

ADP_Red\stats\4.regression.py

```

1
2 # 4.Regression
3 # %% 0. Import Libraries
4 import numpy as np
5 import pandas as pd
6 import seaborn as sns
7 import matplotlib.pyplot as plt
8
9 import scipy.stats as stats
10
11
12 # %% 1. Load Dataset
13 df = pd.read_csv('../ADP_Python/data/insurance.csv')
14
15 # Check dataset and set dependent variable(dv) and independent variable(iv)
16 print(f'{df.info()}')
17 dv = 'charges'
18 iv = 'age'
19
20
21 # %% 2. Data Visualization
22 sns.scatterplot(df, x=iv, y=dv)
23 plt.tight_layout
24 plt.show()
25
26
27 # %% 3. Statistical Test (Statsmodels)
28 import statsmodels.api as sm
29 import statsmodels.formula.api as smf
30
31 lm = smf.ols(
32     formula = f'{dv} ~ 1 + {iv}', data=df
33 ).fit()
34
35 lm.rsquared
36 lm.rsquared_adj
37 lm.params
38 lm.pvalues
39
40 print(lm.summary())
41 print(f'''
42     나이를 사용한 의료비용 예측을 위해 단순 선형 회귀 분석을 시행하였다.
43     모형의 p value는 {lm.f_pvalue:.2f}로 모형이 데이터를 잘 설명한다.
44     모형의 결정계수(R-squared)는 {lm.rsquared:.3f}, 수정결정계수(Adjusted R-squared)는
45     {lm.rsquared_adj:.3f}다.
46     {iv} 변수의 절편은 {lm.params.iloc[1]:.2f}, p value는 {lm.pvalues.iloc[1]:.2f}로
47     통계적으로 유의미하다.
48 ''')
49
50
51 # %% 4. Statistical Test (Scikit-learn SGDRegressor)
52 from sklearn.linear_model import SGDRegressor
53
54 df_iv = np.array(df[iv]).reshape(-1,1)
55 df_dv = np.array(df[dv]).reshape(-1,1)
56

```

```
57 lm = SGDRegressor(max_iter=1000, random_state=34)
58 lm.fit(df_iv, df_dv)
59
60 print(f'''
61     {lm.intercept_}
62     {lm.coef_}
63 ''')
64
65 # %% 5. Evaluation
66
67
68
```