

DOMAIN:

Automobile

CONTEXT:

The data concerns city-cycle fuel consumption in miles per gallon, to be predicted in terms of 3 multivalued discrete and 5 continuous attributes

DATA DESCRIPTION:

The data concerns city-cycle fuel consumption in miles per gallon Attribute Information:

1. mpg: continuous
2. cylinders(cyl): multi-valued discrete
3. displacement(displ): continuous
4. horsepower(hp): continuous
5. weight(wt): continuous
6. acceleration(acc): continuous
7. model year(yr): multi-valued discrete
8. origin: multi-valued discrete
9. car name: string (unique for each instance)

PROJECT OBJECTIVE:

Goal is to cluster the data and treat them as individual datasets to train Regression models to predict 'mpg'

In [1]:

```
1 #importing required packages
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.linear_model import LinearRegression
7 from scipy import stats
8 from scipy.stats import zscore
9 from scipy.cluster.hierarchy import dendrogram, linkage
10 from scipy.cluster.hierarchy import fcluster
11 from sklearn.cluster import KMeans
12 from sklearn.metrics import silhouette_samples, silhouette_score
13 %matplotlib inline
14 sns.set(color_codes=True)
15 import warnings
16 warnings.filterwarnings('ignore')
17 from sklearn.model_selection import train_test_split
18
```

In [2]:

```
1 #Load data files
2 ca=pd.read_json('car.json')
3 ca1=pd.read_csv('car name.csv')
4 car=pd.concat([ca,ca1],axis=1)
5
6 row,column=car.shape
7 print('the dataset contain ',row,'row and',column,'columns')
8 car.head()
```

the dataset contain 398 row and 9 columns

Out[2]:

	mpg	cyl	disp	hp	wt	acc	yr	origin	car_name
0	18.0	8	307.0	130	3504	12.0	70	1	chevrolet chevelle malibu
1	15.0	8	350.0	165	3693	11.5	70	1	buick skylark 320
2	18.0	8	318.0	150	3436	11.0	70	1	plymouth satellite
3	16.0	8	304.0	150	3433	12.0	70	1	amc rebel sst
4	17.0	8	302.0	140	3449	10.5	70	1	ford torino

In [3]:

```
1 #save this data
2 car.to_csv('mpg.csv',index=False)
3 car.to_excel('mpg.xlsx',index=False)
4 car.to_json('mpg.json',orient='split',compression='infer',index='true')
```

In [4]:

```
1 car=car.drop('car_name',axis=1)
2 car['origin']=car['origin'].replace({1:'america',2:'europe',3:'asia'})
3 car.head()
```

Out[4]:

	mpg	cyl	disp	hp	wt	acc	yr	origin
0	18.0	8	307.0	130	3504	12.0	70	america
1	15.0	8	350.0	165	3693	11.5	70	america
2	18.0	8	318.0	150	3436	11.0	70	america
3	16.0	8	304.0	150	3433	12.0	70	america
4	17.0	8	302.0	140	3449	10.5	70	america

In [5]:

```
1 car.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 398 entries, 0 to 397
Data columns (total 8 columns):
 #   Column  Non-Null Count  Dtype  
---  -
 0   mpg     398 non-null    float64
 1   cyl     398 non-null    int64   
 2   disp    398 non-null    float64
 3   hp      398 non-null    object  
 4   wt      398 non-null    int64   
 5   acc     398 non-null    float64
 6   yr      398 non-null    int64   
 7   origin  398 non-null    object  
dtypes: float64(3), int64(3), object(2)
memory usage: 25.0+ KB
```

In [6]:

```
1 car.describe()
```

Out[6]:

	mpg	cyl	disp	wt	acc	yr
count	398.000000	398.000000	398.000000	398.000000	398.000000	398.000000
mean	23.514573	5.454774	193.425879	2970.424623	15.568090	76.010050
std	7.815984	1.701004	104.269838	846.841774	2.757689	3.697627
min	9.000000	3.000000	68.000000	1613.000000	8.000000	70.000000
25%	17.500000	4.000000	104.250000	2223.750000	13.825000	73.000000
50%	23.000000	4.000000	148.500000	2803.500000	15.500000	76.000000
75%	29.000000	8.000000	262.000000	3608.000000	17.175000	79.000000
max	46.600000	8.000000	455.000000	5140.000000	24.800000	82.000000

In [7]:

```
1 #horse power
2 hpISDigit=pd.DataFrame(car.hp.str.isdigit())
```

In [8]:

```
1 car[hpISDigit['hp']==False]
```

Out[8]:

	mpg	cyl	disp	hp	wt	acc	yr	origin
32	25.0	4	98.0	?	2046	19.0	71	america
126	21.0	6	200.0	?	2875	17.0	74	america
330	40.9	4	85.0	?	1835	17.3	80	europe
336	23.6	4	140.0	?	2905	14.3	80	america
354	34.5	4	100.0	?	2320	15.8	81	europe
374	23.0	4	151.0	?	3035	20.5	82	america

In [9]:

```
1 car=car.replace('?',np.nan)
2 car[hpISDigit['hp']==False]
```

Out[9]:

	mpg	cyl	disp	hp	wt	acc	yr	origin
32	25.0	4	98.0	NaN	2046	19.0	71	america
126	21.0	6	200.0	NaN	2875	17.0	74	america
330	40.9	4	85.0	NaN	1835	17.3	80	europe
336	23.6	4	140.0	NaN	2905	14.3	80	america
354	34.5	4	100.0	NaN	2320	15.8	81	europe
374	23.0	4	151.0	NaN	3035	20.5	82	america

In [10]:

```
1 car.median()
```

Out[10]:

```
mpg      23.0
cyl       4.0
disp     148.5
hp        93.5
wt       2803.5
acc       15.5
yr        76.0
dtype: float64
```

In [11]:

```
1 car['hp'].fillna((car['hp'].median()),inplace=True)
2 car.isnull()
```

Out[11]:

	mpg	cyl	disp	hp	wt	acc	yr	origin
0	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False
...
393	False	False	False	False	False	False	False	False
394	False	False	False	False	False	False	False	False
395	False	False	False	False	False	False	False	False
396	False	False	False	False	False	False	False	False
397	False	False	False	False	False	False	False	False

398 rows × 8 columns

In [12]:

```
1 car.isnull().sum()
```

Out[12]:

```
mpg      0
cyl      0
disp     0
hp       0
wt       0
acc      0
yr       0
origin   0
dtype: int64
```

In [13]:

```
1 car['mpg_level']=car['mpg'].apply(lambda X:'low' if X<17 else 'high' if X>29 else 'me
2 car.head()
```

Out[13]:

	mpg	cyl	disp	hp	wt	acc	yr	origin	mpg_level
0	18.0	8	307.0	130.0	3504	12.0	70	america	median
1	15.0	8	350.0	165.0	3693	11.5	70	america	low
2	18.0	8	318.0	150.0	3436	11.0	70	america	median
3	16.0	8	304.0	150.0	3433	12.0	70	america	low
4	17.0	8	302.0	140.0	3449	10.5	70	america	median

In [14]:

```
1 car_cat=car.iloc[:,[1,6,7,8]]
2 car_cat.head()
```

Out[14]:

	cyl	yr	origin	mpg_level
0	8	70	america	median
1	8	70	america	low
2	8	70	america	median
3	8	70	america	low
4	8	70	america	median

In [15]:

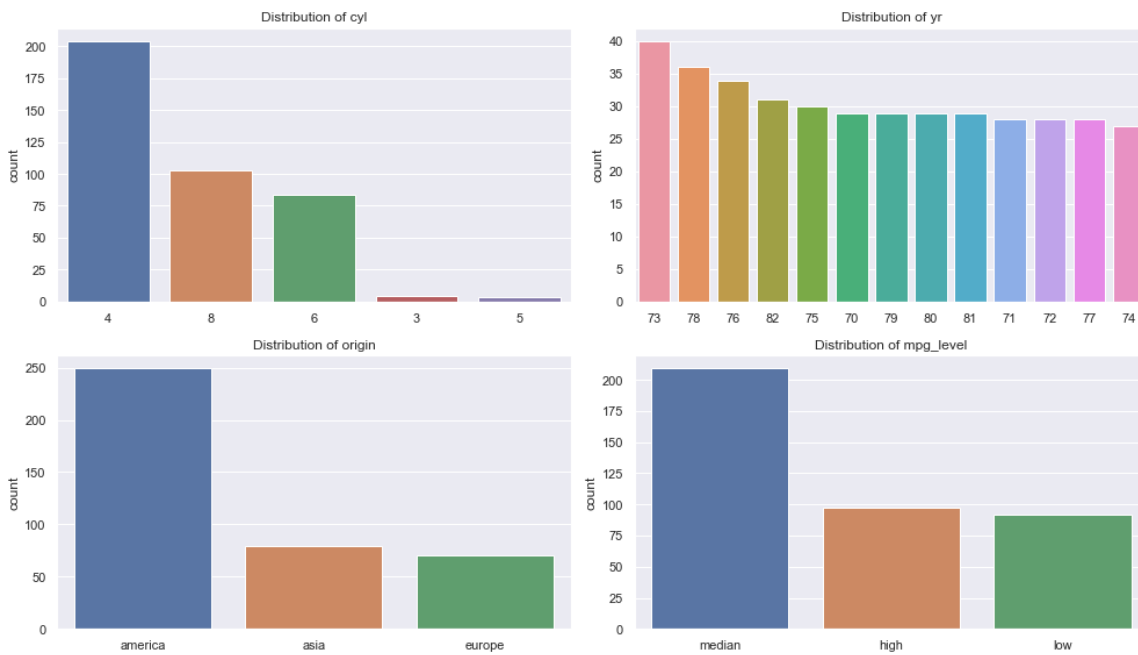
```
1 car_num=car.drop(['cyl','yr','origin','mpg_level'],axis=1)
2 car_num.head()
```

Out[15]:

	mpg	disp	hp	wt	acc
0	18.0	307.0	130.0	3504	12.0
1	15.0	350.0	165.0	3693	11.5
2	18.0	318.0	150.0	3436	11.0
3	16.0	304.0	150.0	3433	12.0
4	17.0	302.0	140.0	3449	10.5

In [16]:

```
1 #plot cat variable
2 fig=plt.figure(1,(14,8))
3
4 for i,car in enumerate(car_cat.columns):
5     ax=plt.subplot(2,2,i+1)
6     sns.countplot(car_cat[car],order=car_cat[car].value_counts().index)
7     ax.set_xlabel(None)
8     ax.set_title(f'Distribution of {car}')
9     plt.tight_layout()
10
11
12 plt.show()
13
```



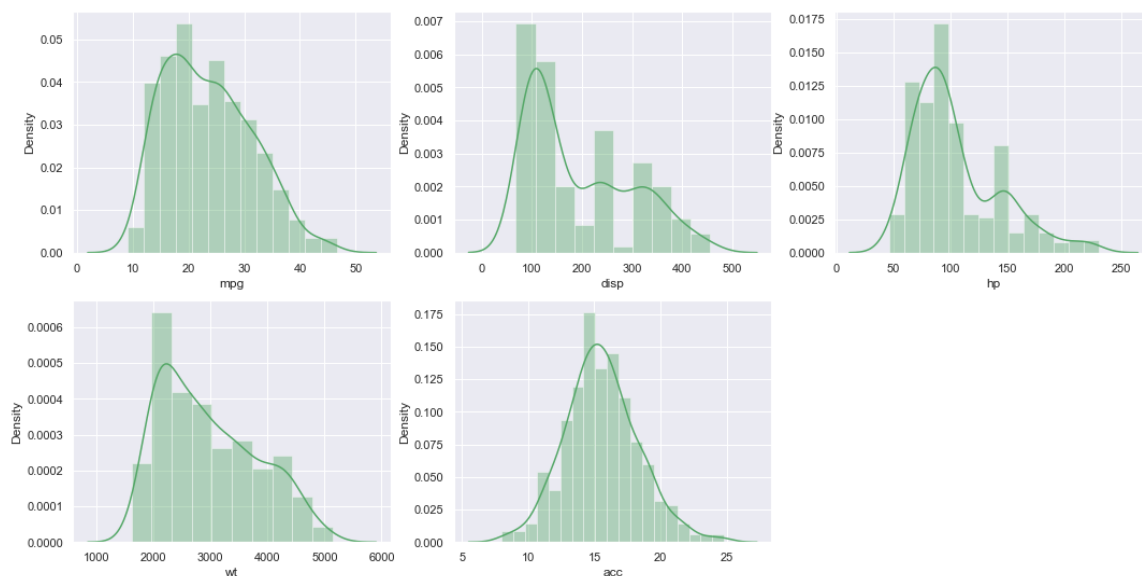
In [17]:

```
1 car_num.hist(bins=20,figsize=(10,8),color='pink')
2 plt.show()
```



In [18]:

```
1 plt.figure(figsize=(18,14))
2 col=1
3 for i in car_num.columns:
4     plt.subplot(3,3,col)
5     sns.distplot(car_num[i],color='g')
6     col+=1
```



In [19]:

```
1 car=pd.concat([car_cat,car_num],axis=1)
```

In [20]:

```
1 car.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 398 entries, 0 to 397
Data columns (total 9 columns):
#   Column      Non-Null Count  Dtype
---  -
0   cyl         398 non-null    int64
1   yr          398 non-null    int64
2   origin      398 non-null    object
3   mpg_level   398 non-null    object
4   mpg         398 non-null    float64
5   disp        398 non-null    float64
6   hp          398 non-null    float64
7   wt          398 non-null    int64
8   acc         398 non-null    float64
dtypes: float64(4), int64(3), object(2)
memory usage: 28.1+ KB
```

In [21]:

```
1 car=pd.get_dummies(car,columns=['origin'])
```

In [22]:

```
1 car=pd.get_dummies(car,columns=['mpg_level'])
```

In [23]:

```
1 car.head()
```

Out[23]:

	cyl	yr	mpg	disp	hp	wt	acc	origin_america	origin_asia	origin_europe	mpg_level
0	8	70	18.0	307.0	130.0	3504	12.0	1	0	0	
1	8	70	15.0	350.0	165.0	3693	11.5	1	0	0	
2	8	70	18.0	318.0	150.0	3436	11.0	1	0	0	
3	8	70	16.0	304.0	150.0	3433	12.0	1	0	0	
4	8	70	17.0	302.0	140.0	3449	10.5	1	0	0	

In [24]:

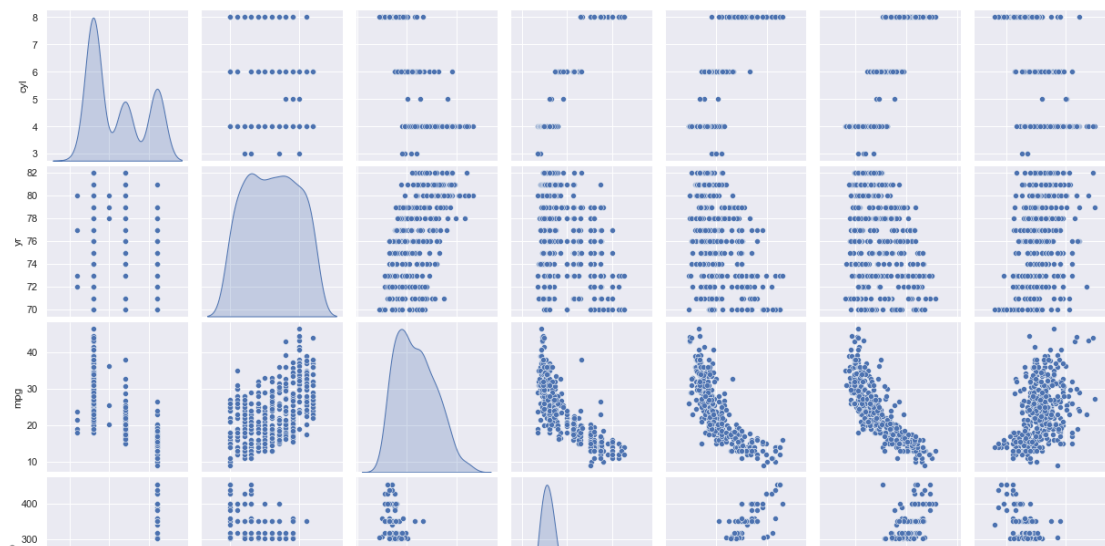
```
1 Hcar=car.copy()
2 Kcar=car.copy()
```

In [25]:

```
1 #pair plot for the numeric attributes
2 car_attr=car.iloc[:,0:7]
3 sns.pairplot(car_attr,diag_kind='kde')
4 sns.pairplot(car_attr,diag_kind='hist')
```

Out[25]:

<seaborn.axisgrid.PairGrid at 0x235138b9b80>

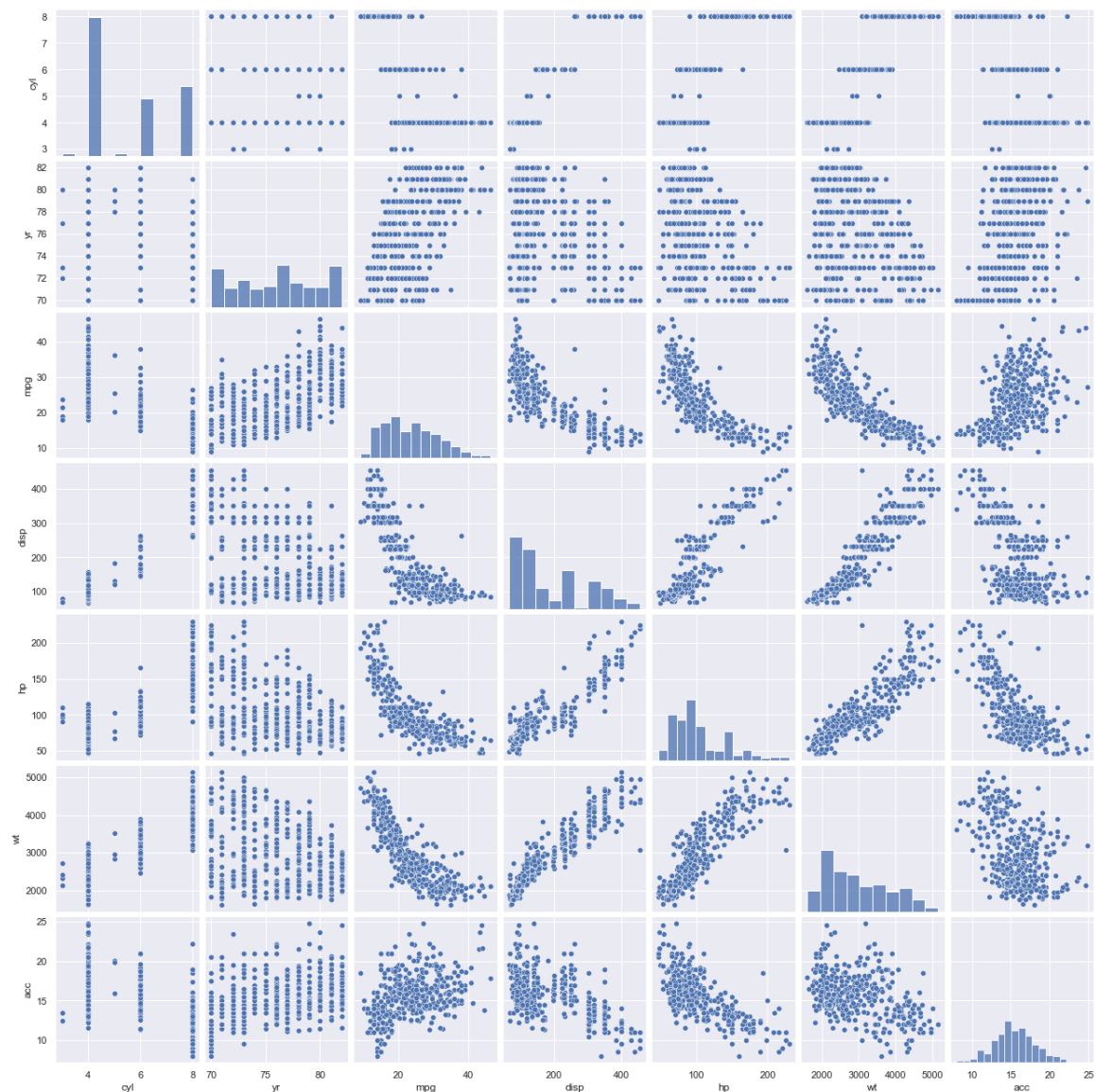


In [26]:

```
1 sns.pairplot(car_attr,diag_kind='auto')
```

Out[26]:

<seaborn.axisgrid.PairGrid at 0x2351632a790>

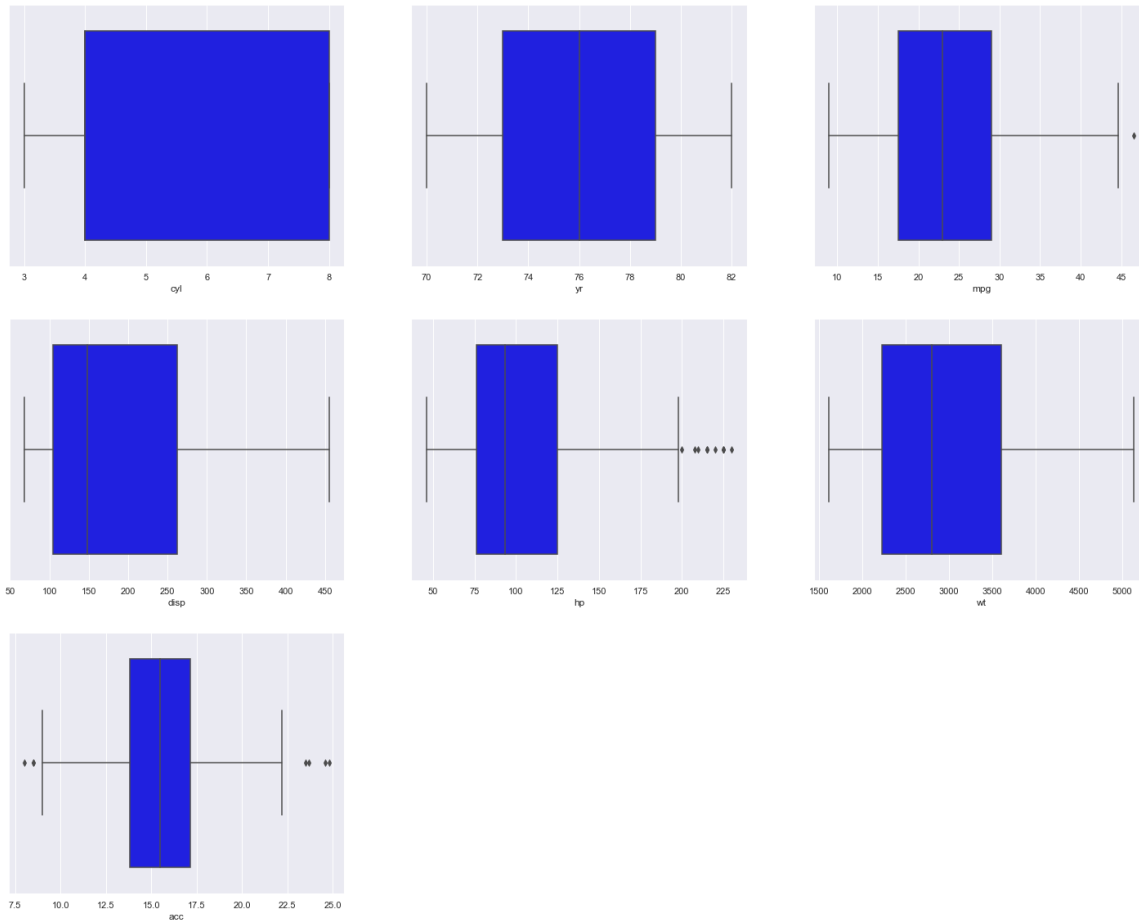


In [27]:

```
1 #dropping the created dummy variable
2 car2=car.drop(['origin_america','origin_asia','origin_europe','mpg_level_high','mpg_
```

In [28]:

```
1 plt.figure(figsize=(25,20))
2 col = 1
3 for i in car2.columns:
4     plt.subplot(3,3,col)
5     sns.boxplot(car2[i],color='blue')
6     col +=1
```

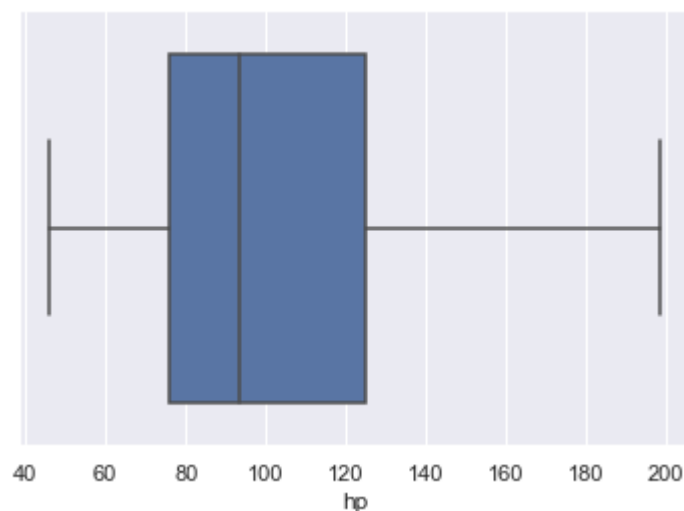


In [29]:

```
1 #replacing outliers with IQR (Q1 and Q3 +-1.5*IQR)
2 IQR1 = stats.iqr(car2['hp'], interpolation = 'midpoint')
3 IQR2 = stats.iqr(car2['acc'], interpolation = 'midpoint')
```

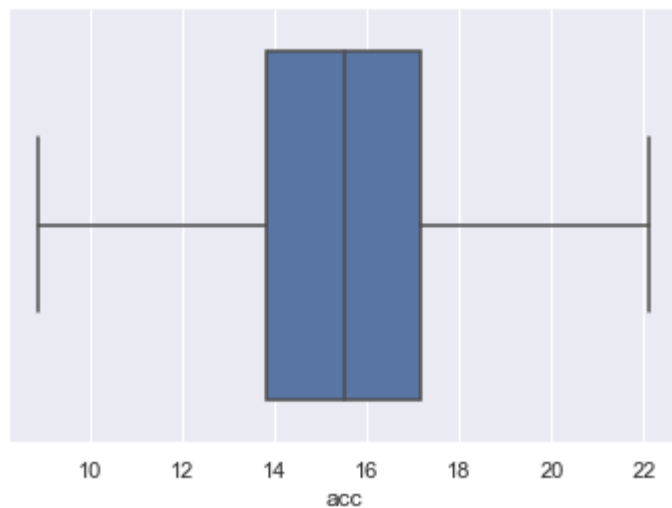
In [30]:

```
1 Q3 = car2['hp'].quantile(0.75)
2 car2['hp'] = np.where(car2["hp"] > (Q3+1.5*IQR1), 198.5, car2['hp'])
3 sns.boxplot(car2['hp']);
```



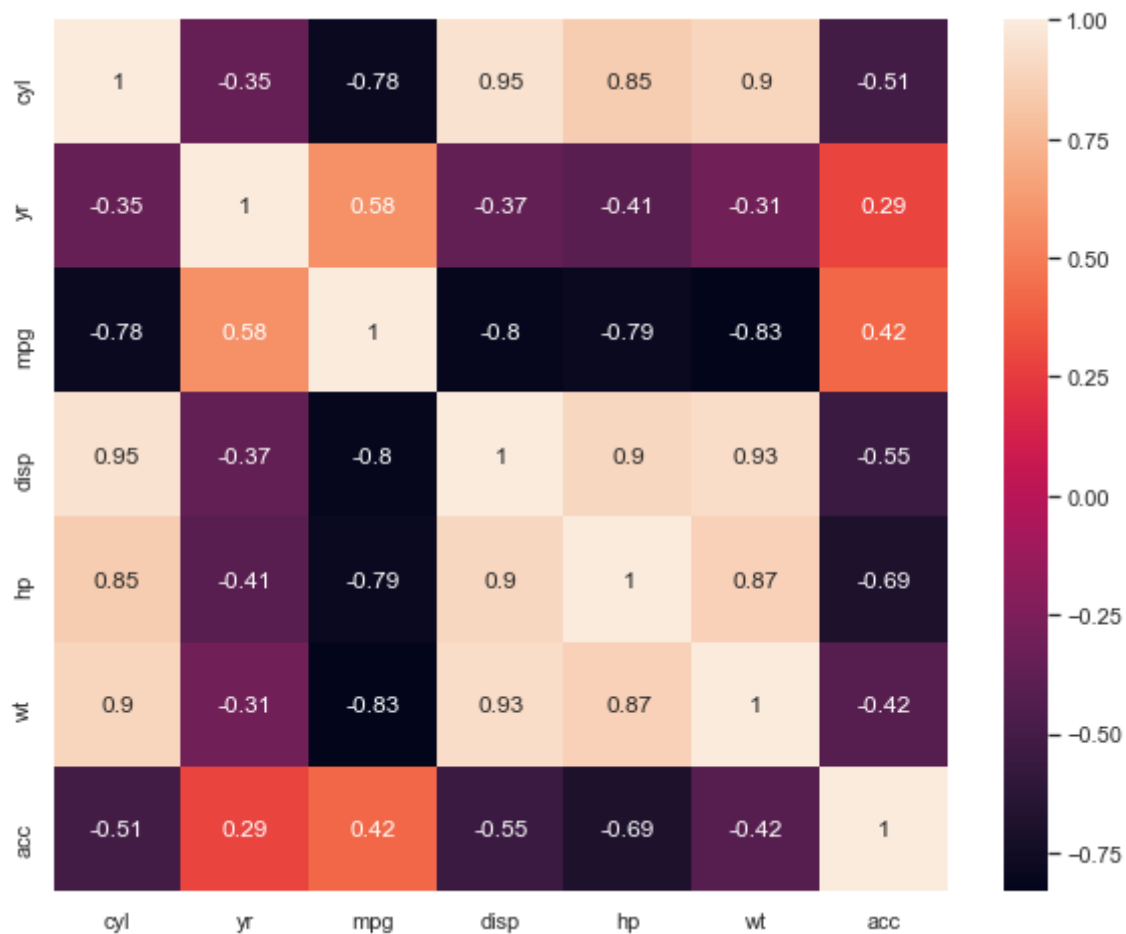
In [31]:

```
1 Q1 = car2['acc'].quantile(0.25)
2 Q31=car2['acc'].quantile(0.75)
3 car2['acc'] = np.where(car2["acc"] > (Q31+1.5*IQR2), 22.10 , car2['acc'])
4 car2['acc'] = np.where(car2["acc"] < (Q1-1.5*IQR2), (Q1-
5 1.5*IQR2), car2['acc'])
6 sns.boxplot(car2['acc']);
```



In [32]:

```
1 #checking for correlation
2 plt.figure(figsize=(10,8))
3 corr=car2.corr()
4 sns.heatmap(corr,annot=True);
```



Heirarchical Clustering

In [33]:

```
1 #separating numeric variables
2 cc = car.iloc[:,0:7]
3 cc.head()
```

Out[33]:

	cyl	yr	mpg	disp	hp	wt	acc
0	8	70	18.0	307.0	130.0	3504	12.0
1	8	70	15.0	350.0	165.0	3693	11.5
2	8	70	18.0	318.0	150.0	3436	11.0
3	8	70	16.0	304.0	150.0	3433	12.0
4	8	70	17.0	302.0	140.0	3449	10.5

In [34]:

```
1 cc_z = cc.apply(zscore)
2 cc_z.head()
```

Out[34]:

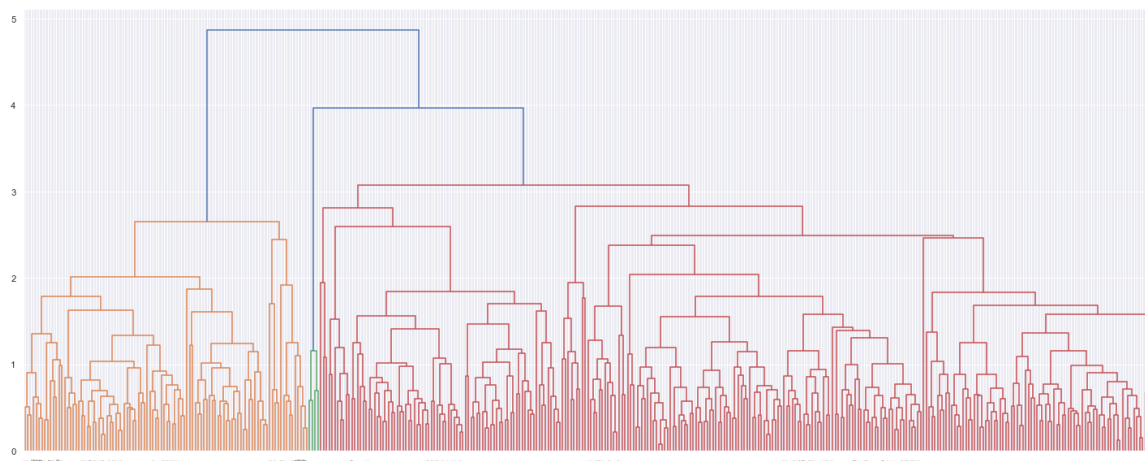
	cyl	yr	mpg	disp	hp	wt	acc
0	1.498191	-1.627426	-0.706439	1.090604	0.673118	0.630870	-1.295498
1	1.498191	-1.627426	-1.090751	1.503514	1.589958	0.854333	-1.477038
2	1.498191	-1.627426	-0.706439	1.196232	1.197027	0.550470	-1.658577
3	1.498191	-1.627426	-0.962647	1.061796	1.197027	0.546923	-1.295498
4	1.498191	-1.627426	-0.834543	1.042591	0.935072	0.565841	-1.840117

In [35]:

```
1 link_method = linkage(cc_z.iloc[:,0:7], method = 'average')
```

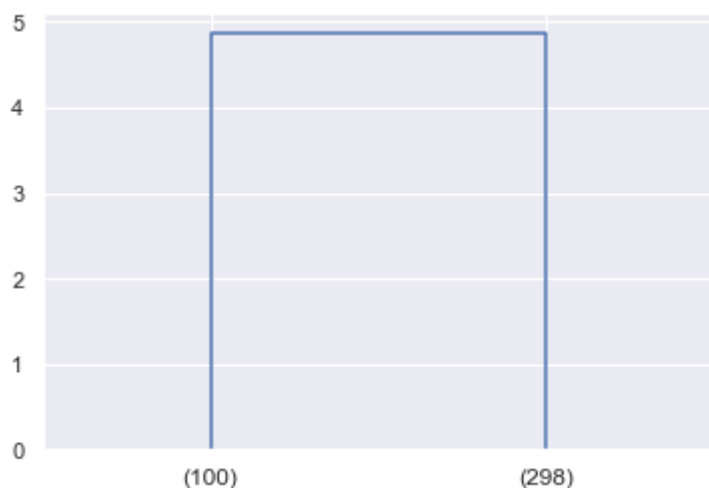
In [36]:

```
1 #plotting the H-cluster
2 plt.figure(figsize=(25, 10))
3 dendrogram(link_method)
4 plt.show()
```



In [37]:

```
1 # dendrogram function to arrive at dendrogram
2 dendrogram(
3     link_method,
4     truncate_mode='lastp',
5     p=2,
6 )
7 plt.show()
```



In [38]:

```
1 #vieweing the clusters formed
2 clusters = fcluster(link_method, 2, criterion='maxclust')
3 clusters
```

Out[38]:

[illegible]

In [39]:

```
1 cc_z['clusters_H'] = clusters
2 cc_z.head()
```

Out[39]:

	cyl	yr	mpg	disp	hp	wt	acc	clusters_H
0	1.498191	-1.627426	-0.706439	1.090604	0.673118	0.630870	-1.295498	1
1	1.498191	-1.627426	-1.090751	1.503514	1.589958	0.854333	-1.477038	1
2	1.498191	-1.627426	-0.706439	1.196232	1.197027	0.550470	-1.658577	1
3	1.498191	-1.627426	-0.962647	1.061796	1.197027	0.546923	-1.295498	1
4	1.498191	-1.627426	-0.834543	1.042591	0.935072	0.565841	-1.840117	1

In [40]:

```
1 #attaching the clusters formed to the scales data
2 cc_z.clusters_H.value_counts().sort_index()
```

Out[40]:

```
1    100
2    298
Name: clusters_H, dtype: int64
```

In [41]:

```
1 cc['clusters_H']=clusters
2 Hcar['clusters_H']=clusters
3 cc.head()
```

Out[41]:

	cyl	yr	mpg	disp	hp	wt	acc	clusters_H
0	8	70	18.0	307.0	130.0	3504	12.0	1
1	8	70	15.0	350.0	165.0	3693	11.5	1
2	8	70	18.0	318.0	150.0	3436	11.0	1
3	8	70	16.0	304.0	150.0	3433	12.0	1
4	8	70	17.0	302.0	140.0	3449	10.5	1

In [42]:

```
1 #create a new data set named Hclus
2 Hclus=cc
3 Hclus.head()
```

Out[42]:

	cyl	yr	mpg	disp	hp	wt	acc	clusters_H
0	8	70	18.0	307.0	130.0	3504	12.0	1
1	8	70	15.0	350.0	165.0	3693	11.5	1
2	8	70	18.0	318.0	150.0	3436	11.0	1
3	8	70	16.0	304.0	150.0	3433	12.0	1
4	8	70	17.0	302.0	140.0	3449	10.5	1

In [43]:

```
1 aggdata=cc.iloc[:,0:8].groupby('clusters_H').mean()
2 aggdata['Freq']=cc.clusters_H.value_counts().sort_index()
3 aggdata
```

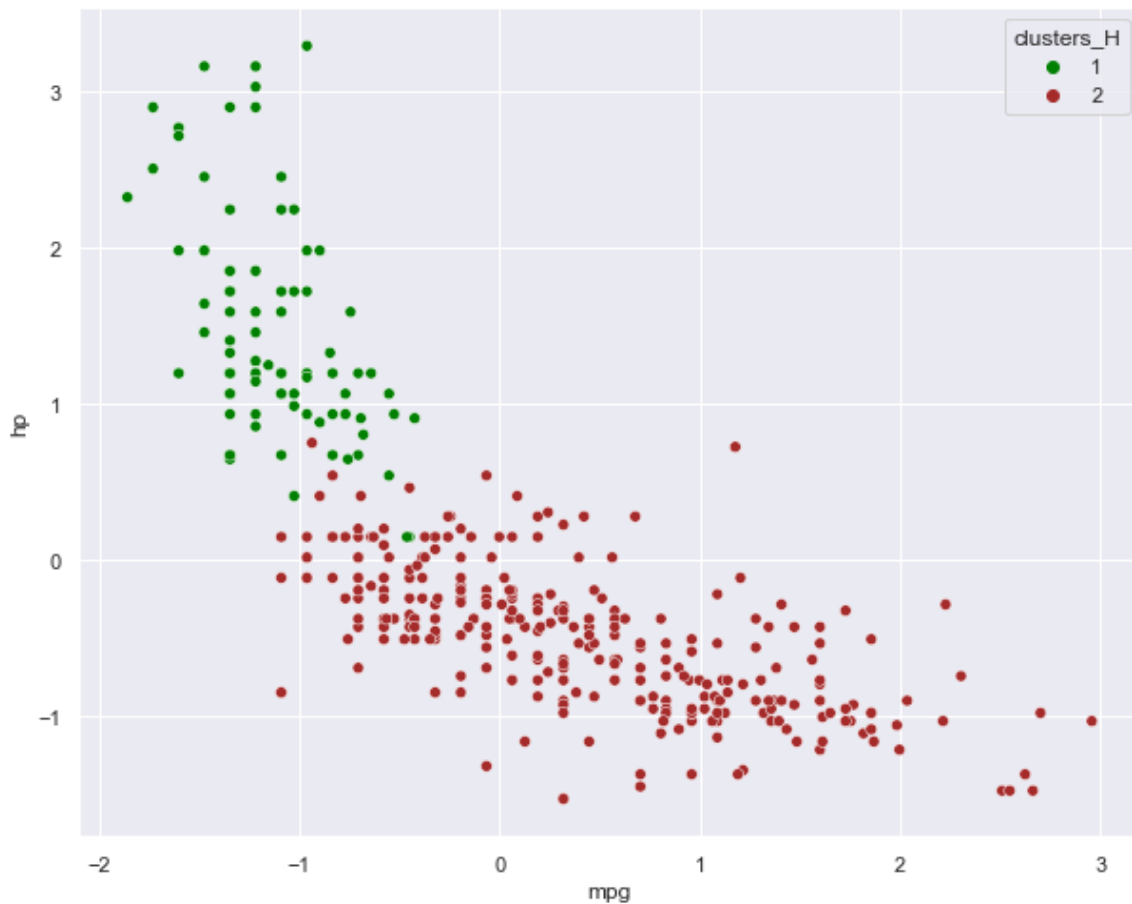
Out[43]:

	cyl	yr	mpg	disp	hp	wt	acc	l
clusters_H								
1	7.980000	73.740000	14.684000	345.470000	160.400000	4121.560000	12.702000	
2	4.607383	76.771812	26.477852	142.404362	85.479866	2584.137584	16.529866	



In [44]:

```
1 #plotting the clusters formed
2 plt.figure(figsize=(10, 8))
3 sns.scatterplot(x="mpg", y="hp", hue="clusters_H", data=cc_z, palette=['green', 'brown'])
```



K-Means Clustering

In [45]:

```
1 #seperating the numeric values
2 cc = car.iloc[:,0:7]
3 cc_z1 = cc.apply(zscore)
4 cc_z1.head()
```

Out[45]:

	cyl	yr	mpg	disp	hp	wt	acc
0	1.498191	-1.627426	-0.706439	1.090604	0.673118	0.630870	-1.295498
1	1.498191	-1.627426	-1.090751	1.503514	1.589958	0.854333	-1.477038
2	1.498191	-1.627426	-0.706439	1.196232	1.197027	0.550470	-1.658577
3	1.498191	-1.627426	-0.962647	1.061796	1.197027	0.546923	-1.295498
4	1.498191	-1.627426	-0.834543	1.042591	0.935072	0.565841	-1.840117

In [46]:

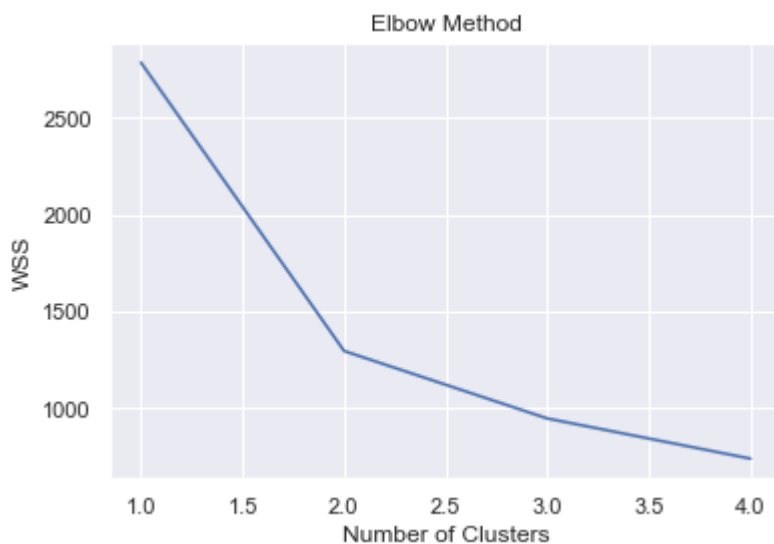
```
1 #calculatint the within sum of squares
2 wss =[]
3 for i in range(1,5):
4     KM = KMeans(n_clusters=i)
5     KM.fit(cc_z1)
6     wss.append(KM.inertia_)
7 wss
```

Out[46]:

```
[2785.9999999999995, 1294.8418950727323, 946.019790855379, 738.39232815273
18]
```

In [47]:

```
1 #plotting the WSS against the number of cluster to come up with optimal number of clu
2 plt.plot(range(1,5), wss);
3 plt.title('Elbow Method');
4 plt.xlabel("Number of Clusters")
5 plt.ylabel("WSS");
```



In [48]:

```
1 k_means = KMeans(n_clusters = 2)
2 k_means.fit(cc_z1)
3 labels = k_means.labels_
```

In [49]:

```
1 silhouette_score(cc_z1, labels)
```

Out[49]:

```
0.48235946103916116
```

In [50]:

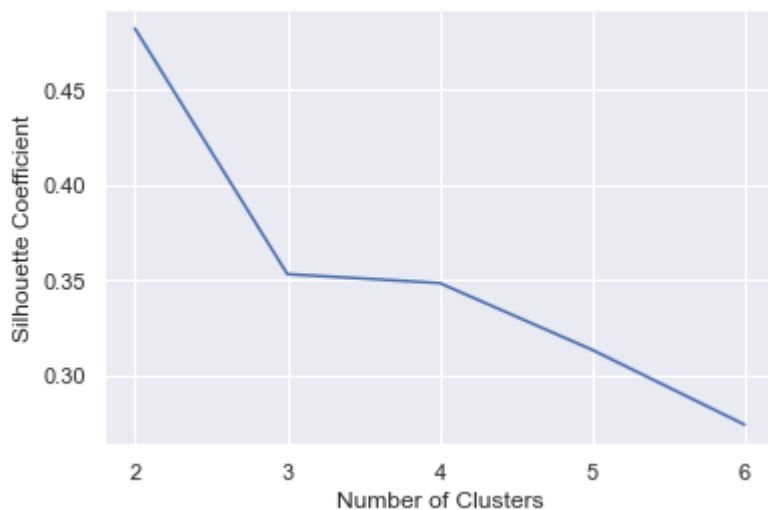
```
1 #calculating silhouette score for different centroids
2 kmeans_kwargs = {
3     "init": "random",
4     "n_init": 10,
5     "max_iter": 300,
6     "random_state": 42,
7 }
8
9
10 silhouette_coefficients = []
```

In [51]:

```
1 for k in range(2, 7):
2     kmeans = KMeans(n_clusters=k, **kmeans_kwargs)
3     kmeans.fit(cc_z1)
4     score = silhouette_score(cc_z1, kmeans.labels_)
5     silhouette_coefficients.append(score)
```

In [52]:

```
1 #plotting silhouette score for different centroids
2 plt.plot(range(2, 7), silhouette_coefficients)
3 plt.xticks(range(2, 7))
4 plt.xlabel("Number of Clusters")
5 plt.ylabel("Silhouette Coefficient")
6 plt.show()
```



In [53]:

```
1 #attaching the labels to the datasets
2 cc["cluster_K"] = labels
3 Kcar['cluster_K']=labels
4 Kclus=cc
5 Kclus.head()
```

Out[53]:

	cyl	yr	mpg	disp	hp	wt	acc	cluster_K
0	8	70	18.0	307.0	130.0	3504	12.0	0
1	8	70	15.0	350.0	165.0	3693	11.5	0
2	8	70	18.0	318.0	150.0	3436	11.0	0
3	8	70	16.0	304.0	150.0	3433	12.0	0
4	8	70	17.0	302.0	140.0	3449	10.5	0

In [54]:

```
1 cc.cluster_K.value_counts().sort_index()
```

Out[54]:

```
0    105
1    293
Name: cluster_K, dtype: int64
```

In [55]:

```
1 cc_z1["cluster_K"] = labels
2 cc_z1.head()
```

Out[55]:

	cyl	yr	mpg	disp	hp	wt	acc	cluster_K
0	1.498191	-1.627426	-0.706439	1.090604	0.673118	0.630870	-1.295498	0
1	1.498191	-1.627426	-1.090751	1.503514	1.589958	0.854333	-1.477038	0
2	1.498191	-1.627426	-0.706439	1.196232	1.197027	0.550470	-1.658577	0
3	1.498191	-1.627426	-0.962647	1.061796	1.197027	0.546923	-1.295498	0
4	1.498191	-1.627426	-0.834543	1.042591	0.935072	0.565841	-1.840117	0

In [56]:

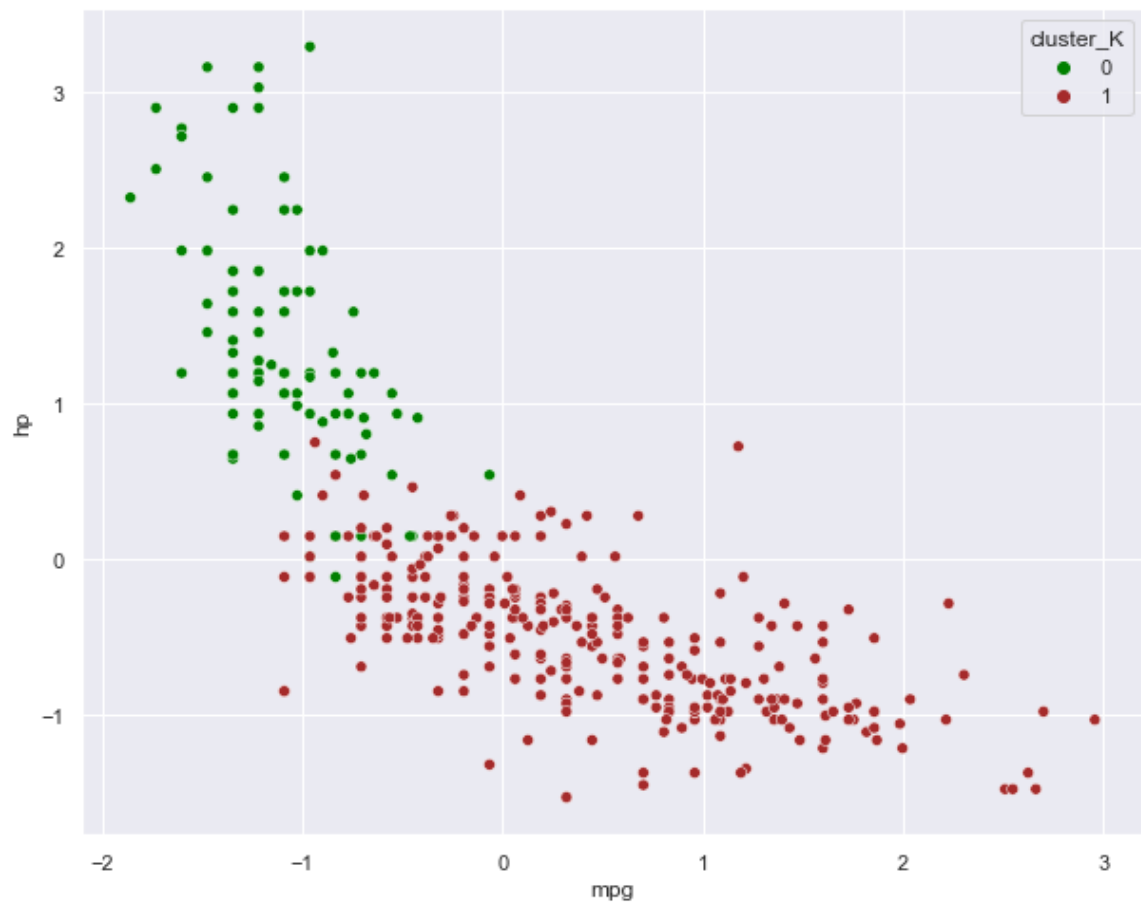
```
1 aggdata=cc.iloc[:,0:8].groupby('cluster_K').mean()
2 aggdata['Freq']=cc.cluster_K.value_counts().sort_index()
3 aggdata
```

Out[56]:

	cyl	yr	mpg	disp	hp	wt	acc	F
cluster_K								
0	7.923810	73.742857	14.851429	341.809524	158.000000	4093.771429	12.867619	✓
1	4.569966	76.822526	26.619113	140.250853	85.061433	2567.860068	16.535836	✓

In [57]:

```
1 #plotting the clusters
2 plt.figure(figsize=(10, 8))
3 sns.scatterplot(x="mpg", y="hp", hue="cluster_K",
4 data=cc_z1,
5 palette=['green', 'brown']);
```



In [58]:

```
1 Hcar.clusters_H.value_counts().sort_index()
```

Out[58]:

```
1    100
2    298
Name: clusters_H, dtype: int64
```

In [59]:

```
1 Kcar.cluster_K.value_counts().sort_index()
```

Out[59]:

```
0    105
1    293
Name: cluster_K, dtype: int64
```

In [60]:

```
1 Hcar.shape
```

Out[60]:

```
(398, 14)
```

In [61]:

```
1 Kcar.shape
```

Out[61]:

```
(398, 14)
```

In [62]:

```
1 car.head()
```

Out[62]:

	cyl	yr	mpg	disp	hp	wt	acc	origin_america	origin_asia	origin_europe	mpg_level
0	8	70	18.0	307.0	130.0	3504	12.0	1	0	0	
1	8	70	15.0	350.0	165.0	3693	11.5	1	0	0	
2	8	70	18.0	318.0	150.0	3436	11.0	1	0	0	
3	8	70	16.0	304.0	150.0	3433	12.0	1	0	0	
4	8	70	17.0	302.0	140.0	3449	10.5	1	0	0	

Linear regression on the original dataset

In [63]:

```
1 X = car.drop(['mpg', 'origin_europe', 'mpg_level_low'], axis=1)
```

In [64]:

```
1 y = car[['mpg']]
```


In [65]:

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y,
2 test_size=0.30, random_state=15)
3 regression_model = LinearRegression()
4 regression_model.fit(X_train, y_train)
5 LinearRegression()
```

Out[65]:

LinearRegression()

In [66]:

```
1 for idx, col_name in enumerate(X_train.columns):
2     print("The coefficient for {} is {}".format(col_name,
3     regression_model.coef_[0][idx]))
```

The coefficient for cyl is -0.5134441386218385
The coefficient for yr is 0.44346504291168337
The coefficient for disp is 0.010688858394646908
The coefficient for hp is 0.010315514536314019
The coefficient for wt is -0.004538788568737129
The coefficient for acc is 0.19183425608862725
The coefficient for origin_america is -1.7306209513688977
The coefficient for origin_asia is -0.8976724344009391
The coefficient for mpg_level_high is 8.552374663817035
The coefficient for mpg_level_median is 1.5941218694850414

In [67]:

```
1 intercept = regression_model.intercept_[0]
2 print("The intercept for our model is {}".format(intercept))
```

The intercept for our model is -1.663571756865423

In [68]:

```
1 V=regression_model.score(X_train, y_train)
2 V
```

Out[68]:

0.8967703023839787

Linear regression on data with K means cluster

In [69]:

```
1 Kcar['cluster_K']=Kcar['cluster_K'].astype('category')
2 Kcar['cluster_K'] = Kcar['cluster_K'].replace({1: 'heavy', 0:
3 'light'})
4 Kcar = pd.get_dummies(Kcar, columns=['cluster_K'])
```

In [70]:

```
1 Kcar.head()
```

Out[70]:

	cyl	yr	mpg	disp	hp	wt	acc	origin_america	origin_asia	origin_europe	mpg_level
0	8	70	18.0	307.0	130.0	3504	12.0	1	0	0	
1	8	70	15.0	350.0	165.0	3693	11.5	1	0	0	
2	8	70	18.0	318.0	150.0	3436	11.0	1	0	0	
3	8	70	16.0	304.0	150.0	3433	12.0	1	0	0	
4	8	70	17.0	302.0	140.0	3449	10.5	1	0	0	

In [71]:

```
1 X = Kcar.drop(['mpg', 'origin_europe', 'mpg_level_low', 'cluster_K_light'],  
2 axis=1)
```

In [72]:

```
1 y = Kcar[['mpg']]
```

In [73]:

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y,  
2 test_size=0.30, random_state=12)  
3 regression_model = LinearRegression()  
4 regression_model.fit(X_train, y_train)  
5 LinearRegression()
```

Out[73]:

LinearRegression()

In [74]:

```
1 for idx, col_name in enumerate(X_train.columns):  
2     print("The coefficient for {} is {}".format(col_name, regression_model.coef_[0][idx]))
```

The coefficient for cyl is -1.1945995644778
The coefficient for yr is 0.4318651041505994
The coefficient for disp is 0.01747749627911041
The coefficient for hp is -0.010138045835905619
The coefficient for wt is -0.0040684301693864056
The coefficient for acc is 0.1856482874625008
The coefficient for origin_america is -1.6918315494304075
The coefficient for origin_asia is -0.7407779192303029
The coefficient for mpg_level_high is 9.283120939156875
The coefficient for mpg_level_median is 2.2500017142312454
The coefficient for cluster_K_heavy is -2.5115140143384753

In [75]:

```
1 intercept = regression_model.intercept_[0]
2 print("The intercept for our model is {}".format(intercept))
```

The intercept for our model is 3.715660821055735

In [76]:

```
1 regression_model.score(X_train, y_train)
```

Out[76]:

0.8942370456543635

In [77]:

```
1 K=regression_model.score(X_test, y_test)
```

In [78]:

```
1 K
```

Out[78]:

0.9117893808052381

Linear regression on data with H-clusters

In [80]:

```
1 Hcar['clusters_H']=Hcar['clusters_H'].astype('category')
2 Hcar['clusters_H']=Hcar['clusters_H'].replace({1:'heavy', 2:'light'})
3 Hcar = pd.get_dummies(Hcar,columns=['clusters_H'])
4 Hcar.head()
```

Out[80]:

	cyl	yr	mpg	disp	hp	wt	acc	origin_america	origin_asia	origin_europe	mpg_level
0	8	70	18.0	307.0	130.0	3504	12.0	1	0	0	heavy
1	8	70	15.0	350.0	165.0	3693	11.5	1	0	0	heavy
2	8	70	18.0	318.0	150.0	3436	11.0	1	0	0	heavy
3	8	70	16.0	304.0	150.0	3433	12.0	1	0	0	heavy
4	8	70	17.0	302.0	140.0	3449	10.5	1	0	0	heavy

In [81]:

```
1 X = Hcar.drop(['mpg', 'origin_europe', 'mpg_level_low', 'clusters_H_light'],axis=1)
```

In [82]:

```
1 y = Hcar[['mpg']]
```

In [83]:

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y,
2 test_size=0.30, random_state=10)
3 regression_model = LinearRegression()
4 regression_model.fit(X_train, y_train)
5 LinearRegression()
```

Out[83]:

LinearRegression()

In [84]:

```
1 for idx, col_name in enumerate(X_train.columns):
2     print("The coefficient for {} is {}".format(col_name,
3 regression_model.coef_[0][idx]))
```

The coefficient for cyl is -1.0104832432577142
The coefficient for yr is 0.4475417357550146
The coefficient for disp is 0.01511520052461407
The coefficient for hp is -0.013301584387234512
The coefficient for wt is -0.004264179780672395
The coefficient for acc is 0.11805139164484812
The coefficient for origin_america is -2.1174569315391127
The coefficient for origin_asia is -1.3974915348558072
The coefficient for mpg_level_high is 8.565948239298272
The coefficient for mpg_level_median is 1.6577250698582815
The coefficient for clusters_H_heavy is 2.038974468807401

In [85]:

```
1 intercept = regression_model.intercept_[0]
2 print("The intercept for our model is {}".format(intercept))
```

The intercept for our model is 2.5727293182332573

In [86]:

```
1 regression_model.score(X_train, y_train)
```

Out[86]:

0.8988409890950728

In [87]:

```
1 H=regression_model.score(X_test, y_test)
2 H
```

Out[87]:

0.9010238373846695

In [88]:

```
1 modellists = []
2 modellists.append(['Linear Regression on Original Data set', V*100])
3 modellists.append(['Linear Regression with K means clusters', K*100])
4 modellists.append(['Linear Regression with Heirarchical clusters',H*100])
5 mdl_df = pd.DataFrame(modellists, columns = ['Model','r^2 on Test'])
6 mdl_df
```

Out[88]:

	Model	r^2 on Test
0	Linear Regression on Original Data set	89.677030
1	Linear Regression with K means clusters	91.178938
2	Linear Regression with Heirarchical clusters	90.102384