



# Titanic

## Prediction for Survived

EDA / Feature Engineering / Machine Learning

# Titanic\_Kaggle

Learning

## Titanic with Youhan Lee



게시자: Youhan Lee · 재생목록 · 동영상 15개 · 조회수 24,852회

### 1 - 14강 강의 수강

- 강의 수강 기록지 작성
- 강의 에러 기록지 작성

1

캐글 타이타닉 Titanic - 1. Dataset check  
Youhan Lee · 조회수 3.1만회 · 6년 전  
30:30

2

캐글 타이타닉 Titanic - 2. EDA - Pclass  
Youhan Lee · 조회수 9.1천회 · 6년 전  
36:30

3

캐글 타이타닉 Titanic - 3. EDA - Sex(성별)  
Youhan Lee · 조회수 4.6천회 · 6년 전  
15:31

4

캐글 타이타닉 Titanic - 4. EDA - Age  
Youhan Lee · 조회수 4.3천회 · 6년 전  
32:21

5

캐글 타이타닉 Titanic - 5. Age, Sex, Pclass (violinplot)  
Youhan Lee · 조회수 2.7천회 · 6년 전  
10:33

6

캐글 타이타닉 Titanic - 6. EDA - Embarked  
Youhan Lee · 조회수 2.2천회 · 6년 전  
12:22



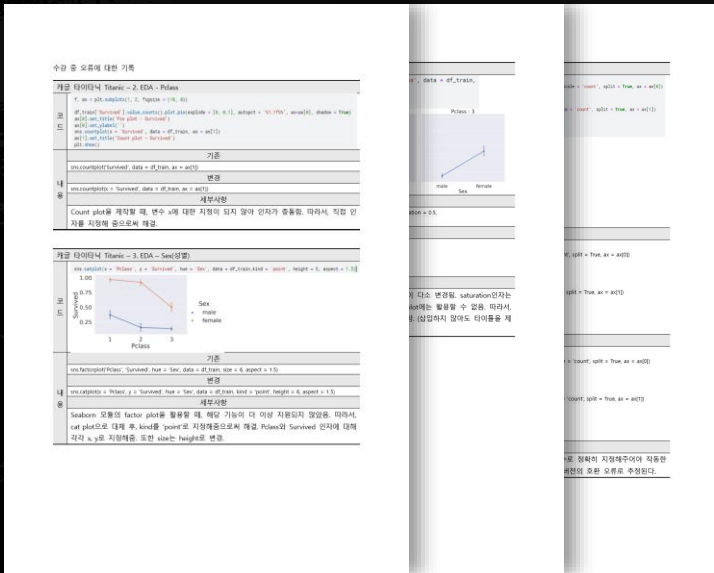


# Titanic\_Kaggle

## Learning

### 1 - 14강 강의 수강

- 강의 수강 기록지 작성
- 강의 에러 기록지 작성



경설검 / 경희대학교 서울캠퍼스 컴퓨터과학 동아리		@khu_sc.com
경설검 기초 프로젝트 강의 기록 (1~8강)		
소속	성명	
경희대학교 경영대학 경영학과	최석훈	
강의 주제		
캐글 타이타닉 (데이터 사이언스) by Youhan Lee in YouTube URL: <a href="https://www.youtube.com/playlist?list=PLC_wC_PMBLSMnqmgTLqDgu4tO8mQakuF">https://www.youtube.com/playlist?list=PLC_wC_PMBLSMnqmgTLqDgu4tO8mQakuF</a>		
강의 목차	내용 요약 및 추가할 점	
① Dataset Check	데이터셋 구성 확인 및 결측치 시각화 1. df_train에 pandas로 csv파일을 읽도록 할당한다. 2. Pandas 자체 기능 (describe, head)을 활용해 수치 확인 3. Missingno를 활용한 결측치 시각화	
② EDA - Pclass	Pclass와 Survived 데이터 활용에 시각화 1. Pandas 활용에 데이터 수치 표 제작 및 색상 추가 2. 파이 차트, 바 차트 활용 시각화 3. 차트에 대한 색상 및 구조 설정	
③ EDA - Sex	Pclass와 Sex(성별) 간 상관관계 분석 및 시각화 1. 이전 강의와 대체적으로 동일하게 진행함. 2. Point plot을 활용에 Pclass와 Survived 인자의 상관관계 분석을 진행하고 차트와 함.	
④ EDA - Age	Age(나이) 기준, Pclass 분포 및 생존 확률 시각화 1. KDE plot을 활용한 분포도 확인 2. 히스토그램과 KDE plot의 차이점 확인 3. Age와 Survived와의 상관관계 분석 (▶연령이 낮을수록 생존확률이 높다.)	
⑤ Age, Sex, Pclass (Violinplot)	Violin plot을 활용에 Age와 Sex가 생존확률과 연관성을 갖는지 Pclass로 나누어 분석 1. Violin plot의 시각화 (3가지 인자로 구성됨) 2. 각각의 인자와 생존의 상관관계를 파악	
⑥ EDA - Embarked	Count plot을 활용에 Embarked, Sex, Survived, Pclass 네가지 인자를 2x2 구조 (figsize = (2, 2))로 시각화 1. 하나의 결과 값에 2*2, 총 4개의 그래프 삽입 2. Embarked 기준, Sex, Survived, Pclass 세가지 변수 표현	
⑦ EDA - Family Size	Count plot을 활용에 Family size별 Survived 인자와 연관성 확인 1. Count plot을 1*3구조로 형성 2. Family size에 대한 재정의 (데이터 프레임 재구성) 3. Family size별 생존확률에 대한 분석	
⑧ EDA - Fare, Cabin, Ticket	Distplot을 활용에 Fare(요금)의 분포를 파악 1. Fare를 기준으로 분포도를 형성 2. Skewness(왜도): 데이터가 왜곡된 정도에 정보의 비대칭성을 의미. 과도한 왜도를 조정해주는 수단 → 로그 3. Lambda로 Log를 위해 Skewness 4.78 → 0.44 조정	
- 24학년도 2학기 기초 프로젝트팀 -		

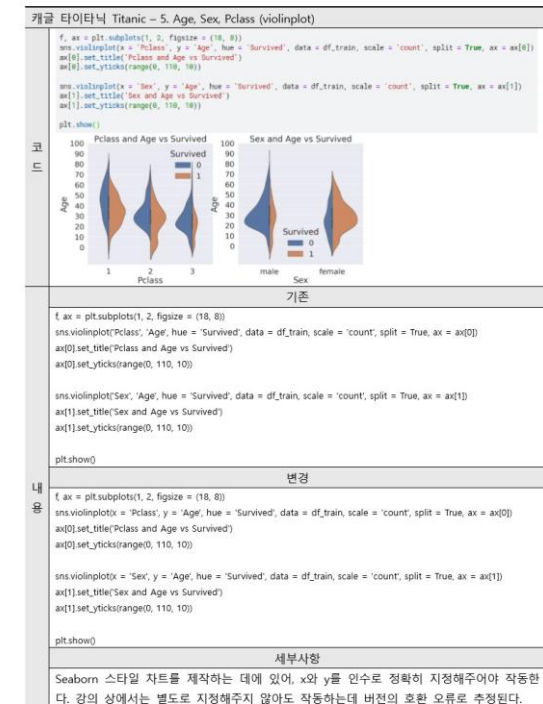
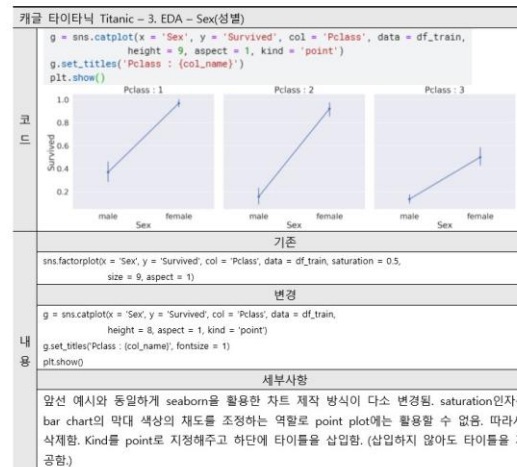
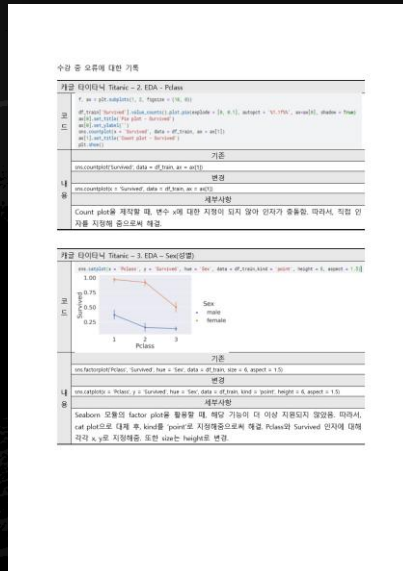
경설검 / 경희대학교 서울캠퍼스 컴퓨터과학 동아리		@khu_sc.com
경설검 기초 프로젝트 강의 기록		
소속	성명	
경희대학교 경영대학 경영학과	최석훈	
강의 주제		
캐글 타이타닉 (데이터 사이언스) by Youhan Lee in YouTube URL: <a href="https://www.youtube.com/playlist?list=PLC_wC_PMBLSMnqmgTLqDgu4tO8mQakuF">https://www.youtube.com/playlist?list=PLC_wC_PMBLSMnqmgTLqDgu4tO8mQakuF</a>		
강의 목차	내용 요약 및 추가할 점	
⑨ Feature Engineering - Fill Null in Age	Feature engineering - Age의 결측치 보완 1. Name에서 Initial을 추출하여 승객의 성격을 분류한다. 2. Crosstab 시각화를 통해 Initial 구성을 확보. 3. replace함수를 통해 5가지 분류로 재구성 4. loc(로케이션) 함수로 NaN값에 대한 대체 값을 지정 5. Age에 대한 결측치를 채워준다. (결측치는 임의 지정)	
⑩ Feature Engineering - Fill Null in Embarked	Feature engineering - Embarked 결측치 보완 및 Age 카테고리화 1. Embarked 항목에 대해 fillna함수를 통해 결측치 보완 2. Age 항목에 대해 0~7까지 지수로 연령대 구분 (2가지) A. loc(로케이션) 함수를 사용하여 직접 채움 B. apply함수 활용하여 내부 함수를 순환시켜 적용 3. Age 카테고리화 적용법의 등장성 여부 판단 (call함수) 4. dropna함수를 통해 필요 없는 항목 제거	
⑪ Feature Engineering - Change string to categorical Person	Feature engineering - Initial, Embarked 문자열 항목 수치화 1. Initial과 Embarked의 문자열을 숫자로 변환 2. seaborn-타입 히트맵을 통한 데이터 시각화 3. Pearson Correlation Coefficient 피어슨 상관 계수 활용 A. C1: 종의 상관관계, 1: 양의 상관관계	
⑫ Feature Engineering - One-hot encoding on the initial and Embarked	Feature engineering - Initial, Embarked 항목에 대한 합성 분포 1. 기존: Initial, Embarked 항목의 값이 각각 0-4, 0-2로 지정되어 있었음 2. 변환 pandas의 get_dummies 함수를 서로 다른 columns으로 분리함 (예: initial0, initial1, ...) → 직접 하나하나 지정해주지 않아도 자동으로 작업할 수 있어 더 데이터 활용에 용이함 3. 불필요한 항목(passengerid)을 drop으로 제거함	
⑬ Model development - Machine learning (Randomforest)	Sklearn ensemble 중 RandomForest를 활용하여 모델 구축 1. 데이터 전처리 A. X_train에 Survived를 제외한 나머지 항목에 대한 데이터프레임을 Numpy 배열로 변환 (입력변수) B. target_label에 Survived 값을 인자(출력변수)로 활용하기 위해 Numpy 배열로 변환 (종속변수) C. X_test에 test.csv의 모든 결측치를 Numpy 배열로 변환 (입력변수) 2. 훈련 세트와 검증 세트 분리 A. X_train, y_train: 각각 훈련 세트와 검증 세트의 입력 데이터 B. X_test, y_test: 각각 훈련 세트와 검증 세트의 출력 데이터 (타겟) C. Test_size = 0.3: 검증 세트를 30%, 훈련 세트 70% 할당 D. Random_state: 시드 값을 고정하여 동일한 랜덤 (기본은 랜덤) 3. Random Forest Classifier 모델 지정 및 활용 A. fit: 훈련 세트 입력하여 training 실시 B. predict: 훈련된 모델을 바탕으로 X_test 배열로 검증 C. metrics의 accuracy_score함수를 통해 정확도(성능) 평가 4. 최종적으로 test.csv파일을 기반으로 테스트 실시 A. Test.csv파일을 대상으로 model을 검증한다. 2. 검증된 모델로 2가지 submission을 제출한다. 3. 수행한 결과를 제출에서 확인한다.	
⑭ Machine learning prediction - feature importance and prediction on test set	Feature importance and prediction on test set	
- 24학년도 2학기 기초 프로젝트팀 -		

# Titanic\_Kaggle

## Learning

### 1 - 14강 강의 수강

- 강의 수강 기록지 작성
- 강의 에러 기록지 작성



# Titanic\_Kaggle

## Process for Machine Learning

---

EDA

### EDA

탐색적 데이터 분석으로, 가공하지 않은 일차적 Raw Data를 통해, 주어진 데이터의 구성과 형태를 확인하고 ML까지의 과정을 계획한다.

Feature  
Engineering

### Feature Engineering

주어진 데이터를 ML에 적용하기 위해 적절한 형태로 가공, 변환, 통합, 축소 등 데이터 전처리를 수행한다.

Machine  
Learning

### Machine Learning with Random Forest

Feature Engineering을 수행한 정제된 데이터를 통해, 선택한 ML모델에 적용한다. 결과를 수치로 도출하고 정확도를 측정해 모델의 완성도를 확인한다. 결과적으로, Test파일을 통해 Target을 예측한다.



# Titanic\_Kaggle

## Structure of Dataset

---

**\* Data Dictionary** Columns: 10, Row: 891 (train.csv)

**Target:** Survived

---

### Variable

- Survival
- Pclass
- Sex
- Age
- SibSp
- Parch
- Ticket
- Fare
- Cabin
- Embarked

### Definition

- Survival
- Ticket class
- Sex
- Age in years
- # of siblings / spouses aboard the Titanic
- # of parents / children aboard the Titanic
- Ticket number
- Passenger fare
- Cabin number
- Port of embarked

### Key

- 0 = No, 1 = Yes
- 1 = 1st, 2 = 2nd, 3 = 3rd
- C = Cherbourg, Q = Queenstown, S = Southampton

# Titanic\_Kaggle

## Before Beginning

---

```
1. import numpy as np
2. import pandas as pd
3. import matplotlib.pyplot as plt
4. import seaborn as sns
5. plt.style.use('seaborn')
6. sns.set(font_scale = 2.5)
7.
8. import missingno as msno
9. import warnings
10. warnings.filterwarnings('ignore')
11.
12. %matplotlib inline
```

### Numpy & Pandas

---

데이터 분석을 위한 도구로서 Pearson Correlation Coefficient 등에 활용.

### Matplotlib & Seaborn

---

분석한 데이터 셋을 바탕으로 시각화에 도움을 준다. Heatmap 등에 활용.

### Missingno

---

데이터 분석을 위한 도구로서 Pearson Correlation Coefficient 등에 활용.

### Warnings

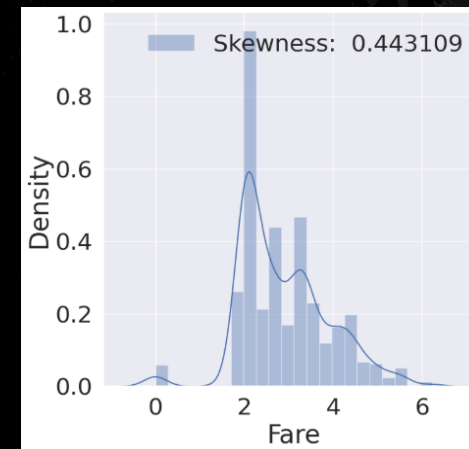
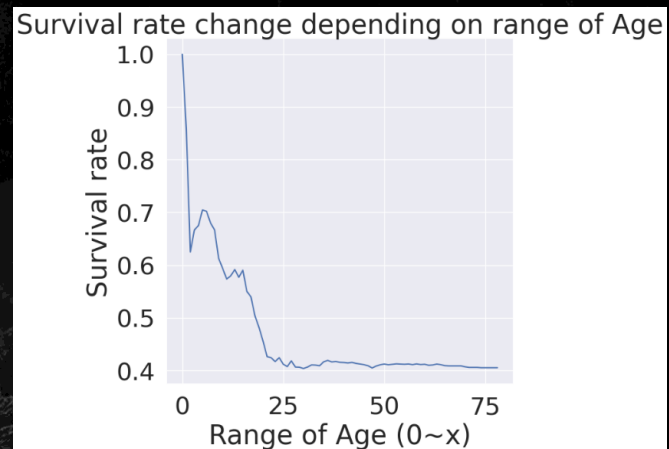
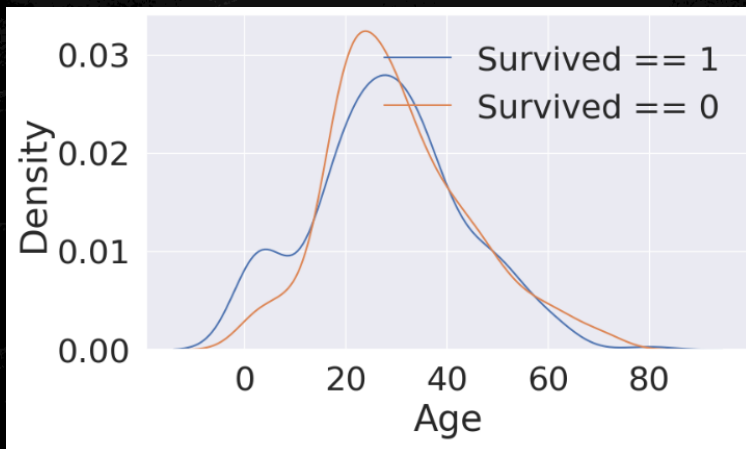
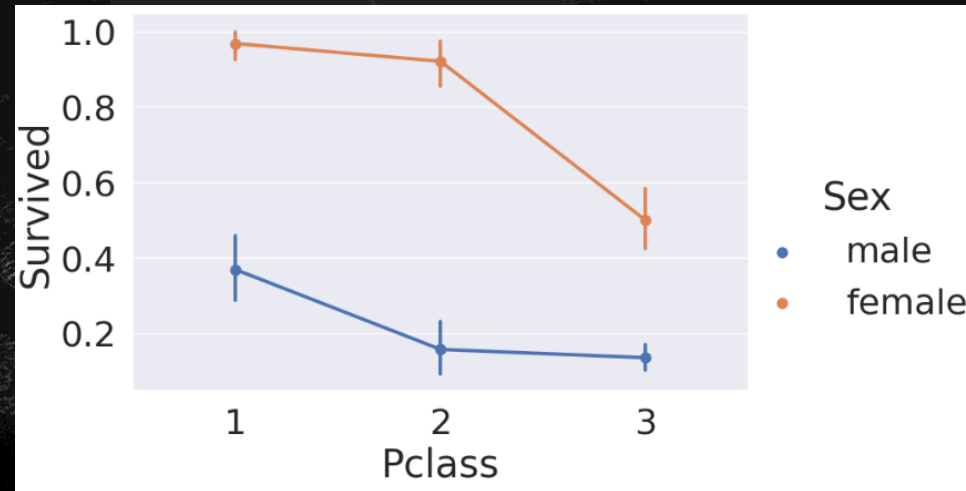
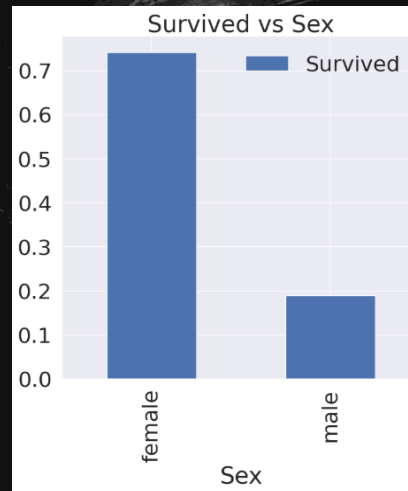
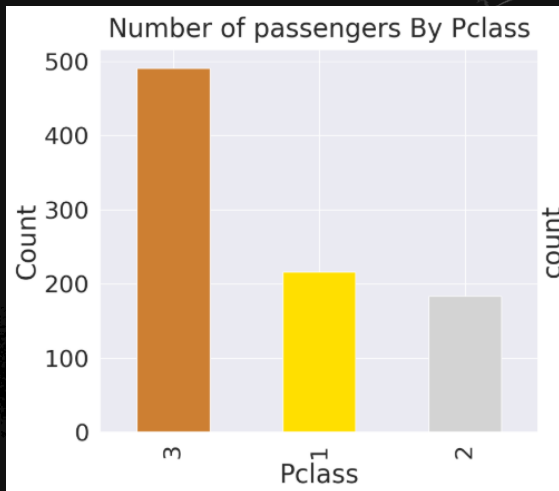
---

코드 작성 시 발생하는 경고에 대해 무시하는 명령.



# Titanic\_Kaggle

## EDA



Survived	0	1	All
Pclass			
1	80	136	216
2	97	87	184
3	372	119	491
All	549	342	891

# Titanic\_Kaggle

## Preprocess

```
#train, test파일 불러오기
df_train = pd.read_csv('/kaggle/input/titanic/train.csv')
df_test = pd.read_csv('/kaggle/input/titanic/test.csv')

#레이블 변형
#Familysize 생성
df_train['Familysize'] = df_train['SibSp'] + df_train['Parch'] + 1
df_test['Familysize'] = df_test['SibSp'] + df_test['Parch'] + 1

#test.csv 컬럼에 추가, Fare
#Pclass를 값에 따라 Age
fare_by_pclass = df_test.groupby('Pclass')['Fare'].transform('mean')
df_test['Fare'] = df_test['Fare'].fillna(fare_by_pclass, inplace=True)

#로그 변환
df_train['Fare'] = df_train['Fare'].map(lambda i: np.log(1 + i) if i > 0 else 0)
df_test['Fare'] = df_test['Fare'].map(lambda i: np.log(1 + i) if i > 0 else 0)

#Initial 생성
df_train['Initial'] = df_train.Name.str.extract('([A-Za-z]+)\.')
df_test['Initial'] = df_test.Name.str.extract('([A-Za-z]+)\.')

#Initial 수정 (5가지)
df_train['Initial'].replace(['Mlle', 'Mme', 'Ms', 'Dr', 'Major', 'Lady', 'Countess', 'Jonkheer', 'Col', 'Rev', 'Capt', 'Sir', 'Don', 'Dona'],
                           ['Miss', 'Miss', 'Miss', 'Mr', 'Mr', 'Mrs', 'Mrs', 'Other', 'Other', 'Other', 'Mr', 'Mr', 'Mr', 'Mr'], inplace=True)

df_test['Initial'].replace(['Mlle', 'Mme', 'Ms', 'Dr', 'Major', 'Lady', 'Countess', 'Jonkheer', 'Col', 'Rev', 'Capt', 'Sir', 'Don', 'Dona'],
                           ['Miss', 'Miss', 'Miss', 'Mr', 'Mr', 'Mrs', 'Mrs', 'Other', 'Other', 'Other', 'Mr', 'Mr', 'Mr', 'Mr'], inplace=True)

#Age의 Null값 처리
df_train.loc[(df_train.Age.isnull()) & (df_train.Initial=='Mr'), 'Age'] = 33
df_train.loc[(df_train.Age.isnull()) & (df_train.Initial=='Mrs'), 'Age'] = 36
df_train.loc[(df_train.Age.isnull()) & (df_train.Initial=='Master'), 'Age'] = 5
df_train.loc[(df_train.Age.isnull()) & (df_train.Initial=='Miss'), 'Age'] = 22
df_train.loc[(df_train.Age.isnull()) & (df_train.Initial=='Other'), 'Age'] = 46

df_test.loc[(df_test.Age.isnull()) & (df_test.Initial=='Mr'), 'Age'] = 33
df_test.loc[(df_test.Age.isnull()) & (df_test.Initial=='Mrs'), 'Age'] = 36
df_test.loc[(df_test.Age.isnull()) & (df_test.Initial=='Master'), 'Age'] = 5
df_test.loc[(df_test.Age.isnull()) & (df_test.Initial=='Miss'), 'Age'] = 22
df_test.loc[(df_test.Age.isnull()) & (df_test.Initial=='Other'), 'Age'] = 46

#Embarked의 Null값을 'S'로 대체
df_train['Embarked'].fillna('S', inplace = True)
```

```
#Embarked의 Null값을 'S'로 대체
df_train['Embarked'].fillna('S', inplace = True)

#Age_cat에 Age를 카테고리화 하여 인코딩 (10살 단위로 구분)
def category_age(x):
    for i in range(7):
        if x < ((i+1)*10):
            return i
    return 7

df_train['Age_cat'] = df_train['Age'].apply(category_age)
df_test['Age_cat'] = df_test['Age'].apply(category_age)
df_train.drop(['Age'], axis = 1, inplace = True) #기존 Age 레이블 삭제
df_test.drop(['Age'], axis = 1, inplace = True)

#Initial을 int값으로 대체 (이름, 직책등을 통한 것과 예측을 위해)
#또한, Initial을 위해서 Age의 Null값 해소와 상관없이 나열에 편리함.
df_train['Initial'] = df_train['Initial'].map({'Master': 0, 'Miss': 1, 'Mr': 2, 'Mrs': 3, 'Other': 4})
df_test['Initial'] = df_test['Initial'].map({'Master': 0, 'Miss': 1, 'Mr': 2, 'Mrs': 3, 'Other': 4})

#Embarked도 동일하게 int값으로 변환
df_train['Embarked'] = df_train['Embarked'].map({'C': 0, 'Q': 1, 'S': 2})
df_test['Embarked'] = df_test['Embarked'].map({'C': 0, 'Q': 1, 'S': 2})

#Sex도 동일하게 int값으로 변환
df_train['Sex'] = df_train['Sex'].map({'female': 0, 'male': 1})
df_test['Sex'] = df_test['Sex'].map({'female': 0, 'male': 1})

#get_dummies 함수를 활용한 Initial 레이블 분리 (Initial0, Initial1, ...)
df_train = pd.get_dummies(df_train, columns=['Initial'], prefix='Initial')
df_test = pd.get_dummies(df_test, columns=['Initial'], prefix='Initial')

#get_dummies 함수를 활용한 Embarked 레이블 분리 (Embarked0, Embarked1, Embarked2)
df_train = pd.get_dummies(df_train, columns=['Embarked'], prefix='Embarked')
df_test = pd.get_dummies(df_test, columns=['Embarked'], prefix='Embarked')

#불필요한 columns 제거
df_train.drop(['PassengerId', 'Name', 'SibSp', 'Parch', 'Ticket', 'Cabin'], axis=1, inplace=True)
df_test.drop(['PassengerId', 'Name', 'SibSp', 'Parch', 'Ticket', 'Cabin'], axis=1, inplace=True)
```

Null data 정리

카테고리화(범주)

레이블 병합

etc.

# Titanic\_Kaggle

## Machine Learning

```
#sklearn 모듈 불러오기
from sklearn.ensemble import RandomForestClassifier
from sklearn import metrics
from sklearn.model_selection import train_test_split

#X_train: Survived를 제외한 입력변수 설정
#target_label: Survived를 출력변수(타겟) 설정
#X_test: 테스트 파일을 Numpy 배열로 변환
X_train = df_train.drop('Survived', axis=1).values
target_label = df_train['Survived'].values
X_test = df_test.values

#훈련 및 검증 데이터를 테스트 세트(70%)와 검증 세트(30%)로 분리
#시드값을 고정하여 동일한 재현
X_tr, X_vld, y_tr, y_vld = train_test_split(X_train, target_label, test_size=0.3, random_state=2018)

#머신러닝 모델 설정 (RandomForestClassifier)
#fit 함수로 데이터 입력
#prediction 변수에 예측값 입력 (변수는 Test의 입력데이터(X_vld))
model = RandomForestClassifier()
model.fit(X_tr, y_tr)
prediction = model.predict(X_vld)

#Pandas를 활용하여 Feature Importance 추출
from pandas import Series

feature_importance = model.feature_importances_
Series_feat_imp = Series(feature_importance, index=df_test.columns)
```

1.  
sklearn

2.  
Submission

3.  
Accuracy

4.  
Feature  
Importance

```
#최종 제출을 위한 Submission.csv 불러오기
submission = pd.read_csv('/kaggle/input/titanic/gender_submission.csv')
```

```
#제작한 모델을 바탕으로 Submission의 Survived 예측값을 출력
prediction = model.predict(X_test)
submission['Survived'] = prediction
```

```
#최종 제출
submission.to_csv('./my_first_submission.csv', index=False)
```

```
print('총 {}명 중 {:.2f}% 정확도로 실적을 맞췄'.format(y_vld.shape[0], 100 * metrics.accuracy_score(prediction, y_vld)))
```

```
plt.figure(figsize=(8, 8))
Series_feat_imp.sort_values(ascending=True).plot.barh()
plt.xlabel('Feature importance')
plt.ylabel('Feature')
plt.show()
```