

Multiple Linear Regression

In [1]:

```
# import libraries
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import statsmodels.formula.api as smf
from statsmodels.graphics.regressionplots import influence_plot
```

In [2]:

```
# Read dataset
cars=pd.read_csv("C:/Users/Ashraf/Documents/Datafiles/Cars.csv")
```

In [3]:

```
cars.head()
```

Out[3]:

	HP	MPG	VOL	SP	WT
0	49	53.700681	89	104.185353	28.762059
1	55	50.013401	92	105.461264	30.466833
2	55	50.013401	92	105.461264	30.193597
3	70	45.696322	92	113.461264	30.632114
4	53	50.504232	92	104.461264	29.889149

In [4]:

```
cars.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 81 entries, 0 to 80
Data columns (total 5 columns):
#   Column  Non-Null Count  Dtype
---  -
0    HP      81 non-null      int64
1    MPG      81 non-null      float64
2    VOL      81 non-null      int64
3    SP       81 non-null      float64
4    WT       81 non-null      float64
dtypes: float64(3), int64(2)
memory usage: 3.3 KB
```

In [5]:

```
# checking for missing values  
cars.isnull().sum()
```

Out[5]:

```
HP      0  
MPG      0  
VOL      0  
SP      0  
WT      0  
dtype: int64
```

Correlation Matrix

In [6]:

```
cars.corr()
```

Out[6]:

	HP	MPG	VOL	SP	WT
HP	1.000000	-0.725038	0.077459	0.973848	0.076513
MPG	-0.725038	1.000000	-0.529057	-0.687125	-0.526759
VOL	0.077459	-0.529057	1.000000	0.102170	0.999203
SP	0.973848	-0.687125	0.102170	1.000000	0.102439
WT	0.076513	-0.526759	0.999203	0.102439	1.000000

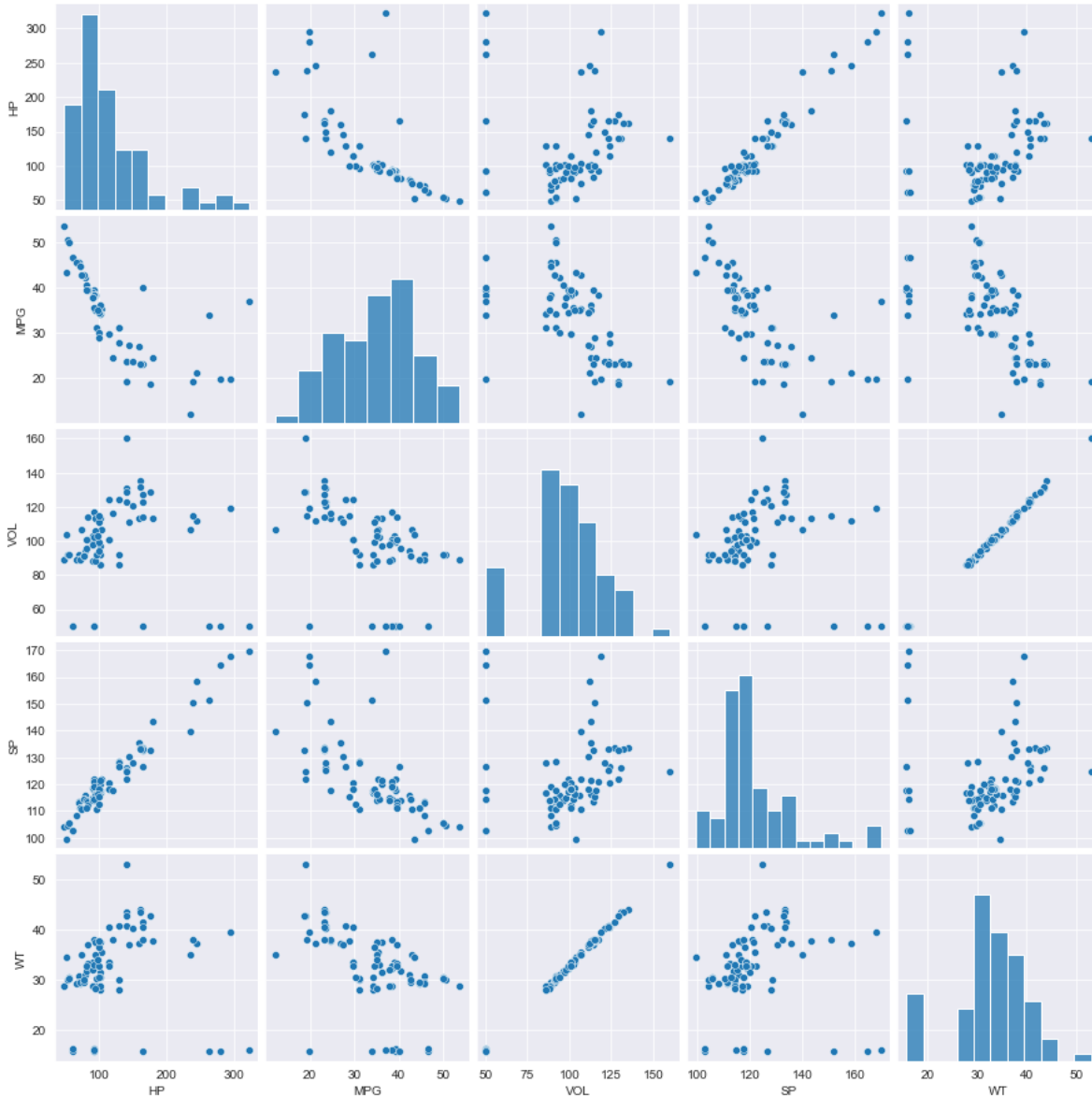
Scatterplot between variables along with histograms

In [7]:

```
#Format the plot background and scatter plots for all the variables  
sns.set_style(style='darkgrid')  
sns.pairplot(cars)
```

Out[7]:

<seaborn.axisgrid.PairGrid at 0xd3786c09a0>



Preparing a model

In [8]:

```
# build model
import statsmodels.formula.api as smf
model=smf.ols("MPG~HP+SP+VOL+WT",data=cars).fit()
```

In [9]:

```
# Coefficients
model.params
```

Out[9]:

```
Intercept    30.677336
HP            -0.205444
SP             0.395627
VOL           -0.336051
WT             0.400574
dtype: float64
```

In [10]:

```
# p-values and t-values
print(model.tvalues, '\n', model.pvalues)
```

```
Intercept    2.058841
HP           -5.238735
SP            2.499880
VOL          -0.590970
WT            0.236541
dtype: float64
Intercept    0.042936
HP            0.000001
SP            0.014579
VOL           0.556294
WT            0.813649
dtype: float64
```

In [11]:

```
# R squared
(model.rsquared, model.rsquared_adj)
```

Out[11]:

```
(0.7705372737359844, 0.7584602881431415)
```

In [12]:

```
# Building slr for VOL,WT
model1=smf.ols("MPG~VOL",data=cars).fit()
```

In [13]:

```
#print t-values and p-values
print(model1.tvalues,'\n',model1.pvalues) # p-value is significant
```

```
Intercept    14.106056
VOL          -5.541400
dtype: float64
Intercept    2.753815e-23
VOL          3.822819e-07
dtype: float64
```

In [14]:

```
model2=smf.ols("MPG~WT",data=cars).fit()
```

In [15]:

```
# print t-values and p-values
(model2.tvalues,'\n',model2.pvalues) # p-values is significant
```

Out[15]:

```
(Intercept    14.248923
WT           -5.508067
dtype: float64,
'\n',
Intercept    1.550788e-23
WT           4.383467e-07
dtype: float64)
```

In [16]:

```
# Building MLR for VOL and WT
model3=smf.ols("MPG~VOL+WT",data=cars).fit()
```

In [17]:

```
# print t-values and p-values
print(model3.tvalues,'\n',model3.pvalues) # p-values of both VOL and WT are insignificant
```

```
Intercept    12.545736
VOL          -0.709604
WT           0.489876
dtype: float64
Intercept    2.141975e-20
VOL          4.800657e-01
WT           6.255966e-01
dtype: float64
```

Calculating VIF

In [18]:

```
rsq_hp=smf.ols("HP~SP+VOL+WT",data=cars).fit().rsquared
vif_hp=1/(1-rsq_hp)

rsq_sp=smf.ols("SP~VOL+WT+HP",data=cars).fit().rsquared
vif_sp=1/(1-rsq_sp)

rsq_vol=smf.ols("VOL~WT+HP+SP",data=cars).fit().rsquared
vif_vol=1/(1-rsq_vol)

rsq_wt=smf.ols("WT~HP+SP+VOL",data=cars).fit().rsquared
vif_wt=1/(1-rsq_wt)

# Storing VIF values in dataframe
d1={'Variables':['HP','SP','VOL','WT'], 'VIF':[vif_hp, vif_sp, vif_vol, vif_wt]}
Vif_frame=pd.DataFrame(d1)
Vif_frame
```

Out[18]:

	Variables	VIF
0	HP	19.926589
1	SP	20.007639
2	VOL	638.806084
3	WT	639.533818

Residual Analysis

Test for Normality of Residuals (Q-Q Plot)

In [19]:

```
import statsmodels.api as sm
qqplot=sm.qqplot(model.resid,line='q')
plt.title("Normal Q-Q plot of residuals")
plt.show()
```

C:\Users\Ashraf\anaconda3\lib\site-packages\statsmodels\graphics\gofplots.py:993: UserWarning: marker is redundantly defined by the 'marker' keyword argument and the fmt string "bo" (-> marker='o'). The keyword argument will take precedence.

```
ax.plot(x, y, fmt, **plot_style)
```



In [20]:

```
list(np.where(model.resid>10))
```

Out[20]:

```
[array([ 0, 76], dtype=int64)]
```

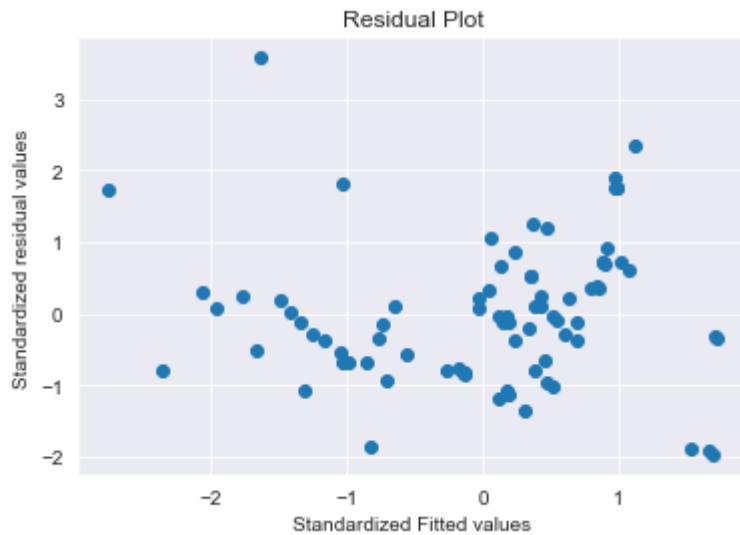
Residual Plot for Homoscedasticity

In [21]:

```
def get_standardized_values( vals ):
    return (vals - vals.mean())/vals.std()
```

In [22]:

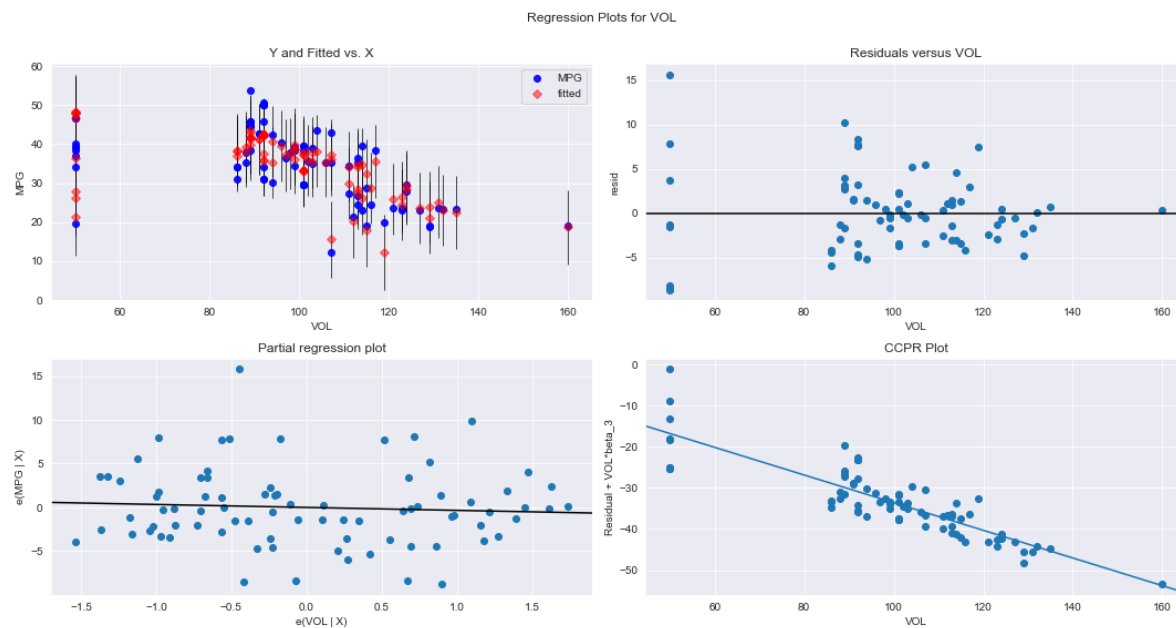
```
plt.scatter(get_standardized_values(model.fittedvalues),  
            get_standardized_values(model.resid))  
  
plt.title('Residual Plot')  
plt.xlabel('Standardized Fitted values')  
plt.ylabel('Standardized residual values')  
plt.show()
```



Residual Vs Regressors

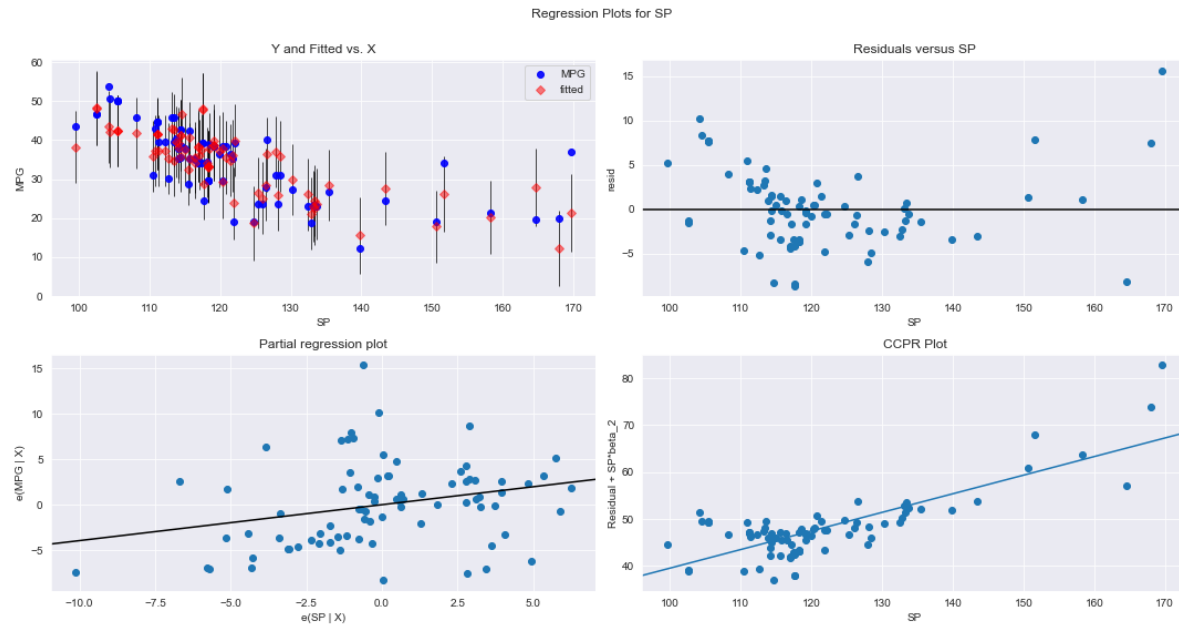
In [23]:

```
fig = plt.figure(figsize=(15,8))  
fig = sm.graphics.plot_regress_exog(model, "VOL", fig=fig)  
plt.show()
```



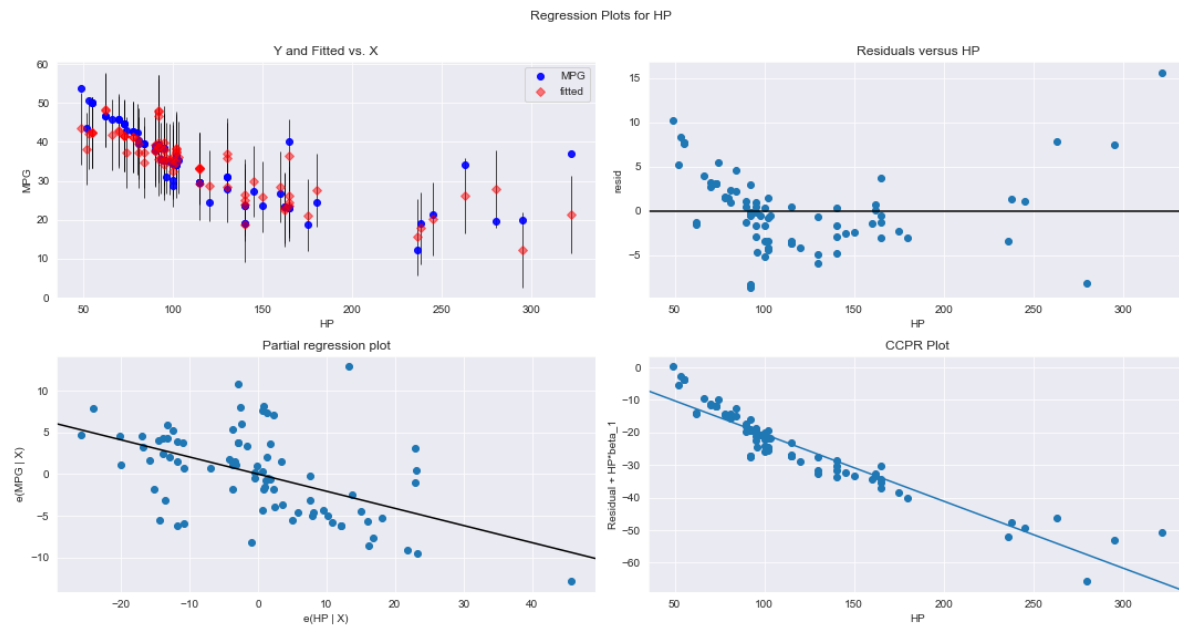
In [24]:

```
fig = plt.figure(figsize=(15,8))
fig = sm.graphics.plot_regress_exog(model, "SP", fig=fig)
plt.show()
```



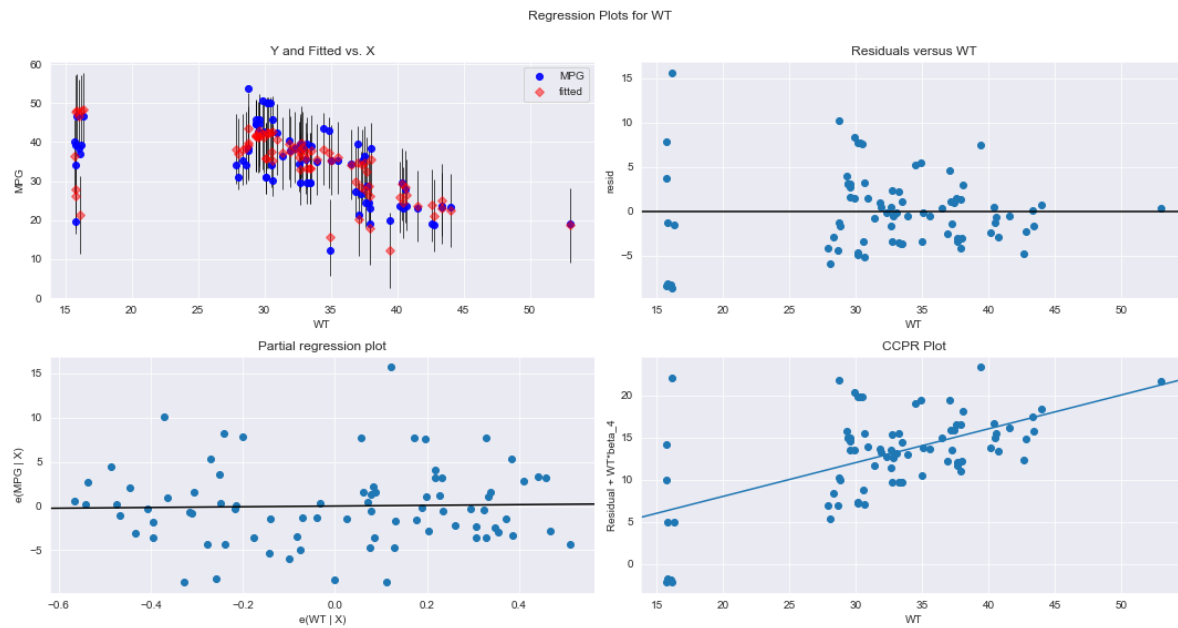
In [25]:

```
fig = plt.figure(figsize=(15,8))
fig = sm.graphics.plot_regress_exog(model, "HP", fig=fig)
plt.show()
```



In [26]:

```
fig = plt.figure(figsize=(15,8))
fig = sm.graphics.plot_regress_exog(model, "WT", fig=fig)
plt.show()
```



Model Deletion Diagnostics

Detecting Influencers/Outliers

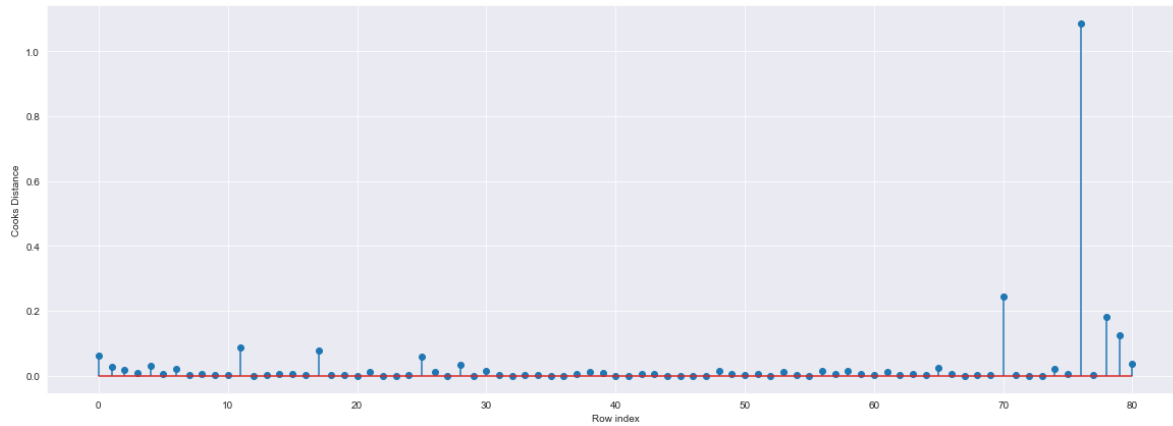
Cook's Distance

In [27]:

```
model_influence = model.get_influence()
(c, _) = model_influence.cooks_distance
```

In [28]:

```
#Plot the influencers values using stem plot  
fig = plt.subplots(figsize=(20, 7))  
plt.stem(np.arange(len(cars)), np.round(c, 3))  
plt.xlabel('Row index')  
plt.ylabel('Cooks Distance')  
plt.show()
```



In [29]:

```
#index and value of influencer where c is more than .5  
(np.argmax(c), np.max(c))
```

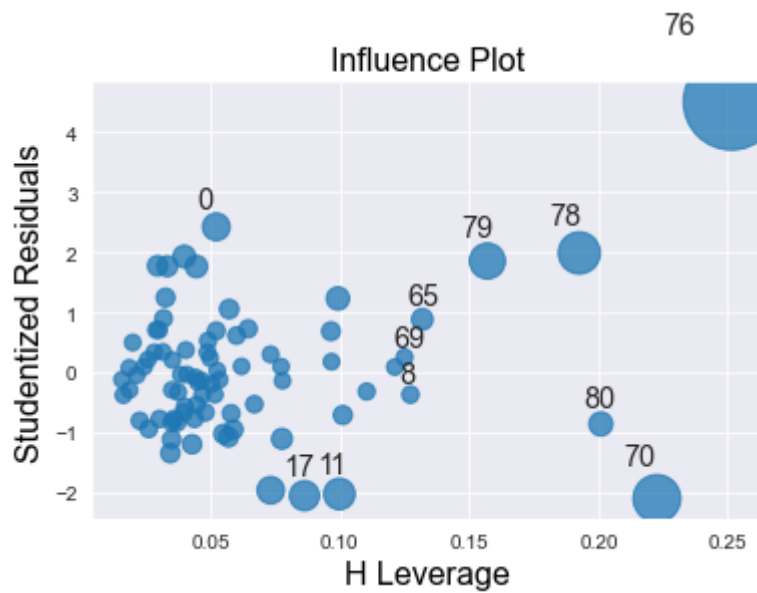
Out[29]:

```
(76, 1.0865193998180045)
```

High Influence points

In [30]:

```
from statsmodels.graphics.regressionplots import influence_plot
influence_plot(model)
plt.show()
```



In [31]:

```
k = cars.shape[1]
n = cars.shape[0]
leverage_cutoff = 3*((k + 1)/n)
```

In [32]:

```
leverage_cutoff
```

Out[32]:

```
0.2222222222222222
```

From the above plot, it is evident that data point 70 and 76 are the influencers

In [33]:

```
cars[cars.index.isin([70, 76])]
```

Out[33]:

	HP	MPG	VOL	SP	WT
70	280	19.678507	50	164.598513	15.823060
76	322	36.900000	50	169.598513	16.132947

In [34]:

```
#See the differences in HP and other variable values  
cars.head()
```

Out[34]:

	HP	MPG	VOL	SP	WT
0	49	53.700681	89	104.185353	28.762059
1	55	50.013401	92	105.461264	30.466833
2	55	50.013401	92	105.461264	30.193597
3	70	45.696322	92	113.461264	30.632114
4	53	50.504232	92	104.461264	29.889149

Improving the model

In [35]:

```
cars_new=pd.read_csv("C:/Users/Ashraf/Documents/Datafiles/Cars.csv")
```

In [36]:

```
#Discard the data points which are influencers and reassign the row number (reset_index())  
car1=cars_new.drop(cars_new.index[[70,76]],axis=0).reset_index()
```

In [37]:

```
car1
```

Out[37]:

	index	HP	MPG	VOL	SP	WT
0	0	49	53.700681	89	104.185353	28.762059
1	1	55	50.013401	92	105.461264	30.466833
2	2	55	50.013401	92	105.461264	30.193597
3	3	70	45.696322	92	113.461264	30.632114
4	4	53	50.504232	92	104.461264	29.889149
...
74	75	175	18.762837	129	132.864163	42.778219
75	77	238	19.197888	115	150.576579	37.923113
76	78	263	34.000000	50	151.598513	15.769625
77	79	295	19.833733	119	167.944460	39.423099
78	80	236	12.101263	107	139.840817	34.948615

79 rows × 6 columns

In [38]:

```
#Drop the original index
car1=car1.drop(['index'],axis=1)
```

In [39]:

car1

Out[39]:

	HP	MPG	VOL	SP	WT
0	49	53.700681	89	104.185353	28.762059
1	55	50.013401	92	105.461264	30.466833
2	55	50.013401	92	105.461264	30.193597
3	70	45.696322	92	113.461264	30.632114
4	53	50.504232	92	104.461264	29.889149
...
74	175	18.762837	129	132.864163	42.778219
75	238	19.197888	115	150.576579	37.923113
76	263	34.000000	50	151.598513	15.769625
77	295	19.833733	119	167.944460	39.423099
78	236	12.101263	107	139.840817	34.948615

79 rows × 5 columns

Build Model

In [40]:

```
#Exclude variable "WT" and generate R-Squared and AIC values
final_ml_V= smf.ols('MPG~VOL+SP+HP',data = car1).fit()
```

In [41]:

```
(final_ml_V.rsquared,final_ml_V.aic)
```

Out[41]:

```
(0.8161692010376007, 446.11722639447726)
```

In [42]:

```
#Exclude variable "VOL" and generate R-Squared and AIC values
final_ml_W= smf.ols('MPG~WT+SP+HP',data = car1).fit()
```

In [43]:

```
(final_ml_W.rsquared,final_ml_W.aic)
```

Out[43]:

```
(0.8160034320495304, 446.1884323575032)
```

Comparing above R-Square and AIC values, model 'final_ml_V' has high R- square and low AIC value hence include variable 'VOL' so that multi collinearity problem would be resolved.

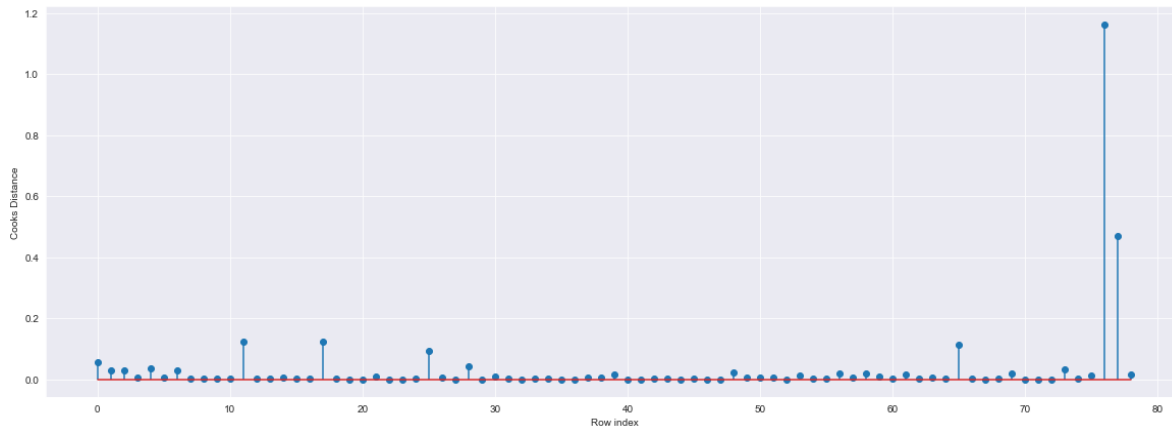
Cook's Distance

In [44]:

```
model_influence_V = final_ml_V.get_influence()  
(c_V, _) = model_influence_V.cooks_distance
```

In [45]:

```
fig= plt.subplots(figsize=(20,7))  
plt.stem(np.arange(len(car1)),np.round(c_V,3));  
plt.xlabel('Row index')  
plt.ylabel('Cooks Distance');
```



In [46]:

```
#index of the data points where c is more than .5  
(np.argmax(c_V),np.max(c_V))
```

Out[46]:

```
(76, 1.1629387469135095)
```

In [47]:

```
#Drop 76 and 77 observations  
car2=car1.drop(car1.index[[76,77]],axis=0)
```

In [48]:

```
car2
```

Out[48]:

	HP	MPG	VOL	SP	WT
0	49	53.700681	89	104.185353	28.762059
1	55	50.013401	92	105.461264	30.466833
2	55	50.013401	92	105.461264	30.193597
3	70	45.696322	92	113.461264	30.632114
4	53	50.504232	92	104.461264	29.889149
...
72	140	19.086341	160	124.715241	52.997752
73	140	19.086341	129	121.864163	42.618698
74	175	18.762837	129	132.864163	42.778219
75	238	19.197888	115	150.576579	37.923113
78	236	12.101263	107	139.840817	34.948615

77 rows × 5 columns

In [49]:

```
#Reset the index and re arrange the row values  
car3=car2.reset_index()
```

In [50]:

```
car3
```

Out[50]:

	index	HP	MPG	VOL	SP	WT
0	0	49	53.700681	89	104.185353	28.762059
1	1	55	50.013401	92	105.461264	30.466833
2	2	55	50.013401	92	105.461264	30.193597
3	3	70	45.696322	92	113.461264	30.632114
4	4	53	50.504232	92	104.461264	29.889149
...
72	72	140	19.086341	160	124.715241	52.997752
73	73	140	19.086341	129	121.864163	42.618698
74	74	175	18.762837	129	132.864163	42.778219
75	75	238	19.197888	115	150.576579	37.923113
76	78	236	12.101263	107	139.840817	34.948615

77 rows × 6 columns

In [51]:

```
car4=car3.drop(['index'],axis=1)
```

In [52]:

```
car4
```

Out[52]:

	HP	MPG	VOL	SP	WT
0	49	53.700681	89	104.185353	28.762059
1	55	50.013401	92	105.461264	30.466833
2	55	50.013401	92	105.461264	30.193597
3	70	45.696322	92	113.461264	30.632114
4	53	50.504232	92	104.461264	29.889149
...
72	140	19.086341	160	124.715241	52.997752
73	140	19.086341	129	121.864163	42.618698
74	175	18.762837	129	132.864163	42.778219
75	238	19.197888	115	150.576579	37.923113
76	236	12.101263	107	139.840817	34.948615

77 rows × 5 columns

In [53]:

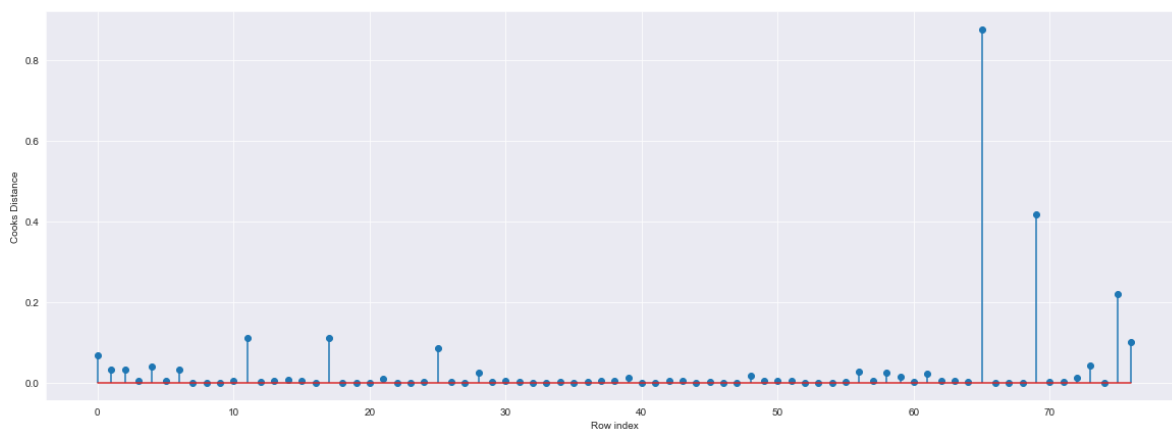
```
#Build the model on the new data
final_ml_V= smf.ols('MPG~VOL+SP+HP',data = car4).fit()
```

In [54]:

```
#Again check for influencers
model_influence_V = final_ml_V.get_influence()
(c_V, _) = model_influence_V.cooks_distance
```

In [55]:

```
fig= plt.subplots(figsize=(20,7))
plt.stem(np.arange(len(car4)),np.round(c_V,3));
plt.xlabel('Row index')
plt.ylabel('Cooks Distance');
```



In [56]:

```
#index of the data points where c is more than .5
(np.argmax(c_V),np.max(c_V))
```

Out[56]:

```
(65, 0.8774556986296826)
```

Since the value is <1 , we can stop the diagnostic process and finalize the model

In [57]:

```
#Check the accuracy of the mode
final_ml_V= smf.ols('MPG~VOL+SP+HP',data = car4).fit()
```

In [58]:

```
(final_ml_V.rsquared,final_ml_V.aic)
```

Out[58]:

```
(0.8669636111859063, 409.41530627195084)
```

Predicting for new data

In [59]:

```
#New data for prediction  
new_data=pd.DataFrame({'HP':40,"VOL":95,"SP":102,"WT":35},index=[1])
```

In [60]:

```
final_ml_V.predict(new_data)
```

Out[60]:

```
1    46.035594  
dtype: float64
```

In [61]:

```
final_ml_V.predict(cars_new.iloc[0:5,])
```

Out[61]:

```
0    45.428872  
1    43.992392  
2    43.992392  
3    43.508150  
4    44.085858  
dtype: float64
```

In [62]:

```
pred_y = final_ml_V.predict(cars_new)
```

In [63]:

```
pred_y
```

Out[63]:

```
0    45.428872  
1    43.992392  
2    43.992392  
3    43.508150  
4    44.085858  
...  
76    7.165876  
77    12.198598  
78    14.908588  
79    4.163958  
80    9.161202  
Length: 81, dtype: float64
```

In []: