

ADP_Red\ml\4.dt_ensemble.py

```

1  # %% 5. Machine Learning - Decision Tree & Ensemble
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  # %% 0. 합성 데이터
10 X = np.sort(5 * np.random.rand(400,1), axis=0)
11 T = np.linspace(0, 5, 500)[: , np.newaxis]
12 y = np.sin(X).ravel()
13
14 y[:,1] += 1 * (0.5 - np.random.rand(400))
15 plt.scatter(X, y, s=20, label='data')
16 plt.show()
17
18 # %% 1. 데이터 수집
19 df = pd.read_csv('../ADP_Python/data/credit_final.csv')
20
21 print(df.shape)
22 print(df.info())
23
24 # Check binary variable
25 for i, var in enumerate(df.columns):
26     print(i, var, len(df[var].unique()))
27
28 # Check data summary
29 print(df.describe())
30
31
32 # %% 2. 데이터 결측치 보정
33 print(df.isna().sum())
34
35 # # 결측치 제거
36 # missing = ['bill_length_mm', 'bill_depth_mm', 'flipper_length_mm', 'body_mass_g']
37
38 # for i in missing:
39 #     df[i] = df[i].fillna(df[i].median())
40 # df['sex'] = df['sex'].fillna('Male')
41
42
43 # %% 3. 라벨 인코딩
44 from sklearn.preprocessing import LabelEncoder
45
46 continuous_vars = ['credit.duration.months', 'credit.amount', 'age']
47 discrete_vars = ['account.balance',
48                 'previous.credit.payment.status', 'credit.purpose',
49                 'savings', 'employment.duration', 'installment.rate', 'marital.status',
50                 'guarantor', 'residence.duration', 'current.assets',
51                 'other.credits', 'apartment.type', 'bank.credits', 'occupation',
52                 'dependents', 'telephone', 'foreign.worker']
53
54 df[discrete_vars] = df[discrete_vars].apply(LabelEncoder().fit_transform)
55 # df['gender'] = np.where(df['class']=='M', 0, 1)
56 # df['class'] = np.where(df['class']=='A', 1, 0)
57

```

```
58 print(df.info())
59 print(df.head())
60
61 # #
62 # from pandas.plotting import scatter_matrix
63
64 # scatter_matrix(df)
65 # plt.show
66
67 # %% 4. 데이터타입, 더미변환 (One-Hot Encoding)
68 # import pandas as pd
69
70 # category = ['gender', 'class']
71 # for i in category:
72 #     df[i] = df[i].astype('category')
73 # df = pd.get_dummies(df)
74 # df.head()
75
76 # %% 5. 파생변수 생성
77 # df['body_mass_g_qcut'] = pd.qcut(df['body_mass_g'], 5, labels=False)
78
79
80 # %% 6. 정규화 또는 스케일 작업
81 from sklearn.preprocessing import StandardScaler, MinMaxScaler
82
83 scaling_vars = continuous_vars
84 # scaler = StandardScaler()
85 scaler = MinMaxScaler()
86 scaler.fit(df[scaling_vars])
87
88 df[scaling_vars] = scaler.transform(df[scaling_vars])
89
90 # Boxplot for scaling check
91 sns.boxplot(df)
92 plt.tight_layout()
93 plt.show()
94
95
96 # %% 7. 데이터 분리
97 from sklearn.model_selection import train_test_split
98
99 # X = df.iloc[:, 1:]
100 # y = df.iloc[:,0]
101
102 X_train, X_test, y_train, y_test = train_test_split(
103     X, y, test_size=0.3, random_state=1,
104     # stratify=y
105 )
106
107 X_train = np.array(X_train)
108 X_test = np.array(X_test)
109 y_train = np.array(y_train)
110 y_test = np.array(y_test)
111
112
113 print('X_train: ', X_train.shape)
114 print('X_test: ', X_test.shape)
115 print('y_train: ', y_train.shape)
116 print('y_test: ', y_test.shape)
117
```

```
118
119 # %% 8. 모델 학습
120 from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor
121 from sklearn.linear_model import LogisticRegression
122
123 # model = DecisionTreeClassifier(max_depth=5) # Decision
124 # Tree Classification
125 model = DecisionTreeRegressor(max_depth=5) # Decision Tree
126 # Regression
127 # model = LogisticRegression() # Logistic
128 # Regression
129 # model = LogisticRegression(multi_class='multinomial', solver='lbfgs') # Softmax
130 # Regression
131 model.fit(X_train, y_train)
132
133 # %% 9. 모델 학습 (2)
134 # (Decision Tree) Feature Importances
135 df_feature = pd.DataFrame([X.columns, model.feature_importances_]).T
136 df_feature.columns = ['feature_nm', 'importances']
137 print(df_feature)
138
139 # %% 10. 앙상블
140
141 # %% 11. 모델 평가
142 from sklearn.metrics import mean_absolute_error, mean_squared_error
143 from sklearn.metrics import r2_score
144 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
145 confusion_matrix, classification_report
146
147 pred = model.predict(X_test)
148 # pred_proba = model.predict_proba(X_test) # Classification
149
150 # X_test_poly = poly_reg.fit_transform(X_test) # Multinomial Regression
151 # pred = model.predict(X_test_poly)
152
153 print(f'MAE {mean_absolute_error(pred, y_test)}')
154 print(f'MSE {mean_squared_error(pred, y_test):.2f}')
155 print(f'RMSE {np.sqrt(mean_squared_error(pred, y_test)):.2f}')
156
157 # Metrics For Regression
158 print(f'R2 Score: {r2_score(pred, y_test):.2f}')
159
160 # # Metrics For Classification
161 # print(f'혼동행렬: {confusion_matrix(pred, y_test)}')
162
163 # print(f'정확도: {accuracy_score(pred, y_test) * 100 :.2f} % ')
164 # print(f'정밀도: {precision_score(pred, y_test) * 100 :.2f} % ')
165 # print(f'재현율: {recall_score(pred, y_test) * 100 :.2f} % ')
166 # print(f'F1 : {f1_score(pred, y_test) * 100 :.2f} % ')
167
168 # # Classification report
169 # report = classification_report(pred, y_test)
170 # print(report)
171
172 # # ROC Curve (For binary classification)
173 # from sklearn.metrics import RocCurveDisplay
174
175 # RocCurveDisplay.from_estimator(model, X_test, y_test)
```

```
174 # plt.show()
175
176
177 # %% 12. 하이퍼파라미터 튜닝
178 # from sklearn.model_selection import GridSearchCV
179
180 # parameters = {'n_estimators':[50,100], 'max_depth':[4,6]}
181 # model4 = RandomForestClassifier()
182 # clf = GridSearchCV(estimator=model4, param_grid=parameters, cv=3)
183 # clf.fit(X_train, y_train)
184
185 # print(f'Best Parameter: {clf.best_params_}')
186
187
188 # %% 13. 예측값 저장
189 # Save Output
190 output = pd.DataFrame({'id': y_test.index, 'pred': pred})
191 output.to_csv('output.csv', index=False)
192
193 # Check Output
194 check = pd.read_csv('output.csv')
195 check.head()
196
197
198
199 # %% References
200 # - [[딥러닝] 로지스틱 회귀](https://circle-square.tistory.com/94)
201 # - [Logistic Regression in Python with statsmodels]
  (https://www.andrewvillazon.com/logistic-regression-python-statsmodels/)
```