**ADP_Red\stats\2.anova.py**

```python
 1  # %% Chapter 6. Statistics - ANOVA
 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/iris.csv')
12
13  print(df.shape)
14  print(df.info())
15
16  # Check binary variable
17  for i, var in enumerate(df.columns):
18      print(i, var, len(df[var].unique()))
19
20  # Check data summary
21  print(df.describe())
22
23  # 2. 데이터 결측치 보정
24  print(df.isna().sum())
25
26  # 3. 라벨 인코딩
27  # 4. 데이터타입, 더미변환 (One-Hot Encoding)
28  # 5. 파생변수 생성
29  # 6. 정규화 또는 스케일 작업
30  # Boxplot for scaling check
31  sns.boxplot(df)
32  plt.tight_layout()
33  plt.show()
34
35  # 7. 데이터 분리
36  # 8. 모델 학습 - ML
37
38
39  # %%
40  mtcars = pd.read_csv('../../ADP_Python/data/mtcars.csv')
41  print(mtcars.shape)
42  print(mtcars.info())
43
44  for i, var in enumerate(mtcars.columns):
45      print(i, var, len(mtcars[var].unique()))
46
47  print(mtcars.describe())
48  print(mtcars.isna().sum())
49
50  df2_sample = mtcars[['mpg', 'am', 'cyl']]
51
52  from pandas.plotting import scatter_matrix
53
54  scatter_matrix(df2_sample)
55  plt.show()
56
57
```

```python
# %% 9. 모델 학습 - Stats
import scipy.stats as stats
import statsmodels.api as sm
import statsmodels.formula.api as smf

from statsmodels.stats.anova import anova_lm

df_A = df[df.target == 'Iris-setosa']['sepal width']
df_B = df[df.target == 'Iris-versicolor']['sepal width']
df_C = df[df.target == 'Iris-virginica']['sepal width']

# ANOVA 정규성 체크 Shapiro (위반 시 Kruskal)
print(stats.shapiro(df_A))
print(stats.shapiro(df_B))
print(stats.shapiro(df_C))

# ANOVA 등분산성 체크 Levene (위반 시 Welch)
print(stats.levene(df_A, df_B, df_C))

# one-way ANOVA using statsmodels
df_sample = df[['sepal width', 'target']]
df_sample.columns = ['sepal_width', 'target']

model = smf.ols(
    formula='sepal_width ~ 1 + target', data=df_sample
).fit()

print(model.summary())

# one-way ANOVA using scipy.stats
anova = stats.f_oneway(df_A, df_B, df_C)
print(anova)

# one-way ANOVA using pinguoin

# %% two-way ANOVA using statsmodels
model2 = smf.ols(
    formula='mpg ~ 1 + am*cyl', data=df2_sample
).fit()

print(model2.summary())
print(anova_lm(model2, typ=2))

sns.lineplot(data=df2_sample, x='cyl', y='mpg', hue='am')
plt.show()

# %% 10. 앙상블
# 11. 모델 평가
# one-way ANOVA post-hoc
from statsmodels.stats.multicomp import pairwise_tukeyhsd, MultiComparison

mc = MultiComparison(data=df['sepal width'], groups=df['target'])
tukeyhsd = mc.tukeyhsd(alpha=0.05)
tukeyhsd.plot_simultaneous()

print(tukeyhsd.summary())

from statsmodels.graphics.factorplots import interaction_plot
interaction_plot(x=df2_sample['cyl'], trace=df2_sample['am'], response=df2_sample['mpg'],
```

```
118                    colors=['red', 'blue'], markers=['D', 'o'])
119
120  plt.show()
121  # 12. 하이퍼파라미터 튜닝
122  # 13. 예측값 저장
123  # References
```