

AICE\DNN_2025_10_25수정.py

```

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.drop(columns=['col1', 'col2'], inplace=True)
30 df = df.drop(['col1', 'col2'], axis=1)
31 '''
32
33 # 결측치 처리(제거, 대치)
34 df.dropna(ignore_index=True, inplace=True)
35 '''
36 # axis = 0 or 1 (행기준, 열기준)
37 결측치 있는 모든 행 제거
38 df.dropna(axis=0, inplace=True) ; how='any'(default) or 'all' 지정가능
39 # 특정 컬럼들만 선택해서 결측치 제거
40 df.dropna(subset=['col1', 'col2'], inplace=True)
41 # 결측치 대치
42 df.fillna({'col_name':df['col_name'].median()}, inplace=True) # 결측치를 중앙값으로 대치.
43 또는
44 df['col_name'].fillna(df['col_name'].mean(), inplace=True)
45 '''
46
47 # 이상치 제거 removing outliers
48 num_cols = df.select_dtypes(include=[np.number]).columns.tolist()
49 cat_cols = [c for c in df.columns if c not in num_cols and c != "survived"]
50
51 for col in num_cols:

```

```

52     Q1 = df[col].quantile(0.25)
53     Q3 = df[col].quantile(0.75)
54     IQR = Q3 - Q1
55     lower_bound = Q1 - 1.5 * IQR
56     upper_bound = Q3 + 1.5 * IQR
57     df = df[ (df[col]>=lower_bound) & (df[col]<=upper_bound) ]
58     '''
59     이상치 제거하는 다른 방법:
60     # index 사용해서 drop
61     df = df.drop(df[df['col_name'] >= 10].index, axis=0)
62     # .clip 사용해서 (lower, upper) 제한(행을 삭제하지 않고 대치)
63     df = df.clip(-10, 200) 또는
64     df['col'] = df['col'].clip(-10, 200) ; or .clip(lower=, upper=)
65     '''
66
67     # feature engineering
68     df['family_size'] = df['sibsp'] + df['parch'] + 1
69
70     # 인코딩 : 레이블 값 대치
71     df['adult_male'] = df['adult_male'].map({True:1, False:0}).astype(int)
72     df['alone'] = df['alone'].map({True:1, False:0}).astype(int)
73     df['sex'] = df['sex'].map({'male':1, 'female':0}).astype(int)
74     '''
75     .map() 과 .replace() 는 다르다. map은 key를 찾지 못하면 NaN 을 반환한다.
76     unit 컬럼에 metres, feet 두 종류의 값이 있다면,
77     df['unit'] = df['unit'].replace({'feet':'metres'}) # 기존 metres 는 유지함.
78     df['unit'] = df['unit'].map({'feet':'metres'})      # 기존 metres 는 NaN 으로 변경됨.
79
80     (연습)
81     # 'shift' 컬럼 직접 인코딩
82     shift_map = {'주간': 0, '야간': 1}
83     defect_pre['shift'] = defect_pre['shift'].map(shift_map)
84     '''
85
86     # 상관관계 출력
87     plt.figure(figsize=(5,5))
88     corr_mat = df[num_cols].corr()
89     sns.heatmap(corr_mat, annot=True)
90     plt.show()
91     '''
92     # heatmap 삼각형으로 출력
93     mask = np.zeros_like(corr_mat, dtype=np.bool)
94     mask[np.triu_indices_from(mask)] = True
95     sns.heatmap(corr_mat, mask=mask, annot=True)
96
97     # 혼동행렬을 히트맵으로 시각화(25년 10월 기출)
98     from sklearn.metrics import confusion_matrix
99     cm = confusion_matrix(y_valid, y_pred)
100    sns.heatmap(cm, annot=True, fmt='d')
101
102    # lmpplot(25년 10월 기출)
103    산점도+회귀직선
104
105    # 클러스터맵

```

```

106 sns.clustermap(corr_matrix, annot=True)
107
108 # boxplot : Show distributions with respect to categories.
109 (25년 8월 기출)
110 sns.boxplot(data=defect_df, x='defect_status', y='temperature')
111
112 # barplot : Show point estimates and errors as rectangular bars.
113 # 범주형 데이터별로 집계(statistical aggregate, default: mean) 값을 막대(bar)로 표현
114 # defaults: estimator='mean', errorbar=(ci, 95)
115 sns.barplot(data=df, x='class', y='fare', hue='sex', errorbar=None, estimator='median')
116 '''
117
118 # 범주형 변수는 자동으로 1/0으로 인코딩된다.
119 df = pd.get_dummies(data=df, drop_first=True).astype(int)
120
121 '''
122 (연습)
123 # 교차표 생성 (비율 기준)
124 cross_table_ratio = pd.crosstab(defect_df['production_line'], defect_df['defect_status'],
125                                 normalize='index')
126
127 # 'shift'로 그룹화하여 주요 측정값들의 평균을 계산합니다.
128 shift_avg = defect_df.groupby('shift')[['measurement_A', 'measurement_B', 'measurement_C']].mean()
129
130
131 # 종속변수(범주형) 수치형으로 인코딩
132 le = LabelEncoder()
133 df['survived'] = le.fit_transform(df['survived'])
134 '''
135 # .map({ key:value,... }) 로 값 대치
136 df['col_name'] = df['col_name'].map({"Y":1, "N":0})
137
138 # to_categorical() 함수는 정수형(integer) 클래스 레이블(label)을 원-핫 인코딩(one-hot encoding) 벡터
139 로 변환
140 from keras.utils import to_categorical
141 df[col] = to_categorical(df[col], num_classes=2) 또는 y = to_categorical(y, num_classes=2)
142 로 이진 분류도 출력층을 2개로 만들 수 있다.
143 --> 이 경우 model.add(Dense(2, activation='softmax')) 로 해야한다.
144
145 ex)
146 labels = [0, 1, 2, 1, 0]
147 one_hot_labels = to_categorical(labels)
148 print(one_hot_labels)
149
150 Out[:
151 [[1. 0. 0.]
152  [0. 1. 0.]
153  [0. 0. 1.]
154  [0. 1. 0.]
155  [1. 0. 0.]]
156
157 '''
158 y = df.pop('survived')

```

```

158 X = df
159
160 # split dataset
161 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
162
163 # 데이터셋 분리 후 스케일링
164 scaler = RobustScaler()
165 X_train = scaler.fit_transform(X_train)
166 X_test = scaler.transform(X_test)
167
168 '''
169 다중클래스 분류 문제라면 출력층을 softmax 함수로 변경하고,
170 y_train, y_test 를 to_categorical() 로 인코딩 해야한다.
171 '''
172
173 from tensorflow.keras.models import Sequential
174 from tensorflow.keras.layers import Dense, Dropout, Input
175 from keras.utils import to_categorical
176
177 model = Sequential()
178 model.add(Input((X_train.shape[1],))) # keras recommend this
179 # model.add(Dense(32, activation='relu', input_shape=(X_train.shape[1],))), # input_dim =
180 # X_train.shape[1]
181 model.add(Dense(32, activation='relu'))
182 model.add(Dense(16, activation='relu'))
183 model.add(Dropout(rate=0.2))
184 model.add(Dense(8, activation='relu'))
185 model.add(Dropout(rate=0.2))
186 model.add(Dense(4, activation='relu'))
187 model.add(Dropout(rate=0.2))
188 # 출력층 시그모이드 함수 사용(이진분류문제)
189 model.add(Dense(1, activation='sigmoid'))
190 # 출력층이 2개이상인 경우 model.add(Dense(n, activation='softmax'))
191 '''
192 또는
193 model = Sequential([
194     Dense(32, activation='relu', input_shape=(X_train.shape[1],)),
195     Dense(16, activation='relu'),
196     Dropout(rate=0.2),
197     ...
198 ])
199
200 keras에서 activation 의 default 값은 없다. activation 함수의 종류:
201     linear :
202         입력과 동일한 값을 출력합니다. 결과적으로 입력에 가중치를 곱하여 모두 더한 값을 그대로 출력
203         합니다.
204     sigmoid :
205         시그모이드 함수를 사용하여 출력값을 0과 1 사이의 값으로 변환합니다. 이진 분류 모델 출력층에
206         많이 사용됩니다.
207     tanh :
208         하이퍼볼릭 탄젠트 함수를 사용하여 출력값을 -1과 1 사이의 값으로 변환합니다. 은닉층에 많이 사
209         용됩니다.
210     relu :

```

```

208         Rectified Linear Unit 함수를 사용하여 출력값을 0 이상의 값으로 변환합니다. 은닉층에 많이 사
        용됩니다.
209         softmax :
210         softmax 함수를 사용하여 출력값을 다중 클래스 분류에 적합한 확률 값으로 변환합니다. 다중 클래
        스 분류 모델 출력층에 많이 사용됩니다.
211         selu :
212         Scaled Exponential Linear Unit 함수를 사용하여 출력값을 0과 1 사이로 스케일링합니다. 자기 정
        규화(Self-Normalizing) 효과로 인해 딥러닝 모델의 성능을 향상시킬 수 있습니다.
213         leaky_relu :
214         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.
215     '''
216
217
218     # 모델 요약
219     model.summary()
220
221     # 모델 컴파일 (클래스가 2개인 2진 분류: 0 or 1)
222     model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
223     '''
224     # loss = 'name' 으로 지정하는 기본값
225     # 회귀문제 : 'mean_squared_error'
226     # 2진 분류 : 'binary_crossentropy'
227     # 다중 클래스 분류
228         - 'categorical_crossentropy' : 원핫인코딩 했을 경우 (pd.get_dummies)
229         - 'sparse_categorical_crossentropy' : 원핫인코딩 안했을 경우 (LabelEncoder 정수 인코딩)
230
231     # metrics 기본 옵션 : 'accuracy', 'binary_accuracy', 'mse'
232     다중 레이블 분류에서는 일반적인 accuracy가 성능을 과소평가 할 수 있음
233         - binary_accuracy: 레이블별 정확도를 평가
234         - precision, recall, f1_score: 클래스 불균형이나 특정 클래스(양성)에 초점을 맞출 때 유용
235         - auc: 클래스 분리 능력을 평가
236     '''
237
238     # 모델 학습
239     # model.fit(X_train, y_train, epochs=30, batch_size=16, verbose=1)
240
241     # 모델 학습을 history 로 저장하려면,
242     from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint      # point 는 소문자임에 주의
243
244     es = EarlyStopping(monitor='val_loss', patience=9, verbose=1)
245     mc = ModelCheckpoint('best_model.keras', monitor='val_loss', save_best_only=True, verbose=1)
246
247     history = model.fit(
248         X_train,
249         y_train,
250         validation_split=0.2,
251         # validation_data=(X_test, y_test),
252         epochs=30,
253         batch_size=16,
254         callbacks=[es, mc]
255     )
256

```

```

257 y_pred = (model.predict(X_test) > 0.5).astype(int)
258 '''
259 시그모이드 출력층과 binary_crossentropy 손실 함수를 사용하는 이진 분류 DNN 모델이다.
260 model.predict(X_test) > 0.5는 각 샘플의 예측 확률이 0.5보다 크면 True
261 .astype(int)는 boolean 값을 0 or 1 로 변환
262 예를 들어 predict 값이 [0.85, 0.11] 이면 [0] 이 반환됨. 즉, Negative 라는 의미로 해석
263
264 softmax 함수로 2개 output 출력되었을 경우,
265 y_pred = np.argmax(model.predict(X_test), axis=1) 로 큰 값을 저장한다. axis=1 은 컬럼방향 최대값 선택의 의미
266 np.argmax() : 행(axis=0) 또는 열(axis=1)을 따라 가장 큰 값(높은 확률)의 index 반환
267 예를 들어 [[0.15, 0.85], [0.89, 0.11] ] 이면 [1, 0] 이 반환됨. 즉, [Positive, Negative] 라는 의미로 해석
268
269 다중 클래스 분류(하나의 클래스 선택)와 달리, 다중 레이블 분류(여러 클래스 동시 선택 가능)에서는
270 softmax 대신 클래스별로 독립적인 sigmoid를 사용하고, binary_crossentropy를 적용
271 '''
272
273 # 점수 출력
274
275 print("Accuracy:", accuracy_score(y_test, y_pred))
276 print("Precision:", precision_score(y_test, y_pred))
277 print("Recall:", recall_score(y_test, y_pred))
278 print("F1 Score:", f1_score(y_test, y_pred))
279 print("Classification Report", classification_report(y_test, y_pred)) # 모아서 출력
280
281 # 모델 성능 시각화
282 def plot_training_history(history):
283     plt.figure(figsize=(10, 4))
284
285     # 손실 값 (Loss)
286     plt.subplot(1, 2, 1)
287     plt.plot(history.history['loss'], label='Training Loss', marker='o')
288     plt.plot(history.history['val_loss'], label='Validation Loss', marker='o')
289     plt.title('Loss Over Epochs')
290     plt.xlabel('Epochs')
291     plt.ylabel('Loss')
292     plt.legend()
293     plt.grid(True)
294
295     # 정확도 (Accuracy)
296     plt.subplot(1, 2, 2)
297     plt.plot(history.history['accuracy'], label='Training Accuracy', marker='o')
298     plt.plot(history.history['val_accuracy'], label='Validation Accuracy', marker='o')
299     plt.title('Accuracy Over Epochs')
300     plt.xlabel('Epochs')
301     plt.ylabel('Accuracy')
302     plt.legend()
303     plt.grid(True)
304
305     plt.tight_layout()
306     plt.show()
307
308 def plot_roc_curve(y_test, y_pred_proba_pos):
309     fpr, tpr, _ = roc_curve(y_test, y_pred_proba_pos)

```

```

309     auc_score = roc_auc_score(y_test, y_pred_proba_pos) # ← 추가
310     plt.plot(fpr, tpr, label=f'ROC curve (AUC = {auc_score:.2f})')
311     plt.plot([0, 1], [0, 1], 'k--', label='Classifier')
312     plt.xlabel('False Positive Rate')
313     plt.ylabel('True Positive Rate')
314     plt.title('Receiver Operating Characteristic (ROC) Curve')
315     plt.legend(loc='lower right')
316     plt.show()
317
318
319 # 모델 학습 결과 시각화
320 plot_training_history(history)
321
322 # plot ROC curve with auc score
323 y_prob = model.predict(X_test).ravel()
324 plot_roc_curve(y_test, y_prob)
325 """
326 sklearn의 분류기(LogisticRegression, RandomForestClassifier 등)에는 .predict_proba()가 있지만,
327 Keras 모델에는 .predict_proba() 메서드는 존재하지 않음. tf.keras 모델은 model.predict가
328 “마지막 층의 출력값”을 그대로 반환한다.
329     y_prob = model.predict(X_test)          # 각 클래스별 확률 (softmax 출력)
330     y_pred = np.argmax(y_prob, axis=1)      # 확률이 가장 큰 클래스 선택
331
332 마지막 층이 sigmoid(이진 분류)면 predict 결과가 곧 양성 확률.
333 마지막 층이 softmax(다중 분류)면 predict 결과가 각 클래스 확률 분포(합=1.0).
334
335 numpy.ravel(a, order='C') : Return a contiguous flattened array.
336 numpy.array.ravel()은 모양(차원)을 (N,1) → (N,)으로 축소함.
337 이진 분류에서 predict가 (샘플수, 1) 형태로 나오므로,
338 y_prob = model.predict(X_test).ravel()로 1차원 벡터로 바꿔
339 roc_auc_score, roc_curve 같은 sklearn 지표에 맞춥니다.
340 """
341
342 ##### y값 Imbalance Handling with SMOTE #####
343 # pip install -U imbalanced-learn
344 from imblearn.over_sampling import SMOTE
345 smote = SMOTE(random_state=42)
346 X_train_ovr, y_train_ovr = smote.fit_resample(X_train, y_train)
347
348 print("SMOTE 적용 전 train/test 데이터셋", X_train.shape, y_train.shape)
349 print("SMOTE 적용 후 train/test 데이터셋", X_train_ovr.shape, y_train_ovr.shape)
350
351 # SMOTE 적용 후 레이블값 분포 확인
352 pd.Series(y_train_ovr).value_counts()
353
354 # 다중클래스로 인코딩(softmax 함수를 출력층에 사용하려고)
355 y_train_ovr = to_categorical(y_train_ovr)
356
357 # build DNN model
358 model = Sequential()
359 model.add(Input((X_train_ovr.shape[1],))) # Keras recommend this
360 # model.add(Dense(32, activation='relu', input_shape=(X_train.shape[1],))), # input_dim =
361 # X_train.shape[1]
361 model.add(Dense(32, activation='relu')),

```

```
362 model.add(Dropout(rate=0.2)),
363 model.add(Dense(16, activation='relu')),
364 model.add(Dropout(rate=0.2)),
365 model.add(Dense(8, activation='relu')),
366 model.add(Dropout(rate=0.2)),
367 model.add(Dense(2, activation='softmax'))
368
369 # model compile
370 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
371 es = EarlyStopping(monitor='val_loss', patience=9, verbose=1)
372 mc = ModelCheckpoint('best_model_smote.keras', monitor='val_loss', save_best_only=True, verbose=1)
373
374 history = model.fit(
375     X_train_ovr,
376     y_train_ovr,
377     validation_split=0.2,
378     # validation_data=(X_test, y_test),
379     epochs=30,
380     batch_size=16,
381     callbacks=[es, mc]
382 )
383
384 # np.argmax() : 행(axis=0) 또는 열(axis=1)을 따라 가장 큰 값(높은 확률)의 index 반환
385 y_pred = np.argmax(model.predict(X_test), axis=1)
386
387 print("Accuracy:", accuracy_score(y_test, y_pred))
388 print("Precision:", precision_score(y_test, y_pred))
389 print("Recall:", recall_score(y_test, y_pred))
390 print("F1 Score:", f1_score(y_test, y_pred))
391 print("Classification Report", classification_report(y_test, y_pred)) # 모아서 출력
392
393 # 모델 학습 결과 시각화
394 plot_training_history(history)
395
396 # plot ROC curve with auc score
397 y_prob_pos = model.predict(X_test)[:,-1]
398 plot_roc_curve(y_test, y_prob_pos)
```