**ADP_Red\stats\4.reg.py**

```python
1   # %% Chapter 6. Statistics - Regression
2   import numpy as np
3   import pandas as pd
4   import seaborn as sns
5   import matplotlib.pyplot as plt
6
7   import scipy.stats as stats
8
9
10  # %% 1. 데이터 수집
11  df = pd.read_csv('../../ADP_Python/data/Cars93.csv')
12  sample = df[['EngineSize', 'RPM', 'Weight', 'Length', 'MPG.city', 'MPG.highway', 'Price']]
13  sample.columns = ['EngineSize', 'RPM', 'Weight', 'Length', 'MPGcity', 'MPGhighway', 'Price'
    ]
14
15  print(sample.shape)
16  print(sample.info())
17
18  # Check binary variable
19  for i, var in enumerate(sample.columns):
20      print(i, var, len(sample[var].unique()))
21
22  # Check data summary
23  print(sample.describe())
24
25  # Check scatterplot
26  from pandas.plotting import scatter_matrix
27
28  scatter_matrix(sample)
29  plt.show()
30
31
32  # %% 2. 데이터 결측치 보정
33  print(sample.isna().sum())
34
35
36  # %% 6. 정규화 또는 스케일 작업
37  # Boxplot for scaling check
38  sns.boxplot(sample)
39  plt.tight_layout()
40  plt.show()
41
42  # Multicollinearity
43  sns.heatmap(sample.corr(), annot=True)
44  plt.show()
45
46  print(sample.corr())
47
48  # %% 9. 모델 학습 - Stats
49  import scipy.stats as stats
50  import statsmodels.api as sm
51  import statsmodels.formula.api as smf
52
53  # Simple Linear Regression
54  model = smf.ols(
55      formula='Price ~ 1 + Length', data=sample
56  ).fit()
```

```
57
58   print(model.summary())
59
60   # Residual Plot
61   sns.scatterplot(model.resid)
62   plt.show()
63
64
65   # Multiple Linear Regression
66   model2 = smf.ols(
67       formula='Price ~ EngineSize + RPM + Weight + Length + MPGcity + MPGhighway', data=
     sample
68   ).fit()
69
70   model2.summary()
```