**ADP_Red\ml\2.classification.py**

```python
1   # %% 5. Machine Learning - Classification
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/bodyPerformance.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  # %% 2. 데이터 결측치 보정
20  print(df.isna().sum())
21
22  # # 결측치 제거
23  # missing = ['bill_length_mm', 'bill_depth_mm', 'flipper_length_mm', 'body_mass_g']
24
25  # for i in missing:
26  #     df[i] = df[i].fillna(df[i].median())
27  # df['sex'] = df['sex'].fillna('Male')
28
29
30  # %% 3. 라벨 인코딩
31  from sklearn.preprocessing import LabelEncoder
32
33  label = ['gender', 'class']
34
35  df[label] = df[label].apply(LabelEncoder().fit_transform)
36  # df['gender'] = np.where(df['class']=='M', 0, 1)
37  # df['class'] = np.where(df['class']=='A', 1, 0)
38
39  print(df.info())
40  print(df.head())
41
42  # # %% data visualization
43  # fig, axes = plt.subplots(nrows=3, ncols=4, constrained_layout=True)
44
45  # for i, var in enumerate(df.columns):
46  #     row, col = i//4, i%4
47  #     sns.regplot(df, x=var, y='class',
48  #                 marker='o', ax=axes[row][col])
49
50  # plt.show()
51
52  # #
53  # from pandas.plotting import scatter_matrix
54
55  # scatter_matrix(df)
56  # plt.show
57
```

```python
# %% 4. 데이터타입, 더미변환 (One-Hot Encoding)
# import pandas as pd

# category  = ['gender', 'class']
# for i in category:
#     df[i] = df[i].astype('category')
# df = pd.get_dummies(df)
# df.head()

# %% 5. 파생변수 생성
# df['body_mass_g_qcut'] = pd.qcut(df['body_mass_g'], 5, labels=False)


# %% 6. 정규화 또는 스케일 작업
# from sklearn.preprocessing import StandardScaler, MinMaxScaler

# scaling_vars = ['age', 'height_cm', 'weight_kg', 'body fat_%', 'diastolic', 'systolic',
'gripForce', 'sit and bend forward_cm', 'sit-ups counts', 'broad jump_cm']
# scaler = StandardScaler()
# # scaler = MinMaxScaler()
# scaler.fit(df[scaling_vars])

# df[scaling_vars] = scaler.transform(df[scaling_vars])


# %% 7. 데이터 분리
from sklearn.model_selection import train_test_split

X = df.iloc[:, :-1]
y = df.iloc[:,-1]

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, stratify=y, random_state=1)

X_train = np.array(X_train)
X_test = np.array(X_test)
y_train = np.array(y_train)
y_test = np.array(y_test)


print('X_train: ', X_train.shape)
print('X_test: ', X_test.shape)
print('y_train: ', y_train.shape)
print('y_test: ', y_test.shape)


# %% 8. 모델 학습
from sklearn.linear_model import LogisticRegression

# model = LogisticRegression()                                    # Logistic Regression
model = LogisticRegression(multi_class='multinomial', solver='lbfgs')   # Softmax Regression
model.fit(X_train, y_train)


# %% 9. 모델 학습 (2)
# import statsmodels.api as sm
# import statsmodels.formula.api as smf
# dv = 'class'
```

```python
116  # model = smf.glm(
117  #     data=df, formula="class~age+height_cm+weight_kg",
118  #     family=sm.families.Binomial()
119  # ).fit()
120
121  # model.summary()
122
123  # %% 10. 앙상블
124
125  # %% 11. 모델 평가
126  from sklearn.metrics import mean_absolute_error, mean_squared_error
127  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
128  from sklearn.metrics import confusion_matrix, classification_report
129
130  pred = model.predict(X_test)
131  pred_proba = model.predict_proba(X_test)
132
133  print(f'MAE {mean_absolute_error(pred, y_test)}')
134  print(f'MSE {mean_squared_error(pred, y_test):.2f}')
135  print(f'RMSE {np.sqrt(mean_squared_error(pred, y_test)):.2f}')
136
137  print(f'혼동행렬: {confusion_matrix(pred, y_test)}')
138
139  print(f'정확도: {accuracy_score(pred, y_test) * 100 :.2f} % ')
140  # print(f'정밀도: {precision_score(pred, y_test) * 100 :.2f} % ')
141  # print(f'재현율: {recall_score(pred, y_test) * 100 :.2f} % ')
142  # print(f'F1   : {f1_score(pred, y_test) * 100 :.2f} % ')
143
144  # ROC Curve
145  # from sklearn.metrics import RocCurveDisplay
146
147  # RocCurveDisplay.from_estimator(model, X_test, y_test)
148  # plt.show()
149
150
151  # %% 12. 하이퍼파라미터 튜닝
152  # from sklearn.model_selection import GridSearchCV
153
154  # parameters = {'n_estimators':[50,100], 'max_depth':[4,6]}
155  # model4 = RandomForestClassifier()
156  # clf = GridSearchCV(estimator=model4, param_grid=parameters, cv=3)
157  # clf.fit(X_train, y_train)
158
159  # print(f'Best Parameter: {clf.best_params_}')
160
161
162  # %% 13. 예측값 저장
163  # Save Output
164  output = pd.DataFrame({'id': y_test.index, 'pred': pred})
165  output.to_csv('output.csv', index=False)
166
167  # Check Output
168  check = pd.read_csv('output.csv')
169  check.head()
170
171
172  # %% References
173  # - [[딥러닝] 로지스틱 회귀](https://circle-square.tistory.com/94)
174  # - [Logistic Regression in Python with statsmodels]
       (https://www.andrewvillazon.com/logistic-regression-python-statsmodels/)
```