

```
In [ ]: import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style="darkgrid")

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
```

```
In [ ]: def concat_df(train_data, test_data):
    # Returns a concatenated df of training and test set on axis 0
    return pd.concat([train_data, test_data], sort=True).reset_index(drop=True)

def divide_df(all_data):
    # Returns divided dfs of training and test set
    return all_data.loc[:890], all_data.loc[891:].drop(['Survived'], axis=1)

df_train = pd.read_csv('train.csv')
df_test = pd.read_csv('test.csv')
df_all = concat_df(df_train, df_test)

df_train.name = 'Training Set'
df_test.name = 'Test Set'
df_all.name = 'All Set'

dfs = [df_train, df_test]

print('Number of Training Examples = {}'.format(df_train.shape[0]))
print('Number of Test Examples = {}\n'.format(df_test.shape[0]))
print('Training X Shape = {}'.format(df_train.shape))
print('Training y Shape = {}\n'.format(df_train['Survived'].shape[0]))
print('Test X Shape = {}'.format(df_test.shape))
print('Test y Shape = {}\n'.format(df_test.shape[0]))
print(df_train.columns)
print(df_test.columns)
```

```
In [ ]: print(df_train.info())
df_train.sample(3)
print(df_test.info())
df_test.sample(3)
```

```

In [ ]: def display_missing(df):
        for col in df.columns.tolist():
            print('{} column missing values: {}'.format(col, df[col].isnull().sum()))
        print('\n')

        for df in dfs:
            print('{}'.format(df.name))
            display_missing(df)

df_all_corr = df_all.corr().abs().unstack().sort_values(kind="quicksort", ascending=False)
df_all_corr.rename(columns={"level_0": "Feature 1", "level_1": "Feature 2", 0: "Correlation"}, inplace=True)
df_all_corr[df_all_corr['Feature 1'] == 'Age']
age_by_pclass_sex = df_all.groupby(['Sex', 'Pclass']).median()['Age']

for pclass in range(1, 4):
    for sex in ['female', 'male']:
        print('Median age of Pclass {} {}s: {}'.format(pclass, sex, age_by_pclass_sex[sex][pclass-1]))
print('Median age of all passengers: {}'.format(df_all['Age'].median()))

# Filling the missing values in Age with the medians of Sex and Pclass groups
df_all['Age'] = df_all.groupby(['Sex', 'Pclass'])['Age'].apply(lambda x: x.fillna(x.median()))

```

```

In [ ]: df_all['Embarked'] = df_all['Embarked'].fillna('S')
med_fare = df_all.groupby(['Pclass', 'Parch', 'SibSp']).Fare.median()[3][0][0]
# Filling the missing value in Fare with the median Fare of 3rd class alone passengers
df_all['Fare'] = df_all['Fare'].fillna(med_fare)

```

```

In [ ]: #there 697 values of caboin missing so we will drop it.
df_all=df_all.drop(['Cabin'],axis=1)
df_train, df_test = divide_df(df_all)
dfs = [df_train, df_test]

```

```
In [ ]:
```

```
df_all['Fare'] = pd.qcut(df_all['Fare'], 13)
df_all['Age'] = pd.qcut(df_all['Age'], 10)
df_all['Family_Size'] = df_all['SibSp'] + df_all['Parch'] + 1
family_map = {1: 'Alone', 2: 'Small', 3: 'Small', 4: 'Small', 5: 'Medium', 6: 'M
df_all['Family_Size_Grouped'] = df_all['Family_Size'].map(family_map)
df_all['Ticket_Frequency'] = df_all.groupby('Ticket')['Ticket'].transform('count
df_all['Title'] = df_all['Name'].str.split(',', expand=True)[1].str.split('.', 
df_all['Title'] = df_all['Title'].replace(['Miss', 'Mrs', 'Ms', 'Mlle', 'Lady', 'I
df_all['Title'] = df_all['Title'].replace(['Dr', 'Col', 'Major', 'Jonkheer', 'Ca

def extract_surname(data):

    families = []

    for i in range(len(data)):
        name = data.iloc[i]

        if '(' in name:
            name_no_bracket = name.split('(')[0]
        else:
            name_no_bracket = name

        family = name_no_bracket.split(',')[0]
        title = name_no_bracket.split(',')[1].strip().split(' ')[0]

        for c in string.punctuation:
            family = family.replace(c, '').strip()

        families.append(family)

    return families
df_all['Family'] = extract_surname(df_all['Name'])
df_train = df_all.iloc[0:890]

df_test = df_all.iloc[891:]

dfs = [df_train, df_test]
non_numeric_features = ['Embarked', 'Sex', 'Title', 'Family_Size_Grouped', 'Age

for df in dfs:
    for feature in non_numeric_features:
        df[feature] = LabelEncoder().fit_transform(df[feature])
cat_features = ['Pclass', 'Sex', 'Embarked', 'Title', 'Family_Size_Grouped']
encoded_features = []

for df in dfs:
    for feature in cat_features:
        encoded_feat = OneHotEncoder().fit_transform(df[feature].values.reshape(
n = df[feature].nunique()
cols = ['{}_{}'.format(feature, n) for n in range(1, n + 1)]
        encoded_df = pd.DataFrame(encoded_feat, columns=cols)
        encoded_df.index = df.index
        encoded_features.append(encoded_df)
df_train = pd.concat([df_train, *encoded_features[:5]], axis=1)
df_test = pd.concat([df_test, *encoded_features[5:]], axis=1)
```

```
In [ ]: df_all = concat_df(df_train, df_test)
drop_cols = ['Embarked', 'Family', 'Family_Size', 'Family_Size_Grouped', 'Survived',
             'Name', 'Parch', 'PassengerId', 'Pclass', 'Sex', 'SibSp', 'Ticket',
             ]

df_all.drop(columns=drop_cols, inplace=True)
```

```
In [ ]: X_train = StandardScaler().fit_transform(df_train.drop(columns=drop_cols))
y_train = df_train['Survived'].values
X_test = StandardScaler().fit_transform(df_test.drop(columns=drop_cols))

print('X_train shape: {}'.format(X_train.shape))
print('y_train shape: {}'.format(y_train.shape))
print('X_test shape: {}'.format(X_test.shape))
```

```
In [ ]: single_best_model = RandomForestClassifier(criterion='gini',
                                                  n_estimators=1100,
                                                  max_depth=5,
                                                  min_samples_split=4,
                                                  min_samples_leaf=5,
                                                  max_features='auto',
                                                  oob_score=True,
                                                  random_state=SEED,
                                                  n_jobs=-1,
                                                  verbose=1)

leaderboard_model = RandomForestClassifier(criterion='gini',
                                          n_estimators=1750,
                                          max_depth=7,
                                          min_samples_split=6,
                                          min_samples_leaf=6,
                                          max_features='auto',
                                          oob_score=True,
                                          random_state=SEED,
                                          n_jobs=-1,
                                          verbose=1)
```

```

In [ ]: N = 5
oob = 0
probs = pd.DataFrame(np.zeros((len(X_test), N * 2)), columns=['Fold_{}_Prob_{}'.format(fold, prob)] for fold in range(N) for prob in [0, 1])
importances = pd.DataFrame(np.zeros((X_train.shape[1], N)), columns=['Fold_{}'.format(fold)] for fold in range(N))
fprs, tprs, scores = [], [], []

skf = StratifiedKFold(n_splits=N, random_state=N, shuffle=True)

for fold, (trn_idx, val_idx) in enumerate(skf.split(X_train, y_train), 1):
    print('Fold {}'.format(fold))

    # Fitting the model
    leaderboard_model.fit(X_train[trn_idx], y_train[trn_idx])

    # Computing Train AUC score
    trn_fpr, trn_tpr, trn_thresholds = roc_curve(y_train[trn_idx], leaderboard_model.predict_proba(X_train[trn_idx])[0])
    trn_auc_score = auc(trn_fpr, trn_tpr)
    # Computing Validation AUC score
    val_fpr, val_tpr, val_thresholds = roc_curve(y_train[val_idx], leaderboard_model.predict_proba(X_train[val_idx])[0])
    val_auc_score = auc(val_fpr, val_tpr)

    scores.append((trn_auc_score, val_auc_score))
    fprs.append(val_fpr)
    tprs.append(val_tpr)
    # X_test probabilities
    probs.loc[:, 'Fold_{}_Prob_0'.format(fold)] = leaderboard_model.predict_proba(X_test)[0]
    probs.loc[:, 'Fold_{}_Prob_1'.format(fold)] = leaderboard_model.predict_proba(X_test)[1]
    importances.iloc[:, fold - 1] = leaderboard_model.feature_importances_

    oob += leaderboard_model.oob_score_ / N
    print('Fold {} OOB Score: {}'.format(fold, leaderboard_model.oob_score_))

print('Average OOB Score: {}'.format(oob))

```

```

In [ ]: importances['Mean_Importance'] = importances.mean(axis=1)
importances.sort_values(by='Mean_Importance', inplace=True, ascending=False)

plt.figure(figsize=(15, 20))
sns.barplot(x='Mean_Importance', y=importances.index, data=importances)

plt.xlabel('')
plt.tick_params(axis='x', labelsize=15)
plt.tick_params(axis='y', labelsize=15)
plt.title('Random Forest Classifier Mean Feature Importance Between Folds', size=16)

plt.show()

```

```

In [ ]: X_test_Survived=leaderboard_model.predict(X_test)
X_ans=pd.DataFrame(X_test_Survived)
X_ans

```

```
In [ ]: final=pd.DataFrame(columns=['Passenger Id','Survived'])
final['Passenger Id']=df_test['PassengerId']
final=final.reset_index()
final['Survived']=X_ans
final
```