

Importing the required libraries for EDA

```
In [3]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
sns.set(color_codes=True)
```

Loading the data into the data frame.

```
In [6]: df = pd.read_csv("teechnook data.csv")
# To display the top 5 rows
df.head(5)
```

	Make	Model	Year	Engine Fuel Type	Engine HP	Engine Cylinders	Transmission Type	Driven_Wheels	Number of Doors	Market Category	Vehicle Size	Vehicle Style	highway MPG	city mpg	Popularity	MSRP
0	BMW	1 Series M	2011	premium unleaded (required)	335.0	6.0	MANUAL	rear wheel drive	2.0	Factory Tuner,Luxury,High-Performance	Compact	Coupe	26	19	3916	46135
1	BMW	1 Series	2011	premium unleaded (required)	300.0	6.0	MANUAL	rear wheel drive	2.0	Luxury,Performance	Compact	Convertible	28	19	3916	40650
2	BMW	1 Series	2011	premium unleaded (required)	300.0	6.0	MANUAL	rear wheel drive	2.0	Luxury,High-Performance	Compact	Coupe	28	20	3916	36350
3	BMW	1 Series	2011	premium unleaded (required)	230.0	6.0	MANUAL	rear wheel drive	2.0	Luxury,Performance	Compact	Coupe	28	18	3916	29450
4	BMW	1 Series	2011	premium unleaded (required)	230.0	6.0	MANUAL	rear wheel drive	2.0	Luxury	Compact	Convertible	28	18	3916	34500

```
In [7]: df.tail(5) # To display the botton 5 rows
```

	Make	Model	Year	Engine Fuel Type	Engine HP	Engine Cylinders	Transmission Type	Driven_Wheels	Number of Doors	Market Category	Vehicle Size	Vehicle Style	highway MPG	city mpg	Popularity	MSRP
11909	Acura	ZDX	2012	premium unleaded (required)	300.0	6.0	AUTOMATIC	all wheel drive	4.0	Crossover,Hatchback,Luxury	Midsize	4dr Hatchback	23	16	204	46120
11910	Acura	ZDX	2012	premium unleaded (required)	300.0	6.0	AUTOMATIC	all wheel drive	4.0	Crossover,Hatchback,Luxury	Midsize	4dr Hatchback	23	16	204	56670
11911	Acura	ZDX	2012	premium unleaded (required)	300.0	6.0	AUTOMATIC	all wheel drive	4.0	Crossover,Hatchback,Luxury	Midsize	4dr Hatchback	23	16	204	50620
11912	Acura	ZDX	2013	premium unleaded (recommended)	300.0	6.0	AUTOMATIC	all wheel drive	4.0	Crossover,Hatchback,Luxury	Midsize	4dr Hatchback	23	16	204	50920
11913	Lincoln	Zephyr	2006	regular unleaded	221.0	6.0	AUTOMATIC	front wheel drive	4.0	Luxury	Midsize	Sedan	26	17	61	28995

Checking the types of data

```
In [8]: df.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
```

Dropping irrelevant columns

```
In [9]: df = df.drop(['Engine Fuel Type', 'Market Category', 'Vehicle Style', 'Popularity', 'Number of Doors', 'Vehicle Size'], axis=1)
df.head(5)
```

	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

Renaming the columns

```
In [10]: df = df.rename(columns={"Engine HP": "HP", "Engine Cylinders": "Cylinders", "Transmission Type": "Transmission", "Driven_Wheels": "Drive Mode","highway MPG": "MPG-H", "city mpg": "MPG-C"})
df.head(5)
```

	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

Dropping the duplicate rows

```
In [11]: df.shape
```

```
Out[11]: (11914, 10)
```

```
In [12]: duplicate_rows_df = df[df.duplicated()]
print("number of duplicate rows: ", duplicate_rows_df.shape)
```

number of duplicate rows: (989, 10)

Now let us remove the duplicate data because it's ok to remove them.

```
In [13]: df.count() # Used to count the number of rows
```

```
Out[13]: 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
```

So seen above there are 11914 rows and we are removing 989 rows of duplicate data

```
In [14]: df = df.drop_duplicates()
df.head(5)
```

	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 [15]: df.count()
```

```
Out[15]: 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
```

Dropping the missing or null values.

```
In [16]: print(df.isnull().sum())
```

```
Make                0
Model               0
Year               0
HP                 0
Cylinders          38
Transmission       0
Drive Mode         0
MPG-H              0
MPG-C              0
Price              0
dtype: int64
```

This is the reason in the above step while counting both Cylinders and Horsepower (HP) had 10856 and 10895 over 10925 rows.

```
In [17]: df = df.dropna() # Dropping the missing values.
df.count()
```

```
Out[17]: 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
```

Now we have removed all the rows which contain the Null or N/A values (Cylinders and Horsepower (HP)).

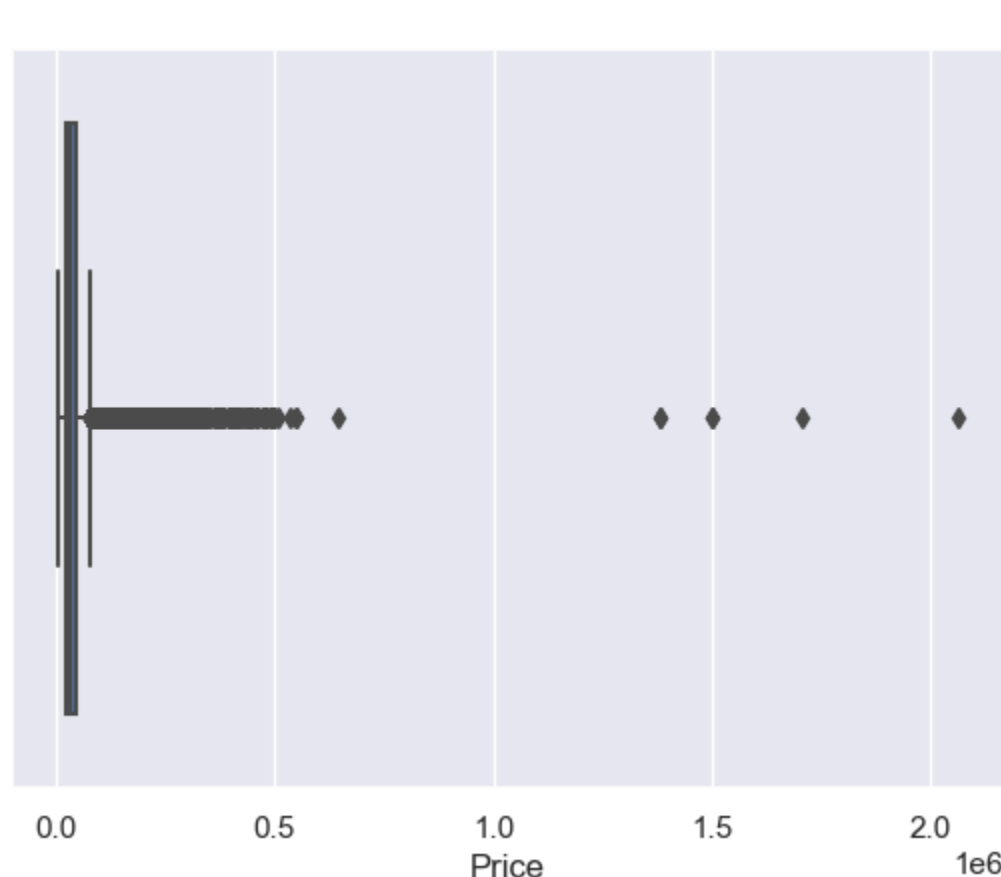
```
In [19]: print(df.isnull().sum()) # After dropping the values
```

```
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
```

Detecting Outliers

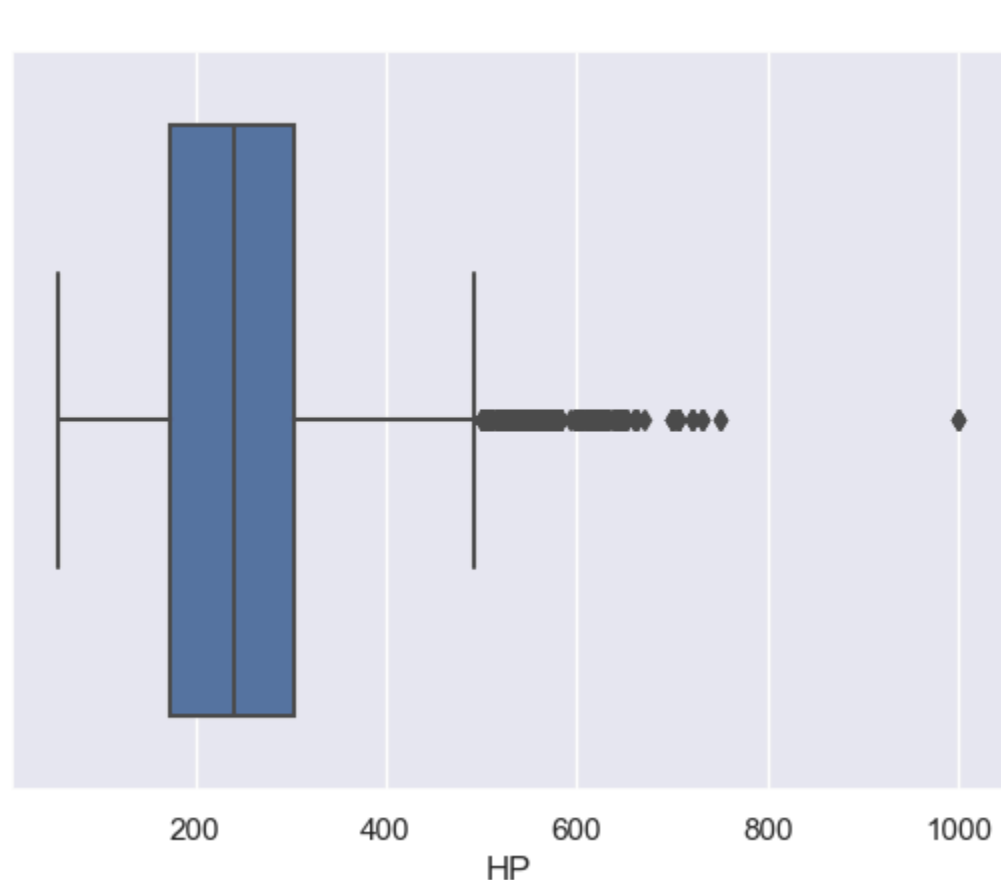
```
In [21]: sns.boxplot(x=df['Price'])
```

```
Out[21]: <Axes: xlabel='Price'>
```



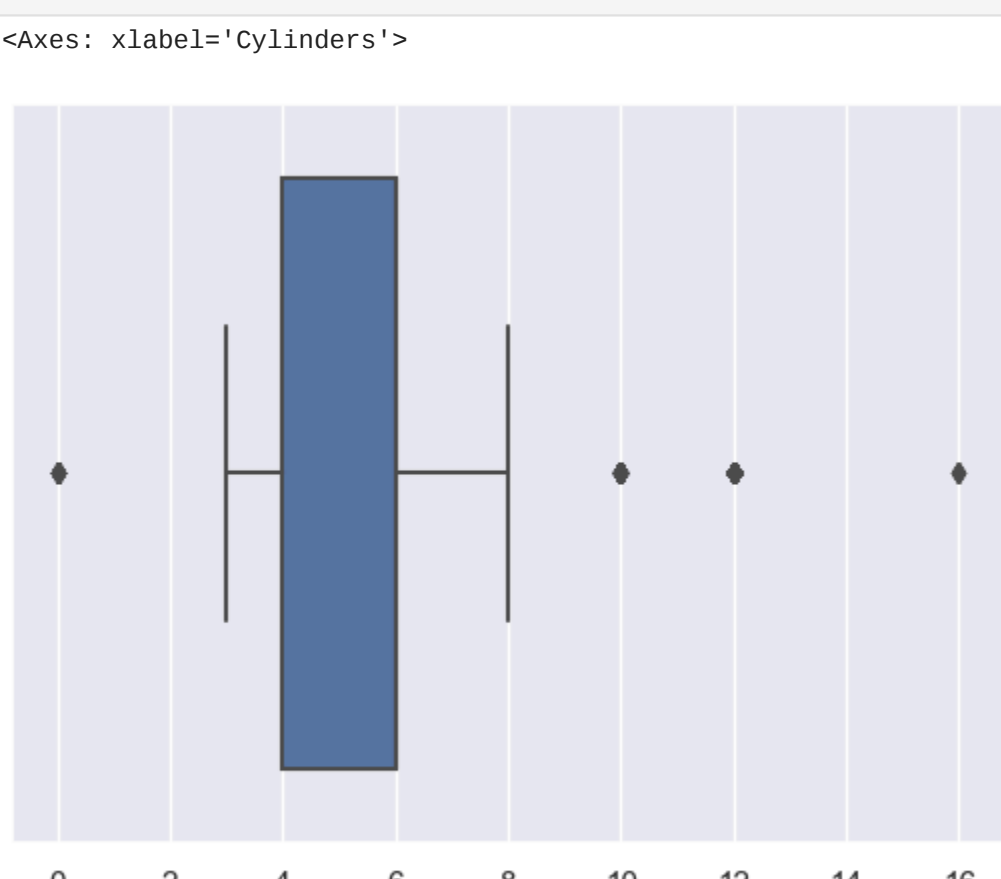
```
In [22]: sns.boxplot(x=df['HP'])
```

```
Out[22]: <Axes: xlabel='HP'>
```



```
In [25]: sns.boxplot(x=df['Cylinders'])
```

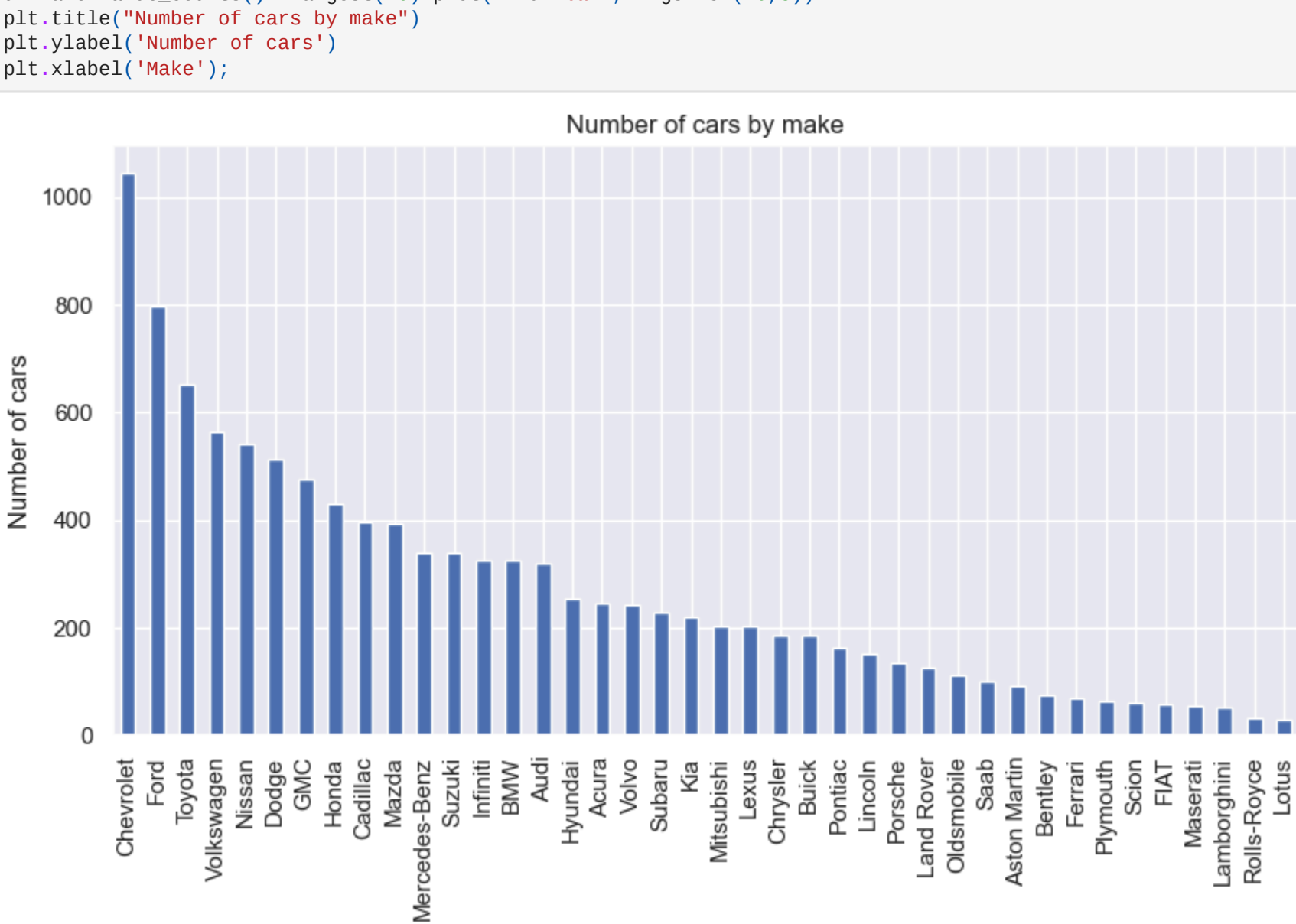
```
Out[25]: <Axes: xlabel='Cylinders'>
```



```
In [ ]:
```

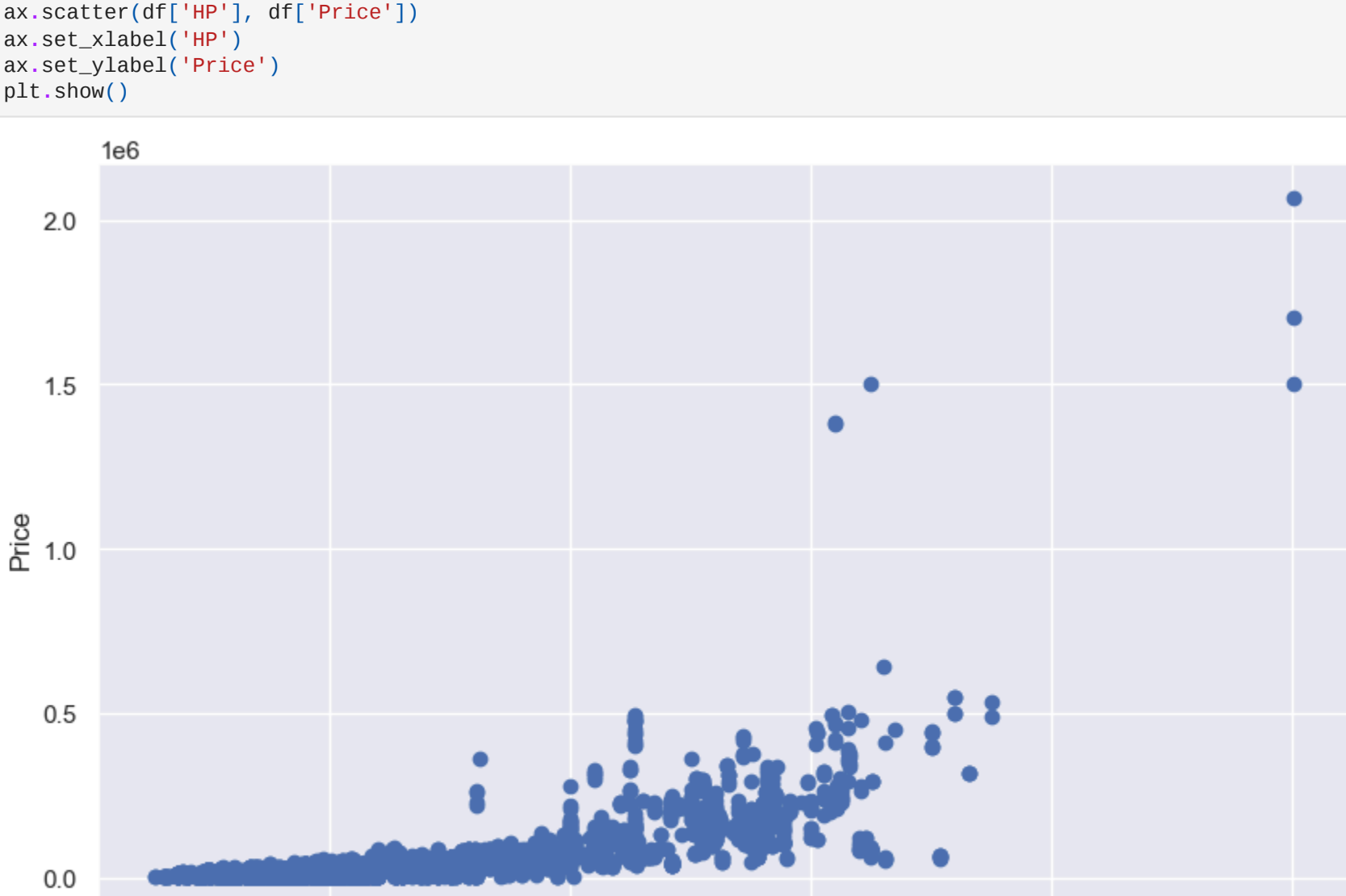
Histogram

```
In [29]: df.Make.value_counts().nlargest(40).plot(kind='bar', figsize=(10,5))
plt.title("Number of cars by make")
plt.ylabel("Number of cars")
plt.xlabel("Make");
```



Scatterplot

```
In [33]: fig, ax = plt.subplots(figsize=(10,6))
ax.scatter(df['HP'], df['Price'])
ax.set_xlabel('HP')
ax.set_ylabel('Price')
plt.show()
```



```
In [ ]:
```