

Ques-1: Perform EDA

```
In [1]: import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
import seaborn as sns

In [2]: df = dfpd.read_csv('Cars.csv')

In [3]: df

Out[2]:
   Name      Location  Year  Kilometers_Driven  Fuel_Type  Transmission  Owner_Type  Mileage  Engine  Power  Colour  Seats  No. of Doors  New_Price  Price
0  Mahindra Scorpio  Pune  2012.0  96000.0  Diesel  Manual  Third  12.05 kmpl  2179 CC  120 bhp  Black/Silver  8.0  5.0  NaN  6.00
1  Maruti Baleno     Kochi  2018.0  18678.0  Petrol  Manual  First  21.1 kmpl  998 CC  100 bhp  Others  5.0  4.0  NaN  8.32
2  Mahindra Xylo     Bangalore  2013.0  197000.0  Diesel  Manual  First  11.66 kmpl  2498 CC  112 bhp  White  7.0  5.0  NaN  4.00
3  Hyundai Grand     Delhi  2014.0  45000.0  Diesel  Manual  First  24.0 kmpl  1120 CC  70 bhp  White  5.0  4.0  NaN  3.49
4  Toyota Innova     Delhi  2011.0  65000.0  Diesel  Manual  First  12.8 kmpl  2494 CC  102 bhp  Others  8.0  5.0  NaN  6.40
...
5956  Honda Civic     Pune  2011.0  47000.0  Petrol  Automatic  Second  13.9 kmpl  1799 CC  130.3 bhp  Others  5.0  4.0  NaN  4.50
5957  Hyundai i20     Delhi  2013.0  63777.0  Petrol  Manual  First  20.4 kmpl  1197 CC  81.80 bhp  Black/Silver  2.0  2.0  NaN  5.54
5958  Maruti Swift     Coimbatore  2016.0  37806.0  Petrol  Manual  First  20.4 kmpl  1197 CC  81.80 bhp  Black/Silver  2.0  2.0  NaN  5.20
5959  Mercedes-Benz SLK-Class  Coimbatore  2016.0  22732.0  Petrol  Automatic  First  18.1 kmpl  3498 CC  306 bhp  Black/Silver  5.0  4.0  NaN  3.60
5960  Hyundai i20     Kolkata  2016.0  7000.0  Petrol  Manual  First  20.36 kmpl  1197 CC  78.9 bhp  White  5.0  4.0  NaN  3.60

5961 rows x 15 columns

In [4]: df.head()

Out[4]:
   Name      Location  Year  Kilometers_Driven  Fuel_Type  Transmission  Owner_Type  Mileage  Engine  Power  Colour  Seats  No. of Doors  New_Price  Price
0  Mahindra Scorpio  Pune  2012.0  96000.0  Diesel  Manual  Third  12.05 kmpl  2179 CC  120 bhp  Black/Silver  8.0  5.0  NaN  6.00
1  Maruti Baleno     Kochi  2018.0  18678.0  Petrol  Manual  First  21.1 kmpl  998 CC  100 bhp  Others  5.0  4.0  NaN  8.32
2  Mahindra Xylo     Bangalore  2013.0  197000.0  Diesel  Manual  First  11.66 kmpl  2498 CC  112 bhp  White  7.0  5.0  NaN  4.00
3  Hyundai Grand     Delhi  2014.0  45000.0  Diesel  Manual  First  24.0 kmpl  1120 CC  70 bhp  White  5.0  4.0  NaN  3.49
4  Toyota Innova     Delhi  2011.0  65000.0  Diesel  Manual  First  12.8 kmpl  2494 CC  102 bhp  Others  8.0  5.0  NaN  6.40

In [5]: df.tail()

Out[5]:
5956  Honda Civic     Pune  2011.0  47000.0  Petrol  Automatic  Second  13.9 kmpl  1799 CC  130.3 bhp  Others  5.0  4.0  NaN  4.50
5957  Hyundai i20     Delhi  2013.0  63777.0  Petrol  Manual  First  20.4 kmpl  1197 CC  81.80 bhp  Black/Silver  2.0  2.0  NaN  5.54
5958  Maruti Swift     Coimbatore  2016.0  37806.0  Petrol  Manual  First  20.4 kmpl  1197 CC  81.80 bhp  Black/Silver  2.0  2.0  NaN  5.20
5959  Mercedes-Benz SLK-Class  Coimbatore  2016.0  22732.0  Petrol  Automatic  First  18.1 kmpl  3498 CC  306 bhp  Black/Silver  5.0  4.0  NaN  3.60
5960  Hyundai i20     Kolkata  2016.0  7000.0  Petrol  Manual  First  20.36 kmpl  1197 CC  78.9 bhp  White  5.0  4.0  NaN  3.60

In [6]: for col in df.columns:
    print(col)

Name
Location
Year
Kilometers_Driven
Fuel_Type
Transmission
Owner_Type
Mileage
Engine
Power
Colour
Seats
No. of Doors
New_Price
Price

In [7]: df.shape
(5961, 15)

In [8]: df['Location'].value_counts()

Out[8]:
Hyderabad    739
Mumbai       639
Coimbatore    639
Delhi         549
Kolkata       526
Chennai       489
Jaipur         406
Bangalore     351
Ahmedabad     222
Name: Location, dtype: int64

In [9]: df['Year'].value_counts()

Out[9]:
2014.0    793
2016.0    740
2015.0    706
2013.0    642
2017.0    586
2012.0    573
2011.0    461
2018.0    338
2018.0    298
2009.0    196
2088.0    170
2094.0    123
2019.0    101
2086.0    75
2085.0    55
2084.0    28
2082.0    14
2083.0    13
2081.0    7
2080.0    4
1995.0    4
1999.0    2
Name: Year, dtype: int64

In [10]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5961 entries, 0 to 5960
Data columns (total 15 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0  Name                5961 non-null   object
 1  Location            5960 non-null   object
 2  Year                5960 non-null   float64
 3  Kilometers_Driven  5963 non-null   float64
 4  Fuel_Type           5961 non-null   object
 5  Transmission        5964 non-null   object
 6  Owner_Type          5946 non-null   object
 7  Mileage             5959 non-null   object
 8  Engine              5944 non-null   object
 9  Power               5929 non-null   object
10  Colour              5950 non-null   object
11  Seats               5956 non-null   float64
12  No. of Doors        5960 non-null   float64
13  New_Price           5961 non-null   float64
14  Price               5961 non-null   float64
memory usage: 698.7+ KB

In [11]: df.isnull()

Out[11]:
Name                0
Location            0
Year                0
Kilometers_Driven  3668
Fuel_Type           0
Transmission        0
Owner_Type          0
Mileage             0
Engine              0
Power               0
Colour              0
Seats               0
No. of Doors        0
New_Price           0
Price               0
dtype: int64

In [12]: df.isnull().sum()

Out[12]:
Name                0
Location            0
Year                0
Kilometers_Driven  3668
Fuel_Type           0
Transmission        0
Owner_Type          0
Mileage             0
Engine              0
Power               0
Colour              0
Seats               0
No. of Doors        0
New_Price           0
Price               0
dtype: int64

In [13]: (df.isnull().sum()/len(df))*100

Out[13]:
Name                0.000000
Location            0.000000
Year                0.000000
Kilometers_Driven  61.533551
Fuel_Type           0.000000
Transmission        0.000000
Owner_Type          0.000000
Mileage             0.000000
Engine              0.000000
Power               0.000000
Colour              0.000000
Seats               0.000000
No. of Doors        0.000000
New_Price           0.000000
Price               0.000000
dtype: float64

In [14]: df['brand'] = df.Name.str.split().str.get(0)
df['Model'] = df.Name.str.split().str.get(1) + df.Name.str.split().str.get(2)
df[['Name', 'brand', 'Model']]

Out[14]:
   Name      brand  Model
0  Mahindra Scorpio  Mahindra  NaN
1  Maruti Baleno     Maruti  NaN
2  Mahindra Xylo     Mahindra  NaN
3  Hyundai Grand     Hyundai  NaN
4  Toyota Innova     Toyota  NaN
...
5956  Honda Civic     Honda  NaN
5957  Hyundai i20     Hyundai  NaN
5958  Maruti Swift     Maruti  NaN
5959  Mercedes-Benz SLK-Class  Mercedes-Benz  NaN
5960  Hyundai i20     Hyundai  NaN

5961 rows x 3 columns

In [15]: print(df.brand.unique())
print(df.Model.unique())

['Mahindra' 'Maruti' 'Hyundai' 'Toyota' 'Honda' 'Chevrolet' 'Audi' 'Skoda'
'Renault' 'Land' 'Datsun' 'Isuzu' 'Jaguar' 'Mercedes-Benz' 'Volvo' 'Audi'
'Tata' 'Mitsubishi' 'Ford' 'Nissan' 'Volvo' 'Fiat' 'Porsche' 'Mini'
'Datsun' 'Jaguar' 'Force' 'Isuzu' 'Jaguar' 'Land Rover' 'Bentley']
30

In [16]: searchfor = ['Isuzu', 'ISUZU', 'Mini', 'Land']
df[df.brand.str.contains('').join(searchfor)]]

Out[16]:
   Name      Location  Year  Kilometers_Driven  Fuel_Type  Transmission  Owner_Type  Mileage  Engine  Power  Colour  Seats  No. of Doors  New_Price  Price  Brand  Model
18  Land Rover Freelander  Hyderabad  2012.0  139000.0  Diesel  Automatic  First  10.9 kmpl  2179 CC  115 bhp  Black/Silver  5.0  4.0  NaN  16.75  Land  RoverFreander
19  ISUZU D-MAX          Jaipur  2017.0  25000.0  Diesel  Manual  First  12.4 kmpl  2499 CC  134 bhp  Others  5.0  4.0  NaN  8.00  ISUZU  NaN
97  Land Rover Range     Delhi  2015.0  37000.0  Diesel  Automatic  First  12.7 kmpl  2179 CC  187.7 bhp  White  5.0  4.0  NaN  36.00  Land  RoverRange
103  Land Rover Range     Coimbatore  2017.0  49275.0  Diesel  Automatic  First  12.7 kmpl  2179 CC  187.7 bhp  Others  5.0  4.0  NaN  45.64  Land  RoverRange
118  Land Rover Range     Coimbatore  2013.0  75995.0  Diesel  Automatic  Second  11.49 kmpl  4367 CC  335.3 bhp  White  5.0  4.0  NaN  65.81  Land  RoverRange

In [17]: df['brand'].replace(['ISUZU','Isuzu','Mini','Mini Cooper','Land','Land Rover'], inplace=True)

In [18]: df.describe().T

Out[18]:
          count      mean      std      min      25%      50%      75%      max
Kilometers_Driven  5953.0  58711.006118  324305.1  1998.00  2011.5  2014.00  2016.0  2019.0
Seats              5956.0  5.269140  0.789048  2.00  5.0  5.00  5.0  10.0
No. of Doors       5960.0  4.114933  0.344757  2.00  4.0  4.00  4.0  5.0
Price              5961.0  9.528103  11.214382  0.44  3.5  5.66  10.0  160.0

In [19]: df.describe(include='all').T

Out[19]:
          count      mean      top  top_freq      mean      std      min      25%      50%      75%      max
Name                5961  5961  Maruti Swift  343      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Location            5950  11  Mumbai  781      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Year                5959.0  NaN      NaN      NaN      2013.389159  324305.1  1998.0  2011.5  2014.0  2016.0  2019.0
Kilometers_Driven  5953.0  NaN      NaN      NaN      58711.006118  324305.1  1998.0  2011.5  2014.0  2016.0  2019.0
Fuel_Type           5961  5  Diesel  3188      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Transmission        5934  2  Manual  4225      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Owner_Type          5946  4  First  4875      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Mileage             5959  420  18.9 kmpl  22      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Engine              5944  143  1197 CC  606      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Power               5929  369  74 hp  233      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Colour              5950  3  White  2115      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Seats               5956.0  NaN      NaN      NaN      5.26914  0.789048  2.0  5.0  5.0  5.0  10.0
No. of Doors        5960.0  NaN      NaN      NaN      4.114933  0.344757  2.0  4.0  4.0  4.0  5.0
New_Price           824  540  4.78 Lakh  6      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Price              5961  29  Maruti  1189      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
Model              57  3  RoverRange  26      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN

In [20]: cat_cols=df.select_types(include=['object']).columns
num_cols = df.select_types(include=np.number).columns.tolist()
print('Categorical Variables:')
print('Numerical Variables:')
print(num_cols)

Categorical Variables:
Index(['Name', 'Location', 'Fuel_Type', 'Transmission', 'Owner_Type',
      'Mileage', 'Engine', 'Power', 'Colour', 'New_Price', 'Brand', 'Model'],
      dtype='object')
Numerical Variables:
['Year', 'Kilometers_Driven', 'Seats', 'No. of doors', 'Price']

In [21]: for col in num_cols:
    print(col)
    print('Skew : ', round(df[col].skew(), 2))
    plt.figure(figsize=(15, 4))
    plt.subplot(2, 2, 1)
    df[col].hist(grid=False)
    plt.ylabel('count')
    plt.subplot(2, 2, 2)
    sns.boxplot(x=df[col])
    plt.show()

Year
Skew : -0.84

Kilometers_Driven
Skew : 58.152

Seats
Skew : 1.85

No. of Doors
Skew : 1.35

Price
Skew : 3.33

fig, axes = plt.subplots(4, 2, figsize=(18, 18))
sns.boxplot(x=axes[0, 0], x = 'Fuel_Type', data = df, color = 'blue',
            order = df['Fuel_Type'].value_counts().index)
sns.boxplot(x=axes[0, 1], x = 'Transmission', data = df, color = 'blue',
            order = df['Transmission'].value_counts().index)
sns.boxplot(x=axes[1, 0], x = 'Owner_Type', data = df, color = 'blue',
            order = df['Owner_Type'].value_counts().index)
sns.boxplot(x=axes[1, 1], x = 'Location', data = df, color = 'blue',
            order = df['Location'].value_counts().index)
sns.boxplot(x=axes[2, 0], x = 'Brand', data = df, color = 'blue',
            order = df['Brand'].head(28).value_counts().index)
sns.boxplot(x=axes[2, 1], x = 'Model', data = df, color = 'blue',
            order = df['Model'].head(28).value_counts().index)
axes[1][1].tick_params(labelrotation=45)
axes[2][1].tick_params(labelrotation=90)
plt.show()

Bar plot for all categorical variables in 'brand'

In [22]: # Function for log transformation of the column
def log_transform(data,col):
    for colname in col:
        if df[colname].isnull().all():
            df[colname] = _log = np.log(df[colname]+1)
        else:
            df[colname] = _log = np.log(df[colname])

In [24]: log_transform(df,['Kilometers_Driven','Price'])

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5961 entries, 0 to 5960
Data columns (total 15 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0  Name                5961 non-null   object
 1  Location            5960 non-null   object
 2  Year                5960 non-null   float64
 3  Kilometers_Driven  5963 non-null   float64
 4  Fuel_Type           5961 non-null   object
 5  Transmission        5964 non-null   object
 6  Owner_Type          5946 non-null   object
 7  Mileage             5959 non-null   object
 8  Engine              5944 non-null   object
 9  Power               5929 non-null   object
10  Colour              5950 non-null   object
11  Seats               5956 non-null   float64
12  No. of Doors        5960 non-null   float64
13  New_Price           5961 non-null   float64
14  Brand              5961 non-null   object
15  Kilometers_Driven_log  5963 non-null   float64
16  Price_log           5961 non-null   float64
17  Price_log           5961 non-null   float64
dtype: object
memory usage: 885.6+ KB

In [25]: #log transformation of the feature 'Kilometers_Driven'
sns.boxplot(df['Kilometers_Driven_log'], axlabel='Kilometers_Driven_log');

C:\Users\Hardi\AppData\Local\Temp\ipykernel_19180\681597622.py(2): UserWarning:
'displot' is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either displot (a figure-level function for histograms)
or histplot (a FacetGrid or histplot on axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://github.com/mwaskom/dea4146ed297487ad693727586be5781

sns.displot(df['Kilometers_Driven_log'], axlabel='Kilometers_Driven_log');

In [26]: plt.figure(figsize=(12, 7))
sns.pairplot(df.drop(['Kilometers_Driven','Price'],axis=1),corr(), annot = True, vsin = -1, vmax = 1)
plt.show()

<Figure size 936x1224 with 0 Axes>

In [27]: fig, axarr = plt.subplots(4, 2, figsize=(12, 18))
df.groupby('Location')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[0][0], fontsize=12)
axarr[0][0].set_title('Location Vs Price', fontsize=18)
df.groupby('Transmission')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[0][1], fontsize=12)
axarr[0][1].set_title('Transmission Vs Price', fontsize=18)
df.groupby('Owner_Type')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[1][0], fontsize=12)
axarr[1][0].set_title('Fuel_Type Vs Price', fontsize=18)
df.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', fontsize=18)
df.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', fontsize=18)
df.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', fontsize=18)
plt.subplots_adjust(wspace=0.5)
sns.despine()

Location Vs Price

Transmission Vs Price

Fuel_Type Vs Price

Owner_Type Vs Price

Brand Vs Price

Seats Vs Price

In [28]: plt.figure(figsize=(12, 7))
plt.barlist(df['Model'].value_counts()[0:2], keys=[1], colors='g')
plt.show()

In [29]: df['Location'].value_counts()[0:16]

Out[29]:
Hyderabad    739
Mumbai       639
Coimbatore    639
Delhi         549
Kolkata       526
Chennai       489
Jaipur         406
Bangalore     351
Ahmedabad     222
Name: Location, dtype: int64

In [30]: df['Brand'].value_counts()[0:16]

Out[30]:
Maruti       1159
Hyundai      602
Toyota       418
Mercedes-Benz 318
Volkswagen   315
Ford         236
Mahindra     226
BMW          267
Audi         236
Name: Brand, dtype: int64

In [31]: df['Model'].value_counts()[0:16]

Out[31]:
RoverRange    28
RoverFreander 17
RoverDiscovery 12
Name: Model, dtype: int64

In [32]: list(df['Price_log'].value_counts()[0:2].keys())

Out[32]:
['RoverRange', 'RoverFreander']

Ques-2: Comparison two model
```

```
In [33]: plt.figure(figsize=(8,5))
plt.barlist(df['Model'].value_counts()[0:2].keys(),list(df['Model'].value_counts()[0:2], colors='g')
plt.show()

In [34]: # Function for log transformation of the column
def log_transform(data,col):
    for colname in col:
        if df[colname].isnull().all():
            df[colname] = _log = np.log(df[colname]+1)
        else:
            df[colname] = _log = np.log(df[colname])

In [24]: log_transform(df,['Kilometers_Driven','Price'])

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5961 entries, 0 to 5960
Data columns (total 15 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0  Name                5961 non-null   object
 1  Location            5960 non-null   object
 2  Year                5960 non-null   float64
 3  Kilometers_Driven  5963 non-null   float64
 4  Fuel_Type           5961 non-null   object
 5  Transmission        5964 non-null   object
 6  Owner_Type          5946 non-null   object
 7  Mileage             5959 non-null   object
 8  Engine              5944 non-null   object
 9  Power               5929 non-null   object
10  Colour              5950 non-null   object
11  Seats               5956 non-null   float64
12  No. of Doors        5960 non-null   float64
13  New_Price           5961 non-null   float64
14  Brand              5961 non-null   object
15  Kilometers_Driven_log  5963 non-null   float64
16  Price_log           5961 non-null   float64
17  Price_log           5961 non-null   float64
dtype: object
memory usage: 885.6+ KB

In [25]: #log transformation of the feature 'Kilometers_Driven'
sns.boxplot(df['Kilometers_Driven_log'], axlabel='Kilometers_Driven_log');

C:\Users\Hardi\AppData\Local\Temp\ipykernel_19180\681597622.py(2): UserWarning:
'displot' is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either displot (a figure-level function for histograms)
or histplot (a FacetGrid or histplot on axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://github.com/mwaskom/dea4146ed297487ad693727586be5781

sns.displot(df['Kilometers_Driven_log'], axlabel='Kilometers_Driven_log');

In [26]: plt.figure(figsize=(12, 7))
sns.pairplot(df.drop(['Kilometers_Driven','Price'],axis=1),corr(), annot = True, vsin = -1, vmax = 1)
plt.show()

<Figure size 936x1224 with 0 Axes>

In [27]: fig, axarr = plt.subplots(4, 2, figsize=(12, 18))
df.groupby('Location')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[0][0], fontsize=12)
axarr[0][0].set_title('Location Vs Price', fontsize=18)
df.groupby('Transmission')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[0][1], fontsize=12)
axarr[0][1].set_title('Transmission Vs Price', fontsize=18)
df.groupby('Owner_Type')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[1][0], fontsize=12)
axarr[1][0].set_title('Fuel_Type Vs Price', fontsize=18)
df.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', fontsize=18)
df.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', fontsize=18)
df.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', fontsize=18)
plt.subplots_adjust(wspace=0.5)
sns.despine()

Location Vs Price

Transmission Vs Price

Fuel_Type Vs Price

Owner_Type Vs Price

Brand Vs Price

Seats Vs Price

In [28]: plt.figure(figsize=(12, 7))
plt.barlist(df['Model'].value_counts()[0:2], keys=[1], colors='g')
plt.show()

In [29]: df['Location'].value_counts()[0:16]

Out[29]:
Hyderabad    739
Mumbai       639
Coimbatore    639
Delhi         549
Kolkata       526
Chennai       489
Jaipur         406
Bangalore     351
Ahmedabad     222
Name: Location, dtype: int64

In [30]: df['Brand'].value_counts()[0:16]

Out[30]:
Maruti       1159
Hyundai      602
Toyota       418
Mercedes-Benz 318
Volkswagen   315
Ford         236
Mahindra     226
BMW          267
Audi         236
Name: Brand, dtype: int64

In [31]: df['Model'].value_counts()[0:16]

Out[31]:
RoverRange    28
RoverFreander 17
RoverDiscovery 12
Name: Model, dtype: int64

In [32]: list(df['Price_log'].value_counts()[0:2].keys())

Out[32]:
['RoverRange', 'RoverFreander']

Ques-2: Comparison two model
```

```
In [33]: plt.figure(figsize=(8,5))
plt.barlist(df['Model'].value_counts()[0:2].keys(),list(df['Model'].value_counts()[0:2], colors='g')
plt.show()

In [34]: # Function for log transformation of the column
def log_transform(data,col):
    for colname in col:
        if df[colname].isnull().all():
            df[colname] = _log = np.log(df[colname]+1)
        else:
            df[colname] = _log = np.log(df[colname])

In [24]: log_transform(df,['Kilometers_Driven','Price'])

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5961 entries, 0 to 5960
Data columns (total 15 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0  Name                5961 non-null   object
 1  Location            5960 non-null   object
 2  Year                5960 non-null   float64
 3  Kilometers_Driven  5963 non-null   float64
 4  Fuel_Type           5961 non-null   object
 5  Transmission        5964 non-null   object
 6  Owner_Type          5946 non-null   object
 7  Mileage             5959 non-null   object
 8  Engine              5944 non-null   object
 9  Power               5929 non-null   object
10  Colour              5950 non-null   object
11  Seats               5956 non-null   float64
12  No. of Doors        5960 non-null   float64
13  New_Price           5961 non-null   float64
14  Brand              5961 non-null   object
15  Kilometers_Driven_log  5963 non-null   float64
16  Price_log           5961 non-null   float64
17  Price_log           5961 non-null   float64
dtype: object
memory usage: 885.6+ KB

In [25]: #log transformation of the feature 'Kilometers_Driven'
sns.boxplot(df['Kilometers_Driven_log'], axlabel='Kilometers_Driven_log');

C:\Users\Hardi\AppData\Local\Temp\ipykernel_19180\681597622.py(2): UserWarning:
'displot' is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either displot (a figure-level function for histograms)
or histplot (a FacetGrid or histplot on axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://github.com/mwaskom/dea4146ed297487ad693727586be5781

sns.displot(df['Kilometers_Driven_log'], axlabel='Kilometers_Driven_log');

In [26]: plt.figure(figsize=(12, 7))
sns.pairplot(df.drop(['Kilometers_Driven','Price'],axis=1),corr(), annot = True, vsin = -1, vmax = 1)
plt.show()

<Figure size 936x1224 with 0 Axes>

In [27]: fig, axarr = plt.subplots(4, 2, figsize=(12, 18))
df.groupby('Location')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[0][0], fontsize=12)
axarr[0][0].set_title('Location Vs Price', fontsize=18)
df.groupby('Transmission')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[0][1], fontsize=12)
axarr[0][1].set_title('Transmission Vs Price', fontsize=18)
df.groupby('Owner_Type')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[1][0], fontsize=12)
axarr[1][0].set_title('Fuel_Type Vs Price', fontsize=18)
df.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', fontsize=18)
df.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', fontsize=18)
df.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', fontsize=18)
plt.subplots_adjust(wspace=0.5)
sns.despine()

Location Vs Price

Transmission Vs Price

Fuel_Type Vs Price

Owner_Type Vs Price

Brand Vs Price

Seats Vs Price

In [28]: plt.figure(figsize=(12, 7))
plt.barlist(df['Model'].value_counts()[0:2], keys=[1], colors='g')
plt.show()

In [29]: df['Location'].value_counts()[0:16]

Out[29]:
Hyderabad    739
Mumbai       639
Coimbatore    639
Delhi         549
Kolkata       526
Chennai       489
Jaipur         406
Bangalore     351
Ahmedabad     222
Name: Location, dtype: int64

In [30]: df['Brand'].value_counts()[0:16]

Out[30]:
Maruti       1159
Hyundai      602
Toyota       418
Mercedes-Benz 318
Volkswagen   315
Ford         236
Mahindra     226
BMW          267
Audi         236
Name: Brand, dtype: int64

In [31]: df['Model'].value_counts()[0:16]

Out[31]:
RoverRange    28
RoverFreander 17
RoverDiscovery 12
Name: Model, dtype: int64

In [32]: list(df['Price_log'].value_counts()[0:2].keys())

Out[32]:
['RoverRange', 'RoverFreander']

Ques-2: Comparison two model
```

```
In [33]: plt.figure(figsize=(8,5))
plt.barlist(df['Model'].value_counts()[0:2].keys(),list(df['Model'].value_counts()[0:2], colors='g')
plt.show()

In [34]: # Function for log transformation of the column
def log_transform(data,col):
    for colname in col:
        if df[colname].isnull().all():
            df[colname] = _log = np.log(df[colname]+1)
        else:
            df[colname] = _log = np.log(df[colname])

In [24]: log_transform(df,['Kilometers_Driven','Price'])

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5961 entries, 0 to 5960
Data columns (total 15 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0  Name                5961 non-null   object
 1  Location            5960 non-null   object
 2  Year                5960 non-null   float64
 3  Kilometers_Driven  5963 non-null   float64
 4  Fuel_Type           5961 non-null   object
 5  Transmission        5964 non-null   object
 6  Owner_Type          5946 non-null   object
 7  Mileage             5959 non-null   object
 8  Engine              5944 non-null   object
 9  Power               5929 non-null   object
10  Colour              5950 non-null   object
11  Seats               5956 non-null   float64
12  No. of Doors        5960 non-null   float64
13  New_Price           5961 non-null   float64
14  Brand              5961 non-null   object
15  Kilometers_Driven_log  5963 non-null   float64
16  Price_log           5961 non-null   float64
17  Price_log           5961 non-null   float64
dtype: object
memory usage: 885.6+ KB

In [25]: #log transformation of the feature 'Kilometers_Driven'
sns.boxplot(df['Kilometers_Driven_log'], axlabel='Kilometers_Driven_log');

C:\Users\Hardi\AppData\Local\Temp\ipykernel_19180\681597622.py(2): UserWarning:
'displot' is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either displot (a figure-level function for histograms)
or histplot (a FacetGrid or histplot on axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://github.com/mwaskom/dea4146ed297487ad693727586be5781

sns.displot(df['Kilometers_Driven_log'], axlabel='Kilometers_Driven_log');

In [26]: plt.figure(figsize=(12, 7))
sns.pairplot(df.drop(['Kilometers_Driven','Price'],axis=1),corr(), annot = True, vsin = -1, vmax = 1)
plt.show()

<Figure size 936x1224 with 0 Axes>

In [27]: fig, axarr = plt.subplots(4, 2, figsize=(12, 18))
df.groupby('Location')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[0][0], fontsize=12)
axarr[0][0].set_title('Location Vs Price', fontsize=18)
df.groupby('Transmission')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[0][1], fontsize=12)
axarr[0][1].set_title('Transmission Vs Price', fontsize=18)
df.groupby('Owner_Type')['Price_log'].mean().sort_values(ascending=False).plot.bar(ax=axarr[1][0], fontsize=12)
axarr[1][0].set_title('Fuel_Type Vs Price', fontsize=18)
df.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', fontsize=18)
df.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', fontsize=18)
df.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', fontsize=18)
plt.subplots_adjust(wspace=0.5)
sns.despine()

Location Vs Price

Transmission Vs Price

Fuel_Type Vs Price

Owner_Type Vs Price

Brand Vs Price

Seats Vs Price

In [28]: plt.figure(figsize=(12, 7))
plt.barlist(df['Model'].value_counts()[0:2], keys=[1], colors='g')
plt.show()

In [29]: df['Location'].value_counts()[0:16]

Out[29]:
Hyderabad    739
Mumbai       639
Coimbatore    639
Delhi         549
Kolkata       526
Chennai       489
Jaipur         406
Bangalore     351
Ahmedabad     222
Name: Location, dtype: int64

In [30]: df['Brand'].value_counts()[0:16]

Out[30]:
Maruti       1159
Hyundai      602
Toyota       418
Mercedes-Benz 318
Volkswagen   315
Ford         236
Mahindra     226
BMW          267
Audi         236
Name: Brand, dtype: int64

In [31]: df['Model'].value_counts()[0:16]

Out[31]:
RoverRange    28
RoverFreander 17
RoverDiscovery 12
Name: Model, dtype: int64

In [32]: list(df['Price_log'].value_counts()[0:2].keys())

Out[32]:
['RoverRange', 'RoverFreander']

Ques-2: Comparison two model
```

```
In [33]: plt.figure(figsize=(8,5))
plt.barlist(df['Model'].value_counts()[0:2].keys(),list(df['Model'].value_counts()[0:2], colors='g')
plt.show()

In [34]: # Function for log transformation of the column
def log_transform(data,col):
    for colname in col:
        if df[colname].isnull().all():
            df[colname] = _log = np.log(df[colname]+1)
        else:
            df[colname] = _log = np.log(df[colname])

In [24]: log_transform(df,['Kilometers_Driven','Price'])

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5961 entries, 0 to 5960
Data columns (total 15 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0  Name                5961 non-null   object
 1  Location            5960 non-null   object
 2  Year                5960 non-null   float64
 3  Kilometers_Driven  5963 non-null   float64
 4  Fuel_Type           5961 non-null   object
 5  Transmission        5964 non-null   object
 6  Owner_Type          5946 non-null   object
 7  Mileage             5959 non-null   object
 8  Engine              5944 non-null   object
 9  Power               5929 non-null  
```