

Titanic

NO.

Library

```
{ import numpy as np
  import pandas as pd
  import seaborn as sns
  sns.set('darkgrid')
  style =
```

```
{ from sklearn.ensemble import RandomForestClassifier
  from sklearn.preprocessing import OneHotEncoder, LabelEncoder, StandardScaler
  from sklearn.metrics import roc_curve, auc
  from sklearn.model_selection import StratifiedKFold
```

```
import string
```

```
import warnings
```

```
warnings.filterwarnings('ignore')
```

```
SEED = 42
```

```
def concat_df(train_data, test_data):
```

```
    return pd.concat([train_data, test_data], sort=True).reset_index(
        drop=True)
```

```
def divide_df(all_data)
```

```
    return all_data.iloc[:890], all_data.iloc[891:].drop(['Survived'],
        axis=1)
```


1.2.1 Age.

```
df-all-corr = df-all-corr.corr().abs().unstack().sort_values(
    kind='quicksort', ascending=False).reset_index()
```

```
df-all-corr.rename(columns={'axel_0' : 'Feature 1', ... 3, inplace=True)
```

```
df-all-corr[df-all-corr['Feature 1'] == 'Age']
```

index 별

median값

찾아 올림.

```
age_by_pclass_sex = df-all.groupby(['Sex', 'pclass']).median()['Age']
```

```
df-all['Age'] = df-all.groupby(['Sex', 'pclass'])['Age'].apply(lambda x:
    x.fillna(x.median()))
```

1.2.2 Embarked.

```
df-all[df-all['Embarked'].isnull()]
```

```
df-all['Embarked'] = df-all['Embarked'].fillna('S')
```

1.2.3 Fare

['Fare']

```
med-fare = df-all.groupby(['pclass', 'parch', 'sibsp'])['Fare'].median()[3][0]
```

['Fare']

1.2.4 Cabin.

```
df-all['Deck'] = df-all['Cabin'].apply(lambda s: s[0] if pd.notnull(s)
    else 'N')
```

```
df-all-decks = df-all.groupby(['Deck', 'Pclass']).count().drop(columns=['...'])
```

```
def get_pass_dist(df):
```

```
    # Creating a dictionary for every passenger class count in every deck.
```

```
    deck_counts = {'A': {}, 'B': {}, ..., 'T': {}}
```

```
    decks = df.columns.levels[0]
```

```
    for deck in decks:
```

```
        for pclass in range(1, 4):
```

```
            try:
```

```
                count = df[deck][pclass][0]
```

```
                deck_counts[deck][pclass] = count
```

```
            except KeyError:
```

```
                deck_counts[deck][pclass] = 0
```

```
    # Percentage.
```

```
    df_decks = pd.DataFrame(deck_counts)
```

```
    deck_percentages = {}
```

```
    for col in df_decks.columns:
```

```
        deck_percentages[col] = [(count / df_decks[col].sum()) * 100
```

```
            for count in df_decks[col]]
```

```
    return deck_counts, deck_percentages
```



```
def display_pclass_dist(percentages):
```

```
    Runtime  
    df_percentages = pd.percentages, transpose)
```

```
    deck_names = ('A', ..., 'T')
```

```
    bar_count = np.arange(len(deck_names))
```

```
    bar_width = 0.85
```

```
    pclass1 = df_percentages[0]
```

```
    pclass2 = df_percentages[1]
```

```
    pclass3 = df_percentages[2]
```

```
    plt.figure(figsize=(20,10))
```

```
    plt.bar(bar_count, pclass1, color='#659494', edgecolor='white',  
            width=bar_width, label='passenger class 1')
```

```
    plt.bar(" " , bottom=pclass1 " ")
```

```
    plt.bar(" " , bottom=pclass1 + pclass2
```

각 줄의 label 사이의 간격

```
    plt.xlabel('Deck', size=15, labelpad=20)
```

```
    plt.ylabel('Passenger Class Percentage', size=15, labelpad=20)
```

```
    plt.xticks(bar_count, deck_names) → X의 label 설정
```

```
    plt.tick_params(axis='x', labelsize=15)
```

```
    " (" " " , " " )
```

```
    plt.legend(loc='upper left', bbox_to_anchor=(1,1), prop={'size':15})
```

```
    plt.title(' ', size=18, y=1.05)
```

```
    plt.show()
```