**ADP_Red\ml\1.regression.py**

```python
1   # %% 5. Machine Learning - 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   # %% 1. 데이터 수집
10  df = pd.read_csv('../../ADP_Python/data/cereal.csv')
11
12  print(df.shape)
13  print(df.info())
14
15  # Check binary variable
16  for i, var in enumerate(df.columns):
17      print(i, var, len(df[var].unique()))
18
19  # Check data summary
20  print(df.describe())
21
22
23  # %% 2. 데이터 결측치 보정
24  print(df.isna().sum())
25
26  # # 결측치 제거
27  # missing = ['bill_length_mm', 'bill_depth_mm', 'flipper_length_mm', 'body_mass_g']
28
29  # for i in missing:
30  #     df[i] = df[i].fillna(df[i].median())
31  # df['sex'] = df['sex'].fillna('Male')
32
33
34  # %% 3. 라벨 인코딩
35  # from sklearn.preprocessing import LabelEncoder
36
37  # label = ['sex', 'smoker', 'region']
38
39  # df[label] = df[label].apply(LabelEncoder().fit_transform)
40  # # df['gender'] = np.where(df['class']=='M', 0, 1)
41  # # df['class'] = np.where(df['class']=='A', 1, 0)
42
43  # print(df.info())
44  # print(df.head())
45
46  # scatter matrix
47  from pandas.plotting import scatter_matrix
48
49  scatter_matrix(df)
50  plt.show
51
52  # %% 4. 데이터타입, 더미변환 (One-Hot Encoding)
53  # import pandas as pd
54
55  # category  = ['gender', 'class']
56  # for i in category:
57  #     df[i] = df[i].astype('category')
```

```python
58  # df = pd.get_dummies(df)
59  # df.head()
60
61
62  # %% 5. 파생변수 생성
63  # df['body_mass_g_qcut'] = pd.qcut(df['body_mass_g'], 5, labels=False)
64
65
66  # %% 6. 정규화 또는 스케일 작업
67  # from sklearn.preprocessing import StandardScaler, MinMaxScaler
68
69  # scaling_vars = ['age', 'bmi', 'children', 'charges']
70  # # scaler = StandardScaler()
71  # scaler = MinMaxScaler()
72  # scaler.fit(df[scaling_vars])
73
74  # df[scaling_vars] = scaler.transform(df[scaling_vars])
75
76  # Boxplot for scaling check
77  sns.boxplot(df)
78  plt.show()
79
80
81  # %% 7. 데이터 분리
82  from sklearn.model_selection import train_test_split
83
84  X = df.iloc[:, :-1]
85  y = df.iloc[:,-1]
86
87  X_train, X_test, y_train, y_test = train_test_split(
88      X, y, test_size=0.3, random_state=1,
89      # stratify=y
90      )
91
92  X_train = np.array(X_train)
93  X_test = np.array(X_test)
94  y_train = np.array(y_train)
95  y_test = np.array(y_test)
96
97
98  print('X_train: ', X_train.shape)
99  print('X_test: ', X_test.shape)
100 print('y_train: ', y_train.shape)
101 print('y_test: ', y_test.shape)
102
103
104 # %% 8. 모델 학습
105 from sklearn.linear_model import LinearRegression, LogisticRegression
106 from sklearn.preprocessing import PolynomialFeatures
107
108 model = LinearRegression()
109 # model = LogisticRegression()                                   # Logistic
    Regression
110 # model = LogisticRegression(multi_class='multinomial', solver='lbfgs')   # Softmax
    Regression
111 # model.fit(X_train, y_train)
112
113 # Polynomial Regresion
114 poly_reg = PolynomialFeatures(degree=2)
115 X_train_poly = poly_reg.fit_transform(X_train)
```

```python
116  model.fit(X_train_poly, y_train)
117
118  print(f'절편: {model.intercept_}, 기울기: {model.coef_}')
119
120  # %% 9. 모델 학습 (2)
121  from sklearn.linear_model import SGDRegressor
122
123  model2 = SGDRegressor(max_iter=1000)
124  model2.fit(X_train, y_train)
125
126  # %% 10. 앙상블
127
128  # %% 11. 모델 평가
129  from sklearn.metrics import mean_absolute_error, mean_squared_error
130  from sklearn.metrics import r2_score
131  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
     confusion_matrix, classification_report
132
133  # pred = model.predict(X_test)
134  # pred_proba = model.predict_proba(X_test)
135
136  X_test_poly = poly_reg.fit_transform(X_test)
137  pred = model.predict(X_test_poly)
138
139  print(f'MAE {mean_absolute_error(pred, y_test)}')
140  print(f'MSE {mean_squared_error(pred, y_test):.2f}')
141  print(f'RMSE {np.sqrt(mean_squared_error(pred, y_test)):.2f}')
142
143  # Metrics For Regression
144  print(f'R2 Score: {r2_score(pred, y_test):.2f}')
145
146  # Metrics For Classification
147  # print(f'혼동행렬: {confusion_matrix(pred, y_test)}')
148
149  # print(f'정확도: {accuracy_score(pred, y_test) * 100 :.2f} % ')
150  # print(f'정밀도: {precision_score(pred, y_test) * 100 :.2f} % ')
151  # print(f'재현율: {recall_score(pred, y_test) * 100 :.2f} % ')
152  # print(f'F1   : {f1_score(pred, y_test) * 100 :.2f} % ')
153
154  # ROC Curve
155  # from sklearn.metrics import RocCurveDisplay
156
157  # RocCurveDisplay.from_estimator(model, X_test, y_test)
158  # plt.show()
159
160
161  # %% 12. 하이퍼파라미터 튜닝
162  # from sklearn.model_selection import GridSearchCV
163
164  # parameters = {'n_estimators':[50,100], 'max_depth':[4,6]}
165  # model4 = RandomForestClassifier()
166  # clf = GridSearchCV(estimator=model4, param_grid=parameters, cv=3)
167  # clf.fit(X_train, y_train)
168
169  # print(f'Best Parameter: {clf.best_params_}')
170
171
172  # %% 13. 예측값 저장
173  # Save Output
174  output = pd.DataFrame({'id': y_test.index, 'pred': pred})
```

```python
175  output.to_csv('output.csv', index=False)
176
177  # Check Output
178  check = pd.read_csv('output.csv')
179  check.head()
180
181
182  # %% References
183  # - [[딥러닝] 로지스틱 회귀](https://circle-square.tistory.com/94)
184  # - [Logistic Regression in Python with statsmodels]
     (https://www.andrewvillazon.com/logistic-regression-python-statsmodels/)
```