

# Data Science with Python Career Program - Capstone Project

- By Sumit S. Chaure  
(Batch – 10)



- Data Exploration (Using Excel & Python)
- Data insights (Excel & Python)
- EDA Graphs
- Graphical Analysis and conclusion on Data
- Data Cleaning & Pre-Processing Steps
- ML Modeling
- Model Test Evaluation & Prediction Analysis
- Deployment of ML Models using Streamlit

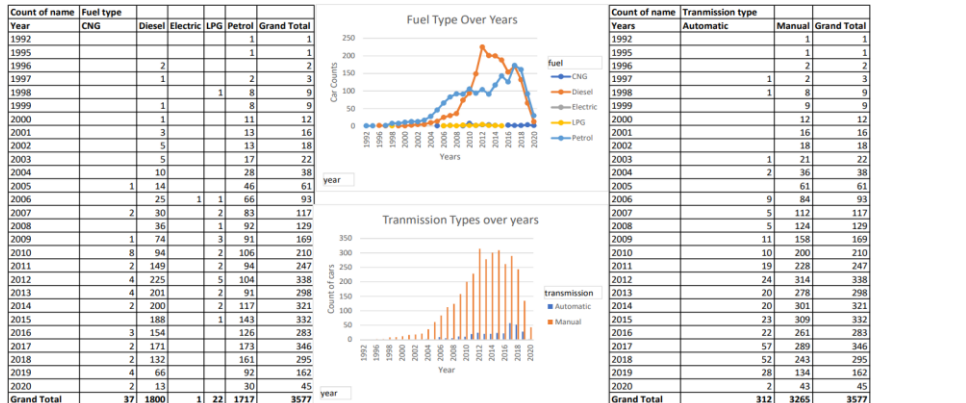
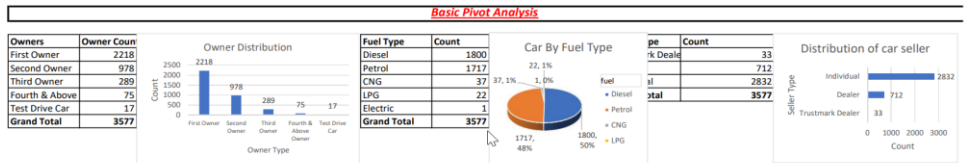


## Basic Analysis Using Excel Sheet

*Basic Data Analysis Using Excel*

	A	B	C	D	E	F	G	H	I	J	K	L
1												
2												
3	Total Count			4340								
4	Blanks		FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	
5	Blanks Count		0	0	0	0	0	0	0	0	0	
6	% Blanks		0%	0%	0%	0%	0%	0%	0%	0%	0%	
7	Mean	#DIV/0!	2013.1	504127.3118	66215.7	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	
8	Mode	#N/A	2017	300000	70000	#N/A	#N/A	#N/A	#N/A	#N/A	#N/A	
9	Median	#NUM!	2014	350000	60000	#NUM!	#NUM!	#NUM!	#NUM!	#NUM!	#NUM!	
10												
11		name	year	selling_price	km_driven	fuel	seller_type	transmission	owner			
12		Maruti 800 AC	2007	60000	70000	Petrol	Individual	Manual	First Owner			
13		Maruti Wagon R LXI Minor	2007	135000	50000	Petrol	Individual	Manual	First Owner			
14		Hyundai Verano 1.6 SX	2012	600000	100000	Diesel	Individual	Manual	First Owner			
15		Datsun RediGO T Option	2017	250000	46000	Petrol	Individual	Manual	First Owner			
16		Honda Amaze VX i-DTEC	2014	450000	141000	Diesel	Individual	Manual	Second Owner			
17		Maruti Alto LX BSIII	2007	140000	125000	Petrol	Individual	Manual	First Owner			
18		Hyundai Xcent 1.2 Kappa S	2016	550000	25000	Petrol	Individual	Manual	First Owner			
19		Tata Indigo Grand Petrol	2014	240000	60000	Petrol	Individual	Manual	Second Owner			
20		Hyundai Creta 1.6 VIVT S	2015	850000	25000	Petrol	Individual	Manual	First Owner			
21		Maruti Celerio Green VXI	2017	365000	78000	CNG	Individual	Manual	First Owner			

- The Dataset contains **8 rows & 4240** columns showing the information of used cars dataset.
- The Basic analysis shows that there is *no null value* present.
- The Analysis further shows that there are **763 duplicate** values.



Microsoft Excel

763 duplicate values found and removed; 3577 unique values remain.

OK

Sumit S. C

## Step 1 – Import Library & Read File

```
2. Module Imports

Click here to ask Blackbox to help you code faster
1 import pandas as pd # for data cleaning and data pre-processing, CSV file I/O,etc
2 import numpy as np # linear algebra & for mathematical computation
3 import matplotlib.pyplot as plt # for visualization
4 %matplotlib inline
5 import seaborn as sns # for visualization
6 from collections import Counter # to count occurrences
7 from tabulate import tabulate # to make tables for results
8
9 import warnings # for warning removals in code output
10 warnings.filterwarnings('ignore')
```

```
2.1) Importing the dataset (With error handling)

Click here to ask Blackbox to help you code faster
1 # 2.1) Importing the dataset (With error handling)
2 # If you want to upload the dataset directly (Since on Google Colab it will be lost
  on re-run) - uncomment the below 2 line code and run
3 # from google.colab import files
4 # uploaded = files.upload()
5
6 file_path = "../data/raw/CAR DETAILS.csv"
7 file_name = file_path.split("/")[-1]
8
9 try:
10 # Reading the CSV file into a Pandas DataFrame
11 df = pd.read_csv(file_path)
12 # Store the filename as an attribute in the DataFrame
13 df.file_name = file_name
14 print(f"\n '{df.file_name}' loaded successfully.")
15
16 # Exception to check if the file has some error like no file at the path, etc.
17 except FileNotFoundError:
18 print(f"Error: '{file_name}' not found at the specified location (
  {file_path}.")
19 except Exception as e:
20 print(f"An unexpected error occurred: {e}")
21
Python
'CAR DETAILS.csv' loaded successfully.
```

## Step 2 – Basic Data Lookup

```
Click here to ask Blackbox to help you code faster
1 df

name year selling_price km_driven fuel seller_type transmission owner
0 Maruti 800 AC 2007 60000 70000 Petrol Individual Manual First Owner
1 Maruti Wagon R LXI Minor 2007 135000 50000 Petrol Individual Manual First Owner
2 Hyundai Verma 1.6 SX 2012 600000 100000 Diesel Individual Manual First Owner
```

*df* – loads the data frame we gave while loading the csv & displays data.

```
2.2) Showing Basic Dataset Information

Click here to ask Blackbox to help you code faster
1 # Check the shape of the DataFrame
2 print("\nShape of the DataFrame:")
3 print(df.shape)
4 print(df.size)
5 num_rows, num_columns = df.shape
6 print(f"Rows: {num_rows}, Columns: {num_columns}")
7
8 # Display information about the dataset
9 print("\nDataset information for (df.file_name):")
10 df.head(3)
11 df.tail(3)
12 print("\nDataset information:")
13 print(df.info())

Shape of the DataFrame:
(4340, 8)
34720
Rows: 4340, Columns: 8

Dataset information for CAR DETAILS.csv:

Dataset information:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4340 entries, 0 to 4339
Data columns (total 8 columns):
# Column Non-Null Count Dtype
---
0 name 4340 non-null object
1 year 4340 non-null int64
2 selling_price 4340 non-null int64
3 km_driven 4340 non-null int64
4 fuel 4340 non-null object
5 seller_type 4340 non-null object
6 transmission 4340 non-null object
7 owner 4340 non-null object
dtypes: int64(3), object(5)
memory usage: 271.4+ KB
```

*df.shape* & *df.size* – gives the rows & column info  
*df.head* & *df.tail* – displays 5 values from top & bottom of dataset  
*df.info* – gives the information about the column its datatype and null values.

## Step 3 – Data Checks

```
1 # Null values
2 print("Missing values in the dataset:\n")
3 print(df.isnull().sum())
4
5 # NA value calculation
6 nullval = df.isna().sum()
7 nullval = nullval[nullval > 0]
8 print("\nSum of Missing values:\n", nullval)
9
10 # Calculating the number of missing values or null values in df
11 total_missing_values = df_dup_dropped.isnull().sum().sum()
12 print("The number of missing values/NA in dataframe :", total_missing_values)
13
14 # Calculate the total number of values in df (excluding the missing values)
15 total_values = df_dup_dropped.size
16 print("Total number of values in dataframe :", total_values)
17 # percentage of missing values or null values in df
18 percentage_missing_values = (total_missing_values / total_values) * 100
19 print("Percentage of missing values or null values in df :",
20       percentage_missing_values)
21
22 # Duplicate values
23 print("\nChecking for duplicated values:\n")
24 print(df.duplicated())
25 print("\nSum of Duplicated Values in Dataframe :", df.duplicated().sum())
26 # Dropping Duplicate values from dataframe
27 df_dup_dropped = df.drop_duplicates().reset_index(drop=True)
28 print("\nSum of Duplicated Values in Dataframe :",
29       df_dup_dropped.duplicated().sum())
30 print(f"Original Dataframe Shape : (df.shape)\nNew Dataframe Shape after Dropping Duplicates : (
31       df_dup_dropped.shape)
32       \nNo of duplicate rows removed : (df.duplicated().sum())")
33
34 # Value check
35 print("Value counts of the dataset by datatypes")
36 df_dup_dropped.dtypes.value_counts()
37 # Unique values
38 print("Unique Value counts inside each columns")
39 df_dup_dropped.nunique()
40
41 # Datatype check
42 print("Data types of each column:")
43 df_dup_dropped.dtypes
44
45 # Category Column Check
46 df + df_dup_dropped.select_dtypes(include="category")
47 categorical_columns = df_dup_dropped.select_dtypes(include="object").columns
48 print("\nCategorical columns:")
49 print(categorical_columns)
50
51 # Summary Stats
52 print("\nSummary statistics:\n")
53 print(df_dup_dropped.describe())
54
55 # Descriptive stats
56 print(f"The descriptive Stats for the {file_name} dataset:")
57 df_dup_dropped.describe()
```

## Step 4 – Data Pre-Processing

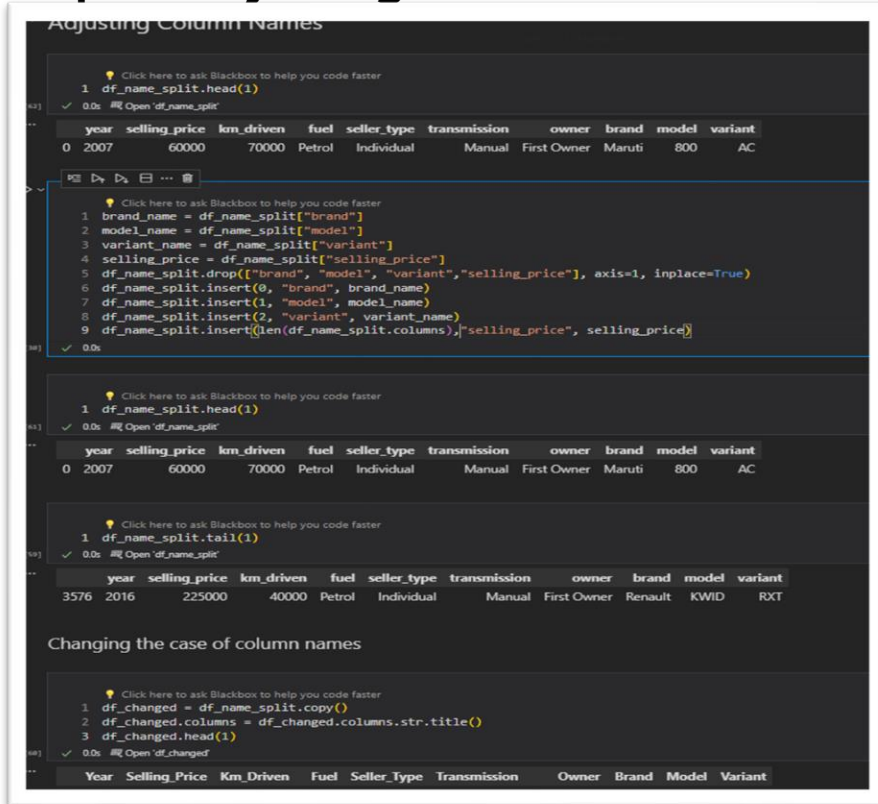
Extracting Car Manufacturer, Model & Variant Name from Car Name (To do more In-depth analysis)

```
1 # Unique values in name column
2 print(f"Unique Values in 'name' column : {df_dup_dropped['name'].nunique()}") #unique: Unknown word.
3
4 Unique Values in 'name' column : 1491
5
6 # Creating a copy to perform the splitting
7 df_name_split = df_dup_dropped.copy()
8 # name of the column to be split
9 car_names = df_dup_dropped["name"]
10 # Functions to split car names into brand,model & variant
11 def split_car_name(name):
12     words = name.split()
13     if len(words) == 3:
14         brand, model, variant = words
15     else:
16         brand = words[0]
17         model = ""
18         variant = ""
19         for i, word in enumerate(words[1:]):
20             if word[0].isdigit():
21                 if i == 0: # Number at the second position
22                     model = " ".join(words[1:2])
23                     variant = " ".join(words[2:])
24                 else: # Number at the third position or later
25                     model = " ".join(words[1 : i + 1])
26                     variant = " ".join(words[i + 1 :])
27                     break
28             else:
29                 model = " ".join(words[1:-1])
30                 variant = words[-1]
31     return brand, model, variant
32 # Apply the split_car_name function to create new columns
33 df_name_split[["brand", "model", "variant"]] = pd.DataFrame(car_names.apply(split_car_name).tolist())
34 # Drop the name column
35 df_name_split.drop(columns=["name"], inplace=True)
36
```

year	selling_price	km_driven	fuel	seller_type	transmission	owner	brand	model	variant
0	2007	60000	70000	Petrol	Individual	Manual	First Owner	Maruti	800 AC
1	2007	126000	60000	Petrol	Individual	Manual	First Owner	Maruti	800 AC

- **Step 3** – We do basic checks to find *null*, *duplicates*, *NA*, *unique values* in dataset.
- We also find the categorical columns & statistics to get info about numerical data.
- **Step 4** – We *Pre-process the data* and handle the features like splitting name into model, brand, variant to do analysis.
- Dropped the name column after split.

## Step 5 – Adjusting Data



```
Adjusting Column Names

Click here to ask Blackbox to help you code faster
1 df_name_split.head(1)
0.0s Open 'df_name_split'

year selling_price km_driven fuel seller_type transmission owner brand model variant
0 2007 60000 70000 Petrol Individual Manual First Owner Maruti 800 AC

Click here to ask Blackbox to help you code faster
1 brand_name = df_name_split["brand"]
2 model_name = df_name_split["model"]
3 variant_name = df_name_split["variant"]
4 selling_price = df_name_split["selling_price"]
5 df_name_split.drop(["brand", "model", "variant"], axis=1, inplace=True)
6 df_name_split.insert(0, "brand", brand_name)
7 df_name_split.insert(1, "model", model_name)
8 df_name_split.insert(2, "variant", variant_name)
9 df_name_split.insert((len(df_name_split.columns)), "selling_price", selling_price)
0.0s

Click here to ask Blackbox to help you code faster
1 df_name_split.head(1)
0.0s Open 'df_name_split'

year selling_price km_driven fuel seller_type transmission owner brand model variant
0 2007 60000 70000 Petrol Individual Manual First Owner Maruti 800 AC

Click here to ask Blackbox to help you code faster
1 df_name_split.tail(1)
0.0s Open 'df_name_split'

year selling_price km_driven fuel seller_type transmission owner brand model variant
3576 2016 225000 40000 Petrol Individual Manual First Owner Renault KWID RXT

Changing the case of column names

Click here to ask Blackbox to help you code faster
1 df_changed = df_name_split.copy()
2 df_changed.columns = df_changed.columns.str.title()
3 df_changed.head(1)
0.0s Open 'df_changed'

Year Selling_Price Km_Driven Fuel Seller_Type Transmission Owner Brand Model Variant
```

## Insights :

- Data Exploration is very useful step in the process of analytics and model building as it helps us to remove any redundant data & also construct new features to enhance the analysis and model building.
- In the data exploration part we eye-ball the data to look at features that help us in analysis and model formation and help us to remove any unnecessary junk from our data.
- Most times the data from sources is not cleaned so data cleaning is an important part of the analysis work.
- Gathering and exploring data and finding meaningful insights helps us to analyse and train the model more efficiently.
- In the step given on the side we adjusted the column names such that the newly added feature at the end are given index value & adjusted likewise selling price being target we move it to end.

## Basic Info On Dataset

```
Shape of the DataFrame:
(4340, 8)
34720
Rows: 4340, Columns: 8

Dataset information for CAR DETAILS.csv:

Dataset information:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4340 entries, 0 to 4339
Data columns (total 8 columns):
#   Column          Non-Null Count  Dtype
---  ---            -
0   name            4340 non-null   object
1   year            4340 non-null   int64
2   selling_price   4340 non-null   int64
3   km_driven       4340 non-null   int64
4   fuel            4340 non-null   object
5   seller_type     4340 non-null   object
6   transmission    4340 non-null   object
7   owner           4340 non-null   object
dtypes: int64(3), object(5)
memory usage: 271.4+ KB
None
```

The dataset contains **4340 rows & 8 columns**, all are non-null values **3 are numerical** while **5 are categorical** type.

## Dataset after cleaning & adding features

```
Shape of the DataFrame:
(3577, 10)
35770
Rows: 3577, Columns: 10

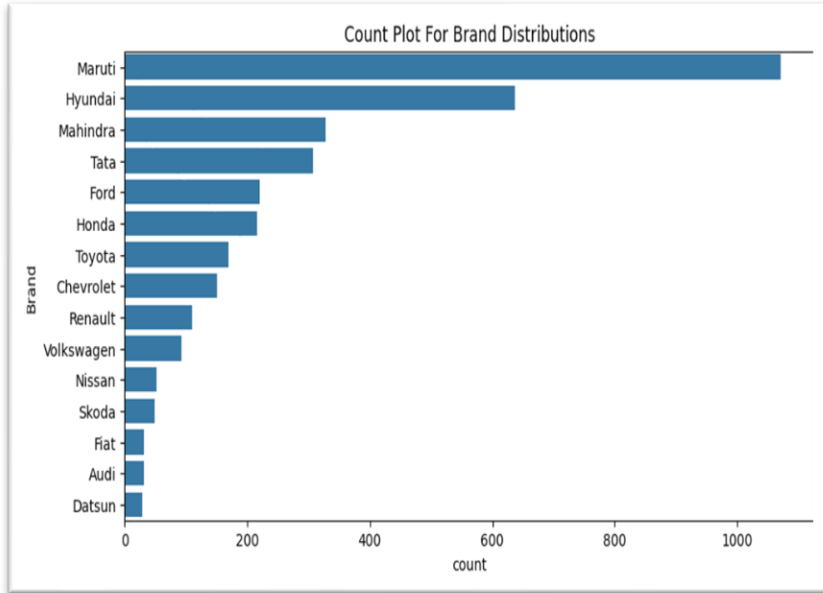
Dataset information for CAR DETAILS.csv:

Dataset information:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3577 entries, 0 to 3576
Data columns (total 10 columns):
#   Column          Non-Null Count  Dtype
---  ---            -
0   Year            3577 non-null   int64
1   Selling_Price   3577 non-null   int64
2   Km_Driven       3577 non-null   int64
3   Fuel            3577 non-null   object
4   Seller_Type     3577 non-null   object
5   Transmission    3577 non-null   object
6   Owner           3577 non-null   object
7   Brand           3577 non-null   object
8   Model           3577 non-null   object
9   Variant         3577 non-null   object
dtypes: int64(3), object(7)
memory usage: 279.6+ KB
None
```

After cleaning the data, removing duplicates and adding removing features from dataset the new processed dataset contains **3577 rows & 10 columns**. **3 are numerical** while **7 are categorical** columns.

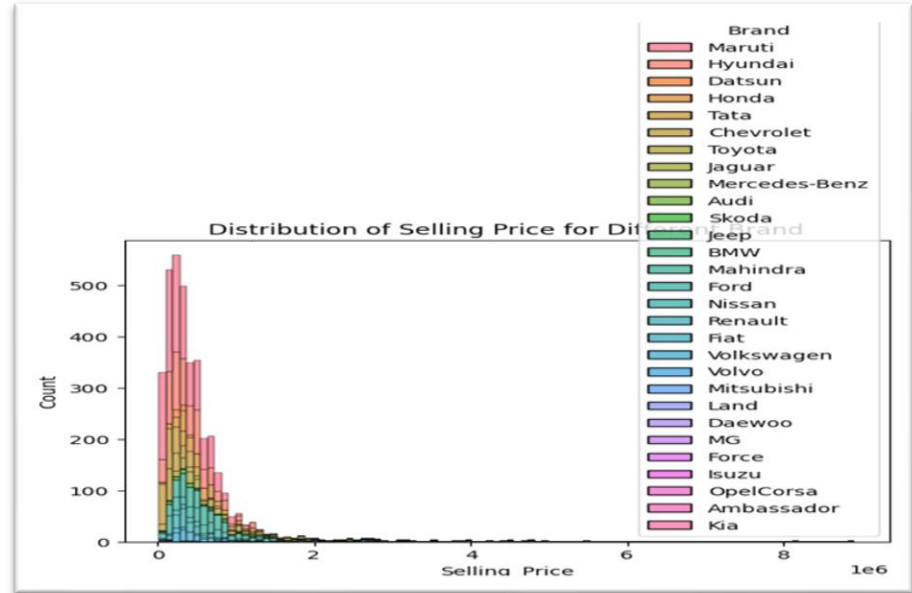
The graphs below are made using matplotlib & Seaborn library (code in notebook file)

## Count Plot



Count-plot showing the distribution of cars according to brand

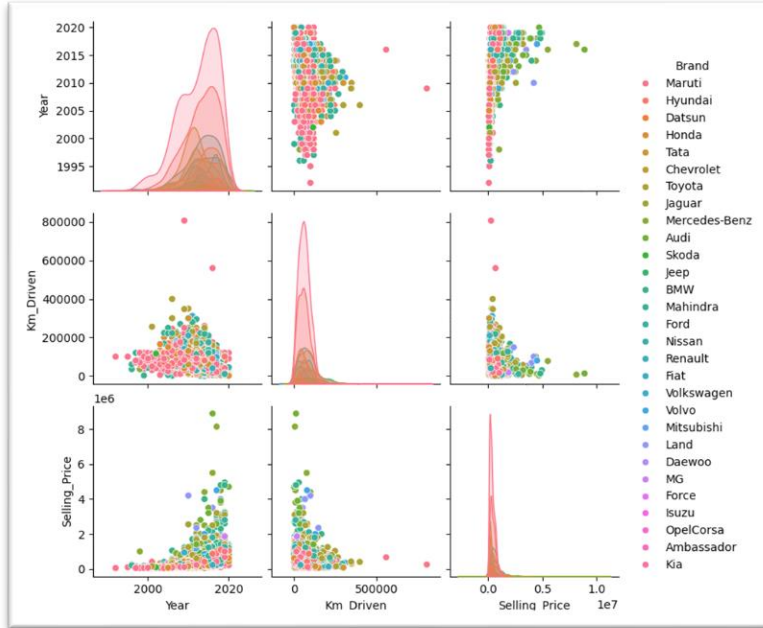
## Hist Plot



Hist-plot to showcase the selling price of various brands

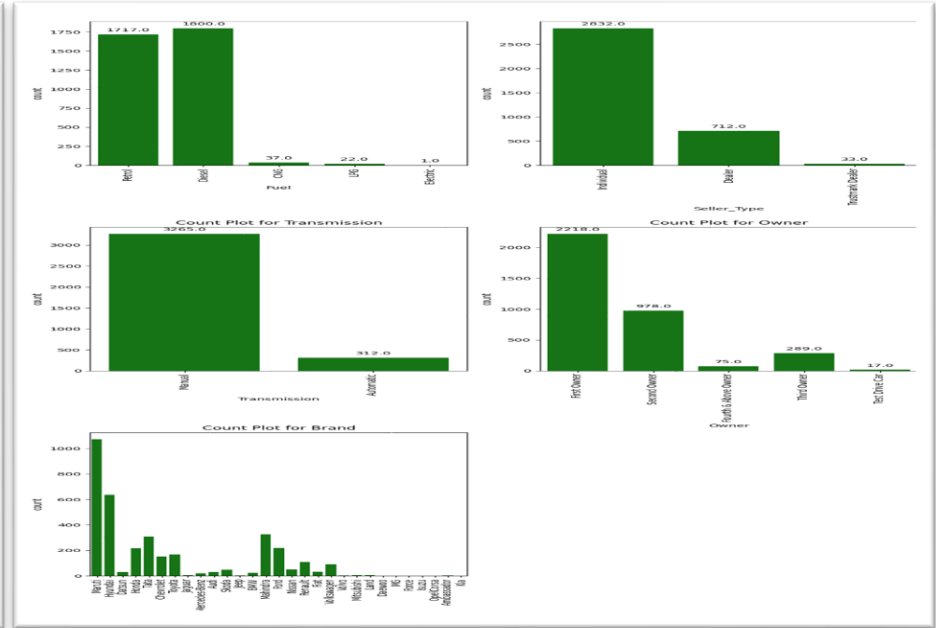


## Pair Plot (Numerical Data)



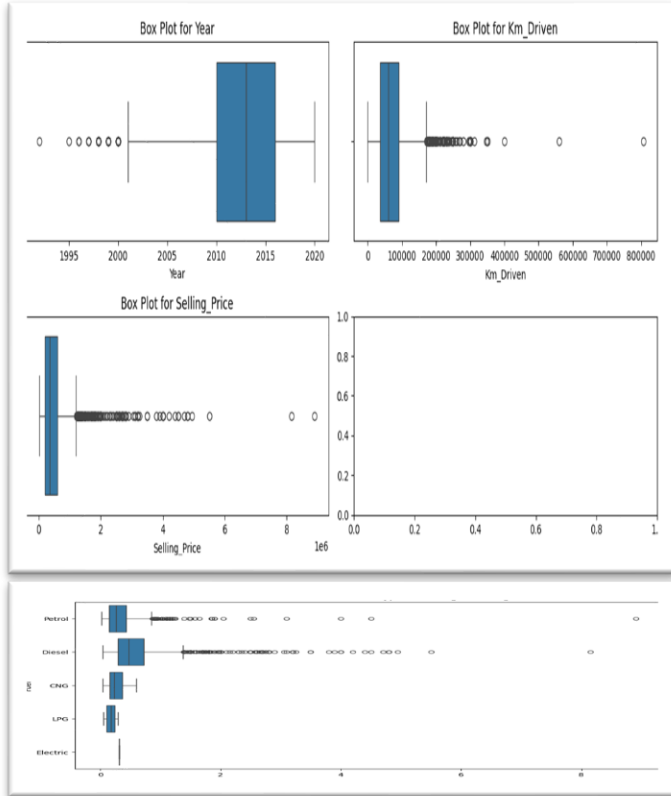
Pair-plot showing the distribution of numerical data (selling price, km\_driven, years) according to brand

## Count Plot (Categorical Data)

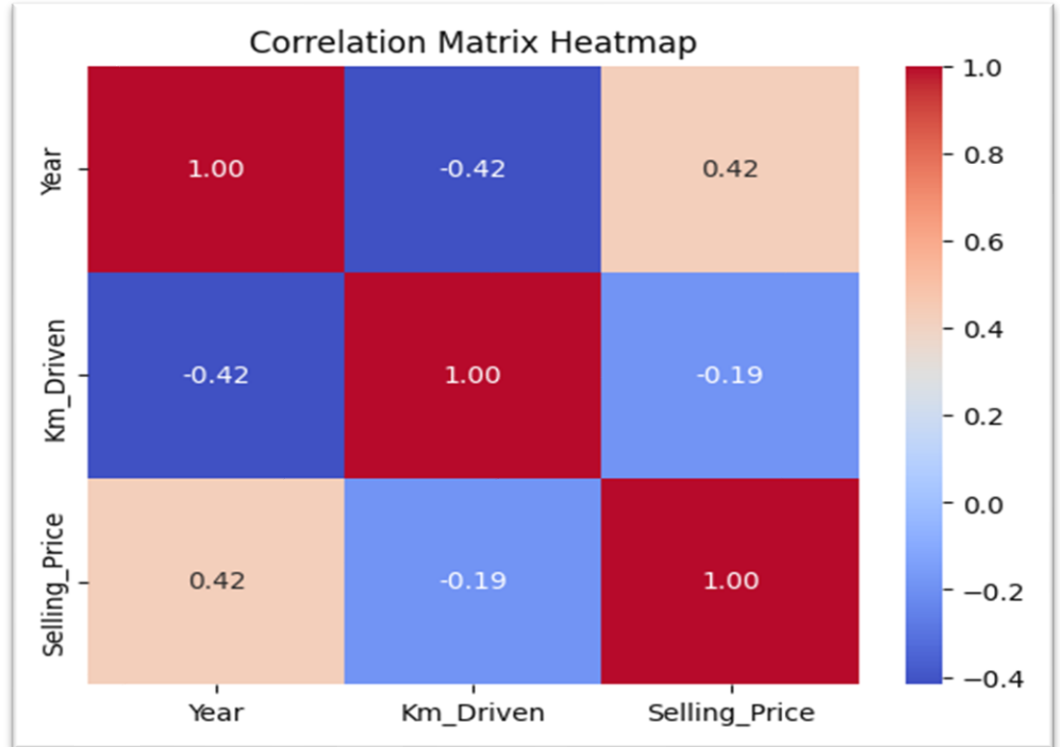


Count Plot of various categorical data representation.

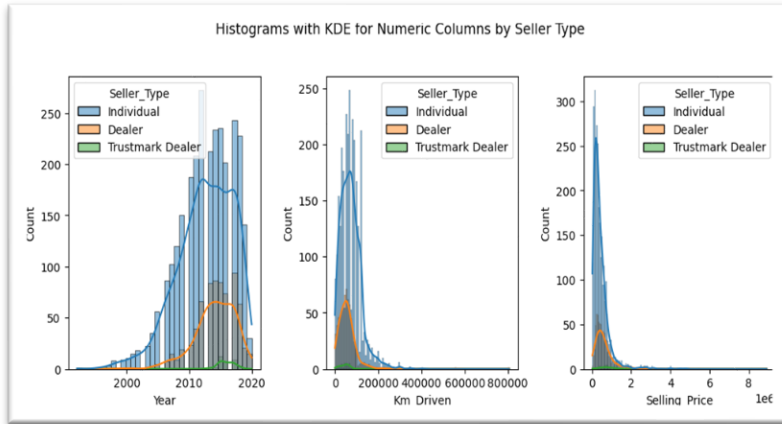
**Box Plot** to find the Outlier data



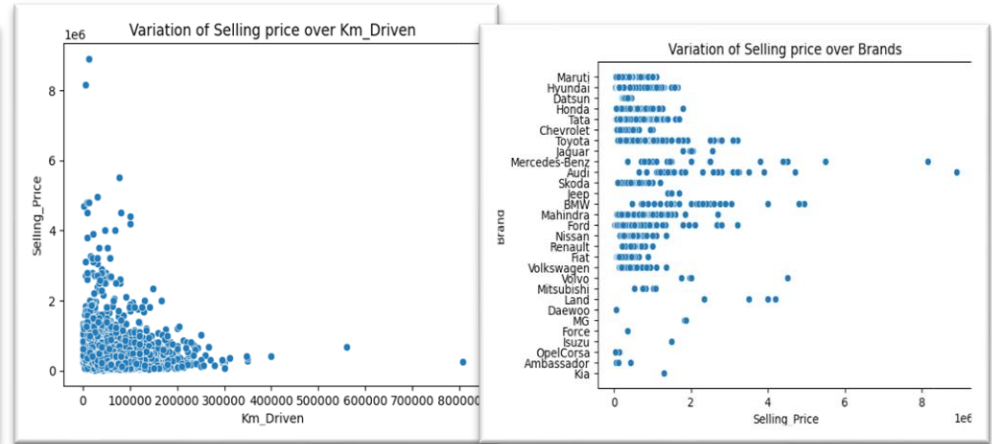
**HeatMap** to find the Correlation in Numerical



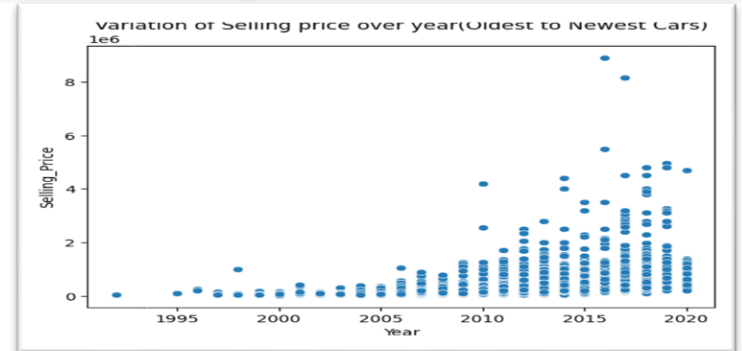
## Hist Plot on continuous data



## Scatter Plot of Selling price over km, brand & Year



- The **Histplot** shows the distribution of continuous numerical.
- The various **Scatter plot** shows the variation of selling price indicating the impact of year(age), kms (use), brand (popularity) on the price of resale value of a car.

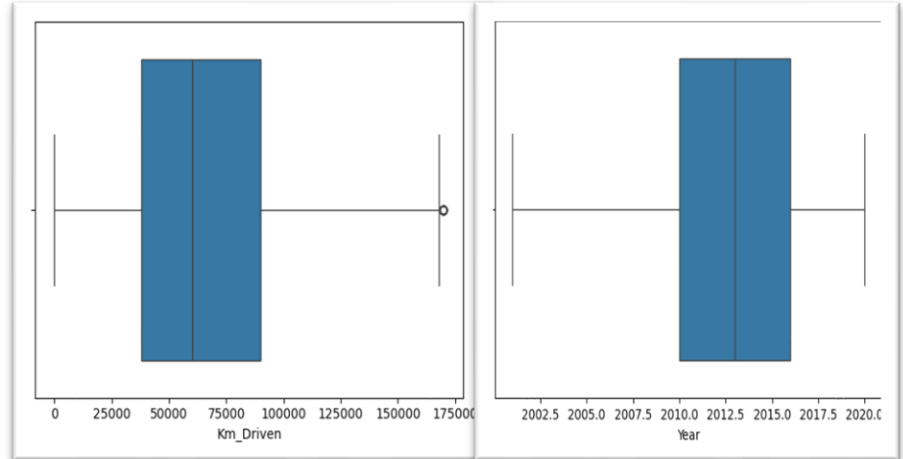


Few data cleaning steps are discussed in previous slides.

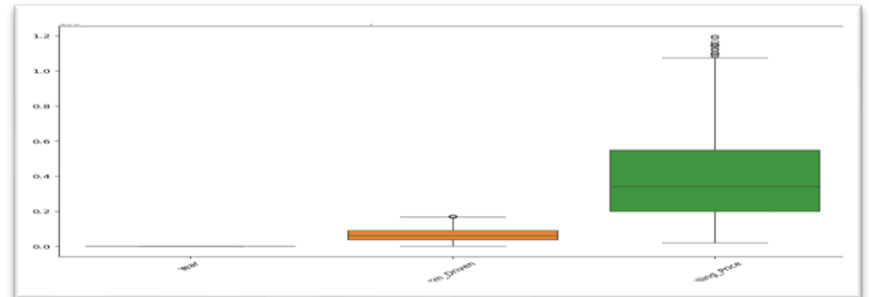
## Treating Outliers (Data Capping)

## ScatterPlot (After Capping Data)

```
1 df_cap = df_changed
2 def capping_data(df_cap):
3     for i in df_cap.columns:
4         print("Check : Capped the", i, "Column")
5         if df_cap[i].dtype in ["float64", "int64"]:
6             Q1 = df_cap[i].quantile(0.25)
7             Q3 = df_cap[i].quantile(0.75)
8             IQR = Q3 - Q1
9             df_cap = df_cap[
10                 ~((df_cap[i] < (Q1 - 1.5 * IQR)) |
11                  (df_cap[i] > (Q3 + 1.5 * IQR)))
12             ]
13         else:
14             df_cap[i] = df_cap[i]
15     return df_cap
16 data = capping_data(df_cap)
```



- The above code caps the data value of numerical data columns using the IQR formulas.
- The scatter plot on the right shows the capped effect on outlier.
- Capping helps us to remove the data which may change the model prediction due to outlier data.



## Step 1 – Import Necessary Library and dataset

```
Module Imports

Click here to ask Blackbox to help you code faster
1 import pandas as pd # for data cleaning and data pre-processing, CSV file I/O, etc
2 import numpy as np # linear algebra & for mathematical computation
3 import matplotlib.pyplot as plt # for visualization
4 %matplotlib inline
5 import seaborn as sns # for visualization
6 from collections import Counter # to count occurrences
7 from tabulate import tabulate # to make tables for results
8
9 import warnings # for warning removals in code output
10 warnings.filterwarnings('ignore')
11
12 # Scalers & Encoders
13 from sklearn.preprocessing import StandardScaler, LabelEncoder
14 #train-test split
15 from sklearn.model_selection import train_test_split
16 # Metrics
17 from sklearn.metrics import (mean_squared_error, r2_score)
18 # Model Libraries
19 from sklearn.linear_model import LinearRegression, Ridge, Lasso
20 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor, AdaBoostRegressor
21 from sklearn.neighbors import KNeighborsRegressor
22 from sklearn.tree import DecisionTreeRegressor
23
24 import pickle #to save and load model files as pickle file
```

- Import the basic library for data handling & the model handling ones from sklearn library.
- The next step is similar to previous to import the dataset but we will import the processed dataset instead of raw which we handled in previous steps in graphical analysis
- To save the dataset in last file at end we did - **`data.to_csv("<file-path/file-name>.csv", index=False)`**

## Step 2 – Encoding dataset

```
Click here to ask Blackbox to help you code faster
1 # Using Label Encoder to encode categorical data
2 label_encoder = LabelEncoder() # instance of encoder
3 # Loop to encode data in df
4 for feature in category_col:
5     df_one[feature] = label_encoder.fit_transform(df[feature])
6 df_one.head()

Brand  Model  Variant  Year  Km_Driven  Fuel  Seller_Type  Transmission  Owner  Selling_Price
0  12  547  657  2007  50000  4  1  1  0  135000
1  10  517  311  2012  100000  1  1  1  0  600000
2  5  399  659  2017  40000  4  1  1  0  250000
3  9  35  729  2014  141000  1  1  1  2  450000

Click here to ask Blackbox to help you code faster
1 corr = df_one.corr()
2 sns.heatmap(corr, cmap="coolwarm", annot=True, fmt=".2f", linewidths=0.5) # c
3 plt.title("Correlation Matrix for Used Car Dataset (Encoded)", fontsize=16) # t
4 plt.savefig("../src/visualization/Correlation Graph - After Label Encoding") # s
5 plt.show()

Correlation Matrix for Used Car Dataset(Encoded)

Brand  1.00  0.09  0.02  0.01  0.10  0.13  0.06  0.03  0.03  0.00
Model  -0.09  1.00  0.17  -0.02  0.09  -0.11  0.04  0.00  0.02  0.06
Variant -0.02  0.17  1.00  -0.11  0.08  -0.04  0.03  0.05  0.06  -0.12
Year    -0.01  -0.03  -0.11  1.00  -0.48  -0.08  -0.11  -0.08  -0.42  0.64
Km_Driven -0.10  0.09  0.08  -0.48  1.00  0.34  0.14  0.09  0.33  0.27
Fuel     -0.11  -0.11  0.04  -0.08  0.34  1.00  0.04  0.07  -0.02  0.22
Seller_Type -0.06  0.04  0.03  -0.11  0.14  0.03  1.00  0.09  0.18  -0.15
Transmission -0.01  -0.00  0.05  -0.08  0.09  -0.07  0.09  1.00  0.03  0.00
Owner    0.01  0.02  0.06  -0.13  0.33  0.02  0.18  0.03  1.00  0.30
Selling_Price -0.00  0.06  0.12  0.64  -0.27  -0.12  -0.14  -0.00  -0.10  1.00

Saving the encoded dataset

Click here to ask Blackbox to help you code faster
1 cleaned_data = df_one.to_csv("../data/processed/cleaned_data.csv", index=False)
```

- The above code encodes the loaded dataset using label encoder so that we can use the categorical values in model training as it needs number to build model.
- The correlation matrix shows us that the new encoded values are in numeric.

## Step 3 – Splitting data into Dependant & Independent variables for train-test split

```
Dependant (y) & independent (x) Features
1. Dropping dependant feature from dataset
Click here to ask Blackbox to help you code faster
1 x = df_one.drop(["Selling_Price"], axis=1)
2 y = df_one["Selling_Price"]
3
4 print(type(x))
5 print(type(y))
6 print(x.shape)
7 print(y.shape)
...
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
(3250, 9)
(3250,)

2. Splitting The dataset
Click here to ask Blackbox to help you code faster
1 # splitting the dataset into 70% training data and 30% test data
2 X_train, X_test, y_train, y_test = train_test_split(
3     x, y, test_size=0.3, random_state=42
4 )
5 print(f"Split Check Test values : {3194 * 0.3} & Train values : {3194 * 0.7}")
6 # rows, columns
7 print(X_train.shape)
8 print(X_test.shape)
9 print(y_train.shape)
10 print(y_test.shape)
...
Split Check Test values : 958.1999999999999 & Train values : 2235.7999999999997
(2275, 9)
(975, 9)
(2275,)
(975,)

4. Standarizing the dataset
Click here to ask Blackbox to help you code faster
1 sc = StandardScaler()
2 sc.fit_transform(X_train)
3 sc.transform(X_test)
...
array([[ -0.00522286,  0.38325919, -0.213698, ...,  0.44793592,
         0.26187172, -0.71016833],
       [ 1.65825717, -0.23064065,  0.93736498, ...,  0.44793592,
         0.26187172,  2.45929306],
       [-0.48050286,  1.12202081, -0.8306776, ...,  0.44793592,
         0.26187172, -0.71016833]])
```

- The code shows us the basic step of model building where in we split the dataset into training data & testing data.
- This step is very important as without doing the split if we train the model it will know the results and if unknown values are introduced it might not function as intended.
- The practice of splitting data helps us to find various parameters such as accuracy, precision, recall values, r2 score, mse values etc which helps us to determine if our model is fit or not or should we tune it more or use some other modelling technique.
- After these step we move on to actual model building & selection step.

*I Will add the code in loop as writing the same piece of code repeatedly changing variable names and remembering them will be not good idea so instead I will do looping.*

```
1 # Create a list of models
2 models = [
3     LinearRegression(),
4     Ridge(),
5     Lasso(),
6     RandomForestRegressor(),
7     KNeighborsRegressor(),
8     DecisionTreeRegressor(),
9     GradientBoostingRegressor(),
10    AdaBoostRegressor(),
11 ]
12
13 # Define a list of model names
14 model_names = [
15     "Linear Regression",
16     "Ridge Regression",
17     "Lasso Regression",
18     "Random Forest",
19     "k-Nearest Neighbors",
20     "Decision Tree",
21     "Gradient Boosting",
22     "Ada Boost",
23 ]
24
25 # Initialize variables to keep track of the best model
26 best_model_name = None
27 best_r2_score = -float("inf")
28 best_model = None
29
30 # Initialize an empty list to store results
31 all_results = []
32
33 # Function to evaluate regression models
34
35
36 def evaluate_regression_model(model, model_name, X_test, y_test):
37     y_pred = model.predict(X_test)
38
39     # Mean Squared Error and R-squared Score
40     mse = mean_squared_error(y_test, y_pred)
41     r2 = r2_score(y_test, y_pred)
42
43     # Display results in tabular format
44     results_table = [
45         ["Model", model_name],
46         ["Mean Squared Error", mse],
47         ["R-squared Score", r2],
48     ]
49
50     print(tabulate(results_table, headers=[
51         "Metric", "Value"], tablefmt="heavy_grid"))
52
53     # Store results in the dictionary
54     return {
55         "Model": model_name,
56         "Mean Squared Error": mse,
57         "R-squared Score": r2,
58     }
59
```

```
1 # Iterate over the models
2 for model, model_name in zip(models, model_names):
3     print(f"\n({'-' * 40})\n{model_name}\n({'-' * 40}")
4
5     # Train the model
6     model.fit(X_train, y_train)
7
8     # Evaluate and store results for regression model
9     results = evaluate_regression_model(model, model_name, X_test, y_test)
10    all_results.append(results)
11
12    # Update the best model if needed
13    if results["R-squared Score"] > best_r2_score:
14        best_r2_score = results["R-squared Score"]
15        best_model_name = model_name
16        best_model = model
17
18 # Create a DataFrame from the results
19 results_df = pd.DataFrame(all_results)
20
21 # Sorting the dataframe by 'R-squared Score' in descending order
22 sorted_results_df = results_df.sort_values(
23     by="R-squared Score", ascending=False)
24
25 # Displaying the sorted dataframe
26 print("\nSorted Models by R-squared Score:")
27 print(sorted_results_df)
28
29 # Print the best model based on R-squared score
30 print(f"\nBest Model based on R-squared Score: {best_model_name}")
31 best_model_name
```

In the previous slide we saw the code for model building – scoring ,parameters , function to check them on various passed models , save the results and iterate to next model. Here we see few glimpse of the code o/p.

**Linear Regression**

Metric	Value
Model	Linear Regression
Mean Squared Error	28427461410.760574
R-squared Score	0.5370094957944969

**Ridge Regression**

Metric	Value
Model	Ridge Regression
Mean Squared Error	28429012702.72386
R-squared Score	0.53698423031483

**Lasso Regression**

Metric	Value
Model	Lasso Regression
Mean Squared Error	28427489885.780876
R-squared Score	0.5370090320294136

**Random Forest**

Metric	Value
Model	Random Forest
Mean Squared Error	13871188171.954166
R-squared Score	0.774083647045899

**k-Nearest Neighbors**

Metric	Value
Model	k-Nearest Neighbors
Mean Squared Error	54114480862.42737
R-squared Score	0.11865184100362591

**Decision Tree**

Metric	Value
Model	Decision Tree
Mean Squared Error	26104259077.00128
R-squared Score	0.5748468743889662

**Gradient Boosting**

Metric	Value
Model	Gradient Boosting
Mean Squared Error	14369733662.247465
R-squared Score	0.7659639692250402

**Ada Boost**

Metric	Value
Model	Ada Boost
Mean Squared Error	29315494521.837402
R-squared Score	0.522546336671431

Sorted Models by R-squared Score:

	Model	Mean Squared Error	R-squared Score
3	Random Forest	1.387119e+10	0.774084
6	Gradient Boosting	1.436973e+10	0.765964
5	Decision Tree	2.610426e+10	0.574847
0	Linear Regression	2.842746e+10	0.537009
2	Lasso Regression	2.842749e+10	0.537009
1	Ridge Regression	2.842901e+10	0.536984
7	Ada Boost	2.931549e+10	0.522546
4	k-Nearest Neighbors	5.411448e+10	0.118652

Best Model based on R-squared Score: Random Forest  
'Random Forest'

Model performance table (Basic)

```
[ ] model_performance = pd.DataFrame(results_df)
model_performance
```

	Model	Mean Squared Error	R-squared Score
0	Linear Regression	2.842746e+10	0.537009
1	Ridge Regression	2.842901e+10	0.536984
2	Lasso Regression	2.842749e+10	0.537009
3	Random Forest	1.387119e+10	0.774084
4	k-Nearest Neighbors	5.411448e+10	0.118652
5	Decision Tree	2.610426e+10	0.574847
6	Gradient Boosting	1.436973e+10	0.765964
7	Ada Boost	2.931549e+10	0.522546



In the previous slide we saw the model scores o/p and best model the next step is to save the model(will skip ss can look into notebook) now we see the testing part of model against test data that we split.

## Generating sample data from cleaned df to test on the trained model.

```
[ ] random_datasample = df_one.sample(20)
random_datasample_df = random_datasample.drop("Selling_Price", axis=1)
print(random_datasample_df.shape)
random_datasample_df.head()
```

	Brand	Model	Variant	Year	Km_Driven	Fuel	Seller_Type	Transmission	Owner
(20, 9)	2308	12	588	37	2008	80000	4	1	1
...	...	...	...	...	...	...	...	...	...

## Making Predictions on sample dataset against the trained model

```
# making prediction on random data
predicted_data = load_model.predict(testsample_df)
print(f"The predicted data from {f_modelname} model:\n", predicted_data)
```

The predicted data from Random Forest model:  
[126650. 132452.66666667 572970. 523030.]

- The code above generates 20 random sample data from original file which we used for train test split and test them to predict value against model.
- The code on side checks the actual value and predicted value ( I tweaked it a bit so that if difference of selling price is say under certain limit its safe prediction.

## Comparison of Actual and Predicted values by the model

```
# Compare the actual data and predicted data
prediction_data = random_datasample.copy()
prediction_data["predicted_target"] = predicted_data
# Calculate the absolute percentage difference
prediction_data["percentage_difference"] = abs(
    (random_datasample["Selling_Price"] - predicted_data) /
    random_datasample["Selling_Price"]) * 100
# Print the actual and predicted data
print(f"Actual Data and Predicted Data Comparison based on {
    f_modelname} model:\n")
# Display the results where the absolute percentage difference is less than or equal to 20%
safe_predictions = prediction_data[prediction_data["percentage_difference"] <= 20]
# Print the safe predictions
print("Safe Predictions:")
print(safe_predictions[["Selling_Price",
    "predicted_target", "percentage_difference"]])
# Print the count and percentage of safe predictions
safe_percentage = (len(safe_predictions) / len(prediction_data)) * 100
print(f"\nPercentage of Safe Predictions: {safe_percentage:.2f}%")
if safe_percentage >= 90:
    print(
        f"\nOur model based on '{f_modelname}' is well trained, with {
            safe_percentage:.2f}% safe predictions."
    )
else:
    print(f"\nOur model based on '{f_modelname}' needs more training to improve safety, currently at {
        safe_percentage:.2f}% safe predictions.")
# Save the results as a DataFrame
final_results_df = prediction_data[
    ["Selling_Price", "predicted_target", "percentage_difference"]]
final_results_df.to_csv("../reports/model_predicted_results.csv", index=False)
```

### Actual Data and Predicted Data Comparison based on Random Forest model:

Safe Predictions:	Selling_Price	predicted_target	percentage_difference
2308	126650	126650.000000	2.576923
379	580000	572970.000000	14.594600
3116	535000	523030.000000	2.227383
2143	490000	509022.940000	3.883661
1971	836000	768859.580000	8.031183
1509	530000	504139.990000	4.879247
2672	280000	266950.000000	4.982143
692	400000	369299.980000	7.675805
1268	70000	73140.000000	4.485714
1014	430000	410189.970000	4.606984
480	537000	637940.000000	18.456875
2279	280000	238370.000000	14.867857
389	320000	318026.666667	0.616667
2800	280000	295100.000000	5.392857
720	940000	951100.000000	1.189362
1613	190000	213750.000000	12.500000
1750	280000	267481.101111	4.471835

Percentage of Safe Predictions: 85.00%  
Our model based on "Random Forest" needs more training to improve safety, currently at 85.00% safe predictions

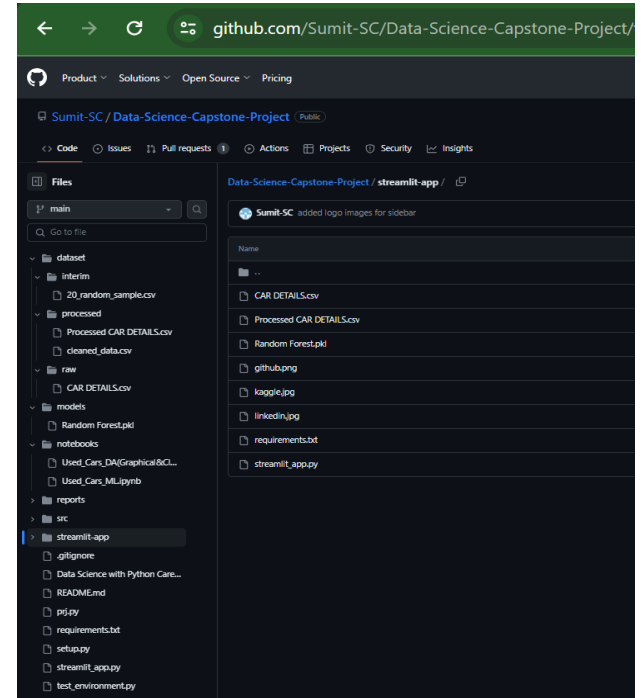
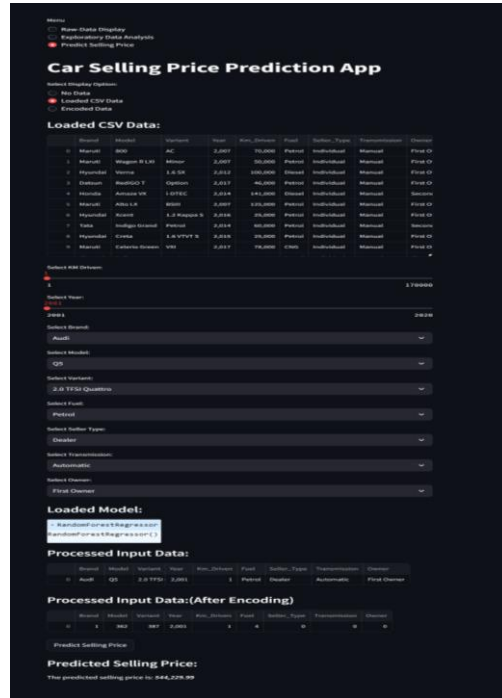
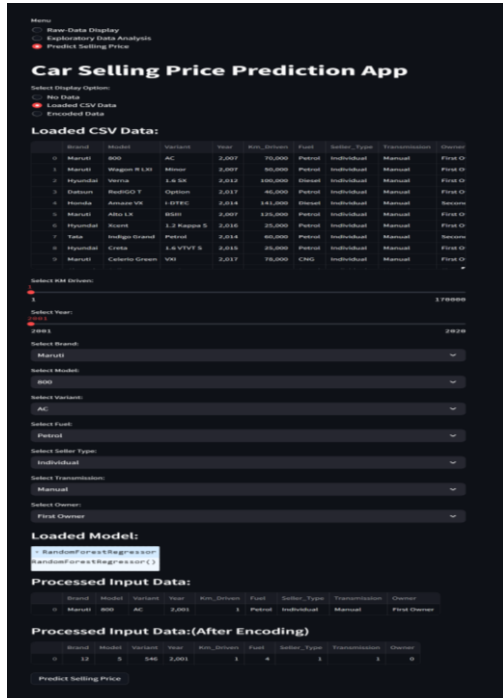
# Deployment of ML Models using Streamlit.

Please visit the [link](#) to view the actual working code of the project (won't paste long code)

**Without any selection**

**Prediction of Selling Price**

**Github Code Structure**



## Reference Links:-

- GitHub [Repo Link](#)
- Streamlit App [Weblink](#)
- Streamlit App [Alternate](#)
- Google Drive [Folder](#)
- Zip [file](#)
- Google Colab Links (Run & View code Directly)
  - EDA & Graphical Analysis [Colab](#)
  - Model Training & Evaluation [Colab](#)
  - Readme [File](#) (To get the project structure & all links)

## About Me :

- Name – Sumit S. Chaure
- Batch – 10
- Course – Data Science With Python
- GitHub - [Sumit-SC](#)
- Kaggle – [Sumit](#)
- Email - [@](#)

END

