question

Problem Statement: Data Analysis The dataset contains more of 10, 000 rows and more than 10 columns which contains features of the car and its (MSRP) manufacturer's suggested retail price. Clean the data and analyse it making it ready for modelling.

solution

```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```python
dt=pd.read_csv('Documents/github/JeetPython_tops/assignment/data analisis with pyth
```

```python
dt
```

| | Make | Model | Year | Engine Fuel Type | Engine HP | Engine Cylinders | Transmission Type | Driven_Wheels | Nu |
|---|---|---|---|---|---|---|---|---|---|
| 0 | BMW | 1 Series M | 2011 | premium unleaded (required) | 335.0 | 6.0 | MANUAL | rear wheel drive | |
| 1 | BMW | 1 Series | 2011 | premium unleaded (required) | 300.0 | 6.0 | MANUAL | rear wheel drive | |
| 2 | BMW | 1 Series | 2011 | premium unleaded (required) | 300.0 | 6.0 | MANUAL | rear wheel drive | |
| 3 | BMW | 1 Series | 2011 | premium unleaded (required) | 230.0 | 6.0 | MANUAL | rear wheel drive | |
| 4 | BMW | 1 Series | 2011 | premium unleaded (required) | 230.0 | 6.0 | MANUAL | rear wheel drive | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 11909 | Acura | ZDX | 2012 | premium unleaded (required) | 300.0 | 6.0 | AUTOMATIC | all wheel drive | |
| 11910 | Acura | ZDX | 2012 | premium unleaded (required) | 300.0 | 6.0 | AUTOMATIC | all wheel drive | |
| 11911 | Acura | ZDX | 2012 | premium unleaded (required) | 300.0 | 6.0 | AUTOMATIC | all wheel drive | |
| 11912 | Acura | ZDX | 2013 | premium unleaded (recommended) | 300.0 | 6.0 | AUTOMATIC | all wheel drive | |
| 11913 | Lincoln | Zephyr | 2006 | regular unleaded | 221.0 | 6.0 | AUTOMATIC | front wheel drive | |

11914 rows × 16 columns

```
dt.head(5)
```

| | Make | Model | Year | Engine Fuel Type | Engine HP | Engine Cylinders | Transmission Type | Driven_Wheels | Number of Doors | Ma |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | BMW | 1 Series M | 2011 | premium unleaded (required) | 335.0 | 6.0 | MANUAL | rear wheel drive | 2.0 | Tun |
| 1 | BMW | 1 Series | 2011 | premium unleaded (required) | 300.0 | 6.0 | MANUAL | rear wheel drive | 2.0 | Luxu |
| 2 | BMW | 1 Series | 2011 | premium unleaded (required) | 300.0 | 6.0 | MANUAL | rear wheel drive | 2.0 | |
| 3 | BMW | 1 Series | 2011 | premium unleaded (required) | 230.0 | 6.0 | MANUAL | rear wheel drive | 2.0 | Luxu |
| 4 | BMW | 1 Series | 2011 | premium unleaded (required) | 230.0 | 6.0 | MANUAL | rear wheel drive | 2.0 | |

In [7]: `dt.tail(5)`

| | Make | Model | Year | Engine Fuel Type | Engine HP | Engine Cylinders | Transmission Type | Driven_Wheels | Nu |
|---|---|---|---|---|---|---|---|---|---|
| 11909 | Acura | ZDX | 2012 | premium unleaded (required) | 300.0 | 6.0 | AUTOMATIC | all wheel drive | |
| 11910 | Acura | ZDX | 2012 | premium unleaded (required) | 300.0 | 6.0 | AUTOMATIC | all wheel drive | |
| 11911 | Acura | ZDX | 2012 | premium unleaded (required) | 300.0 | 6.0 | AUTOMATIC | all wheel drive | |
| 11912 | Acura | ZDX | 2013 | premium unleaded (recommended) | 300.0 | 6.0 | AUTOMATIC | all wheel drive | |
| 11913 | Lincoln | Zephyr | 2006 | regular unleaded | 221.0 | 6.0 | AUTOMATIC | front wheel drive | |

In [8]: `dt.dtypes`

```
Out[8]: Make                  object
        Model                 object
        Year                   int64
        Engine Fuel Type      object
        Engine HP            float64
        Engine Cylinders     float64
        Transmission Type     object
        Driven_Wheels         object
        Number of Doors      float64
        Market Category       object
        Vehicle Size          object
        Vehicle Style         object
        highway MPG            int64
        city mpg              int64
        Popularity            int64
        MSRP                  int64
        dtype: object
```

In [9]: 
```
dt = dt.drop(['Engine Fuel Type', 'Market Category', 'Vehicle Style', 'Popularity'
dt.head(5)
```

Out[9]:

| | Make | Model | Year | Engine HP | Engine Cylinders | Transmission Type | Driven_Wheels | highway MPG | city mpg | MSRP |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | BMW | 1 Series M | 2011 | 335.0 | 6.0 | MANUAL | rear wheel drive | 26 | 19 | 46135 |
| 1 | BMW | 1 Series | 2011 | 300.0 | 6.0 | MANUAL | rear wheel drive | 28 | 19 | 40650 |
| 2 | BMW | 1 Series | 2011 | 300.0 | 6.0 | MANUAL | rear wheel drive | 28 | 20 | 36350 |
| 3 | BMW | 1 Series | 2011 | 230.0 | 6.0 | MANUAL | rear wheel drive | 28 | 18 | 29450 |
| 4 | BMW | 1 Series | 2011 | 230.0 | 6.0 | MANUAL | rear wheel drive | 28 | 18 | 34500 |

In [10]: 
```
dt = dt.rename(columns={"Engine HP": "HP", "Engine Cylinders": "Cylinders", "Transi
dt.head(5)
```

Out[10]:

| | Make | Model | Year | HP | Cylinders | Transmission | Drive Mode | MPG-H | MPG-C | Price |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | BMW | 1 Series M | 2011 | 335.0 | 6.0 | MANUAL | rear wheel drive | 26 | 19 | 46135 |
| 1 | BMW | 1 Series | 2011 | 300.0 | 6.0 | MANUAL | rear wheel drive | 28 | 19 | 40650 |
| 2 | BMW | 1 Series | 2011 | 300.0 | 6.0 | MANUAL | rear wheel drive | 28 | 20 | 36350 |
| 3 | BMW | 1 Series | 2011 | 230.0 | 6.0 | MANUAL | rear wheel drive | 28 | 18 | 29450 |
| 4 | BMW | 1 Series | 2011 | 230.0 | 6.0 | MANUAL | rear wheel drive | 28 | 18 | 34500 |

In [13]: 
```
duplicate_rows_dt = dt[dt.duplicated()]
print("number of duplicate rows: ", duplicate_rows_dt.shape)
```

```
number of duplicate rows:  (989, 10)
```

In [15]: `dt.count()`

Out[15]:
```
Make            11914
Model           11914
Year            11914
HP              11845
Cylinders       11884
Transmission    11914
Drive Mode      11914
MPG-H           11914
MPG-C           11914
Price           11914
dtype: int64
```

In [16]:
```python
dt = dt.drop_duplicates()
dt.head(5)
```

Out[16]:

| | Make | Model | Year | HP | Cylinders | Transmission | Drive Mode | MPG-H | MPG-C | Price |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | BMW | 1 Series M | 2011 | 335.0 | 6.0 | MANUAL | rear wheel drive | 26 | 19 | 46135 |
| 1 | BMW | 1 Series | 2011 | 300.0 | 6.0 | MANUAL | rear wheel drive | 28 | 19 | 40650 |
| 2 | BMW | 1 Series | 2011 | 300.0 | 6.0 | MANUAL | rear wheel drive | 28 | 20 | 36350 |
| 3 | BMW | 1 Series | 2011 | 230.0 | 6.0 | MANUAL | rear wheel drive | 28 | 18 | 29450 |
| 4 | BMW | 1 Series | 2011 | 230.0 | 6.0 | MANUAL | rear wheel drive | 28 | 18 | 34500 |

In [17]: `dt.count()`

Out[17]:
```
Make            10925
Model           10925
Year            10925
HP              10856
Cylinders       10895
Transmission    10925
Drive Mode      10925
MPG-H           10925
MPG-C           10925
Price           10925
dtype: int64
```

In [18]: `print(dt.isnull().sum())`

```
Make             0
Model            0
Year             0
HP              69
Cylinders       30
Transmission     0
Drive Mode       0
MPG-H            0
MPG-C            0
Price            0
dtype: int64
```

```
In [20]: dt = dt.dropna()
         dt.count()
```
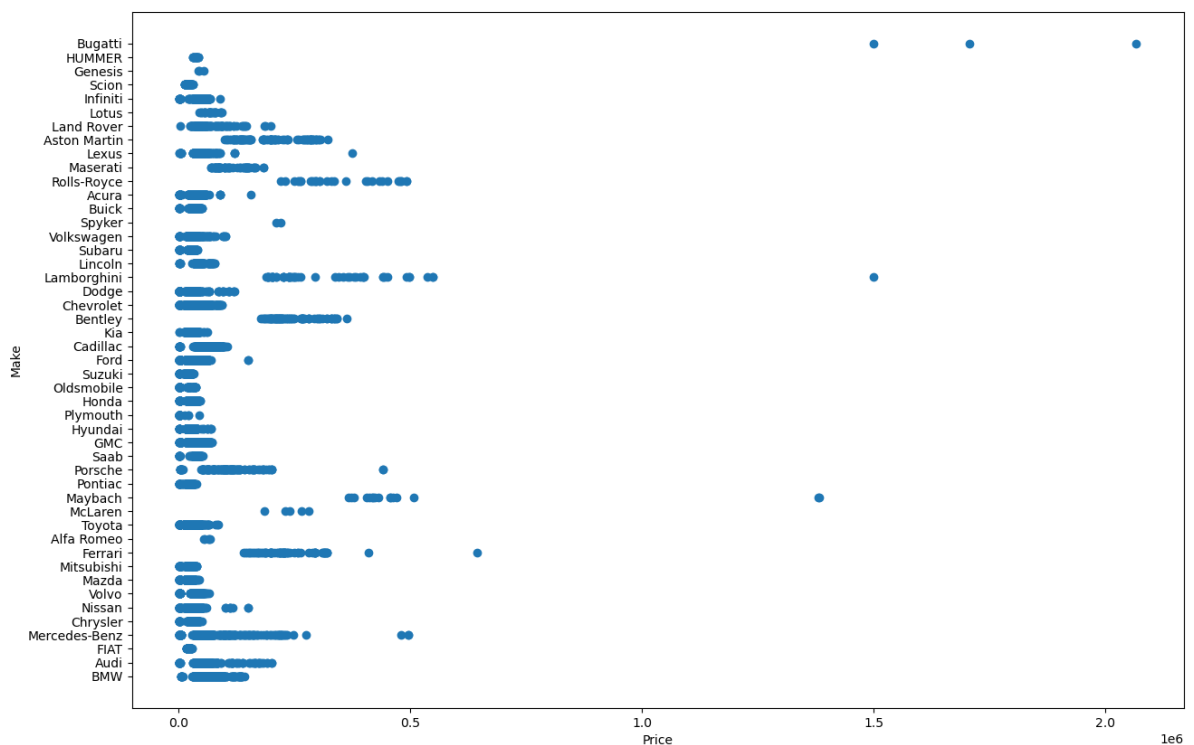
```
Out[20]: Make            10827
         Model           10827
         Year            10827
         HP              10827
         Cylinders       10827
         Transmission    10827
         Drive Mode      10827
         MPG-H           10827
         MPG-C           10827
         Price           10827
         dtype: int64
```

```
In [22]: print(dt.isnull().sum())
```

```
         Make            0
         Model           0
         Year            0
         HP              0
         Cylinders       0
         Transmission    0
         Drive Mode      0
         MPG-H           0
         MPG-C           0
         Price           0
         dtype: int64
```

```
In [40]: plt.subplots(figsize=(15,10))
         plt.scatter(dt['Price'],dt['Make'])
         plt.xlabel('Price')
         plt.ylabel('Make')
         plt.show()
```
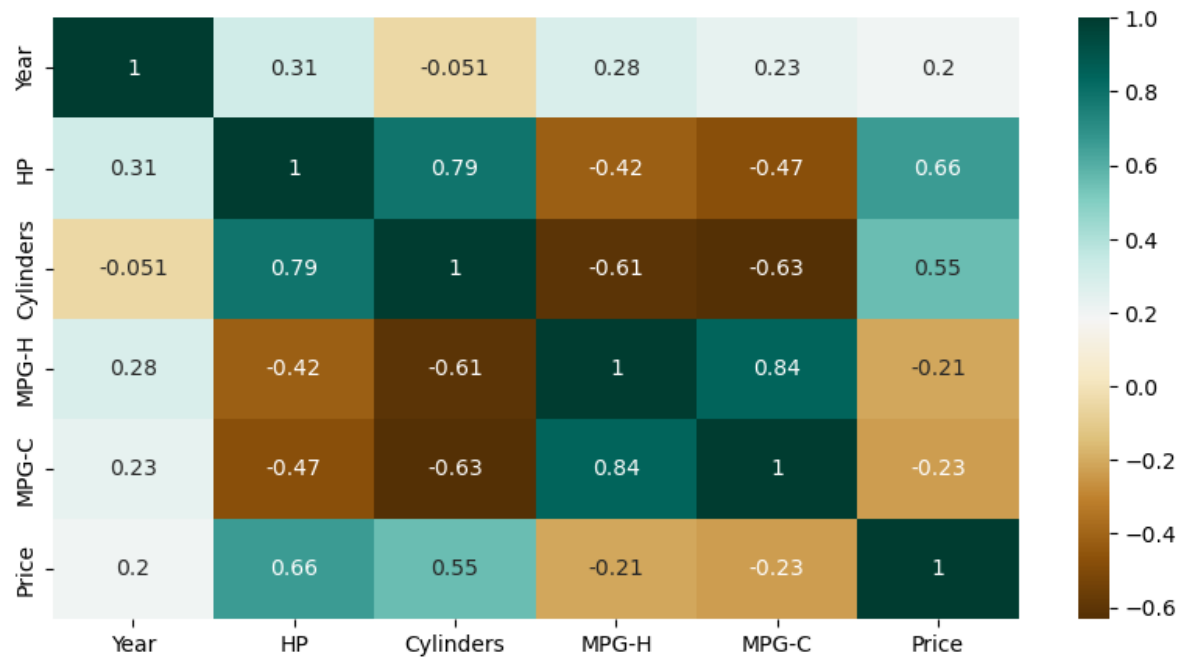


```
In [38]: plt.figure(figsize=(10,5))
         c= dt.corr()
         sns.heatmap(c,cmap="BrBG",annot=True)
         c
```

Out[38]:

|  | Year | HP | Cylinders | MPG-H | MPG-C | Price |
|---|---|---|---|---|---|---|
| **Year** | 1.000000 | 0.314971 | -0.050598 | 0.284237 | 0.234135 | 0.196789 |
| **HP** | 0.314971 | 1.000000 | 0.788007 | -0.420281 | -0.473551 | 0.659835 |
| **Cylinders** | -0.050598 | 0.788007 | 1.000000 | -0.611576 | -0.632407 | 0.554740 |
| **MPG-H** | 0.284237 | -0.420281 | -0.611576 | 1.000000 | 0.841229 | -0.209150 |
| **MPG-C** | 0.234135 | -0.473551 | -0.632407 | 0.841229 | 1.000000 | -0.234050 |
| **Price** | 0.196789 | 0.659835 | 0.554740 | -0.209150 | -0.234050 | 1.000000 |



In [ ]: