

## ✓ 2.6장 실습 - 타이타닉 생존자 예측

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 %matplotlib inline
6
7 df = pd.read_csv('/content/train.csv')
8 df.head(3)

```



	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282

Next steps: [View recommended plots](#)

```

1 # 데이터 칼럼 타입 확인
2 print("학습 데이터 정보 \n")
3 df.info()

```

학습 데이터 정보

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
#   Column          Non-Null Count  Dtype
---  -
0   PassengerId     891 non-null   int64
1   Survived        891 non-null   int64
2   Pclass          891 non-null   int64
3   Name            891 non-null   object
4   Sex             891 non-null   object
5   Age            714 non-null   float64
6   SibSp           891 non-null   int64
7   Parch          891 non-null   int64
8   Ticket         891 non-null   object

```

```

9   Fare      891 non-null   float64
10  Cabin     204 non-null   object
11  Embarked  889 non-null   object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB

```

```

1 # Null 값 처리
2 df['Age'].fillna(df['Age'].mean(), inplace = True)
3 df['Cabin'].fillna('N', inplace = True)
4 df['Embarked'].fillna('N', inplace = True)
5
6 print('데이터 세트 Null 값 개수: ', df.isnull().sum().sum())

```

데이터 세트 Null 값 개수: 0

```

1 # 피처의 값 분류 살펴보기
2 print('Sex 값 분포: \n', df['Sex'].value_counts())
3 print('\n Cabin 값 분포: \n', df['Cabin'].value_counts())
4 print('\n Embarked 값 분포: \n', df['Embarked'].value_counts())

```

```

Sex 값 분포:
male      577
female    314
Name: Sex, dtype: int64

```

```

Cabin 값 분포:
N      687
C23 C25 C27      4
G6      4
B96 B98      4
C22 C26      3
...
E34      1
C7      1
C54      1
E36      1
C148     1
Name: Cabin, Length: 148, dtype: int64

```

```

Embarked 값 분포:
S      644
C      168
Q       77
N        2
Name: Embarked, dtype: int64

```

```

1 # Cabin 속성 앞 글자만 추출
2 df['Cabin'] = df['Cabin'].str[:1]
3 print(df['Cabin'].head(3))

```

```

0   N
1   C
2   N
Name: Cabin, dtype: object

```

```
1 # 성별에 따른 생존자의 수 비교
```

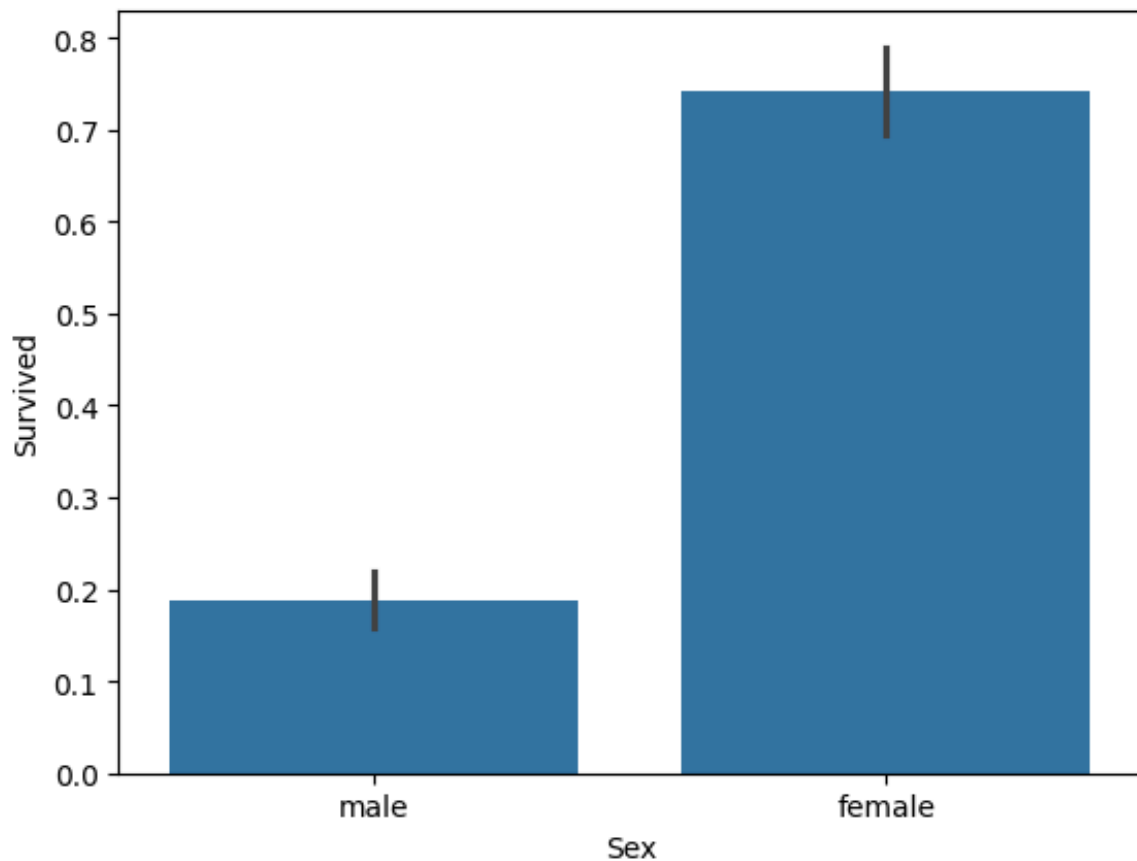
```
2 df.groupby(['Sex', 'Survived'])['Survived'].count()
```

```
Sex      Survived
female  0           81
        1          233
male    0          468
        1          109
Name: Survived, dtype: int64
```

```
1 # 분석 결과 시각화 - 성별
```

```
2 sns.barplot(x = 'Sex', y = 'Survived', data = df)
```

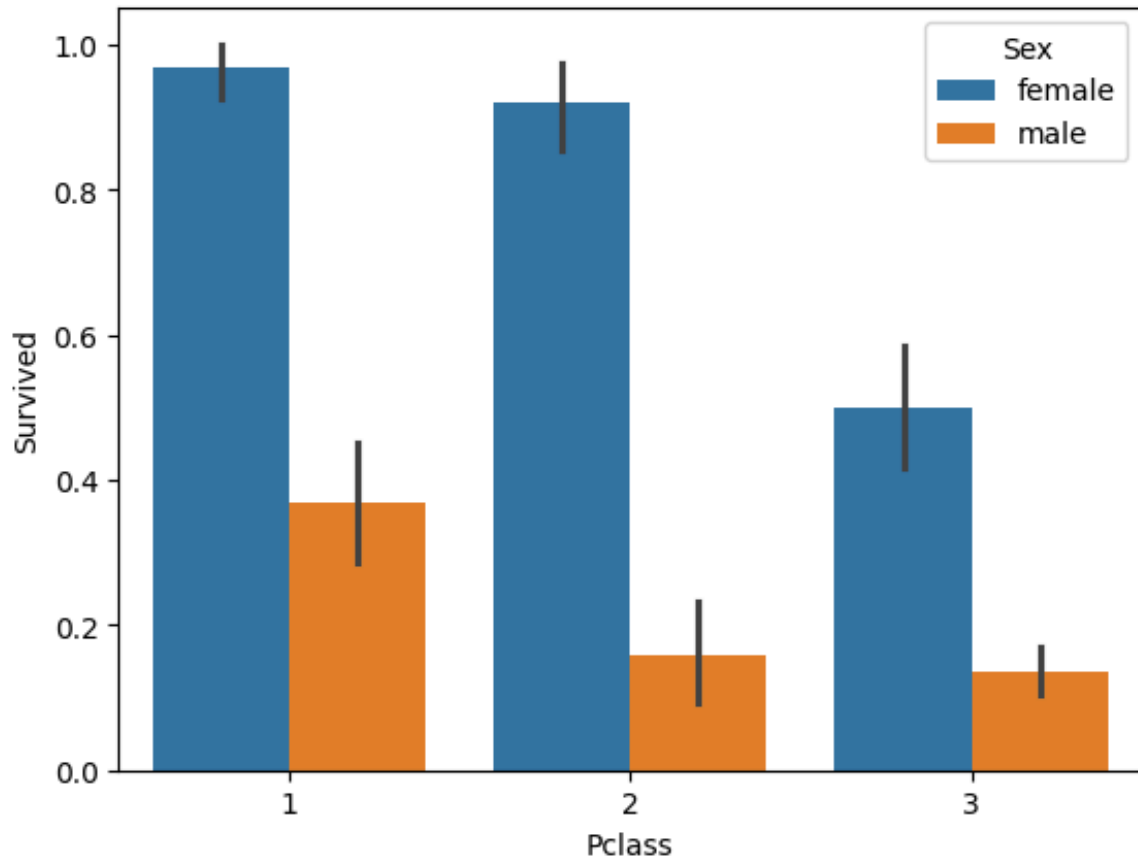
```
<Axes: xlabel='Sex', ylabel='Survived'>
```



```
1 # 분석 결과 시각화 - 부자 + 성별
```

```
2 sns.barplot(x = 'Pclass', y = 'Survived', hue = 'Sex', data = df)
```

<Axes: xlabel='Pclass', ylabel='Survived'>



```

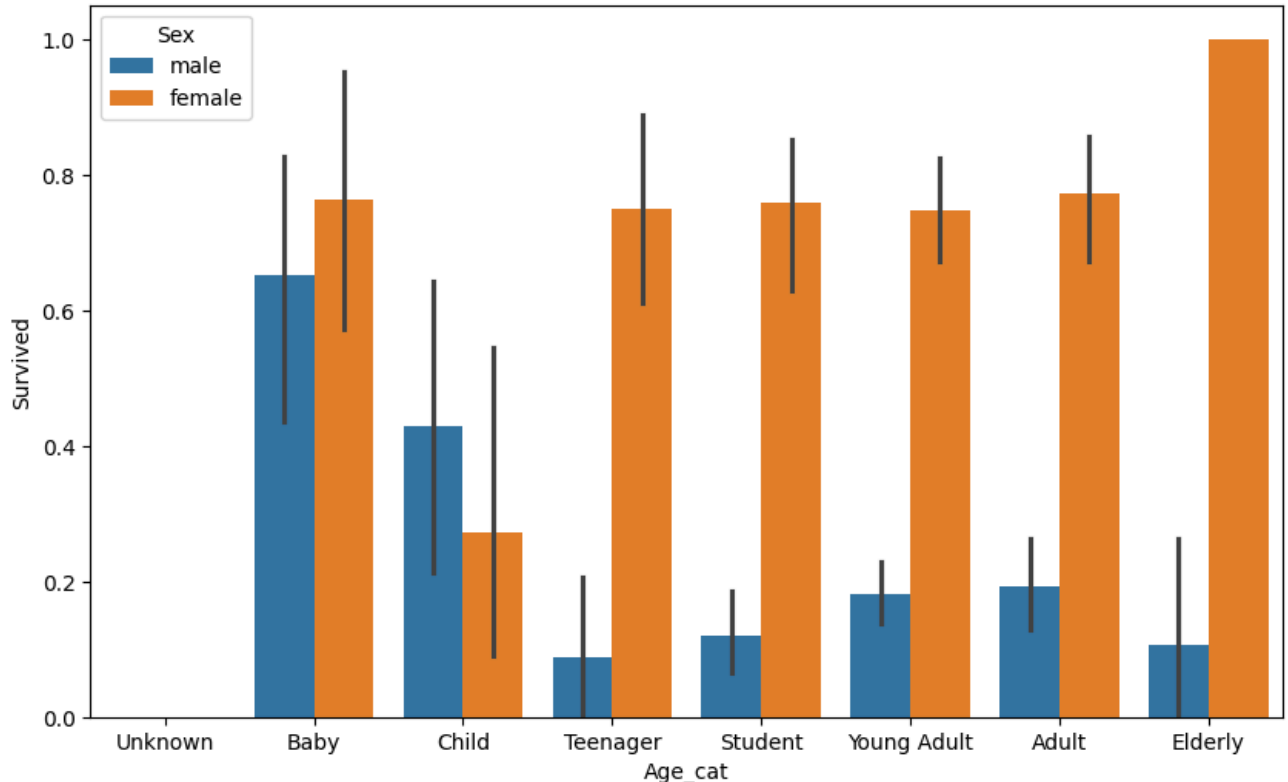
1 # Age에 따른 생존 확률 비교
2 # 값을 범위별로 분류, 카테고리 값 할당
3 def get_category(age):
4     cat = ''
5     if age <= -1:
6         cat = 'Unknown'
7     elif age <= 5:
8         cat = 'Baby'
9     elif age <= 12:
10        cat = 'Child'
11    elif age <= 18:
12        cat = 'Teenager'
13    elif age <= 25:
14        cat = 'Student'
15    elif age <= 35:
16        cat = 'Young Adult'
17    elif age <= 60:
18        cat = 'Adult'
19    else:
20        cat = 'Elderly'
21    return cat

```

```

1 plt.figure(figsize = (10, 6))
2
3 group_names = ['Unknown', 'Baby', 'Child', 'Teenager', 'Student', 'Young Adu
4
5 df['Age_cat'] = df['Age'].apply(lambda x: get_category(x))
6 sns.barplot(x = 'Age_cat', y = 'Survived', hue = 'Sex', data = df, order = c
7 df.drop('Age_cat', axis = 1, inplace = True)

```



여자 Baby의 경우 비교적 생존 확률이 높았다.

여자 Child의 경우 비교적 생존 확률이 낮았다.

여자 Elderly의 경우 생존 확률이 매우 높았다.

=> Sex, Age, PClass가 생존 여부를 결정하는 중요한 요소로 작용함을 확인.

```

1 # 문자열 피쳐 -> 숫자형 카테고리 피쳐로 변환
2 # 인코딩 - 레이블 인코더 객체 사용
3 from sklearn import preprocessing
4
5 def encode_features(dataDF):
6     features = ['Cabin', 'Sex', 'Embarked']
7     for feature in features:
8         le = preprocessing.LabelEncoder()
9         le = le.fit(dataDF[feature])
10        dataDF[feature] = le.transform(dataDF[feature])
11    return dataDF
12
13 df = encode_features(df)
14 df.head()

```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
0	1	0	3	Braund, Mr. Owen Harris	1	22.0	1	0	A/5 21171
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	0	38.0	1	0	PC 17599 7
2	3	1	3	Heikkinen, Miss. Laina	0	26.0	0	0	STON/O2. 3101282
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	0	35.0	1	0	113803 5
4	5	0	3	Allen, Mr. William Henry	1	35.0	0	0	373450 8

Next steps: [View recommended plots](#)

**transform\_features( )**: 데이터의 전처리를 전체적으로 호출하는 함수

- Null 처리
- 포매팅
- 인코딩

```
1 # Null 처리 함수
2 def fillna(df):
3     df['Age'].fillna(df['Age'].mean(), inplace=True)
4     df['Cabin'].fillna('N', inplace=True)
5     df['Embarked'].fillna('N', inplace=True)
6     df['Fare'].fillna(0, inplace=True)
7     return df
8
9 # 머신러닝 알고리즘에 불필요한 속성 제거
10 def drop_features(df):
11     df.drop(['PassengerId', 'Name', 'Ticket'], axis=1, inplace=True)
12     return df
13
14 # 레이블 인코딩 수행.
15 def format_features(df):
16     df['Cabin'] = df['Cabin'].str[:1]
17     features = ['Cabin', 'Sex', 'Embarked']
18     for feature in features:
19         le = preprocessing.LabelEncoder()
20         le = le.fit(df[feature])
21         df[feature] = le.transform(df[feature])
22     return df
23
24 # 앞에서 설정한 데이터 전처리 함수 호출
25 def transform_features(df):
26     df = fillna(df)
27     df = drop_features(df)
28     df = format_features(df)
29     return df


1 df = pd.read_csv('/content/train.csv')
2 y_df = df['Survived']
3 x_df = df.drop('Survived', axis=1)
4
5 x_df = transform_features(x_df)


1 # 테스트 데이터 추출
2 from sklearn.model_selection import train_test_split
3 x_train, x_test, y_train, y_test = train_test_split(x_df, y_df, test_size=0.

1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score
```

```

1 # 결정 트리, RandomForest, 로지스틱 회귀를 위한 사이킷런 Classifier 클래스 생성
2 dt_clf = DecisionTreeClassifier(random_state=11)
3 rf_clf = RandomForestClassifier(random_state=11)
4 lr_clf = LogisticRegression()
5
6 # DecisionTreeClassifier 학습/예측/평가
7 dt_clf.fit(x_train, y_train)
8 dt_pred = dt_clf.predict(x_test)
9 print('DecisionTreeClassifier 정확도: {0:.4f}'.format(accuracy_score(y_test, dt_pred)))
10
11 # RandomForestClassifier 학습/예측/평가
12 rf_clf.fit(x_train, y_train)
13 rf_pred = rf_clf.predict(x_test)
14 print('RandomForestClassifier 정확도: {0:.4f}'.format(accuracy_score(y_test, rf_pred)))
15
16 # LogisticRegression 학습/예측/평가
17 lr_clf.fit(x_train, y_train)
18 lr_pred = lr_clf.predict(x_test)
19 print('LogisticRegression 정확도 : {0:.4f}'.format(accuracy_score(y_test, lr_pred)))

```

DecisionTreeClassifier 정확도: 0.7877

RandomForestClassifier 정확도: 0.8547

LogisticRegression 정확도 : 0.8492

/usr/local/lib/python3.10/dist-packages/sklearn/linear\_model/\_logistic.py:4

STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

n\_iter\_i = \_check\_optimize\_result(

Logistic Regression이 높은 정확도를 보이고 있다.

그러나 아직 최적화 작업을 수행하지 않았고 데이터 양도 충분하지 않기에 평가하기에는 조금 이릅니다.



```

1 # 교차 검증으로 결정 트리 모델 평가하기
2 from sklearn.model_selection import KFold
3 def exec_kfold(clf, folds=5):
4     # 폴드 세트를 5개인 KFold 객체를 생성, 폴드 수만큼 예측 결과 저장을 위한 리스트 객체 생성.
5     kfold = KFold(n_splits=folds)
6     scores = []
7
8     # KFold 교차검증 수행.
9     for iter_count, (train_index, test_index) in enumerate(kfold.split(x_df)):
10    # x_df 데이터에서 교차검증별로 학습과 검증 데이터를 가리키는 index 생성
11        x_train, x_test = x_df.values[train_index], x_df.values[test_index]
12        y_train, y_test = y_df.values[train_index], y_df.values[test_index]
13        # Classifier 학습, 예측, 정확도 계산
14        clf.fit(x_train, y_train)
15        predictions = clf.predict(x_test)
16        accuracy = accuracy_score(y_test, predictions)
17        scores.append(accuracy)
18        print('교차 검증 {0} 정확도 : {1:.4f}'.format(iter_count, accuracy))
19 # 5개 fold에서의 평균 정확도 계산
20 mean_score = np.mean(scores)
21 print('평균 정확도: {0:.4f}'.format(mean_score))
22 # exec_kfold 호출
23 exec_kfold(dt_clf, folds = 5)

```

```

교차 검증 0 정확도 : 0.7542
교차 검증 1 정확도 : 0.7809
교차 검증 2 정확도 : 0.7865
교차 검증 3 정확도 : 0.7697
교차 검증 4 정확도 : 0.8202
평균 정확도: 0.7823

```

평균 정확도는 약 78.23%

```

1 # cross_val_score()를 이용해 교차 검증 수행하기
2 from sklearn.model_selection import cross_val_score
3 scores = cross_val_score(dt_clf, x_df, y_df, cv=5)
4 for iter_count, accuracy in enumerate(scores):
5     print('교차 검증 {0} 정확도: {1:.4f}'.format(iter_count, accuracy))
6 print('평균 정확도: {0:.4f}'.format(np.mean(scores)))

```

```

교차 검증 0 정확도: 0.7430
교차 검증 1 정확도: 0.7753
교차 검증 2 정확도: 0.7921
교차 검증 3 정확도: 0.7865
교차 검증 4 정확도: 0.8427
평균 정확도: 0.7879

```

K 폴드와 cross\_val\_score()의 수치가 다른 이유:

cross\_val\_score()가 StratifiedKFold를 이용해 폴드 세트를 분할하기 때문.

```

1 # GridSearchCV를 이용해 DecisionTreeClassifier의 최적 하이퍼 파라미터를 찾고, 이를 예측
2 from sklearn.model_selection import GridSearchCV
3 parameters = {'max_depth':[2, 3, 5, 10],
4               'min_samples_split': [2, 3, 5],
5               'min_samples_leaf': [1, 5, 8]}
6 grid_dclf = GridSearchCV(dt_clf, param_grid=parameters, scoring='accuracy',
7 grid_dclf.fit(x_train, y_train)
8
9 print('GridSearchCV 최적 하이퍼 파라미터: ', grid_dclf.best_params_)
10 print('GridSearchCV 최고 정확도: {0:.4f}'.format(grid_dclf.best_score_))
11 best_dclf = grid_dclf.best_estimator_
12
13 # GridSearchCV의 최적 하이퍼 파라미터로 학습된 Estimator로 예측 및 평가 수행.
14 dpredictions = best_dclf.predict(x_test)
15 accuracy = accuracy_score(y_test, dpredictions)
16 print('테스트 세트에서의 DecisionTreeClassifier 정확도:{0:.4f}'.format(accuracy))

```

```

GridSearchCV 최적 하이퍼 파라미터: {'max_depth': 3, 'min_samples_leaf': 5, 'min
GridSearchCV 최고 정확도: 0.7992
테스트 세트에서의 DecisionTreeClassifier 정확도:0.8715

```

학습 후 예측 정확도가 87.15%로 향상됨.

## ✓ 3.6장 실습 - 피마 인디언 당뇨병 예측

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import accuracy_score, precision_score, recall_score, r
8 from sklearn.metrics import f1_score, confusion_matrix, precision_recall_cur
9 from sklearn.preprocessing import StandardScaler
10 from sklearn.linear_model import LogisticRegression
11
12 dd = pd.read_csv('/content/diabetes.csv')
13 print(dd['Outcome'].value_counts())
14 dd.head(3)

```

```
0    500
```

```
1    268
```

```
Name: Outcome, dtype: int64
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Diabetes
0	6	148	72	35	0	33.6	
1	1	85	66	29	0	26.6	
2	8	183	64	0	0	23.3	

Next steps:

[View recommended plots](#)

```
1 # feature값과 Null 개수 살펴보기
```

```
2 dd.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 768 entries, 0 to 767
```

```
Data columns (total 9 columns):
```

#	Column	Non-Null Count	Dtype
0	Pregnancies	768 non-null	int64
1	Glucose	768 non-null	int64
2	BloodPressure	768 non-null	int64
3	SkinThickness	768 non-null	int64
4	Insulin	768 non-null	int64
5	BMI	768 non-null	float64
6	DiabetesPedigreeFunction	768 non-null	float64
7	Age	768 non-null	int64
8	Outcome	768 non-null	int64

```
dtypes: float64(2), int64(7)
```

```
memory usage: 54.1 KB
```

Null값은 없으며, 피쳐 타입은 모두 숫자형이다.

- 별도의 피쳐 인코딩은 필요없음

```
1 # 사전 작업 - 함수 정의 - 정확도
```

```
2 from sklearn.base import BaseEstimator
```

```
3 from sklearn.datasets import load_digits
```

```
4 from sklearn.model_selection import train_test_split
```

```
5 from sklearn.base import BaseEstimator
```

```
6
```

```
7 from sklearn.metrics import accuracy_score, precision_score, recall_score, c
```

```
8
```

```
9 import numpy as np
```

```
10 import pandas as pd
```

```

1 def get_clf_eval(y_test, pred=None, pred_proba=None):
2     confusion = confusion_matrix(y_test, pred)
3     accuracy = accuracy_score(y_test, pred)
4     precision = precision_score(y_test, pred)
5     recall = recall_score(y_test, pred)
6     f1 = f1_score(y_test, pred)
7     roc_auc = roc_auc_score(y_test, pred_proba)
8     print ('오차행렬')
9     print(confusion)
10    print('정확도:{0:.4f}, 정밀도:{1:.4f}, 재현율:{2:.4f}, F1:{3:.4f}, AUC:{4:.4f}')

```

```

1 # 로지스틱 회귀를 이용한 예측 모델 생성
2
3 # 피쳐 데이터 세트 X, 레이블 데이터 세트 y를 추출.
4 # 맨 끝이 outcome 칼럼으로 레이블값. 칼럼 위치 -1을 이용해 추출
5 X = dd.iloc[:, :-1]
6 y = dd.iloc[:, -1]
7
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, r
9
10 # 로지스틱 회귀로 학습, 예측 및 평가 수행.
11 lr_clf = LogisticRegression()
12 lr_clf.fit(X_train, y_train)
13 pred = lr_clf.predict(X_test)
14 pred_proba = lr_clf.predict_proba(X_test)[: , 1]
15 get_clf_eval(y_test, pred, pred_proba)

```

오차행렬

[[88 12]

[23 31]]

정확도:0.7727, 정밀도:0.7209, 재현율:0.5741, F1:0.639175, AUC:0.7919

/usr/local/lib/python3.10/dist-packages/sklearn/linear\_model/\_logistic.py:4

STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

n\_iter\_i = \_check\_optimize\_result(

```

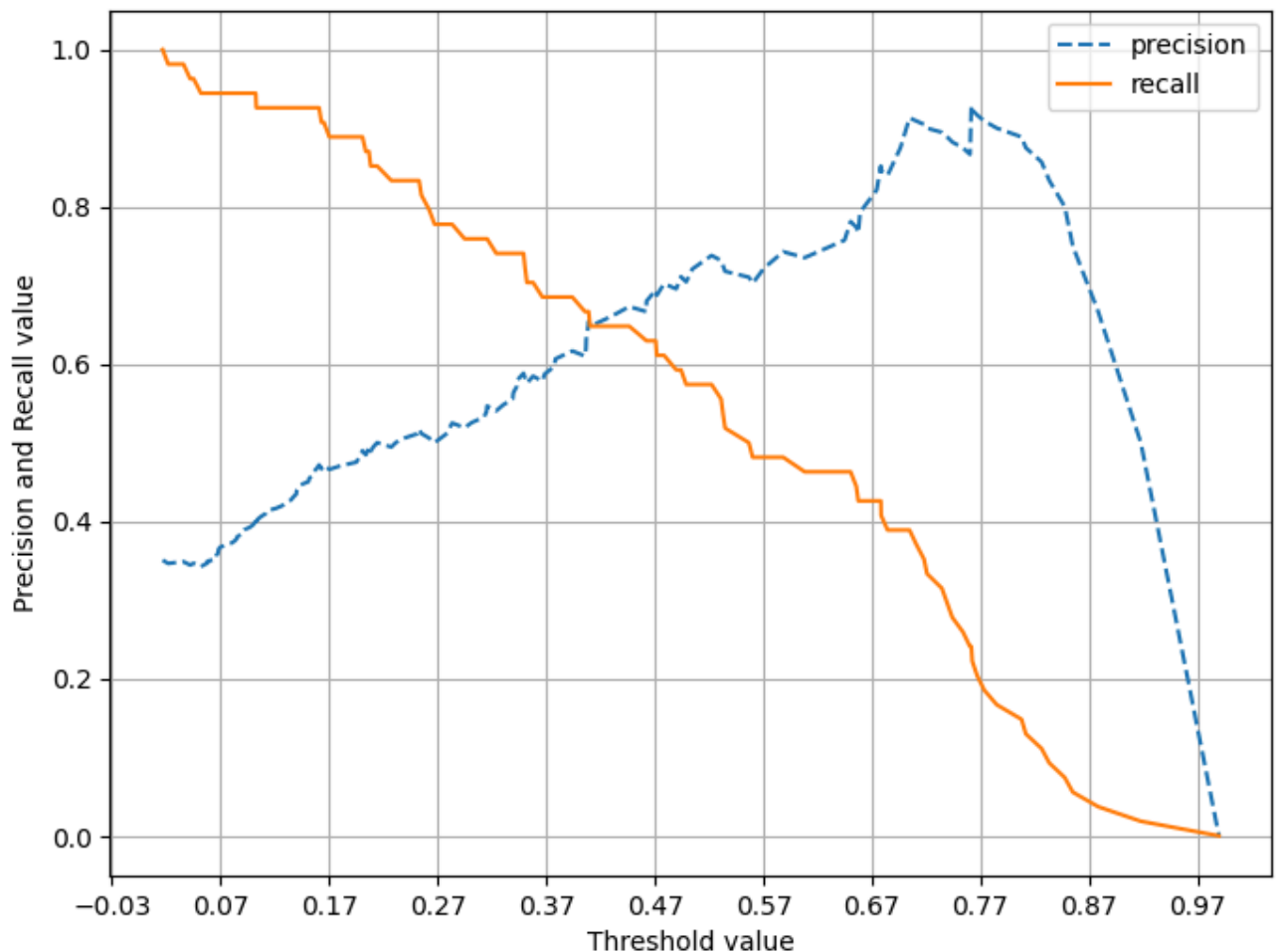
1 import matplotlib.pyplot as plt
2 import matplotlib.ticker as ticker
3 %matplotlib inline
4
5 def precision_recall_curve_plot(y_test, pred_proba_c1):
6     precisions, recalls, thresholds = precision_recall_curve(y_test, pred_proba_c1)
7     plt.figure(figsize=(8, 6))
8     threshold_boundary = thresholds.shape[0]
9     plt.plot(thresholds, precisions[0: threshold_boundary], linestyle='--',
10            plt.plot(thresholds, recalls[0: threshold_boundary], label='recall')
11
12     start, end = plt.xlim()
13     plt.xticks(np.round(np.arange(start, end, 0.1), 2))
14     plt.xlabel('Threshold value'); plt.ylabel('Precision and Recall value')
15     plt.legend(); plt.grid()
16     plt.show()

```

```

1 pred_proba_c1 = lr_clf.predict_proba(X_test)[: , 1]
2 precision_recall_curve_plot(y_test, pred_proba_c1)

```



```

1 # 피쳐 값 분포도 살펴보기
2 dd.describe()

```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.0
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.0
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.3
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.0
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.6

```
1 plt.hist(dd['Glucose'], bins=10)
```

```
2 # 0값이 일정 수준 존재함
```

```
(array([ 5.,  0.,  4., 32., 156., 211., 163., 95., 56., 46.]),
 array([ 0., 19.9, 39.8, 59.7, 79.6, 99.5, 119.4, 139.3, 159.2,
        179.1, 199. ]),
 <BarContainer object of 10 artists>)
```

