

## ADP\_Red\stats\5.classification.py

```

1  # 5. Classification
2  # - Logistic/Softmax Regression
3  # - Sigmoid Function (Logistic Function)
4  # - Odds
5
6  # %% 5-1. Logistic/Softmax Regression (Scikit-Learn)
7  import numpy as np
8  from sklearn import datasets
9
10 # Load Dataset
11 iris = datasets.load_iris()
12 X = iris.data[:, [2,3]]
13 y = iris.target
14
15 # Train-Test Split
16 from sklearn.model_selection import train_test_split
17
18 X_train, X_test, y_train, y_test = train_test_split(
19     X, y, test_size=0.3, random_state=1, stratify=y    # stratification: 훈련데이터와 테스트
    # 데이터의 클래스 레이블 비율 동일하게
20 )
21
22 # Scaling
23 from sklearn.preprocessing import StandardScaler
24
25 sc = StandardScaler()
26 sc.fit(X_train)
27
28 X_train_std = sc.transform(X_train)
29 X_test_std = sc.transform(X_test)    # 훈련데이터의 mu와 sigma로 scaling
30
31 # Model Fitting
32 from sklearn.linear_model import LogisticRegression
33
34 lr = LogisticRegression(C=100.0, random_state=1)
35 lr.fit(X_train_std, y_train)
36
37 lr.predict_proba(X_test_std[:, :])
38 lr.predict(X_test_std[:, :])
39
40
41
42
43
44
45
46
47 # %%
48 # %% 1. 데이터 수집
49 import numpy as np
50 import pandas as pd
51 import seaborn as sns
52 import matplotlib.pyplot as plt
53
54 import scipy.stats as stats
55
56 # Load Dataset

```

```
57 df = pd.read_csv('../ADP_Python/data/bodyPerformance.csv')
58
59 print(df.shape)
60 print(df.info())
61
62 # %% 2. 데이터 결측치 보정
63 print(df.isna().sum())
64
65 # # 결측치 제거
66 # missing = ['bill_length_mm', 'bill_depth_mm', 'flipper_length_mm', 'body_mass_g']
67
68 # for i in missing:
69 #     df[i] = df[i].fillna(df[i].median())
70 # df['sex'] = df['sex'].fillna('Male')
71
72
73 # %% 3. 라벨 인코딩
74 from sklearn.preprocessing import LabelEncoder
75
76 label = ['gender', 'class']
77
78 df[label] = df[label].apply(LabelEncoder().fit_transform)
79
80
81 # %% 4. 데이터타입, 더미변환 (One-Hot Encoding)
82 # import pandas as pd
83
84 # category = ['gender', 'class']
85 # for i in category:
86 #     df[i] = df[i].astype('category')
87 # df = pd.get_dummies(df)
88
89
90 # %% 5. 파생변수 생성
91 # df['body_mass_g_qcut'] = pd.qcut(df['body_mass_g'], 5, labels=False)
92
93
94 # %% 6. 정규화 또는 스케일 작업
95 # from sklearn.preprocessing import StandardScaler, MinMaxScaler
96
97 # scaling_vars = ['bill_length_mm', 'bill_depth_mm', 'flipper_length_mm', 'body_mass_g']
98 # scaler = StandardScaler()
99 # # scaler = MinMaxScaler()
100 # scaler.fit(df[scaling_vars])
101
102 # df[scaling_vars] = scaler.transform(df[scaling_vars])
103
104
105 # %% 7. 데이터 분리
106 from sklearn.model_selection import train_test_split
107
108 X_train, X_test, y_train, y_test = train_test_split(
109     df.iloc[:, :-1], df.iloc[:, -1], test_size=0.3, stratify=df.iloc[:, -1], random_state=1)
110
111 X_train = np.array(X_train)
112 X_test = np.array(X_test)
113 y_train = np.array(y_train).reshape(-1,1)
114 y_test = np.array(y_test).reshape(-1,1)
115
116
```

```
117 print('X_train: ', X_train.shape)
118 print('X_test: ', X_test.shape)
119 print('y_train: ', y_train.shape)
120 print('y_test: ', y_test.shape)
121
122
123 # %% 8. 모델 학습
124 from sklearn.linear_model import LogisticRegression
125
126 lr = LogisticRegression(C=100.0, random_state=1)
127 lr.fit(X_train, y_train)
128
129 lr.predict_proba(X_test[:3, :])
130 lr.predict(X_test[:3, :])
131
132
133 # %% 11. 모델 평가
134 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
135 from sklearn.metrics import confusion_matrix, classification_report,
136
137 pred = lr.predict(X_test)
138
139 print(f'Model Accuracy {accuracy_score(y_test, pred)}')
140 print()
141
142 # %% 12. 하이퍼파라미터 튜닝
143 from sklearn.model_selection import GridSearchCV
144
145 parameters = {'n_estimators':[50,100], 'max_depth':[4,6]}
146 model4 = RandomForestClassifier()
147 clf = GridSearchCV(estimator=model4, param_grid=parameters, cv=3)
148 clf.fit(X_train, y_train)
149
150 print(f'Best Parameter: {clf.best_params_}')
151
152
153 # %% 13. 예측값 저장
154 # Save Output
155 output = pd.DataFrame({'id': y_test.index, 'pred': pred3})
156 output.to_csv('00300.csv', index=False)
157
158 # Check Output
159 check = pd.read_csv('00300.csv')
160 check.head()
161
```