

In [1]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import statsmodels.api as sm
from numpy.polynomial.polynomial import polyfit
from sklearn.linear_model import LinearRegression
import seaborn as sns
import statsmodels.stats.tests.test_influence
from sklearn.feature_selection import RFE
from statsmodels.stats.outliers_influence import variance_inflation_factor
```

executed in 32.2s, finished 18:54:02 2022-01-05

In [2]:

```
df = pd.read_csv("ToyotaCorolla.csv", encoding='latin1')
df1 = df [['Price', 'Age_08_04', 'KM', 'HP', 'cc', 'Doors', 'Gears', 'Quarterly_Tax', 'Weight']]
df1
```

executed in 251ms, finished 18:54:07 2022-01-05

Out[2]:

	Price	Age_08_04	KM	HP	cc	Doors	Gears	Quarterly_Tax	Weight
0	13500	23	46986	90	2000	3	5	210	1165
1	13750	23	72937	90	2000	3	5	210	1165
2	13950	24	41711	90	2000	3	5	210	1165
3	14950	26	48000	90	2000	3	5	210	1165
4	13750	30	38500	90	2000	3	5	210	1170
...	...	...	...	...	...	...	...	...	...
1431	7500	69	20544	86	1300	3	5	69	1025
1432	10845	72	19000	86	1300	3	5	69	1015
1433	8500	71	17016	86	1300	3	5	69	1015
1434	7250	70	16916	86	1300	3	5	69	1015
1435	6950	76	1	110	1600	5	5	19	1114

1436 rows × 9 columns

In [5]:

df1.describe

executed in 24ms, finished 13:53:20 2021-11-27

Out[5]:

```
<bound method NDFrame.describe of
Doors  Gears  Quarterly_Tax  Weight  Price  Age_08_04  KM  HP  cc
0      13500      23  46986   90  2000      3      5      210  116
5
1      13750      23  72937   90  2000      3      5      210  116
5
2      13950      24  41711   90  2000      3      5      210  116
5
3      14950      26  48000   90  2000      3      5      210  116
5
4      13750      30  38500   90  2000      3      5      210  117
0
...      ...      ...      ...  ...      ...      ...      ...
...
1431   7500      69  20544   86  1300      3      5      69  102
5
1432  10845      72  19000   86  1300      3      5      69  101
5
1433   8500      71  17016   86  1300      3      5      69  101
5
1434   7250      70  16916   86  1300      3      5      69  101
5
1435   6950      76      1  110  1600      5      5      19  111
4
```

[1436 rows x 9 columns]&gt;

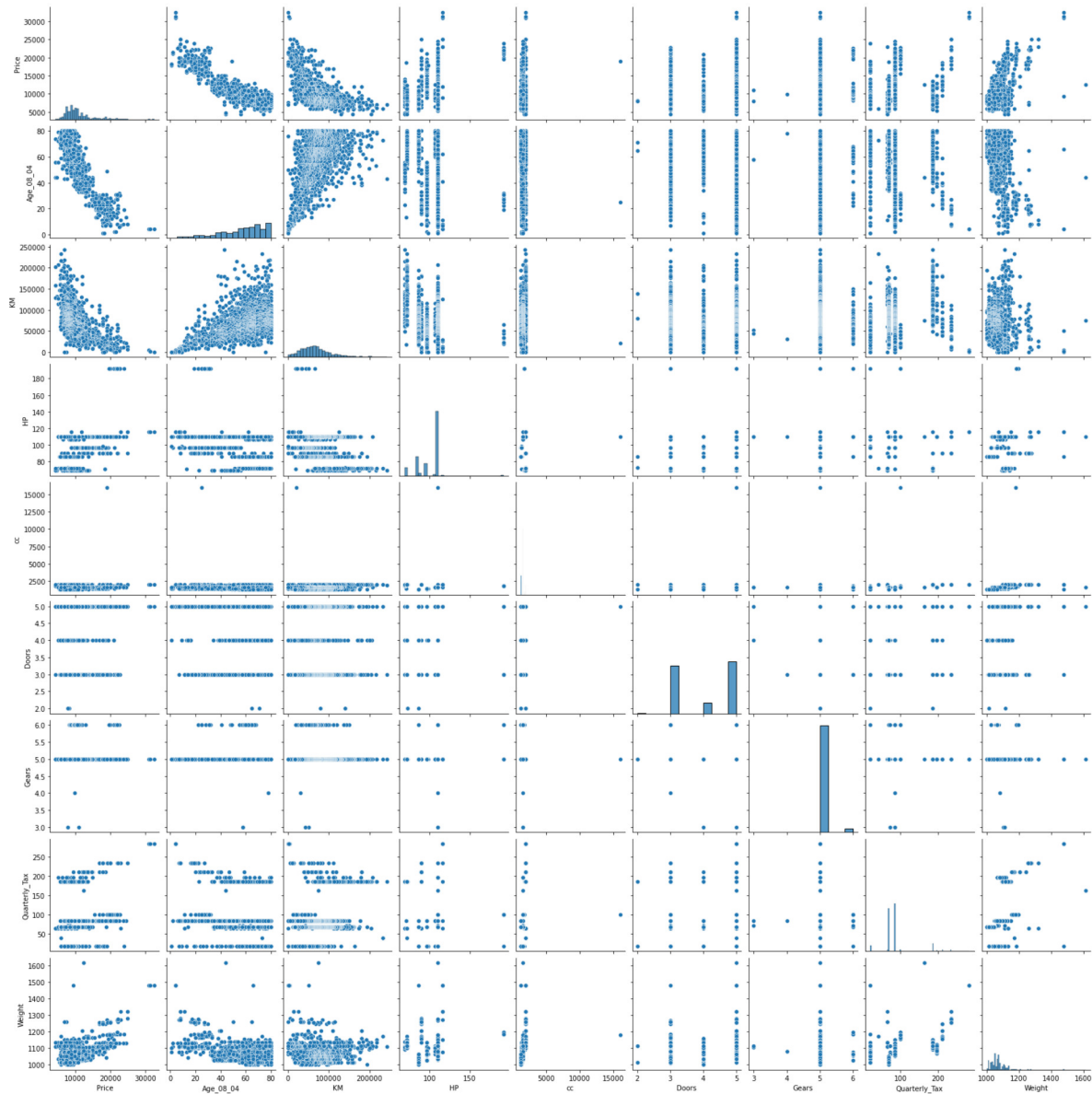
In [6]:

```
sns.pairplot(df1)
```

executed in 27.5s, finished 13:54:03 2021-11-27

Out[6]:

&lt;seaborn.axisgrid.PairGrid at 0x2d407fa1c70&gt;



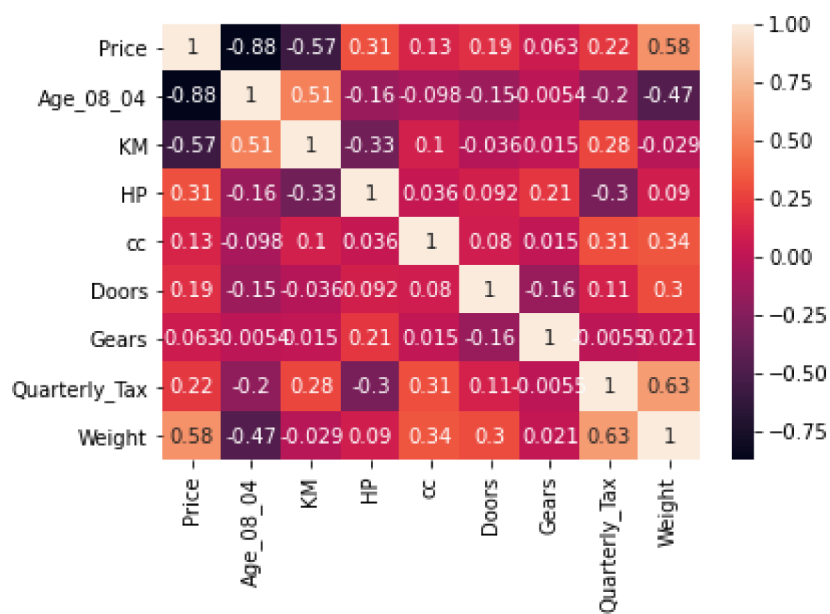
In [7]:

```
corrMatrix = df1.corr()
sns.heatmap(corrMatrix, annot=True)
```

executed in 2.31s, finished 13:55:12 2021-11-27

Out[7]:

&lt;AxesSubplot:&gt;



In [8]:

```
X = df1[['Age_08_04', 'KM', 'HP', 'cc', 'Doors', 'Gears', 'Quarterly_Tax', 'Weight']]
Y = df1[['Price']]
model = sm.OLS(Y, X).fit()
predictions = model.predict(X)
model.summary()
```

executed in 89ms, finished 13:55:53 2021-11-27

Out[8]:

OLS Regression Results

<b>Dep. Variable:</b>	Price	<b>R-squared (uncentered):</b>	0.986
<b>Model:</b>	OLS	<b>Adj. R-squared (uncentered):</b>	0.986
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1.247e+04
<b>Date:</b>	Sat, 27 Nov 2021	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	13:55:53	<b>Log-Likelihood:</b>	-12383.
<b>No. Observations:</b>	1436	<b>AIC:</b>	2.478e+04
<b>Df Residuals:</b>	1428	<b>BIC:</b>	2.482e+04
<b>Df Model:</b>	8		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>Age_08_04</b>	-125.4510	2.445	-51.303	0.000	-130.248	-120.654
<b>KM</b>	-0.0205	0.001	-16.305	0.000	-0.023	-0.018
<b>HP</b>	33.4737	2.796	11.973	0.000	27.990	38.958
<b>cc</b>	-0.1032	0.090	-1.141	0.254	-0.281	0.074
<b>Doors</b>	-7.2494	40.184	-0.180	0.857	-86.075	71.576
<b>Gears</b>	78.3780	148.258	0.529	0.597	-212.449	369.205
<b>Quarterly_Tax</b>	5.8258	1.227	4.748	0.000	3.419	8.233
<b>Weight</b>	14.0322	0.773	18.157	0.000	12.516	15.548

<b>Omnibus:</b>	108.641	<b>Durbin-Watson:</b>	1.509
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	562.996
<b>Skew:</b>	0.019	<b>Prob(JB):</b>	5.59e-123
<b>Kurtosis:</b>	6.067	<b>Cond. No.</b>	3.26e+05

Notes:

- [1]  $R^2$  is computed without centering (uncentered) since the model does not contain a constant.
- [2] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [3] The condition number is large, 3.26e+05. This might indicate that there are strong multicollinearity or other numerical problems.

In [9]:

```
infl = model.get_influence()
summ_df = infl.summary_frame()
summ_df.sort_values('cooks_d', ascending=False)
```

executed in 3.58s, finished 13:56:18 2021-11-27

Out[9]:

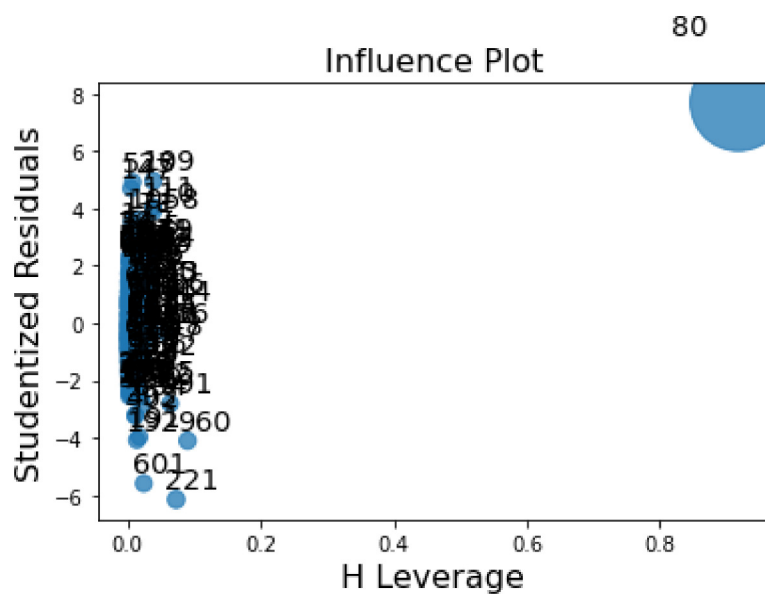
	dfb_Age_08_04	dfb_KM	dfb_HP	dfb_cc	dfb_Doors	dfb_Gears	dfb_Quarl
<b>80</b>	-0.289686	-2.363023e+00	-2.781734	2.615263e+01	0.645518	0.905398	-
<b>221</b>	-0.255051	-2.607849e-02	0.174427	4.179288e-01	0.553762	1.572311	
<b>960</b>	-0.208597	3.740370e-02	0.443719	1.977218e-01	0.539659	0.995730	
<b>109</b>	0.104572	-2.934313e-01	0.114571	-1.902240e-01	-0.220145	-0.606577	
<b>601</b>	-0.213100	3.186680e-01	0.385097	-1.107466e-01	0.404397	0.551158	
...	...	...	...	...	...	...	
<b>1167</b>	-0.000069	-5.512316e-05	0.000142	6.568459e-05	0.000134	-0.000053	
<b>482</b>	-0.000005	2.712110e-08	-0.000085	2.916558e-07	-0.000162	0.000013	-
<b>1433</b>	-0.000136	1.856487e-04	0.000104	1.209905e-05	0.000052	-0.000037	-
<b>397</b>	-0.000015	4.001993e-05	0.000028	1.344847e-06	-0.000098	-0.000034	
<b>922</b>	-0.000027	1.632787e-05	-0.000027	-8.283992e-06	0.000052	0.000017	

1436 rows × 14 columns



```
infl.plot_influence()
```

Out[10]:



In [11]:

```
vif = pd.DataFrame()
vif["VIF Factor"] = [variance_inflation_factor(X.values, i) for i in range(X.shape[1])]
vif["features"] = X.columns
vif.round(1)
```

executed in 64ms, finished 13:57:16 2021-11-27

Out[11]:

	VIF Factor	features
0	16.4	Age_08_04
1	7.6	KM
2	64.9	HP
3	17.2	cc
4	21.9	Doors
5	438.6	Gears
6	11.0	Quarterly_Tax
7	543.1	Weight



In [12]:

```
# Removed Weight from the dataframe and Tested the model
new_X = df1[['Age_08_04', 'KM', 'HP', 'cc', 'Doors', 'Gears', 'Quarterly_Tax']]
new_model = sm.OLS(Y, new_X).fit()
new_predictions = new_model.predict(new_X)
new_model.summary()
```

executed in 72ms, finished 13:57:32 2021-11-27

Out[12]:

OLS Regression Results

<b>Dep. Variable:</b>	Price	<b>R-squared (uncentered):</b>	0.983
<b>Model:</b>	OLS	<b>Adj. R-squared (uncentered):</b>	0.983
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1.155e+04
<b>Date:</b>	Sat, 27 Nov 2021	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	13:57:32	<b>Log-Likelihood:</b>	-12532.
<b>No. Observations:</b>	1436	<b>AIC:</b>	2.508e+04
<b>Df Residuals:</b>	1429	<b>BIC:</b>	2.512e+04
<b>Df Model:</b>	7		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>Age_08_04</b>	-132.0628	2.682	-49.245	0.000	-137.323	-126.802
<b>KM</b>	-0.0208	0.001	-14.947	0.000	-0.024	-0.018
<b>HP</b>	44.4711	3.027	14.692	0.000	38.533	50.409
<b>cc</b>	0.1805	0.099	1.827	0.068	-0.013	0.374
<b>Doors</b>	272.5298	41.159	6.621	0.000	191.791	353.269
<b>Gears</b>	2417.9083	81.331	29.729	0.000	2258.368	2577.449
<b>Quarterly_Tax</b>	17.0169	1.177	14.462	0.000	14.709	19.325

<b>Omnibus:</b>	184.883	<b>Durbin-Watson:</b>	1.396
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	583.580
<b>Skew:</b>	0.640	<b>Prob(JB):</b>	1.89e-127
<b>Kurtosis:</b>	5.849	<b>Cond. No.</b>	1.63e+05

Notes:

- [1]  $R^2$  is computed without centering (uncentered) since the model does not contain a constant.
- [2] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [3] The condition number is large, 1.63e+05. This might indicate that there are strong multicollinearity or other numerical problems.

In [13]:

```
# Removed Index with highest Cook's distance to remove the highest influential ROW
new_df = df1.drop(df1.index[80])
new_X = new_df[['Age_08_04', 'KM', 'HP', 'cc', 'Doors', 'Gears', 'Quarterly_Tax']]
new_Y = new_df[['Price']]
final_model = sm.OLS(new_Y, new_X).fit()
predictions = final_model.predict(new_X)
final_model.summary()
```

executed in 80ms, finished 13:57:52 2021-11-27

Out[13]:

OLS Regression Results

<b>Dep. Variable:</b>	Price	<b>R-squared (uncentered):</b>	0.983
<b>Model:</b>	OLS	<b>Adj. R-squared (uncentered):</b>	0.983
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1.152e+04
<b>Date:</b>	Sat, 27 Nov 2021	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	13:57:52	<b>Log-Likelihood:</b>	-12524.
<b>No. Observations:</b>	1435	<b>AIC:</b>	2.506e+04
<b>Df Residuals:</b>	1428	<b>BIC:</b>	2.510e+04
<b>Df Model:</b>	7		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>Age_08_04</b>	-132.0191	2.684	-49.196	0.000	-137.283	-126.755
<b>KM</b>	-0.0210	0.001	-14.581	0.000	-0.024	-0.018
<b>HP</b>	43.7530	3.287	13.310	0.000	37.304	50.202
<b>cc</b>	0.3468	0.313	1.109	0.268	-0.267	0.960
<b>Doors</b>	270.3889	41.346	6.540	0.000	189.284	351.494
<b>Gears</b>	2394.0486	91.807	26.077	0.000	2213.957	2574.140
<b>Quarterly_Tax</b>	16.4778	1.520	10.843	0.000	13.497	19.459

<b>Omnibus:</b>	183.937	<b>Durbin-Watson:</b>	1.393
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	582.178
<b>Skew:</b>	0.636	<b>Prob(JB):</b>	3.82e-127
<b>Kurtosis:</b>	5.849	<b>Cond. No.</b>	1.82e+05

Notes:

- [1]  $R^2$  is computed without centering (uncentered) since the model does not contain a constant.
- [2] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [3] The condition number is large, 1.82e+05. This might indicate that there are strong multicollinearity or other numerical problems.

