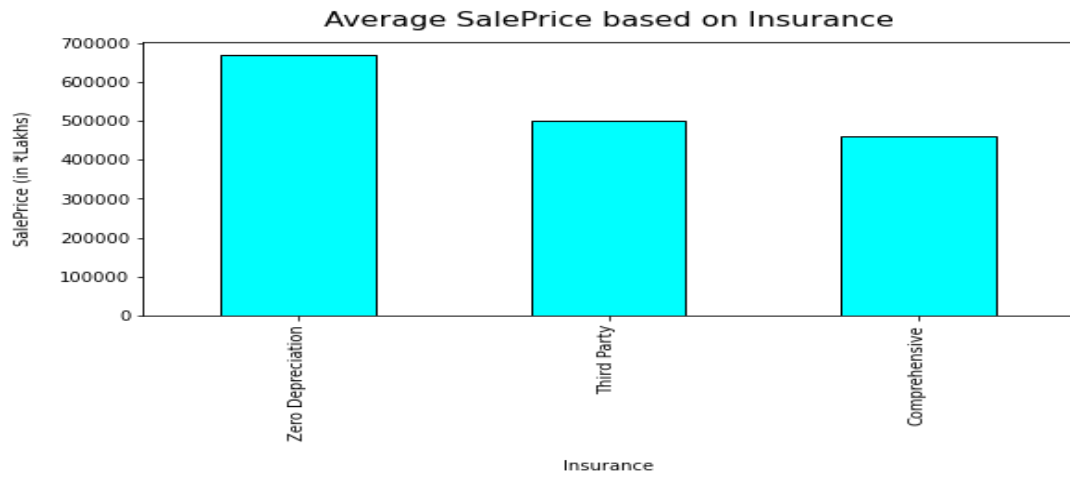


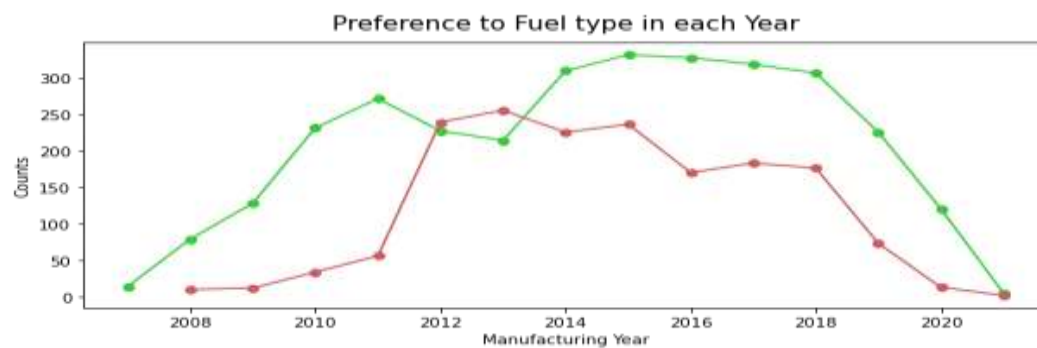
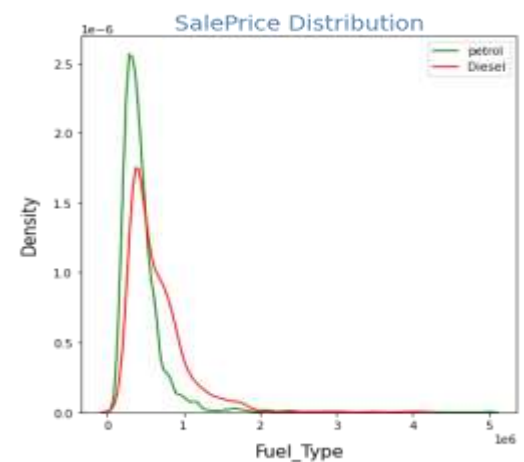
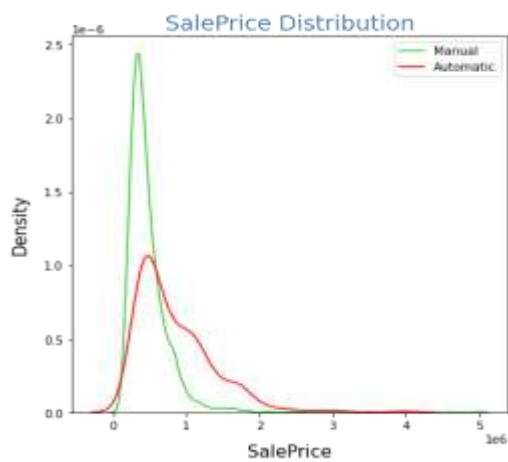
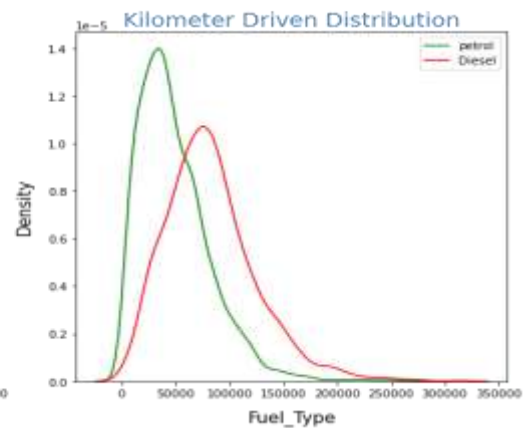
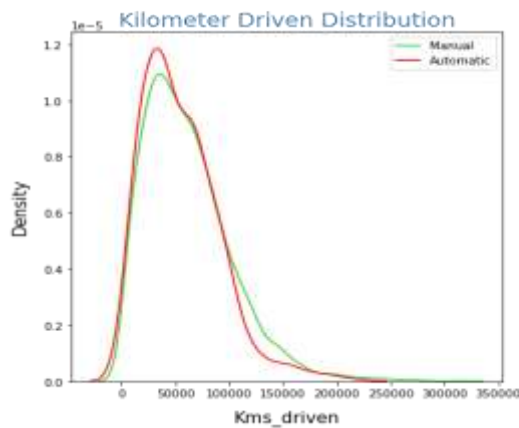
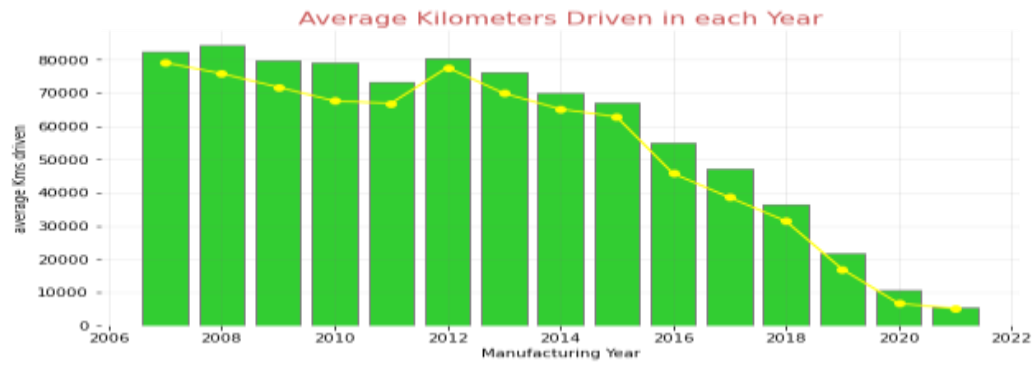


Now checking Average SalePrice based on each variable –









- Interpretation of the Results

1. Cars manufactured in between 2012-2018 are more in number for sale. While a few cars are for resale before 2010.
2. Maruti and Hundai are the most preferred brands for resale in indian market. These brands are sharing more than 50% of resale market share.
3. More likely the cars which are 3-9 years old are for resale.
4. History feature represents there is only "Non-Accidental" values. This feature not providing any valuable information, we can drop it.
5. In car resale market most of the cars for sale are first owners.
6. Minimum kilometers covered by a car is 2`259km`s while maximum kilometers covered is `312882km`. There is huge difference in minimum and maximum kilometers.
7. The distribution plot of Kilometers driven tells that the data gathered is not uniform. Some cars have travelled less while some very high. The difference between the mean and median shows the `presence of outliers` and the `rightskewness` is present in data as `mean is greater than median`.
8. Manual transmission is most preferred.
9. Resale of car in used car market depends on the number of owners present for the car as from the above graph it clearly shows the average saleprice of 1st and 5th owners have a lot of difference.
10. Based on the scatter plot we find that there is slightly negative relation of SalePrice upon Kilometers driven.
11. Another feature which impact the SalePrice is Fuel type. Cars in which Diesel is used are having high average saleprice then Petrol as Deisel cars are more costly than Petrol cars.
12. Based on the transmission ,we find out that cars with automatic transmission are having high SalePrice compared to Manual transmission. Again Automatic transmission are costly than Manual type of transmission.

13. Based on the RTO locations above graph showing that `Andhra Pradesh, West Bengal, Himachal Pradesh` are among the most valued for high Sale Price while Delhi is in least valued RTO locations.
14. In Comprehensive, any damage to vehicle, owner or third person all liabilities are recovered through insurance company. There is no need to bear all the charges own but in third-party insurance as clear with the name only covers for third party liabilities while in Zero Depreciation we have some premium add on.
15. So basically according to the need of vehicle owners Zero Depreciation is more valuable than others and having High sale Price than other types of insurance.
16. Based on the Brands we find out Average Sale Price of some premium brand are very high than other brands present in Indian market like `Jaguar, Landrover` and `Fiat, Datsun, Chevrolet` are the brands with least Sale Price.
17. In covid crisis, whole economy suffered a lot & due to this demand supply for some goods decreases. This decrease in demands sets the sale price at low indexes that's why we are creating a model for our client and this behaviour in Sale Price can be seen in above graph as the constant increase in Sale Price each year suddenly falls in 2020 but in 2021 is again starts to rise.
18. After modelling we find out that Gradient Boosting Regressor algorithm is giving best result than other. So we finalised this algorithm.
19. The selected model gives coefficient of determination value = 0.89 with minimum mean absolute error of 0.039 and cross validation score 0.84 with a least standard deviation of 0.03.

CONCLUSION

- Key Findings and Conclusions of the Study

After working on this project ,we got a lot of insights of how to collect data , preprocess it ,remove unwanted data and how to tackle null values to create a good model.

In this project we find various factors responsible for the sale price of a car. Let's take each variable to describe our observations from this project –

i. Brand Name –

After analysing the brand name variable we find that in Indian car market brands like “Maruti , Hyundai, Honda , Toyota” are covering most of the market while moving towards the Saleprice value we find that the premium brands like “Jaguar, LandRover , Mitsubishi, Mercedes , jeep , Audi , BMW” are among the most valued brands in Re-sale Car market.

ii. Make_yr Variable –

Analysing the Make_yr variable we find that the cars who have been manufactured in 2012-2018 are more in number for sale while at the average graph of SalePrice denotes that there is constant rise in SalePrice each year upto 2020, in 2020 due to covid crisis it comes down but after that again start to rise.

Model Variable –

“Rover , A3 , A6 ,Q3 ,Tuscon, Pajero , Benz’ are among the top model whose average sale price is high.

iii. Owner Variable –

Owner variable shows the number of owners while selling the car. So most of the data contains least owners which is 1 and a few data was present of more than 2 owners. While checking the effect of Owners on SalePrice we find that more the owners less will be the sale price.

iv. Kms_driven Variable –

In this variable we have a right skewness as some of the cars had travelled a lot while some very less. As some people travels more through public transport while some use their own vehicle to travel this cause variance in data. There is very high difference in minimum and maximum kilometers driven in dataset.

v. Fuel_type Variable –

Nowdays there are more than 2 types of fuel used in cars they may be petrol , diesel, electric or hybrid but the data we collected contains only petrol and diesel and based on this we do analysis and found out that people owning petrol cars are more in number but diesel cars are more costlier than petrol cars.

vi. Transmission Variable –

In transmission we have 3 options – Manual, Automatic, Hybrid. But the data we collected contains only Manual and Automatic type of transmission. As Automatic transmission is newly introduced in Indian Market a few years back so there is not much data based on this type of transmission. While finding the saleprice we noticed that Automatic transmission has higher Average SalePrice as compared to Manual.

vii. RTO Variable –

RTO variable define the State in which Road tax Office is present. And we find that Delhi , UttarPardesh are among the States with minimum Avg SalePrice while AndraPradesh and WestBengal are having maximum Avg. salePrice.

viii. Insurance Variable –

Insurance variable contains a lot of dirty data. First we cleaned it using the replace method and then it also have a lot of null values which were treated using the most frequent value.

After this we tried to find that most of the cars having Third Party insurance and the cars with Zero Depreciation are among the most

valued during the sale as in Zero depreication , all liabilities are covered by insurance company.

- **Limitations of this work and Scope for Future Work**

In this project we have to deal with a lot of problems :

- a) Websites are not well maintained as we get almost 10% of duplicated data.
- b) The data present on contains a lot of missing features which create hinderance in analysing the data properly.
- c) We have to face a lot of outliers in SaelPrice and transmission features.

Thankyou