

Import Python Libraries

Pandas and Numpy have been used for Data Manipulation and numerical Calculations

Matplotlib and Seaborn have been used for Data visualizations.

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
```

Reading Dataset

```
In [2]: data = pd.read_csv('C:\Users\vineeta Shrivastava\Downloads\cars4\used_cars_data.csv')
```

```
In [3]: #checking the data
data.head()
```

S.No.	Name	Location	Year	Kilometers_Driven	Fuel_Type	Transmission	Owner_Type	Mileage	Engine	Power	Seats	New_Price	Price
0	Mazd Wagon R LXI CNG	Mumbai	2010	72000	CNG	Manual	First	26.6kmpl	998 CC	56.16 bhp	5.0	NaN	1.75
1	Hyundai Creta 1.6 CRDi SX Option	Pune	2015	43000	Diesel	Manual	First	19.67 kmpl	1592 CC	126.2 bhp	5.0	NaN	1250
2	Hyundai Verna 1.6 CRDi	Chennai	2011	40000	Petrol	Manual	First	18.2 kmpl	1198 CC	88.7 bhp	5.0	8.61 Lakh	430
3	Mazd Ertiga VDi	Chennai	2012	87000	Diesel	Manual	First	20.77 kmpl	1248 CC	86.76 bhp	7.0	NaN	630
4	Audi A4 New 2.0 TDI Multitronic	Coimbatore	2013	40670	Diesel	Automatic	Second	15.2 kmpl	1968 CC	140.8 bhp	5.0	NaN	17.74

```
In [4]: data.tail()
```

S.No.	Name	Location	Year	Kilometers_Driven	Fuel_Type	Transmission	Owner_Type	Mileage	Engine	Power	Seats	New_Price	Price
7248	Volkswagen Vento Diesel Trendline	Hyderabad	2011	94111	Diesel	Manual	First	20.14 kmpl	1096 CC	103.6 bhp	5.0	NaN	NaN
7249	Volkswagen Polo GT TSI	Mumbai	2015	59000	Petrol	Automatic	First	17.23 kmpl	1197 CC	103.6 bhp	5.0	NaN	NaN
7250	Nissan Micra Diesel XV	Kolkata	2012	28000	Diesel	Manual	First	23.08 kmpl	1461 CC	63.1 bhp	5.0	NaN	NaN
7251	Volkswagen Polo GT TSI	Pune	2013	52262	Petrol	Automatic	Third	17.2 kmpl	1197 CC	103.6 bhp	5.0	NaN	NaN
7252	Mercedes-Benz E-Class 2009-2013 E 220 CDI Avani...	Kochi	2014	72443	Diesel	Automatic	First	10.0 kmpl	2148 CC	170 bhp	5.0	NaN	NaN

```
In [5]: <class 'pandas.core.frame.DataFrame'>
RangeIndex: 7252 entries, 0 to 7252
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
#  --  --                --
#   0   Name                  7252 non-null    object
#   1   Location              7252 non-null    object
#   2   Year                  7252 non-null    int64
#   3   Kilometers_Driven     7252 non-null    float64
#   4   Fuel_Type             7252 non-null    object
#   5   Transmission          7252 non-null    object
#   6   Owner_Type           7252 non-null    object
#   7   Mileage               7252 non-null    float64
#   8   Engine               7207 non-null    object
#   9   Power                7207 non-null    float64
#  10   Seats                6619 non-null    float64
#  11   New_Price            1886 non-null    object
#  12   Price                6619 non-null    float64
dtypes: float64(2), int64(2), object(9)
memory usage: 792.4+ KB
```

Check for Duplication

```
In [6]: data.nunique()

Out[6]: S.No.      7253
Name      2641
Location    11
Year        9
Kilometers_Driven  23
Fuel_Type   8
Transmission  2
Owner_Type   4
Mileage     459
Engine      129
Power       396
Seats       66
New_Price   625
Price      1379
dtype: object
```

Missing Values Calculation

```
In [7]: data.isnull().sum()

Out[7]: S.No.      0
Name      0
Location  0
Year      0
Kilometers_Driven  0
Fuel_Type  0
Transmission  0
Owner_Type  0
Mileage    0
Engine     0
Power     46
Seats     66
New_Price  93
Price     6247
dtype: object

In [8]: (data.isnull().sum()/len(data))*100

Out[8]: S.No.      0.00008
Name      0.00008
Location  0.00008
Year      0.00008
Kilometers_Driven  0.00008
Fuel_Type  0.00008
Transmission  0.00008
Owner_Type  0.00008
Mileage    0.02775
Engine     0.02775
Power     0.03228
Seats     0.0091
New_Price  0.08619
Price     0.85369
dtype: float64
```

Data Reduction

```
In [9]: # Remove S.No. column from data
data = data.drop(['S.No.'], axis = 1)
data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7252 entries, 0 to 7252
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
#  --  --                --
#   0   Name                  7252 non-null    object
#   1   Location              7252 non-null    object
#   2   Year                  7252 non-null    int64
#   3   Kilometers_Driven     7252 non-null    float64
#   4   Fuel_Type             7252 non-null    object
#   5   Transmission          7252 non-null    object
#   6   Owner_Type           7252 non-null    object
#   7   Mileage               7252 non-null    float64
#   8   Engine               7207 non-null    object
#   9   Power                7207 non-null    float64
#  10   Seats                6619 non-null    float64
#  11   New_Price            1886 non-null    object
#  12   Price                6619 non-null    float64
dtypes: float64(2), int64(2), object(9)
memory usage: 736.8+ KB
```

Feature Engineering

Feature engineering refers to the process of using domain knowledge to select and transform the most relevant variables from raw data when creating a predictive model using machine learning or statistical modeling. The main goal of Feature engineering is to create meaningful data from raw data.

Creating Features

variables Year and Name in our dataset. If we see the sample data, the column "Year" shows the manufacturing year of the car.It would be difficult to find the car's age if it is in year format as the Age of the car is a contributing factor to Car Price. Introducing a new column, "Car_Age" to know the age of the car

```
In [10]: from datetime import date
date.today()
year = date.today().year
data['Car_Age'] = date.today() - year
data.head()
```

	Name	Location	Year	Kilometers_Driven	Fuel_Type	Transmission	Owner_Type	Mileage	Engine	Power	Seats	New_Price	Price	Car_Age
0	Mazd Wagon R LXI CNG	Mumbai	2010	72000	CNG	Manual	First	26.6 kmpl	998 CC	56.16 bhp	5.0	NaN	1.75	14
1	Hyundai Creta 1.6 CRDi SX Option	Pune	2015	43000	Diesel	Manual	First	19.67 kmpl	1592 CC	126.2 bhp	5.0	NaN	1250	9
2	Hyundai Verna 1.6 CRDi	Chennai	2011	40000	Petrol	Manual	First	18.2 kmpl	1198 CC	88.7 bhp	5.0	8.61 Lakh	430	13
3	Mazd Ertiga VDi	Chennai	2012	87000	Diesel	Manual	First	20.77 kmpl	1248 CC	86.76 bhp	7.0	NaN	630	12
4	Audi A4 New 2.0 TDI Multitronic	Coimbatore	2013	40670	Diesel	Automatic	Second	15.2 kmpl	1968 CC	140.8 bhp	5.0	NaN	17.74	11

Since car names will not be great predictors of the price in our current data. But we can process this column to extract important information using brand and Model names. Let's split the name and introduce new variables "Brand" and "Model"

```
In [11]: data['Brand'] = data.Name.str.split().str.get(0)
data['Model'] = data.Name.str.split().str.get(1) + data.Name.str.split().str.get(2)
data[['Name', 'Brand', 'Model']]

Out[11]:
```

	Name	Brand	Model
0	Mazd Wagon R LXI CNG	Mazd	WagonR
1	Hyundai Creta 1.6 CRDi SX Option	Hyundai	Creta E
2	Honda Jazz V	Honda	Jazz V
3	Mazda Ertiga VDi	Mazd	Ertiga VDi
4	Audi A4 New 2.0 TDI Multitronic	Audi	A4 New
...
7248	Volkswagen Vento Diesel Trendline	Volkswagen	Vento Diesel
7249	Volkswagen Polo GT TSI	Volkswagen	Polo GT
7250	Nissan Micra Diesel XV	Nissan	Micra Diesel
7251	Volkswagen Polo GT TSI	Volkswagen	Polo GT
7252	Mercedes-Benz E-Class 2009-2013 E 220 CDI Avani...	Mercedes-Benz	E-Class 2009-2013

7252 rows x 3 columns

Data Cleaning/Wrangling

```
In [12]: #Some names of the variables are not relevant and not easy to understand. Some data may have data entry errors, and some variables may need data type conversion. we need to fix this issue in the data.In the example, The Brand name

In [13]: print(data.Brand.unique())
print(data.Brand.nunique())

['Maruti', 'Hyundai', 'Honda', 'Audi', 'Nissan', 'Toyota', 'Volkswagen', 'Tata',
'Land', 'Mitsubishi', 'Renault', 'Mercedes-Benz', 'BMW', 'Mahindra', 'Ford',
'Peugeot', 'Datsun', 'Jaguar', 'Volvo', 'Chevrolet', 'Suzuki', 'Mazda', 'Fiat',
'Jeep', 'Skoda', 'Ambassador', 'Isuzu', 'Subaru', 'Force', 'Bentley',
'LandCruiser', 'Mitsubishi', 'Suzuki']
33

In [14]: searchfor = ['Isuzu', 'Subaru', 'Mini', 'Land']
data[data.Brand.str.contains(' ', join='or')].head(5)
```

	Name	Location	Year	Kilometers_Driven	Fuel_Type	Transmission	Owner_Type	Mileage	Engine	Power	Seats	New_Price	Price	Car_Age	Brand	Model
13	Land Rover Range Rover 2.7i Pure	Delhi	2014	72000	Diesel	Automatic	First	12.7 kmpl	2179 CC	187.7 bhp	5.0	NaN	27.00	10	Land	RooverRange
14	Land Rover Freelander 2 Td4 SE	Pune	2012	85000	Diesel	Automatic	Second	0.0 kmpl	2179 CC	115 bhp	5.0	NaN	17.50	12	Land	RooverFreelander
176	Mit Countryman Cooper D	Jaipur	2017	8525	Diesel	Automatic	Second	16.4 kmpl	1998 CC	112 bhp	5.0	NaN	23.00	7	Mit	CountrymanCooper
285	Land Rover Range Rover 2.9 Supercharged	Coimbatore	2008	30001	Diesel	Automatic	First	12.7 kmpl	2179 CC	187.7 bhp	5.0	NaN	55.76	6	Land	RooverRange
228	Mit Cooper Convertible	Kochi	2017	26327	Petrol	Automatic	First	16.82 kmpl	1998 CC	189.08 bhp	4.0	44.28 Lakh	35.67	7	Mit	CooperConvertible

```
In [15]: data['Brand'].replace({'Isuzu': 'Isuzu', 'Mini': 'Mini Cooper', 'Land': 'Land Rover'}, inplace=True)

We have done the fundamental data analysis, Featurang, and data clean-up. Let's move to the EDA process
```

Our Data is ready to perform EDA.

EDA Exploratory Data Analysis

Exploratory Data Analysis refers to the crucial process of performing initial investigations on data to discover patterns to check assumptions with the help of summary statistics and graphical representations.

EDA can be leveraged to check for outliers, patterns, and trends in the given data. EDA helps to find meaningful patterns in data. EDA provides in-depth insights into the data sets to solve our business problems. EDA gives a clue to impute missing values in the dataset

Statistics Summary

```
In [16]: data.describe().T

Out[16]:
```

	count	mean	std	min	25%	50%	75%	max
Year	7252	2014.96366	3.254421	1996.00	2010.0	2014.00	2016.00	2019.0
Kilometers_Driven	7252	58699.962146	84427.726563	171.00	36000.0	53418.00	73000.00	650000.0
Seats	7200	5.279722	0.813660	0.00	5.0	5.00	5.00	10.0
Price	60150	9.479468	11.187917	0.44	3.5	5.64	9.95	160.0
Car_Age	7253.0	10.634864	3.254421	5.00	6.0	10.00	13.00	20.0

```
In [19]: #describe(include='all') provides a statistics summary of all data, include object, category etc
```

```
In [20]: data.describe(include='all').T

Out[20]:
```

	Name	count	unique	top	freq	mean	std	min	25%	50%	75%	max
Location	7253	11	Mumbai	949	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Year	7253	NaN	NaN	NaN	NaN	2013.37	3.25442	1996	2011	2014	2016	2019
Kilometers_Driven	7253	NaN	NaN	NaN	NaN	58699.1	84427.7	171	34000	53418	73000	654+06
Fuel_Type	7253	5	Diesel	3862	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Transmission	7253	2	Manual	3204	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Owner_Type	7253	4	First	1662	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Mileage	7251	450	17.0 kmpl	207	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Engine	7207	150	1197 CC	732	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Power	7207	386	74 bhp	280	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Seats	7200	NaN	NaN	NaN	NaN	5.27972	0.81366	0	5	5	5	10
New_Price	1886	165	63.71 Lakh	0	5	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Price	60150	NaN	NaN	NaN	NaN	9.47947	11.1879	0.44	3.5	5.64	9.95	160
Car_Age	7253	NaN	NaN	NaN	NaN	10.6346	3.25442	5	6	10	13	20
Brand	7253	32	Maruti	1444	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Model	7252	728	SwiftDzire	389	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

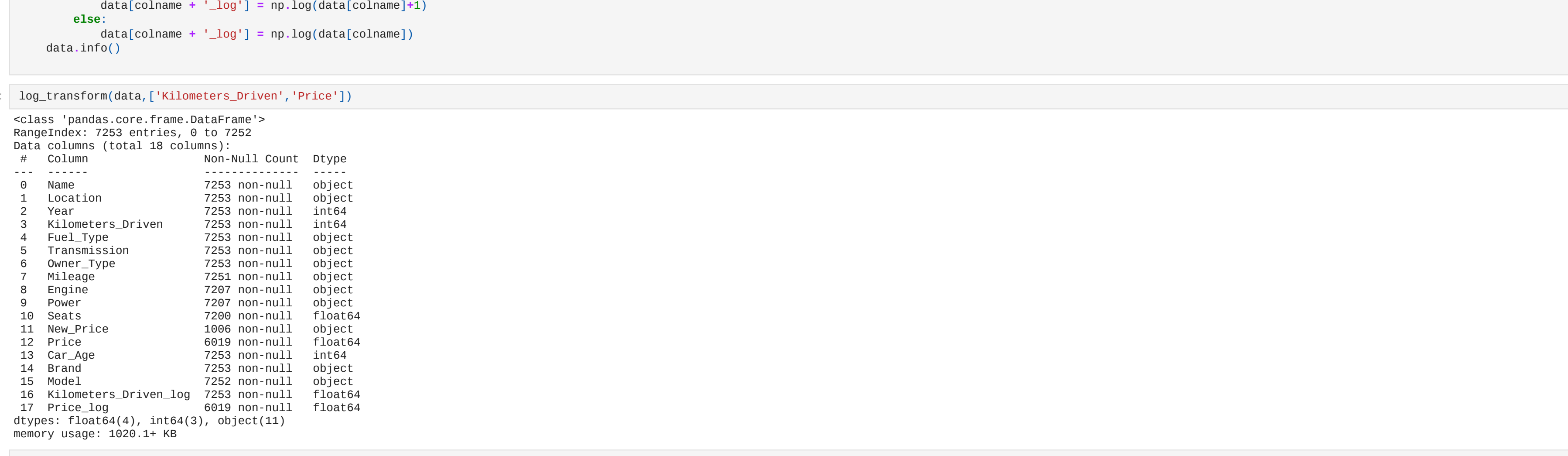
Before we do EDA, lets separate Numerical and categorical variables for easy analysis

```
In [21]: cat_cols = data.select_dtypes(include='object').columns
num_cols = data.select_dtypes(include='number').columns.tolist()
print('Categorical Variables:')
print(cat_cols)
print('Numerical Variables:')
print(num_cols)

Categorical Variables:
Index(['Name', 'Location', 'Fuel_Type', 'Transmission', 'Owner_Type',
'Mileage', 'Engine', 'Power', 'Seats', 'New_Price', 'Brand', 'Model'],
      dtype='object')
Numerical Variables:
['Year', 'Kilometers_Driven', 'Seats', 'Price', 'Car_Age']
```

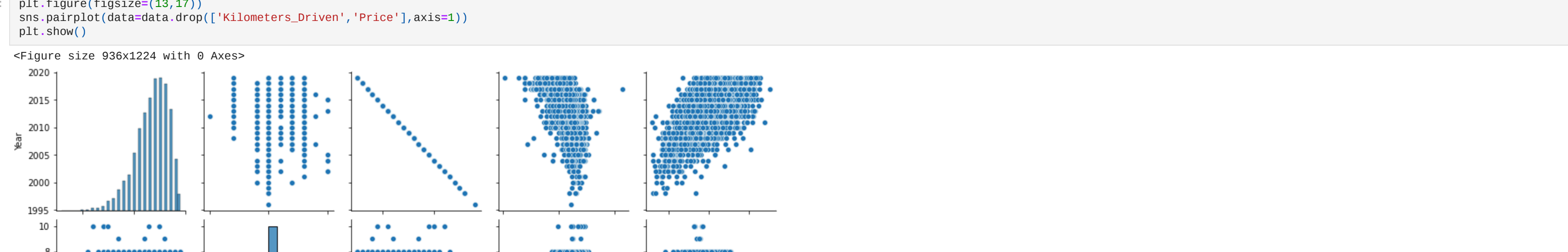
EDA Univariate Analysis

```
In [22]: for col in num_cols:
print(col)
print('Skew : ', round(data[col].skew(), 2))
plt.figure(figsize=(15, 4))
plt.subplot(1, 2, 1)
data[col].hist(grid=False)
plt.xlabel('count')
plt.subplot(1, 2, 2)
sns.boxplot(x=data[col])
plt.show()
```



categorical variables are being visualized using a count plot. Categorical variables provide the pattern of factors influencing car price

```
In [23]: fig, axes = plt.subplots(4, 2, figsize=(15, 10))
fig.suptitle('Bar plot for all categorical variables in the dataset')
sns.countplot(x=axes[0, 0], x='Fuel_Type', data=data, color='blue',
              order = data[Fuel_Type].value_counts().index)
sns.countplot(x=axes[0, 1], x='Transmission', data=data, color='blue',
              order = data[Transmission].value_counts().index)
sns.countplot(x=axes[1, 0], x='Owner_Type', data=data, color='blue',
              order = data[Owner_Type].value_counts().index)
sns.countplot(x=axes[1, 1], x='Location', data=data, color='blue',
              order = data[Location].value_counts().index)
sns.countplot(x=axes[2, 0], x='Brand', data=data, color='blue',
              order = data[Brand].head(20).value_counts().index)
sns.countplot(x=axes[2, 1], x='Model', data=data, color='blue',
              order = data[Model].head(20).value_counts().index)
axes[3][0].tick_params(labelrotation=45)
axes[3][0].set_title('Seats Vs Price')
axes[3][1].tick_params(labelrotation=45)
axes[3][1].set_title('Car_Age Vs Price')
```



Data Transformation

```
In [24]: # Function for log transformation of the column
def log_transform(data, col):
    if data[colname] == 1.0:
        data[colname] = np.log(data[colname]+1)
    else:
        data[colname] = np.log(data[colname]+1)
    data.info()
```

```
In [27]: log_transform(data, 'Kilometers_Driven', 'Price')

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7252 entries, 0 to 7252
Data columns (total 18 columns):
#   Column                Non-Null Count  Dtype
#  --  --                --
#   0   Name                  7252 non-null    object
#   1   Location              7252 non-null    object
#   2   Year                  7252 non-null    int64
#   3   Kilometers_Driven     7252 non-null    float64
#   4   Fuel_Type             7252 non-null    object
#   5   Transmission          7252 non-null    object
#   6   Owner_Type           7252 non-null    object
#   7   Mileage               7252 non-null    float64
#   8   Engine               7207 non-null    object
#   9   Power                7207 non-null    float64
#  10   Seats                6619 non-null    float64
#  11   New_Price            1886 non-null    object
#  12   Car_Age              7253 non-null    int64
#  13   Brand                7253 non-null    object
#  14   Kilometers_Driven_log 7252 non-null    float64
#  15   Price_log             6619 non-null    float64
dtypes: float64(4), int64(3), object(11)
memory usage: 1828+ KB
```

```
In [28]: #log transformation of the feature 'Kilometers_Driven'
sns.distplot(data[['Kilometers_Driven_log']], axlabel='Kilometers_Driven_log');
```



EDA Bivariate Analysis

```
In [29]: plt.figure(figsize=(12, 17))
sns.pairplot(data.drop(['Kilometers_Driven', 'Price'], axis=1), corr=True, annot = True, vmin = -1, vmax = 1)
```



A bar plot can be used to show the relationship between Categorical variables and continuous variables

```
In [30]: fig, axarr = plt.subplots(4, 2, figsize=(12, 10))
data.groupby('Location')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[0][0], fontsize=12)
data.groupby('Transmission')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[0][1], fontsize=12)
axarr[0][0].set_title('Location Vs Price')
axarr[0][1].set_title('Transmission Vs Price')
data.groupby('Owner_Type')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[1][0], fontsize=12)
axarr[1][0].set_title('Owner_Type Vs Price')
data.groupby('Brand')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[1][1], fontsize=12)
axarr[1][1].set_title('Brand Vs Price')
data.groupby('Model')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[2][0], fontsize=12)
axarr[2][0].set_title('Model Vs Price')
data.groupby('Seats')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[2][1], fontsize=12)
axarr[2][1].set_title('Seats Vs Price')
data.groupby('Car_Age')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[3][0], fontsize=12)
axarr[3][0].set_title('Car_Age Vs Price')
plt.subplots_adjust(wspace=0.5)
plt.subplots_adjust(hspace=0.5)
sns.despine()
```



EDA Multivariate Analysis

```
In [31]: #A heat map is widely used for Multivariate Analysis
```

