

---

---

---

---

---



```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
matplotlib inline?

```

7

2018.11.13 일요일

]

titanic\_df = pd.read\_csv('./titanic\_train.csv') → 데이터 불러오기

titanic\_df['Age'].fillna(titanic\_df['Age'].mean(), inplace=True)

titanic\_df['Cabin'].fillna('N', inplace=True)

titanic\_df['Embarked'].fillna('N', inplace=True)

titanic\_df.isnull().sum().sum() → 결측값 개수

titanic\_df['Cabin'] = titanic\_df['Cabin'].str[:1] → cabin 번호 앞자리만 추출

print(titanic\_df['Cabin'].head(3)) → 데이터 3개만 보기

titanic\_df.groupby(['Sex', 'Survived'])['Survived'].count() → 성별에 따른 생존사망 확인

→ 시각화 (Seaborn 라이브러리)

sns.barplot(x='Sex', y='Survived', data=titanic\_df)

sns.barplot(x='Pclass', y='Survived', hue='Sex', data=titanic\_df)

def get\_category(age):

cat = ''

if age <= -1 : cat = 'Unknown'

elif age <= 5 : cat = 'Baby'

elif age <= 12 : cat = 'Child'

elif age <= 18 : cat = 'Teenager'

elif age <= 25 : cat = 'Student'

elif age <= 35 : cat = 'Young Adult'

elif age <= 60 : cat = 'Adult'

else : cat = 'Elderly'

return cat

7

Age 범주화

범주(카테고리) 생성

]

# 막대 그래프의 크기 figure를 더 크게 설정

```
plt.figure(figsize=(10,6))
```

# x축의 값을 순차적으로 표시하기 위한 설정

```
group_names=['Unknown', 'Baby', 'Child', 'Teenager', 'Student', 'Young Adult', 'Adult', 'Elderly']
```

# lambda 식에 의해 생성한 get\_category() 함수를 빈칸으로 지정.

# get\_category(x)는 입력값으로 'Age' 카테고리의 값을 받아서 해당하는 cat 번호

```
titanic_df['Age_cat']=titanic_df['Age'].apply(lambda x: get_category(x))
```

```
sns.barplot(x='Age_cat', y='Survived', hue='Sex', data=titanic_df, order=group_names)
```

```
titanic_df.drop('Age_cat', axis=1, inplace=True) → 데이터프레임의 컬럼 삭제
```

문자열 → 숫자로

입력값 0 ~ (카테고리-1)

```
from sklearn.preprocessing import LabelEncoder
```

```
def encode_features(dataDF):
```

```
    features=['Cabin', 'Sex', 'Embarked']
```

```
    for feature in features:
```

```
        le=LabelEncoder()
```

```
        le=le.fit(dataDF[feature])
```

```
        dataDF[feature]=le.transform(dataDF[feature])
```

```
    return dataDF
```

```
titanic_df=encode_features(titanic_df)
```

```
titanic_df.head() → sn
```

## 테이퍼 진차리 황수 민들기

### 원 Null 처리 방법

det tillrä (dt):

```
df['Age'].fillna(df['Age'].mean(), inplace=True)
```

```
df['Cabin'].fillna('N', inplace=True)
```

```
df['Embarked'].fillna('N', inplace=True)
```

```
df['fare'].fillna(0, inplace=True)
```

return dt

위 알리기식 생각한 데로만 정해지겠지 걱정

def transform\_features(df):

$$df = fillna(df)$$
$$dA = \text{drop\_features}(dA)$$
$$df = \text{format\_features}(df)$$

return fa

위 여성이 남 인리(혹은) 복된도를 피하게

```
def drop_features(df):
```

```
df.drop(['PassengerId', 'Name', 'Ticket'], axis=1, inplace=True)
```

return it

취재이력    인연담    속보

```
def format_features(x):
```

```
features = ['Cabin', 'Sex', 'Embarked']
```

for feature in textures:

```
le = LabelEncoder()
```

```
le = le.fit(df[feature])
```

$$df[\text{texture}] = (e.\text{transform}(df[\text{texture}]))$$

return df

회원분 데이터를 재구성하고, 최적화되어 있는 사이트와 앱을 만들어 새로운

`titanic - df = pd.read_csv('titanic_train.csv')`

```
y_titanic_df = titanic_df['Survived']
```

X-titanic - df = titanic - df . drop('Survived', axis = 1)

$X_{\text{titanic\_dt}} = \text{transform\_features}(X_{\text{titanic\_dt}}) \rightarrow \text{np.ndarray}$

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

↳ 학습, 검증 데이터분할

```
x_train, x_test, y_train, y_test = train_test_split(x_titanic_df, y_titanic_df, test_size=0.2,
```

random\_state = 11)  
↳ 같은 결과 출력 시켜

```

from sklearn.tree import DecisionTreeClassifier # 결정트리
from sklearn.ensemble import RandomForestClassifier # Random Forest
from sklearn.linear_model import LogisticRegression # 로지스틱 회귀
from sklearn.metrics import accuracy_score # 정확도 평가

```

# 결정트리, Random Forest, 로지스틱 회귀를 위한 사이킷런 Classifier 클래스 생성

```
dt_clf = DecisionTreeClassifier(random_state=11)
```

```
rf_clf = RandomForestClassifier(random_state=11)
```

```
lr_clf = LogisticRegression(solver='liblinear') # 최적화 알고리즘 (그라디언트 디센트)
```

# DecisionTreeClassifier 학습 / 예측 / 평가

```
dt_clf.fit(X_train, y_train)
```

```
dt_pred = dt_clf.predict(X_test)
```

```
print('DecisionTreeClassifier 정확도: {:.4f}'.format(accuracy_score(y_test, dt_pred)))
```

# RandomForestClassifier 학습 / 예측 / 평가

```
rf_clf.fit(X_train, y_train)
```

```
rf_pred = rf_clf.predict(X_test)
```

```
print('RandomForestClassifier 정확도: {:.4f}'.format(accuracy_score(y_test, rf_pred)))
```

# LogisticRegression 학습 / 예측 / 평가

```
lr_clf.fit(X_train, y_train)
```

```
lr_pred = lr_clf.predict(X_test)
```

```
print('LogisticRegression 정확도: {:.4f}'.format(accuracy_score(y_test, lr_pred)))
```

model selection 방식  
모든 성능 평가 → 교차검증 (Kfold 클래스, cross\_val\_score(), GridSearchCV)

#1. kfold로 검증하려면 교차검증

```
from sklearn.model_selection import KFold
```

```
def exec_kfold(df, folds=5):
```

# 클래스에 따른 5개의 kfold 객체 생성, fold 순으로 예측값과 지수를 위한 리스트/배열 생성

```
    kfold = KFold(n_splits=folds)
```

```
    scores = []
```

# kfold 교차검증 시작

```
    for iter_count, (train_index, test_index) in enumerate(kfold.split(X_titanic_df)):
```

# X\_titanic\_df 데이터에서 교차검증으로 학습과 검증 데이터를 나누어 index 생성

```
        X_train, X_test = X_titanic_df.values[train_index], X_titanic_df.values[test_index]
```

```
        y_train, y_test = y_titanic_df.values[train_index], y_titanic_df.values[test_index]
```

# Classifier 학습, 예측, 정확도 계산

```
        clf.fit(X_train, y_train)
```

```
        predictions = clf.predict(X_test)
```

```
        accuracy = accuracy_score(y_test, predictions)
```

```
        scores.append(accuracy)
```

```
        print('교차검증 {0} 정확도: {1:.4f}'.format(iter_count, accuracy))
```

# 5fold의 평균 정확도를 계산

```
    mean_score = np.mean(scores)
```

```
    print("평균 정확도: {0:.4f}".format(mean_score))
```

# exec\_kfold 종료

```
exec_kfold(df=clf, folds=5)
```

#2. cross\_val\_score로 검증하려면 교차검증

```
from sklearn.model_selection import cross_val_score
```

```
scores = cross_val_score(df=clf, X=X_titanic_df, y=y_titanic_df, cv=5)
```

```
for iter_count, accuracy in enumerate(scores):
```

```
    print('교차검증 {0}, 정확도 {1:.4f}'.format(iter_count, accuracy))
```

```
print('평균 정확도: {0:.4f}'.format(np.mean(scores)))
```

#3. GridSearchCV를 통한 하이퍼파라미터 최적화

```
from sklearn.model_selection import GridSearchCV
```

```
parameters = {'max_depth': [2, 3, 5, 10],
```

```
               'min_samples_split': [2, 3, 5], 'min_samples_leaf': [1, 5, 6]}
```

```
grid_dclf = GridSearchCV(dt_clf, param_grid=parameters, scoring='accuracy', cv=5)
```

```
grid_dclf.fit(X_train, y_train)
```

```
print('GridSearchCV 최적화 파라미터:', grid_dclf.best_params_)
```

```
print('GridSearchCV 최적화 성능: {:.4f}'.format(grid_dclf.best_score_))
```

```
best_dclf = grid_dclf.best_estimator_
```

# GridSearchCV의 최적 하이퍼파라미터를 활용한 Estimator 생성 및 평가

```
predictions = best_dclf.predict(X_test)
```

```
accuracy = accuracy_score(y_test, predictions)
```

```
print('테스트 세트에 대한 Decision Tree Classifier 정확도: {:.4f}'.format(accuracy))
```