**FLIP ROBO**

# CARPRICE ANALYSIS

Submitted by:

SNEHA SANTRA

# <u>ACKNOWLEDGMENT</u>

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.

# 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.

- ## Conceptual Background of the Domain Problem

  The project is sub-divided following section. These are:

  1. Loading necessary libraries
  2. Loading Dataset from a CSV file
  3. Summarization of Data to understand Dataset (Descriptive Statistics)
  4. Visualization of Data to understand Dataset (Plots, Graphs etc.)
  5. Processing the data for modeling
  6. skewness and outliers detection for better accuracy
  7. Build the model and select the right model and save it

- ## Review of Literature

  There are 14 columns including Name,Location,kilometer driven,fuel-type,engine and main feature is price.
  Most of them are object variable and some int variable and main features is float type.
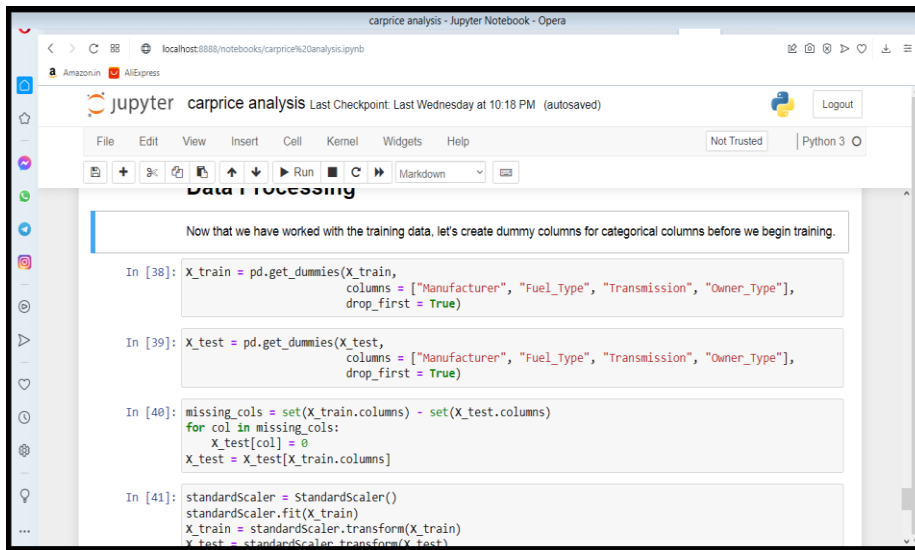
- Motivation for the Problem Undertaken

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
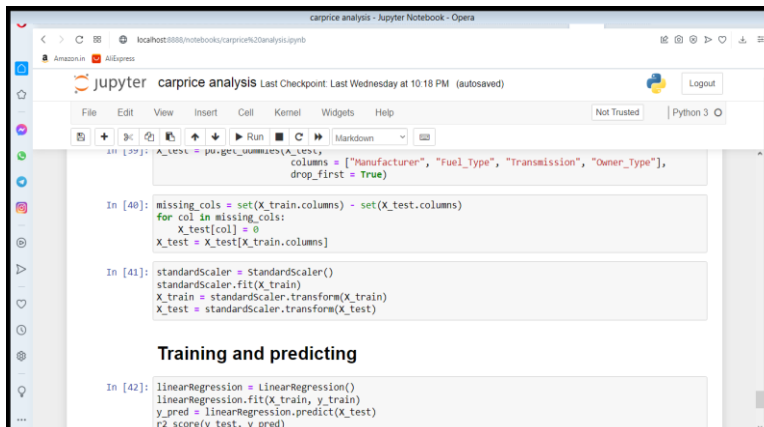4. Model Building
5. Model Evaluation
6. Selecting the best model

# Analytical Problem Framing

- Mathematical/ Analytical Modeling of the Problem
  <u>Data Processing:</u> We have worked with the training data, let's create dummy columns for categorical columns before we begin training<u>.</u>



And then scalling the dataset for model building
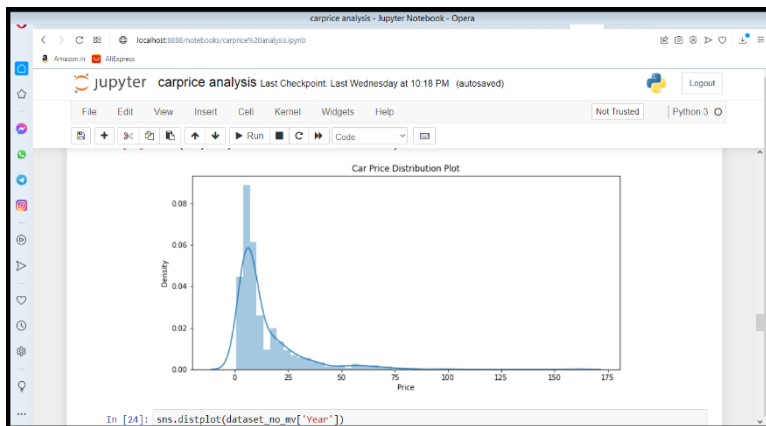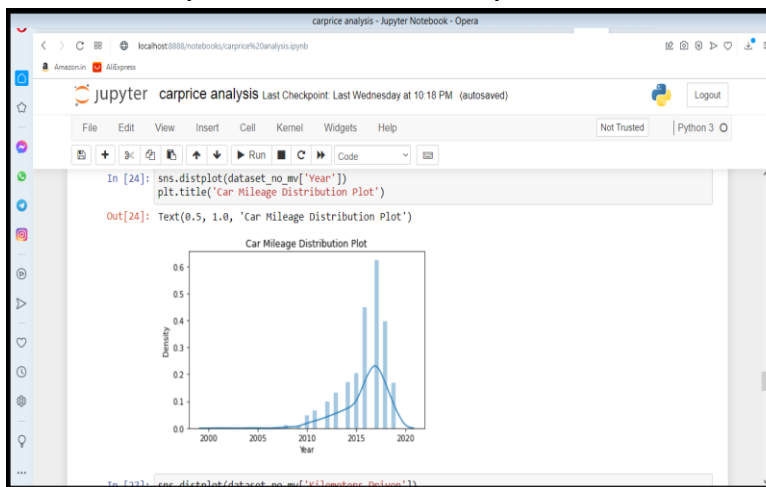


**EDA:**

This is the heatmap of the dataset



This is the relation between manufacturer and the count of cars

This is car price distribution plot



Car Mileage distribution plot

- Data Sources and their formats

The data is collected from various site like cardekho,olx and scrapped the data and make a dataset. There are 14 columns and 6019 rows.

Target variable is price of the cars.

- **Data Inputs- Logic- Output Relationships**

  Describe the relationship behind the data input, its format, the logic in between and the output. Describe how the input affects the output.

- **State the set of assumptions (if any) related to the problem under consideration**

  Here, you can describe any presumptions taken by you.

- **Hardware and Software Requirements and Tools Used**
  Here we use lots of liaberies like pandas,numpy,matplot,seaborn, and we use python language for the coding  purpose and import some other metrics liaberies also for model building like sklearn metrics ,Linear Regression,Random Forest Regressor  ,accuracy report etc.
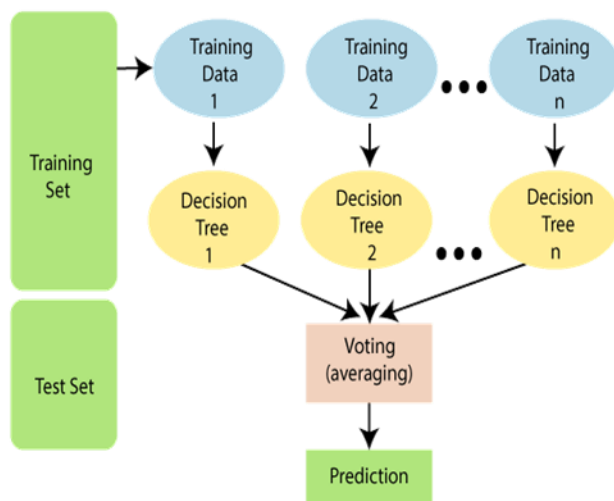
# Model/s Development and Evaluation

- Identification of possible problem-solving approaches
  (methods)
  We use randomforest classifier and Linear Regression for the model
  building as it is regression problem and finally we save the
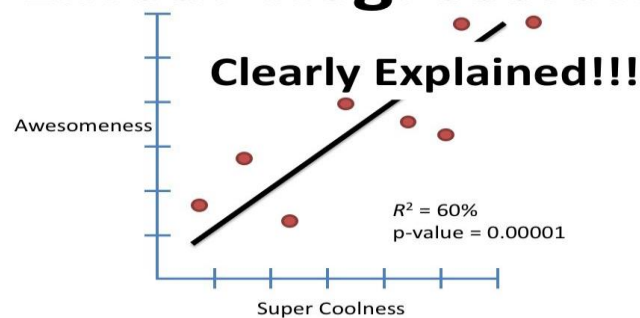  RandomForest Classifier model for gd accuracy.

- Testing of Identified Approaches (Algorithms)
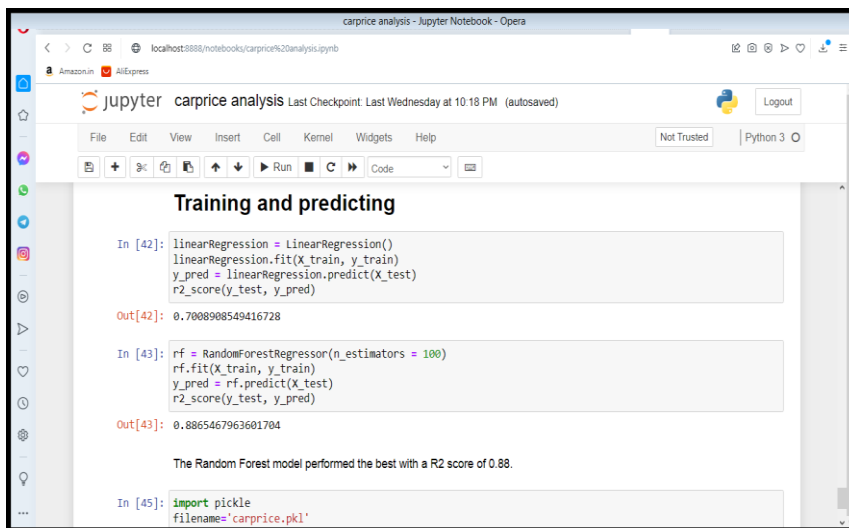  Random Forest Classifier:



Linear Regression:

- Run and Evaluate selected models



We save the best model and this is RandomForest Classifier model.

- Interpretation of the Results



And we save Random Forest Classifier model and accuracy is so good.

# CONCLUSION

- ## Key Findings and Conclusions of the Study

After the Final Submission of test data, my accuracy score was  88%

Feature engineering helped me increase my accuracy.

Amazingly Random Forest Classifier  worked better than all other Ensemble models.

- ## Learning Outcomes of the Study in respect of Data Science

Data science comes as a very important tool to solve problems in the domain to help the companies increase their overall revenue, profits, improving their marketing strategies and focusing on changing trends in car  sales and purchases. Predictive modelling, Market mix modelling, recommendation systems are some of the machine learning techniques used for achieving the business goals for car companies.