

```

1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import random
6  from sklearn.model_selection import train_test_split
7  from sklearn.preprocessing import RobustScaler, LabelEncoder
8  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
   classification_report
9  from sklearn.metrics import roc_curve, roc_auc_score
10 import tensorflow as tf
11 import matplotlib as mpl
12 mpl.rcParams['font.family'] = 'Malgun Gothic'
13 mpl.rcParams['axes.unicode_minus'] = False
14 import warnings
15 warnings.filterwarnings('ignore')
16
17 random.seed(42)
18 np.random.seed(42)
19 tf.random.set_seed(42)
20
21 # load data, 컬럼명 소문자로 일괄 변경
22 df = pd.read_csv("./titanic.csv")
23 df.columns = df.columns.str.lower()
24
25 # 불필요한 컬럼 제거
26 df.drop(columns=["class", "who"], inplace=True)
27 '''
28 # 특정 컬럼 제거(.drop)
29 df = df.drop(['col1', 'col2'], axis=1)
30 df.drop(columns=['col1', 'col2'], inplace=True)
31 '''
32
33 # 결측치 처리(제거, 대치)
34 df.dropna(ignore_index=True, inplace=True)
35 '''
36 # axis = 0 or 1 (행기준, 열기준)
37 df.dropna(axis=0, inplace=True) ; how='any'(default) or 'all' 지정가능
38 # 특정 컬럼들만 선택해서 결측치 제거
39 df.dropna(subset=['col1', 'col2'], inplace=True)
40 # 결측치 대치
41 df.fillna({'col_name':df['col_name'].median()}, inplace=True) # 결측치를 중앙값으로 대치.
42 '''
43
44 # 이상치 제거 removing outliers
45 num_cols = df.select_dtypes(include=[np.number]).columns.tolist()
46 cat_cols = [c for c in df.columns if c not in num_cols and c != "survived"]
47
48 for col in num_cols:
49     Q1 = df[col].quantile(0.25)
50     Q3 = df[col].quantile(0.75)
51     IQR = Q3 - Q1

```

```

52     lower_bound = Q1 - 1.5 * IQR
53     upper_bound = Q3 + 1.5 * IQR
54     df = df[ (df[col]>=lower_bound) & (df[col]<=upper_bound) ]
55     '''
56 이상치 제거하는 다른 방법:
57 # index 사용해서 drop
58 df = df.drop(df[df['col_name'] >= 10].index, axis=0)
59 # .clip 사용해서 (lower, upper) 제한(행을 삭제하지 않고 대치)
60 df = df.clip(-10, 200) 또는
61 df['col'] = df['col'].clip(-10, 200) ; or .clip(lower=, upper=)
62     '''
63
64 # feature engineering
65 df['family_size'] = df['sibsp'] + df['parch'] + 1
66
67 # 인코딩 : 레이블 값 대치
68 df['adult_male'] = df['adult_male'].map({True:1, False:0}).astype(int)
69 df['alone'] = df['alone'].map({True:1, False:0}).astype(int)
70 df['sex'] = df['sex'].map({'male':1, 'female':0}).astype(int)
71     '''
72 .map() 과 .replace() 는 다르다. map은 key를 찾지 못하면 NaN 을 반환한다.
73 unit 컬럼에 metres, feet 두 종류의 값이 있다면,
74 df['unit'] = df['unit'].replace({'feet':'metres'}) # 기존 metres 는 유지함.
75 df['unit'] = df['unit'].map({'feet':'metres'})      # 기존 metres 는 NaN 으로 변경됨.
76     '''
77
78 # 상관관계 출력
79 plt.figure(figsize=(5,5))
80 corr_mat = df[num_cols].corr()
81 sns.heatmap(corr_mat, annot=True)
82 plt.show()
83     '''
84 # heatmap 삼각형으로 출력
85 mask = np.zeros_like(corr_mat, dtype=np.bool)
86 mask[np.triu_indices_from(mask)] = True
87 sns.heatmap(corr_mat, mask=mask, annot=True)
88 # 클러스터맵
89 sns.clustermap(corr_matrix, annot=True)
90     '''
91
92 # 범주형 변수는 자동으로 1/0으로 인코딩
93 df = pd.get_dummies(data=df, drop_first=True).astype(int)
94
95 # 종속변수(범주형) 수치형으로 인코딩
96 le = LabelEncoder()
97 df['survived'] = le.fit_transform(df['survived'])
98     '''
99 # .map({ key:value,... }) 로 값 대치
100 df['col_name'] = df['col_name'].map({"Y":1, "N":0})
101
102 # to_categorical() 함수는 정수형(integer) 클래스 레이블(label)을 원-핫 인코딩(one-hot encoding) 벡터
    로 변환
103 from keras.utils import to_categorical
104 df[col] = to_categorical(df[col], num_classes=2) 또는 y = to_categorical(y, num_classes=2)

```

```

105 로 이진 분류도 출력층을 2개로 만들 수 있다.
106 --> 따라서, model.add(Dense(2, activation='softmax')) 로 해야 한다.
107
108 ex)
109     labels = [0, 1, 2, 1, 0]
110     one_hot_labels = to_categorical(labels)
111     print(one_hot_labels)
112
113     Out[:
114     [[1. 0. 0.]
115     [0. 1. 0.]
116     [0. 0. 1.]
117     [0. 1. 0.]
118     [1. 0. 0.]]
119     '''
120
121 y = df.pop('survived')
122 X = df
123
124 # split dataset
125 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
126
127 # 데이터셋 분리 후 스케일링
128 scaler = RobustScaler()
129 X_train = scaler.fit_transform(X_train)
130 X_test = scaler.transform(X_test)
131
132 from tensorflow.keras.models import Sequential
133 from tensorflow.keras.layers import Dense, Dropout, Input
134 from keras.utils import to_categorical
135
136 model = Sequential()
137 model.add(Input((X_train.shape[1],))) # keras recommend this
138 # model.add(Dense(32, activation='relu', input_shape=(X_train.shape[1],))), # input_dim =
139 # X_train.shape[1]
139 model.add(Dense(32, activation='relu'))
140 model.add(Dense(16, activation='relu'))
141 model.add(Dropout(rate=0.2))
142 model.add(Dense(8, activation='relu'))
143 model.add(Dropout(rate=0.2))
144 model.add(Dense(4, activation='relu'))
145 model.add(Dropout(rate=0.2))
146 # 출력층 시그모이드 함수 사용(이진분류문제)
147 model.add(Dense(1, activation='sigmoid'))
148 # 출력층이 2개이상인 경우 model.add(Dense(n, activation='softmax'))
149 '''
150 keras에서 activation 의 default 값은 없다. activation 함수의 종류:
151     linear :
152         입력과 동일한 값을 출력합니다. 결과적으로 입력에 가중치를 곱하여 모두 더한 값을 그대로 출력
153         합니다.
154     sigmoid :
155         시그모이드 함수를 사용하여 출력값을 0과 1 사이의 값으로 변환합니다. 이진 분류 모델 출력층에
156         많이 사용됩니다.
157     tanh :

```

```

156     하이퍼볼릭 탄젠트 함수를 사용하여 출력값을 -1과 1 사이의 값으로 변환합니다. 은닉층에 많이 사
    용됩니다.
157     relu :
158     Rectified Linear Unit 함수를 사용하여 출력값을 0 이상의 값으로 변환합니다. 은닉층에 많이 사
    용됩니다.
159     softmax :
160     softmax 함수를 사용하여 출력값을 다중 클래스 분류에 적합한 확률 값으로 변환합니다. 다중 클레
    스 분류 모델 출력층에 많이 사용됩니다.
161     selu :
162     Scaled Exponential Linear Unit 함수를 사용하여 출력값을 0과 1 사이로 스케일링합니다. 자기 정
    규화(Self-Normalizing) 효과로 인해 딥러닝 모델의 성능을 향상시킬 수 있습니다.
163     leaky_relu :
164     ReLU 함수의 단점을 보완하기 위해 제안된 활성화 함수입니다. ReLU 함수는 음수 입력에 대해 항상
    0을 출력하는데, 이는 뉴런이 죽는 문제(dying ReLU)를 야기할 수 있습니다. Leaky ReLU는 이러한 문제를
    해결하기 위해 음수 입력에 대해 아주 작은 기울기 (보통 0.01)를 부여합니다.
165     '''
166
167
168 # 모델 요약
169 model.summary()
170
171 # 모델 컴파일 (클래스가 2개인 2진 분류: 0 or 1)
172 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
173 '''
174 # loss = 'name' 으로 지정하는 기본값
175 # 회귀문제 : 'mean_squared_error'
176 # 2진 분류 : 'binary_crossentropy'
177 # 다중 클래스 분류
178     - 'categorical_crossentropy' : 원핫인코딩 했을 경우 (pd.get_dummies)
179     - 'sparse_categorical_crossentropy' : 원핫인코딩 안했을 경우 (LabelEncoder 정수 인코딩)
180 tf.keras.losses._ : 클래스로 지정
181     ex) tf.keras.losses.BinaryCrossentropy(from_logits=True)
182
183 # metrics 기본 옵션 : 'accuracy', 'binary_accuracy', 'mse'
184     다중 레이블 분류에서는 일반적인 accuracy가 성능을 과소평가 할 수 있음
185     - binary_accuracy: 레이블별 정확도를 평가
186     - precision, recall, f1_score: 클래스 불균형이나 특정 클래스(양성)에 초점을 맞출 때 유용
187     - auc: 클래스 분리 능력을 평가
188     tf.keras.metrics._ : precision, recall, f1_score, AUC 등 여러개 지정 가능
189     tf.keras.metrics.AUC(or Accuracy, BinaryAccuracy, CategoricalAccuracy, CosineSimilarity,
190         F1Score, Mean, MeanAbsoluteError, R2Score, Recall,
    RootMeanSquaredError,...)
191 # 클래스로 지정
192     optimizer = tf.keras.optimizers.SGD(learning_rate=0.01)
193     '''
194
195 # 모델 학습
196 # model.fit(X_train, y_train, epochs=30, batch_size=16, verbose=1)
197
198 # 모델 학습을 history 로 저장하려면,
199 from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint # point 는 소문자임에 주의
200
201 es = EarlyStopping(monitor='val_loss', patience=9, verbose=1)
202 mc = ModelCheckpoint('best_model.keras', monitor='val_loss', save_best_only=True, verbose=1)

```

```

203
204 history = model.fit(
205     X_train,
206     y_train,
207     validation_split=0.2,
208     # validation_data=(X_test, y_test),
209     epochs=30,
210     batch_size=16,
211     callbacks=[es, mc]
212 )
213
214 y_pred = (model.predict(X_test) > 0.5).astype(int)
215 # 시그모이드 출력층과 binary_crossentropy 손실 함수를 사용하는 이진 분류 DNN 모델
216 # model.predict(X_test) > 0.5는 각 샘플의 예측 확률이 0.5보다 크면 True
217 # .astype(int)는 boolean 값을 0 or 1 로 변환
218 '''
219 softmax 함수로 2개 output 출력되었을 경우,
220     y_pred = np.argmax(model.predict(X_test), axis=1) 로 큰 값을 저장한다. axis=1 은 컬럼방향 최대
    값 선택의 의미
221     np.argmax() : 행(axis=0) 또는 열(axis=1)을 따라 가장 큰 값(높은 확률)의 index 반환
222
223 다중 클래스 분류(하나의 클래스 선택)와 달리, 다중 레이블 분류(여러 클래스 동시 선택 가능)에서는
224 softmax 대신 클래스별로 독립적인 sigmoid를 사용하고, binary_crossentropy를 적용
225 '''
226
227 # 점수 출력
228
229 print("Accuracy:", accuracy_score(y_test, y_pred))
230 print("Precision:", precision_score(y_test, y_pred))
231 print("Recall:", recall_score(y_test, y_pred))
232 print("F1 Score:", f1_score(y_test, y_pred))
233 print("Classification Report", classification_report(y_test, y_pred)) # 모아서 출력
234
235 # 모델 성능 시각화
236 def plot_training_history(history):
237     plt.figure(figsize=(10, 4))
238
239     # 손실 값 (Loss)
240     plt.subplot(1, 2, 1)
241     plt.plot(history.history['loss'], label='Training Loss', marker='o')
242     plt.plot(history.history['val_loss'], label='Validation Loss', marker='o')
243     plt.title('Loss Over Epochs')
244     plt.xlabel('Epochs')
245     plt.ylabel('Loss')
246     plt.legend()
247     plt.grid(True)
248
249     # 정확도 (Accuracy)
250     plt.subplot(1, 2, 2)
251     plt.plot(history.history['accuracy'], label='Training Accuracy', marker='o')
252     plt.plot(history.history['val_accuracy'], label='Validation Accuracy', marker='o')
253     plt.title('Accuracy Over Epochs')
254     plt.xlabel('Epochs')
255     plt.ylabel('Accuracy')
256     plt.legend()

```

```

256     plt.grid(True)
257
258     plt.tight_layout()
259     plt.show()
260
261 def plot_roc_curve(y_test, y_pred_proba_pos):
262     fpr, tpr, _ = roc_curve(y_test, y_pred_proba_pos)
263     auc_score = roc_auc_score(y_test, y_pred_proba_pos) # ← 추가
264     plt.plot(fpr, tpr, label=f'ROC curve (AUC = {auc_score:.2f})')
265     plt.plot([0, 1], [0, 1], 'k--', label='Classifier')
266     plt.xlabel('False Positive Rate')
267     plt.ylabel('True Positive Rate')
268     plt.title('Receiver Operating Characteristic (ROC) Curve')
269     plt.legend(loc='lower right')
270     plt.show()
271
272
273 # 모델 학습 결과 시각화
274 plot_training_history(history)
275
276 # plot ROC curve with auc score
277 y_prob = model.predict(X_test).ravel()
278 plot_roc_curve(y_test, y_prob)
279 """
280 scikit-learn 추정기(분류기)에는 관례적으로 predict_proba가 있다.
281 tf.keras 모델은 model.predict가 “마지막 층의 출력값”을 그대로 반환한다.
282
283 마지막 층이 sigmoid(이진 분류)면 predict 결과가 곧 양성 확률.
284 마지막 층이 softmax(다중 분류)면 predict 결과가 각 클래스 확률 분포(합=1.0).
285
286 numpy.ravel(a, order='C') : Return a contiguous flattened array.
287 numpy.array.ravel()은 모양(차원)을 (N,1) → (N,)으로 축소함.
288 이진 분류에서 predict가 (샘플수, 1) 형태로 나오므로,
289 y_prob = model.predict(X_test).ravel()로 1차원 벡터로 바꿔
290 roc_auc_score, roc_curve 같은 sklearn 지표에 맞춥니다.
291 """
292
293 ##### y값 Imbalance Handling with SMOTE #####
294 # pip install -U imbalanced-learn
295 from imblearn.over_sampling import SMOTE
296 smote = SMOTE(random_state=42)
297 X_train_ovr, y_train_ovr = smote.fit_resample(X_train, y_train)
298
299 print("SMOTE 적용 전 train/test 데이터셋", X_train.shape, y_train.shape)
300 print("SMOTE 적용 후 train/test 데이터셋", X_train_ovr.shape, y_train_ovr.shape)
301
302 # SMOTE 적용 후 레이블값 분포 확인
303 pd.Series(y_train_ovr).value_counts()
304
305 # 다중클래스로 인코딩(softmax 함수를 출력층에 사용하려고)
306 y_train_ovr = to_categorical(y_train_ovr)
307
308 # build DNN model
309 model = Sequential()

```

```

310 model.add(Input((X_train_ovr.shape[1],))) # Keras recommend this
311 # model.add(Dense(32, activation='relu', input_shape=(X_train.shape[1],))), # input_dim =
    X_train.shape[1]
312 model.add(Dense(32, activation='relu')),
313 model.add(Dropout(rate=0.2)),
314 model.add(Dense(16, activation='relu')),
315 model.add(Dropout(rate=0.2)),
316 model.add(Dense(8, activation='relu')),
317 model.add(Dropout(rate=0.2)),
318 model.add(Dense(2, activation='softmax'))
319
320 # model compile
321 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
322 es = EarlyStopping(monitor='val_loss', patience=9, verbose=1)
323 mc = ModelCheckpoint('best_model_smote.keras', monitor='val_loss', save_best_only=True, verbose=1)
324
325 history = model.fit(
326     X_train_ovr,
327     y_train_ovr,
328     validation_split=0.2,
329     # validation_data=(X_test, y_test),
330     epochs=30,
331     batch_size=16,
332     callbacks=[es, mc]
333 )
334
335 # np.argmax() : 행(axis=0) 또는 열(axis=1)을 따라 가장 큰 값(높은 확률)의 index 반환
336 y_pred = np.argmax(model.predict(X_test), axis=1)
337
338 print("Accuracy:", accuracy_score(y_test, y_pred))
339 print("Precision:", precision_score(y_test, y_pred))
340 print("Recall:", recall_score(y_test, y_pred))
341 print("F1 Score:", f1_score(y_test, y_pred))
342 print("Classification Report", classification_report(y_test, y_pred)) # 모아서 출력
343
344 # 모델 학습 결과 시각화
345 plot_training_history(history)
346
347 # plot ROC curve with auc score
348 y_prob_pos = model.predict(X_test)[: ,1]
349 plot_roc_curve(y_test, y_prob_pos)

```