

0. Introduction

I decided to write this kernel because **Titanic: Machine Learning from Disaster** is one of my favorite competitions on Kaggle. This is a beginner level kernel which focuses on **Exploratory Data Analysis** and **Feature Engineering**. A lot of people start Kaggle with this competition and they get lost in extremely long tutorial kernels. This is a short kernel compared to the other ones. I hope this will be a good guide for starters and inspire them with new feature engineering ideas.

Titanic: Machine Learning from Disaster is a great competition to apply domain knowledge for feature engineering, so I made a research and learned a lot about Titanic. There are many secrets to be revealed beneath the Titanic dataset. I tried to find out some of those secret factors that had affected the survival of passengers when the Titanic was sinking. I believe there are other features still waiting to be discovered.

This kernel has **3** main sections; **Exploratory Data Analysis**, **Feature Engineering** and **Model**, and it can achieve top **2% (0.83732)** public leaderboard score with a tuned Random Forest Classifier. It takes 60 seconds to run whole notebook. If you have any idea that might improve this kernel, please be sure to comment, or fork and experiment as you like. If you didn't understand any part, feel free to ask.

In [1]:

```
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
```

- Training set has **891** rows and test set has **418** rows
- Training set have **12** features and test set have **11** features
- One extra feature in training set is Survived feature, which is the target variable

In [2]:

```
def concat_df(train_data, test_data):
    # Returns a concatenated df of training and test set
    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('../input/train.csv')
df_test = pd.read_csv('../input/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)
```

```
Number of Training Examples = 891
```

```
Number of Test Examples = 418
```

```
Training X Shape = (891, 12)
```

```
Training y Shape = 891
```

```
Test X Shape = (418, 11)
```

```
Test y Shape = 418
```

```
Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',  
       'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],  
      dtype='object')  
Index(['PassengerId', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp', 'Parch',  
       'Ticket', 'Fare', 'Cabin', 'Embarked'],  
      dtype='object')
```

1. Exploratory Data Analysis

1.1 Overview

- `PassengerId` is the unique id of the row and it doesn't have any effect on target
- `Survived` is the target variable we are trying to predict (**0 or 1**):
 - **1 = Survived**
 - **0 = Not Survived**
- `Pclass` (Passenger Class) is the socio-economic status of the passenger and it is a categorical ordinal feature which has **3** unique values (**1, 2 or 3**):
 - **1 = Upper Class**
 - **2 = Middle Class**
 - **3 = Lower Class**
- `Name`, `Sex` and `Age` are self-explanatory
- `SibSp` is the total number of the passengers' siblings and spouse
- `Parch` is the total number of the passengers' parents and children
- `Ticket` is the ticket number of the passenger
- `Fare` is the passenger fare
- `Cabin` is the cabin number of the passenger
- `Embarked` is port of embarkation and it is a categorical feature which has **3** unique values (**C, Q or S**):
 - **C = Cherbourg**
 - **Q = Queenstown**
 - **S = Southampton**

In [3]:

```
print(df_train.info())
df_train.sample(3)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
PassengerId    891 non-null int64
Survived       891 non-null int64
Pclass          891 non-null int64
Name            891 non-null object
Sex             891 non-null object
Age             714 non-null float64
SibSp           891 non-null int64
Parch           891 non-null int64
Ticket          891 non-null object
Fare            891 non-null float64
Cabin           204 non-null object
Embarked        889 non-null object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.6+ KB
None
```

Out[3]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
689	690	1	1	Madill, Miss. Georgette Alexandra	female	15.0	0	1	24160	211.
525	526	0	3	Farrell, Mr. James	male	40.5	0	0	367232	7.75
278	279	0	3	Rice, Master. Eric	male	7.0	4	1	382652	29.1

In [4]:

```
print(df_test.info())
df_test.sample(3)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 418 entries, 0 to 417
Data columns (total 11 columns):
PassengerId    418 non-null int64
Pclass          418 non-null int64
Name            418 non-null object
Sex             418 non-null object
Age             332 non-null float64
SibSp           418 non-null int64
Parch           418 non-null int64
Ticket          418 non-null object
Fare            417 non-null float64
Cabin           91 non-null object
Embarked        418 non-null object
dtypes: float64(2), int64(4), object(5)
memory usage: 36.0+ KB
None
```

Out[4]:

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
395	1287	1	Smith, Mrs. Lucien Philip (Mary Eloise Hughes)	female	18.0	1	0	13695	60.0000	C31
310	1202	3	Cacic, Mr. Jego Grga	male	18.0	0	0	315091	8.6625	NaN
203	1095	2	Quick, Miss. Winifred Vera	female	8.0	1	1	26360	26.0000	NaN

1.2 Missing Values

As seen from below, some columns have missing values. `display_missing` function shows the count of missing values in every column in both training and test set.

- Training set have missing values in `Age` , `Cabin` and `Embarked` columns
- Test set have missing values in `Age` , `Cabin` and `Fare` columns

It is convenient to work on concatenated training and test set while dealing with missing values, otherwise filled data may overfit to training or test set samples. The count of missing values in `Age` , `Embarked` and `Fare` are smaller compared to total sample, but roughly **80%** of the `Cabin` is missing. Missing values in `Age` , `Embarked` and `Fare` can be filled with descriptive statistical measures but that wouldn't work for `Cabin` .

▼ Show hidden code

Training Set

```
PassengerId column missing values: 0
Survived column missing values: 0
Pclass column missing values: 0
Name column missing values: 0
Sex column missing values: 0
Age column missing values: 177
SibSp column missing values: 0
Parch column missing values: 0
Ticket column missing values: 0
Fare column missing values: 0
Cabin column missing values: 687
Embarked column missing values: 2
```

Test Set

```
PassengerId column missing values: 0
Pclass column missing values: 0
Name column missing values: 0
Sex column missing values: 0
Age column missing values: 86
SibSp column missing values: 0
Parch column missing values: 0
Ticket column missing values: 0
Fare column missing values: 1
Cabin column missing values: 327
Embarked column missing values: 0
```

1.2.1 Age

Missing values in `Age` are filled with median age, but using median age of the whole data set is not a good choice. Median age of `Pclass` groups is the best choice because of its high correlation with `Age (0.408106)` and `Survived (0.338481)`. It is also more logical to group ages by passenger classes instead of other features.

▽ Show hidden code

Out[6] :

	Feature 1	Feature 2	Correlation Coefficient
6	Age	Age	1.000000
9	Age	Pclass	0.408106
17	Age	SibSp	0.243699
22	Age	Fare	0.178740
25	Age	Parch	0.150917
29	Age	Survived	0.077221
41	Age	PassengerId	0.028814

In order to be more accurate, `Sex` feature is used as the second level of `groupby` while filling the missing `Age` values. As seen from below, `Pclass` and `Sex` groups have distinct median `Age` values. When passenger class increases, the median age for both males and females also increases. However, females tend to have slightly lower median `Age` than males. The median ages below are used for filling the missing values in `Age` feature.

In [7]:

```
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]))
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()))
```

```
Median age of Pclass 1 females: 36.0
Median age of Pclass 1 males: 42.0
Median age of Pclass 2 females: 28.0
Median age of Pclass 2 males: 29.5
Median age of Pclass 3 females: 22.0
Median age of Pclass 3 males: 25.0
Median age of all passengers: 28.0
```

1.2.2 Embarked

`Embarked` is a categorical feature and there are only **2** missing values in whole data set. Both of those passengers are female, upper class and they have the same ticket number. This means that they know each other and embarked from the same port together. The mode `Embarked` value for an upper class female passenger is **C (Cherbourg)**, but this doesn't necessarily mean that they embarked from that port.

>Show hidden code

Out[8]:

	Age	Cabin	Embarked	Fare	Name	Parch	PassengerId	Pclass	Sex	SibSp	...
61	38.0	B28	NaN	80.0	Icard, Miss. Amelie	0	62	1	female	0	
829	62.0	B28	NaN	80.0	Stone, Mrs. George Nelson (Martha Evelyn)	0	830	1	female	0	

When I googled **Stone, Mrs. George Nelson (Martha Evelyn)**, I found that she embarked from **S (Southampton)** with her maid **Amelie Icard**, in this page [Martha Evelyn Stone: Titanic Survivor](https://www.encyclopedia-titanica.org/titanic-survivor/martha-evelyn-stone.html) (<https://www.encyclopedia-titanica.org/titanic-survivor/martha-evelyn-stone.html>).

Mrs Stone boarded the Titanic in Southampton on 10 April 1912 and was travelling in first class with her maid Amelie Icard. She occupied cabin B-28.

Missing values in `Embarked` are filled with **S** with this information.

In [9]:

```
# Filling the missing values in Embarked with S
df_all['Embarked'] = df_all['Embarked'].fillna('S')
```

1.2.3 Fare

There is only one passenger with missing `Fare` value. We can assume that `Fare` is related to family size (`Parch` and `SibSp`) and `Pclass` features. Median `Fare` value of a male with a third class ticket and no family is a logical choice to fill the missing value.

▼ Show hidden code

Out[10]:

	Age	Cabin	Embarked	Fare	Name	Parch	PassengerId	Pclass	Sex	SibSp	S
1043	60.5	NaN	S	NaN	Storey, Mr. Thomas	0	1044	3	male	0	N

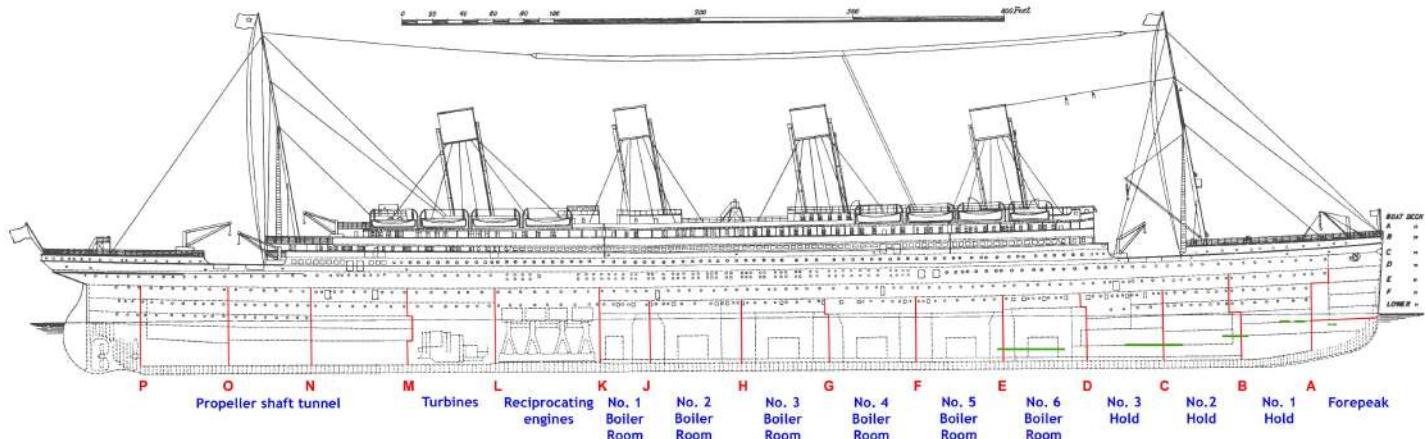
◀ ▶

In [11]:

```
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
# passenger
df_all['Fare'] = df_all['Fare'].fillna(med_fare)
```

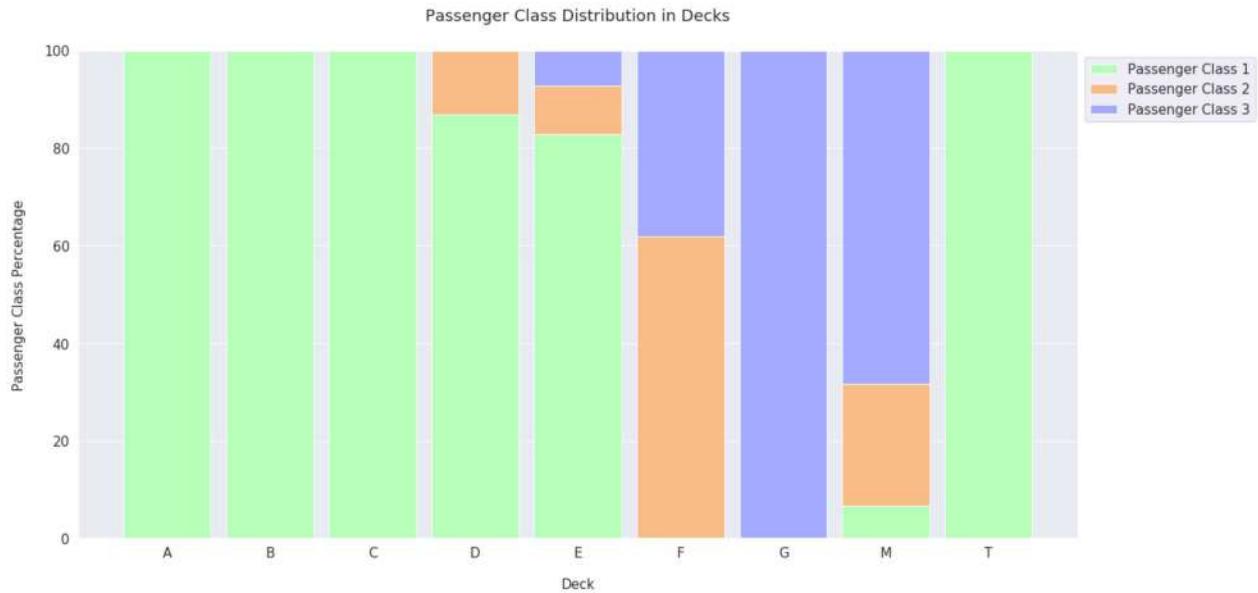
1.2.4 Cabin

Cabin feature is little bit tricky and it needs further exploration. The large portion of the Cabin feature is missing and the feature itself can't be ignored completely because some the cabins might have higher survival rates. It turns out to be the first letter of the Cabin values are the decks in which the cabins are located. Those decks were mainly separated for one passenger class, but some of them were used by multiple passenger classes.



- On the Boat Deck there were **6** rooms labeled as **T, U, W, X, Y, Z** but only the **T** cabin is present in the dataset
- **A, B** and **C** decks were only for 1st class passengers
- **D** and **E** decks were for all classes
- **F** and **G** decks were for both 2nd and 3rd class passengers
- From going **A** to **G**, distance to the staircase increases which might be a factor of survival

▼ Show hidden code

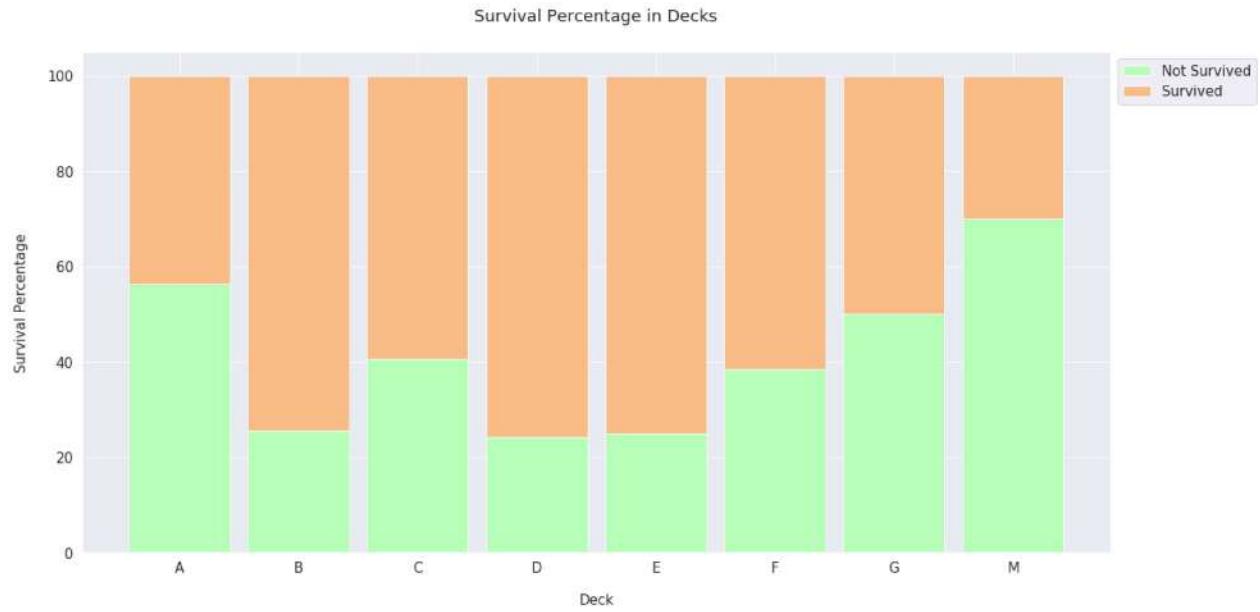


- **100%** of **A, B** and **C** decks are 1st class passengers
- Deck **D** has **87%** 1st class and **13%** 2nd class passengers
- Deck **E** has **83%** 1st class, **10%** 2nd class and **7%** 3rd class passengers
- Deck **F** has **62%** 2nd class and **38%** 3rd class passengers
- **100%** of **G** deck are 3rd class passengers
- There is one person on the boat deck in **T** cabin and he is a 1st class passenger. **T** cabin passenger has the closest resemblance to **A** deck passengers so he is grouped with **A** deck
- Passengers labeled as **M** are the missing values in Cabin feature. I don't think it is possible to find those passengers' real Deck so I decided to use **M** like a deck

In [13]:

```
# Passenger in the T deck is changed to A
idx = df_all[df_all['Deck'] == 'T'].index
df_all.loc[idx, 'Deck'] = 'A'
```

▼ Show hidden code



As I suspected, every deck has different survival rates and that information can't be discarded. Deck **B**, **C**, **D** and **E** have the highest survival rates. Those decks are mostly occupied by 1st class passengers. **M** has the lowest survival rate which is mostly occupied by 2nd and 3rd class passengers. To conclude, cabins used by 1st class passengers have higher survival rates than cabins used by 2nd and 3rd class passengers. In my opinion **M** (Missing Cabin values) has the lowest survival rate because they couldn't retrieve the cabin data of the victims. That's why I believe labeling that group as **M** is a reasonable way to handle the missing data. It is a unique group with shared characteristics. Deck feature has high-cardinality right now so some of the values are grouped with each other based on their similarities.

- **A**, **B** and **C** decks are labeled as **ABC** because all of them have only 1st class passengers
- **D** and **E** decks are labeled as **DE** because both of them have similar passenger class distribution and same survival rate
- **F** and **G** decks are labeled as **FG** because of the same reason above
- **M** deck doesn't need to be grouped with other decks because it is very different from others and has the lowest survival rate.

```
In [15]:
```

```
df_all['Deck'] = df_all['Deck'].replace(['A', 'B', 'C'], 'ABC')
df_all['Deck'] = df_all['Deck'].replace(['D', 'E'], 'DE')
df_all['Deck'] = df_all['Deck'].replace(['F', 'G'], 'FG')

df_all['Deck'].value_counts()
```

```
Out[15]:
```

```
M      1014
ABC    182
DE     87
FG     26
Name: Deck, dtype: int64
```

After filling the missing values in `Age`, `Embarked`, `Fare` and `Deck` features, there is no missing value left in both training and test set. `Cabin` is dropped because `Deck` feature is used instead of it.

▽ Show hidden code

```
Age column missing values: 0
Embarked column missing values: 0
Fare column missing values: 0
Name column missing values: 0
Parch column missing values: 0
PassengerId column missing values: 0
Pclass column missing values: 0
Sex column missing values: 0
SibSp column missing values: 0
Survived column missing values: 0
Ticket column missing values: 0
Deck column missing values: 0
```

```
Age column missing values: 0
Embarked column missing values: 0
Fare column missing values: 0
Name column missing values: 0
Parch column missing values: 0
PassengerId column missing values: 0
Pclass column missing values: 0
Sex column missing values: 0
SibSp column missing values: 0
Ticket column missing values: 0
Deck column missing values: 0
```

1.3 Target Distribution

- **38.38%** (342/891) of training set is **Class 1**
- **61.62%** (549/891) of training set is **Class 0**

▽ Show hidden code

342 of 891 passengers survived and it is the 38.38% of the training set.

549 of 891 passengers didnt survive and it is the 61.62% of the training set.



1.4 Correlations

Features are highly correlated with each other and dependent to each other. The highest correlation between features is **0.549500** in training set and **0.577147** in test set (between `Fare` and `Pclass`). The other features are also highly correlated. There are **9** correlations in training set and **6** correlations in test set that are higher than **0.1**.

>Show hidden code

In [19]:

```
# Training set high correlations
corr = df_train_corr_nd['Correlation Coefficient'] > 0.1
df_train_corr_nd[corr]
```

Out[19]:

	Feature 1	Feature 2	Correlation Coefficient
6	Pclass	Fare	0.549500
8	Pclass	Age	0.417667
10	SibSp	Parch	0.414838
12	Survived	Pclass	0.338481
14	Survived	Fare	0.257307
16	SibSp	Age	0.249747
18	Parch	Fare	0.216225
20	Age	Parch	0.176733
22	SibSp	Fare	0.159651
24	Age	Fare	0.124061

In [20]:

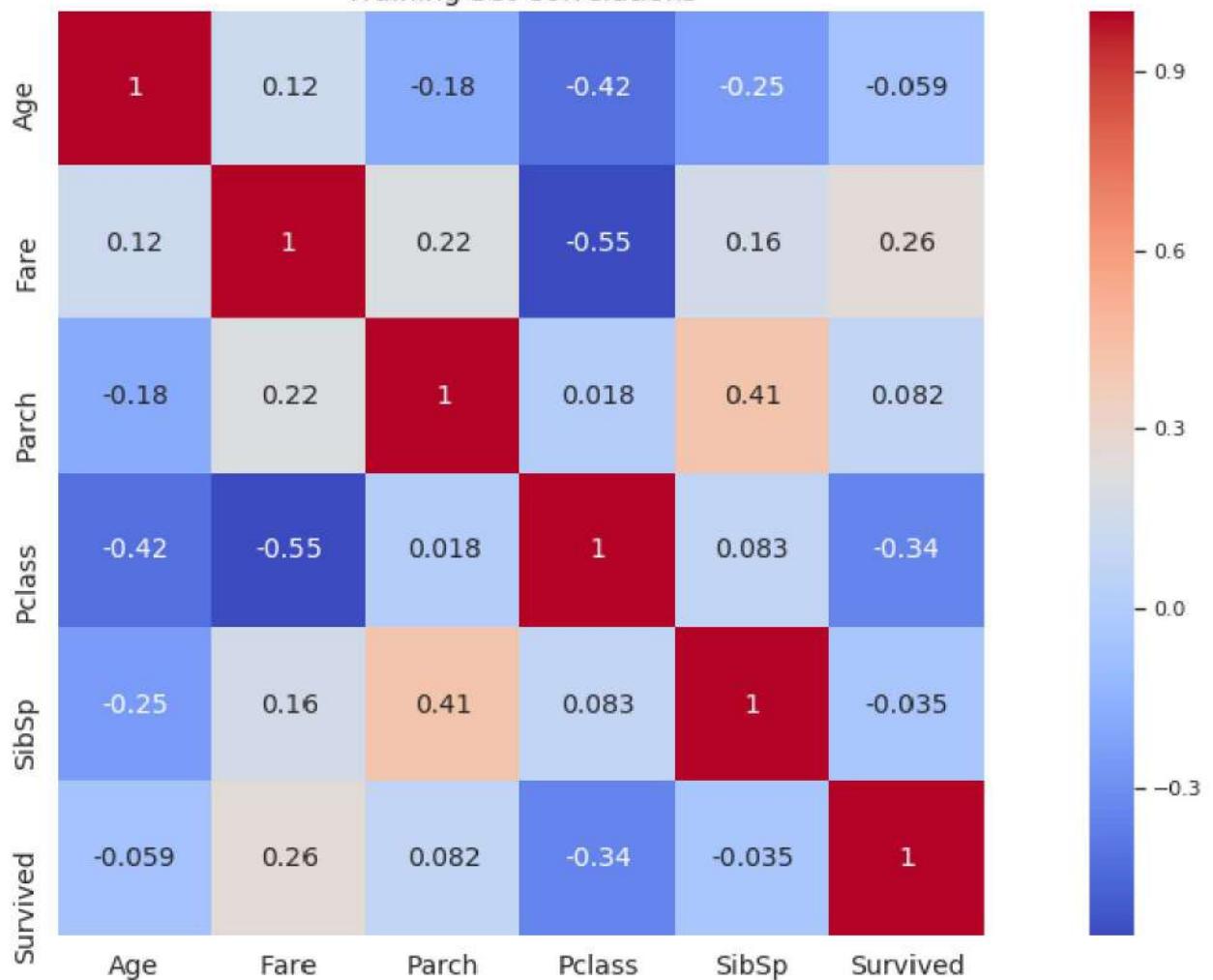
```
# Test set high correlations
corr = df_test_corr_nd['Correlation Coefficient'] > 0.1
df_test_corr_nd[corr]
```

Out[20]:

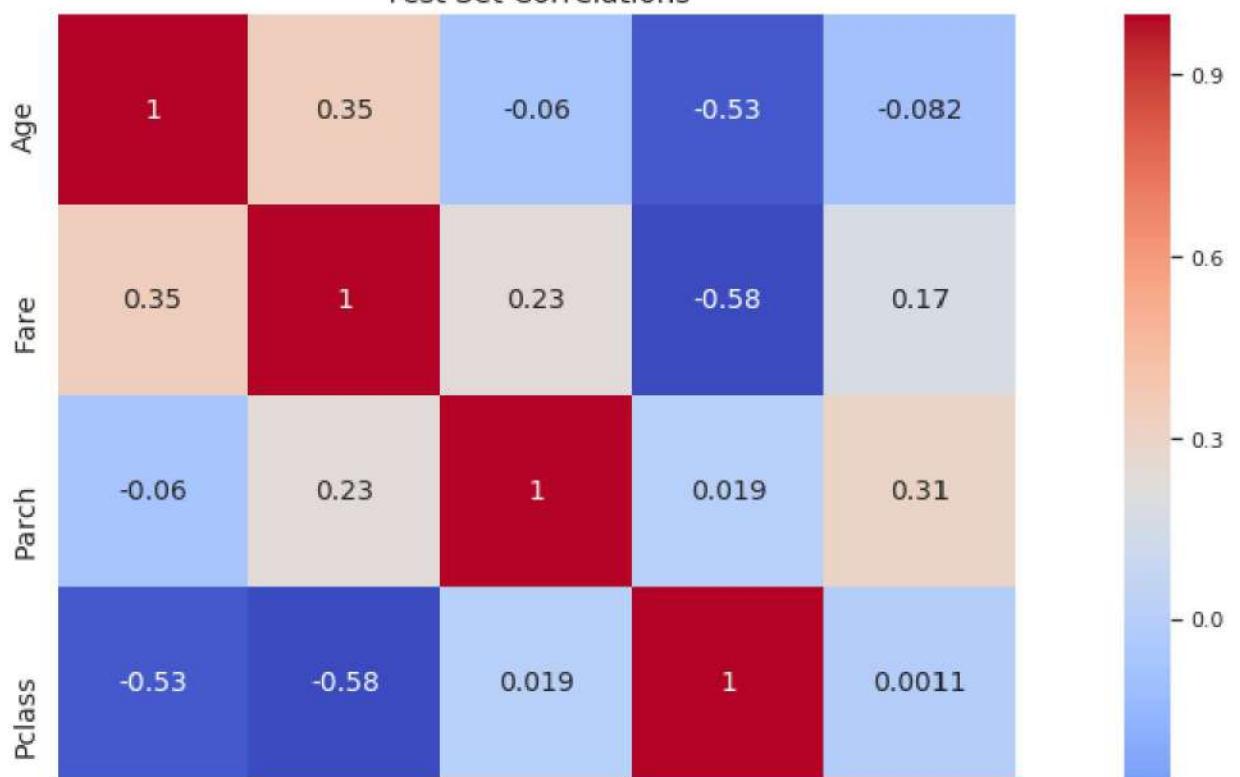
	Feature 1	Feature 2	Correlation Coefficient
6	Fare	Pclass	0.577489
8	Age	Pclass	0.526789
10	Age	Fare	0.345347
12	SibSp	Parch	0.306895
14	Fare	Parch	0.230410
16	SibSp	Fare	0.172032

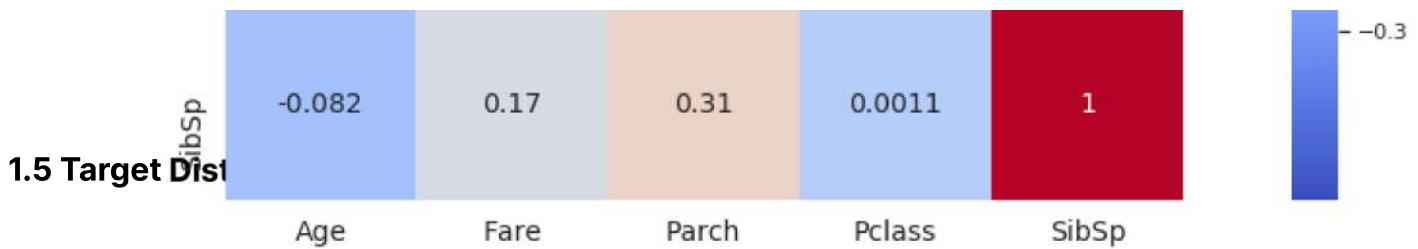
 Show hidden code

Training Set Correlations



Test Set Correlations



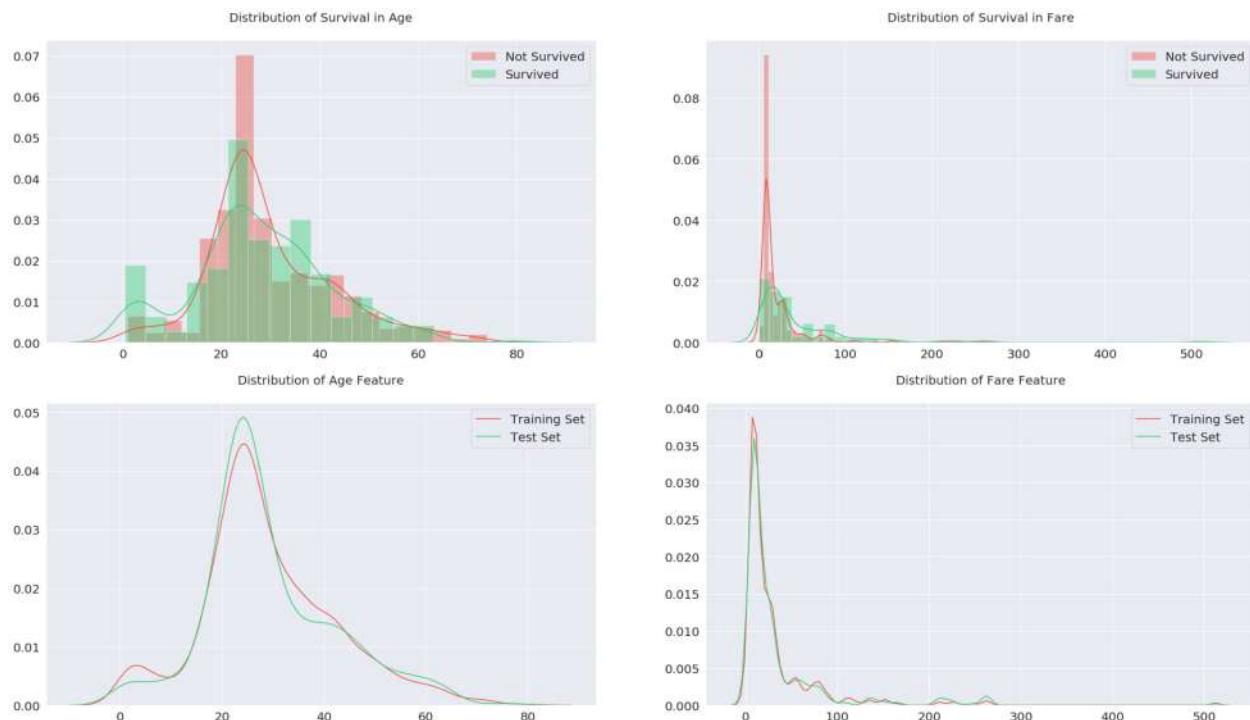


1.5.1 Continuous Features

Both of the continuous features (Age and Fare) have good split points and spikes for a decision tree to learn. One potential problem for both features is, the distribution has more spikes and bumps in training set, but it is smoother in test set. Model may not be able to generalize to test set because of this reason.

- Distribution of Age feature clearly shows that children younger than 15 has a higher survival rate than any of the other age groups
- In distribution of Fare feature, the survival rate is higher on distribution tails. The distribution also has positive skew because of the extremely large outliers

>Show hidden code

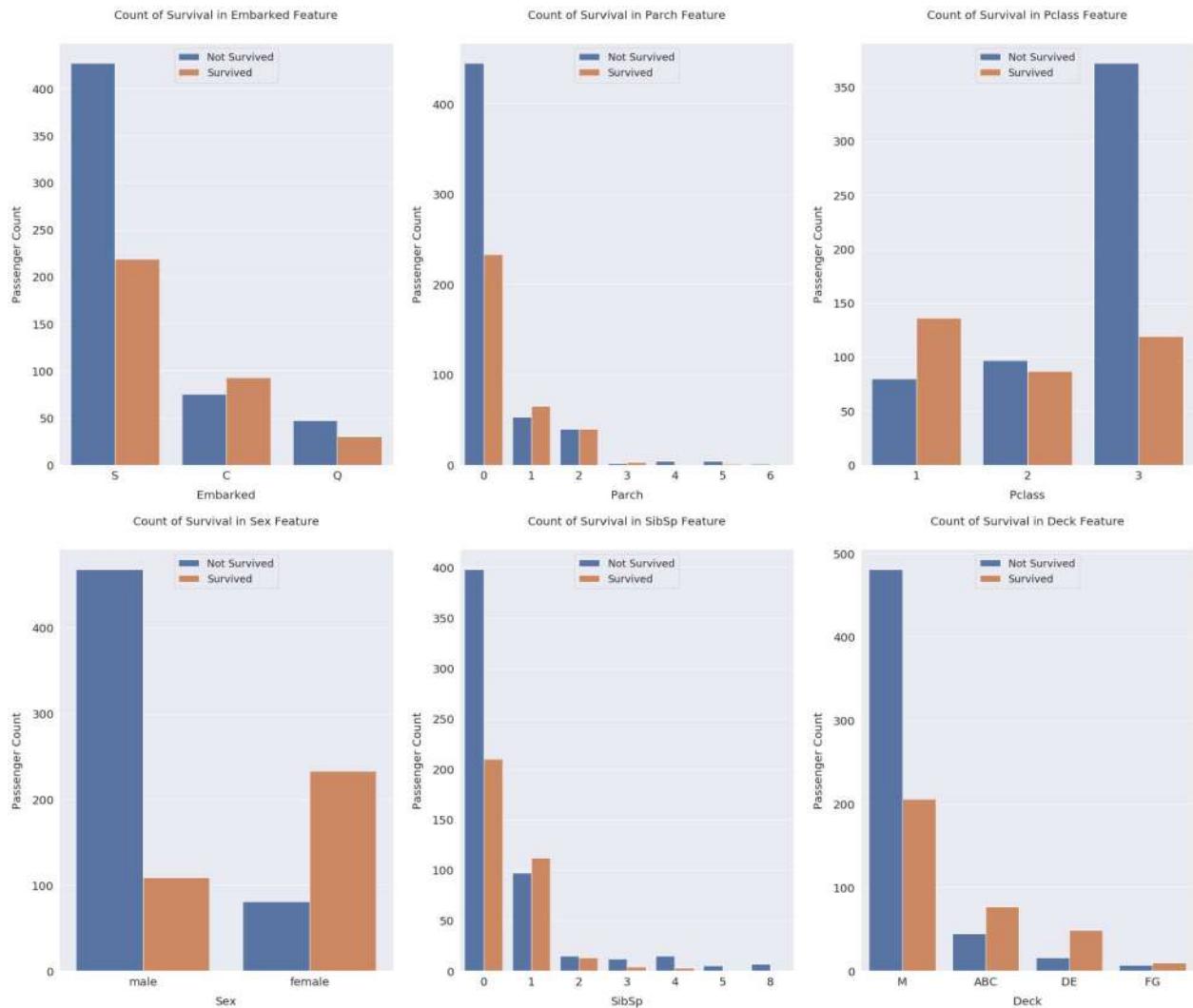


1.5.2 Categorical Features

Every categorical feature has at least one class with high mortality rate. Those classes are very helpful to predict whether the passenger is a survivor or victim. Best categorical features are `Pclass` and `Sex` because they have the most homogenous distributions.

- Passengers boarded from **Southampton** has a lower survival rate unlike other ports. More than half of the passengers boarded from **Cherbourg** had survived. This observation could be related to `Pclass` feature
- `Parch` and `SibSp` features show that passengers with only one family member has a higher survival rate

▼ Show hidden code



1.6 Conclusion

Most of the features are correlated with each other. This relationship can be used to create new features with feature transformation and feature interaction. Target encoding could be very useful as well because of the high correlations with `Survived` feature.

Split points and spikes are visible in continuous features. They can be captured easily with a decision tree model, but linear models may not be able to spot them.

Categorical features have very distinct distributions with different survival rates. Those features can be one-hot encoded. Some of those features may be combined with each other to make new features.

Created a new feature called `Deck` and dropped `Cabin` feature at the **Exploratory Data Analysis** part.

>Show hidden code

Out[24]:

	Age	Deck	Embarked	Fare	Name	Parch	PassengerId	Pclass	Sex	SibSp
0	22.0	M	S	7.2500	Braund, Mr. Owen Harris	0	1	3	male	1
1	38.0	ABC	C	71.2833	Cumings, Mrs. John Bradley (Florence Briggs Th...)	0	2	1	female	1
2	26.0	M	S	7.9250	Heikkinen, Miss. Laina	0	3	3	female	0
3	35.0	ABC	S	53.1000	Futrelle, Mrs. Jacques Heath (Lily May Peel)	0	4	1	female	1
4	35.0	M	S	8.0500	Allen, Mr. William Henry	0	5	3	male	0

2. Feature Engineering

2.1 Binning Continuous Features

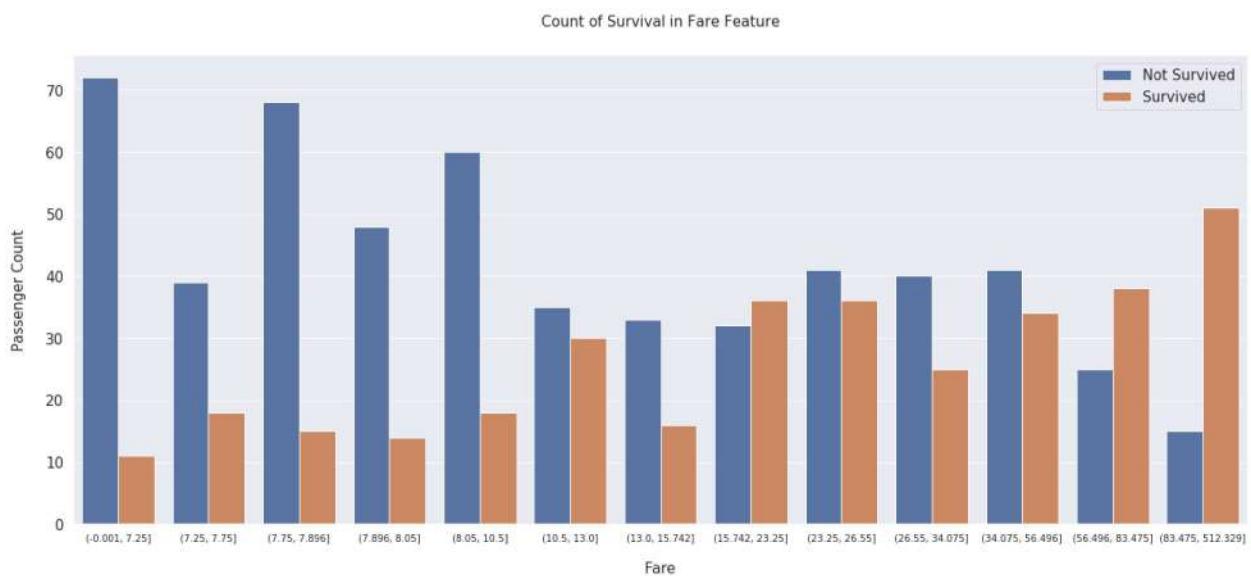
2.1.1 Fare

`Fare` feature is positively skewed and survival rate is extremely high on the right end. **13** quantile based bins are used for `Fare` feature. Even though the bins are too much, they provide decent amount of information gain. The groups at the left side of the graph has the lowest survival rate and the groups at the right side of the graph has the highest survival rate. This high survival rate was not visible in the distribution graph. There is also an unusual group **[15.742, 23.25]** in the middle with high survival rate that is captured in this process.

In [25]:

```
df_all['Fare'] = pd.qcut(df_all['Fare'], 13)
```

>Show hidden code



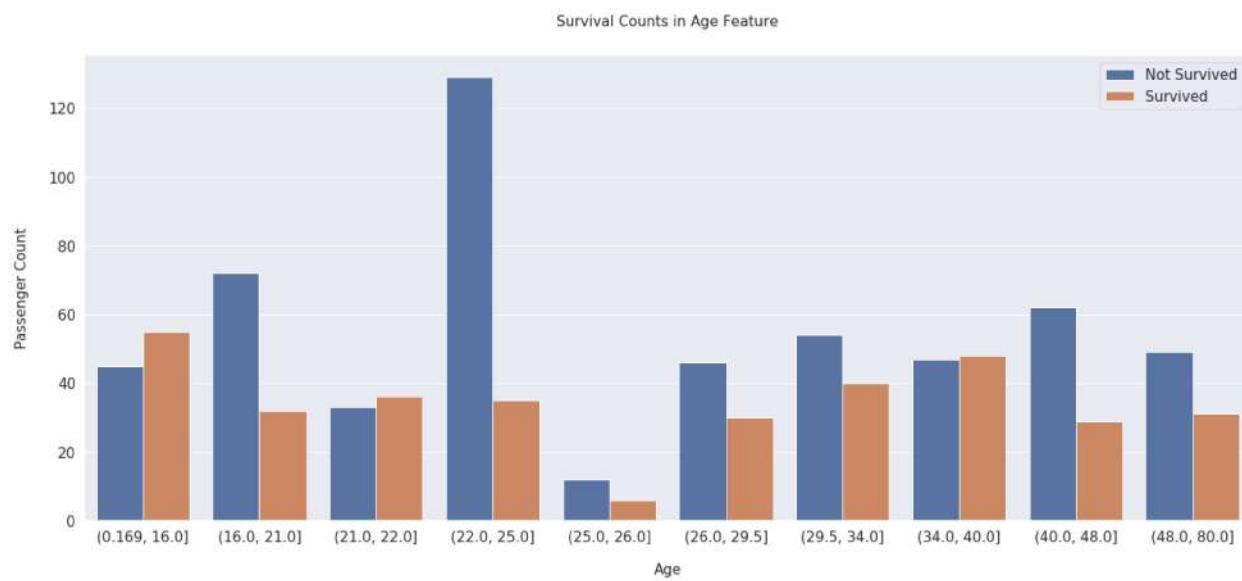
2.1.2 Age

`Age` feature has a normal distribution with some spikes and bumps and **10** quantile based bins are used for `Age`. The first bin has the highest survival rate and 4th bin has the lowest survival rate. Those were the biggest spikes in the distribution. There is also an unusual group (**34.0, 40.0**] with high survival rate that is captured in this process.

In [27]:

```
df_all['Age'] = pd.qcut(df_all['Age'], 10)
```

▽ Show hidden code

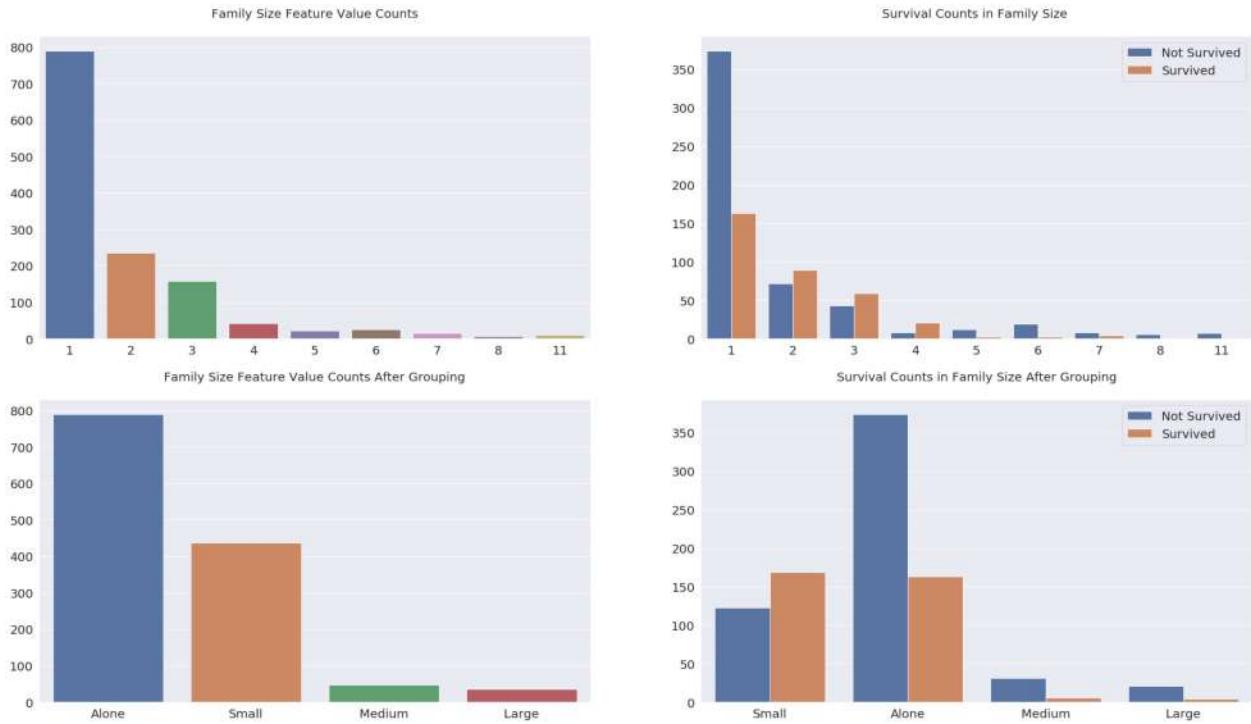


2.2 Frequency Encoding

`Family_Size` is created by adding `SibSp`, `Parch` and **1**. `SibSp` is the count of siblings and spouse, and `Parch` is the count of parents and children. Those columns are added in order to find the total size of families. Adding **1** at the end, is the current passenger. Graphs have clearly shown that family size is a predictor of survival because different values have different survival rates.

- Family Size with **1** are labeled as **Alone**
- Family Size with **2, 3** and **4** are labeled as **Small**
- Family Size with **5** and **6** are labeled as **Medium**
- Family Size with **7, 8** and **11** are labeled as **Large**

▽ Show hidden code



There are too many unique `Ticket` values to analyze, so grouping them up by their frequencies makes things easier.

How is this feature different than `Family_Size`? Many passengers travelled along with groups. Those groups consist of friends, nannies, maids and etc. They weren't counted as family, but they used the same ticket.

Why not grouping tickets by their prefixes? If prefixes in `Ticket` feature has any meaning, then they are already captured in `Pclass` or `Embarked` features because that could be the only logical information which can be derived from the `Ticket` feature.

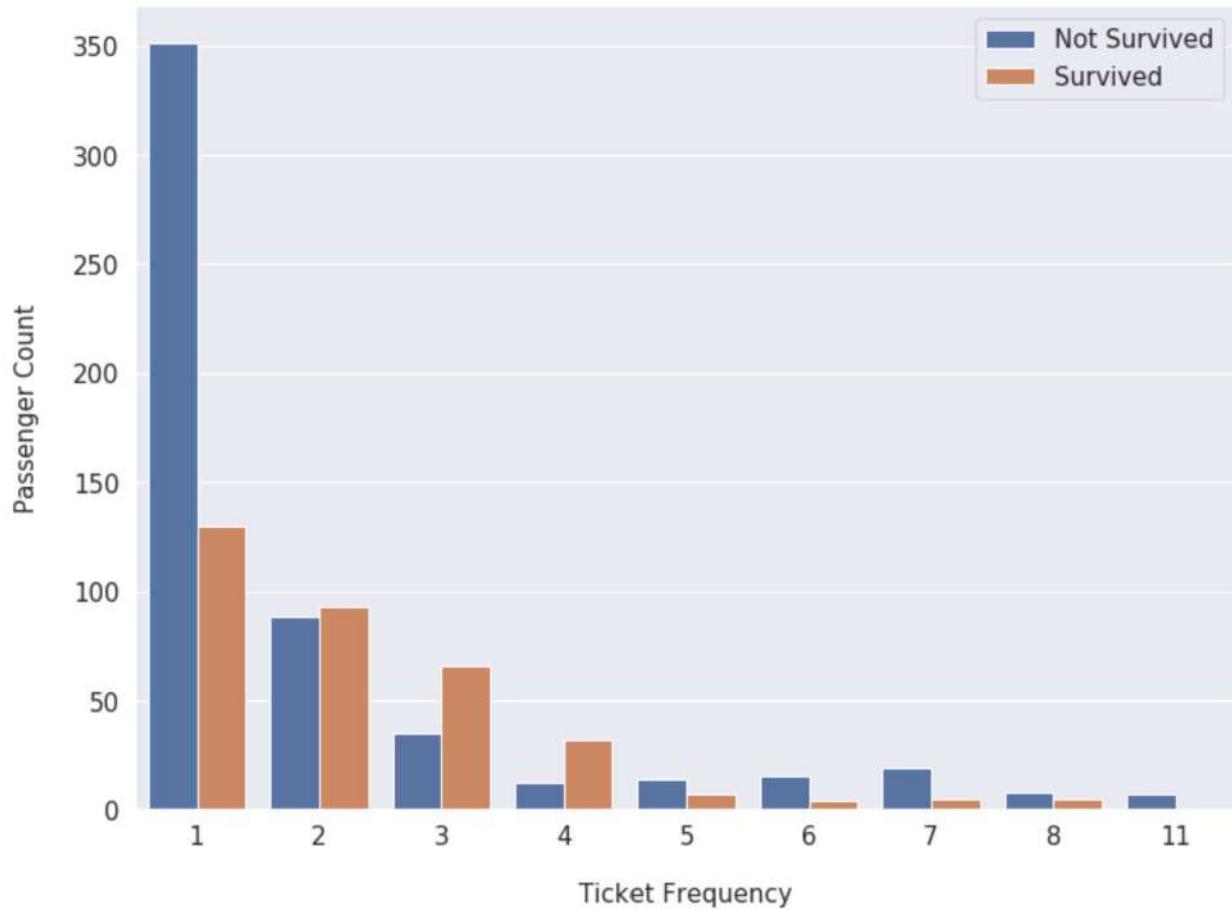
According to the graph below, groups with **2,3** and **4** members had a higher survival rate. Passengers who travel alone has the lowest survival rate. After **4** group members, survival rate decreases drastically. This pattern is very similar to `Family_Size` feature but there are minor differences. `Ticket_Frequency` values are not grouped like `Family_Size` because that would basically create the same feature with perfect correlation. This kind of feature wouldn't provide any additional information gain.

In [30]:

```
df_all['Ticket_Frequency'] = df_all.groupby('Ticket')['Ticket'].transform('count')
```

▽ Show hidden code

Count of Survival in Ticket Frequency Feature



2.3 Title & Is Married

`Title` is created by extracting the prefix before `Name` feature. According to graph below, there are many titles that are occurring very few times. Some of those titles doesn't seem correct and they need to be replaced. **Miss, Mrs, Ms, Mlle, Lady, Mme, the Countess, Dona** titles are replaced with **Miss/Mrs/Ms** because all of them are female. Values like **Mlle, Mme** and **Dona** are actually the name of the passengers, but they are classified as titles because `Name` feature is split by comma. **Dr, Col, Major, Jonkheer, Capt, Sir, Don** and **Rev** titles are replaced with **Dr/Military/Noble/Clergy** because those passengers have similar characteristics. **Master** is a unique title. It is given to male passengers below age **26**. They have the highest survival rate among all males.

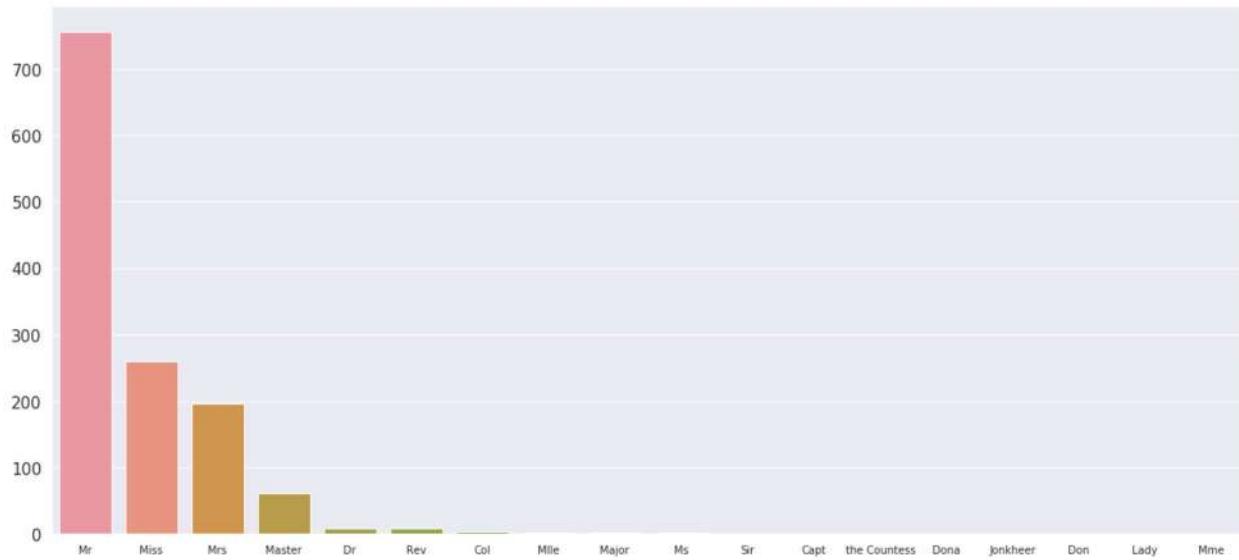
`Is_Married` is a binary feature based on the **Mrs** title. **Mrs** title has the highest survival rate among other female titles. This title needs to be a feature because all female titles are grouped with each other.

In [32]:

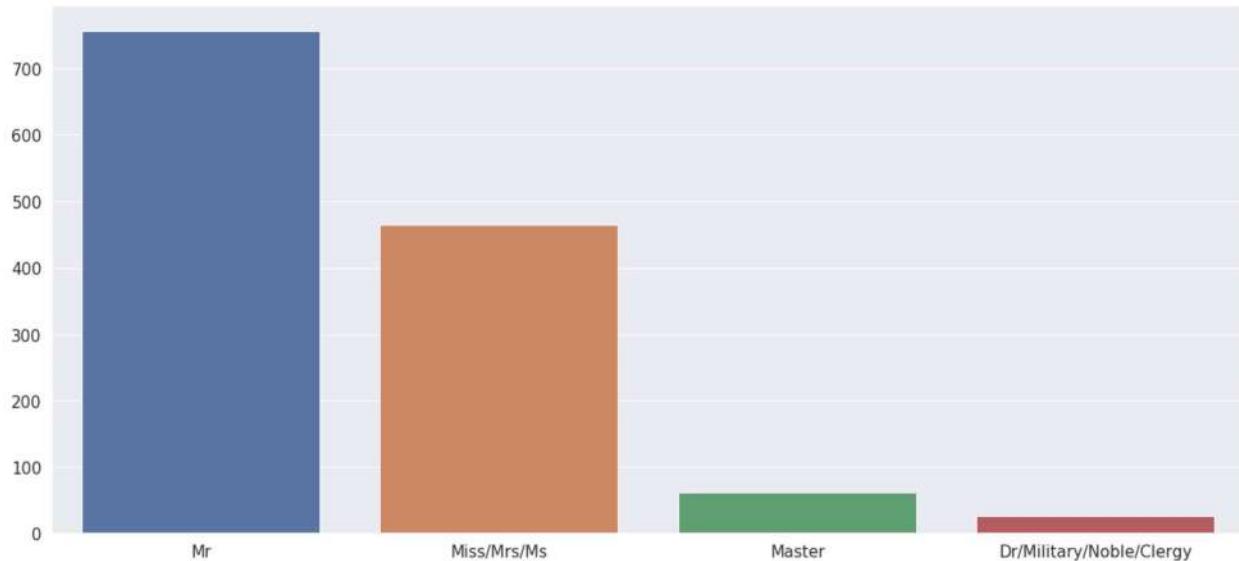
```
df_all['Title'] = df_all['Name'].str.split(', ', expand=True)[1].str.split('. ', expand=True)[0]
df_all['Is_Married'] = 0
df_all['Is_Married'].loc[df_all['Title'] == 'Mrs'] = 1
```

>Show hidden code

Title Feature Value Counts



Title Feature Value Counts After Grouping



2.4 Target Encoding

`extract_surname` function is used for extracting surnames of passengers from the `Name` feature. `Family` feature is created with the extracted surname. This is necessary for grouping passengers in the same family.

In [34]:

```
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.loc[:890]
df_test = df_all.loc[891:]
dfs = [df_train, df_test]
```

`Family_Survival_Rate` is calculated from families in training set since there is no `Survived` feature in test set. A list of family names that are occurring in both training and test set (`non_unique_families`), is created. The survival rate is calculated for families with more than 1 members in that list, and stored in `Family_Survival_Rate` feature.

An extra binary feature `Family_Survival_Rate_NA` is created for families that are unique to the test set. This feature is also necessary because there is no way to calculate those families' survival rate. This feature implies that family survival rate is not applicable to those passengers because there is no way to retrieve their survival rate.

`Ticket_Survival_Rate` and `Ticket_Survival_Rate_NA` features are also created with the same method. `Ticket_Survival_Rate` and `Family_Survival_Rate` are averaged and become `Survival_Rate`, and `Ticket_Survival_Rate_NA` and `Family_Survival_Rate_NA` are also averaged and become `Survival_Rate_NA`.

In [35]:

```
# Creating a list of families and tickets that are occurring in both training and test set
non_unique_families = [x for x in df_train['Family'].unique() if x in df_test['Family'].unique()]
non_unique_tickets = [x for x in df_train['Ticket'].unique() if x in df_test['Ticket'].unique()]

df_family_survival_rate = df_train.groupby('Family')[['Survived', 'Family', 'Family_Size']].median()
df_ticket_survival_rate = df_train.groupby('Ticket')[['Survived', 'Ticket', 'Ticket_Frequency']].median()

family_rates = {}
ticket_rates = {}

for i in range(len(df_family_survival_rate)):
    # Checking a family exists in both training and test set, and has members more than 1
    if df_family_survival_rate.index[i] in non_unique_families and df_family_survival_rate.iloc[i, 1] > 1:
        family_rates[df_family_survival_rate.index[i]] = df_family_survival_rate.iloc[i, 0]

for i in range(len(df_ticket_survival_rate)):
    # Checking a ticket exists in both training and test set, and has members more than 1
    if df_ticket_survival_rate.index[i] in non_unique_tickets and df_ticket_survival_rate.iloc[i, 1] > 1:
        ticket_rates[df_ticket_survival_rate.index[i]] = df_ticket_survival_rate.iloc[i, 0]
```

In [36]:

```
mean_survival_rate = np.mean(df_train['Survived'])

train_family_survival_rate = []
train_family_survival_rate_NA = []
test_family_survival_rate = []
test_family_survival_rate_NA = []

for i in range(len(df_train)):
    if df_train['Family'][i] in family_rates:
        train_family_survival_rate.append(family_rates[df_train['Family']
[i]])
        train_family_survival_rate_NA.append(1)
    else:
        train_family_survival_rate.append(mean_survival_rate)
        train_family_survival_rate_NA.append(0)

for i in range(len(df_test)):
    if df_test['Family'].iloc[i] in family_rates:
        test_family_survival_rate.append(family_rates[df_test['Family'].i
loc[i]])
        test_family_survival_rate_NA.append(1)
    else:
        test_family_survival_rate.append(mean_survival_rate)
        test_family_survival_rate_NA.append(0)

df_train['Family_Survival_Rate'] = train_family_survival_rate
df_train['Family_Survival_Rate_NA'] = train_family_survival_rate_NA
df_test['Family_Survival_Rate'] = test_family_survival_rate
df_test['Family_Survival_Rate_NA'] = test_family_survival_rate_NA

train_ticket_survival_rate = []
train_ticket_survival_rate_NA = []
test_ticket_survival_rate = []
test_ticket_survival_rate_NA = []

for i in range(len(df_train)):
    if df_train['Ticket'][i] in ticket_rates:
        train_ticket_survival_rate.append(ticket_rates[df_train['Ticket']
[i]])
        train_ticket_survival_rate_NA.append(1)
    else:
        train_ticket_survival_rate.append(mean_survival_rate)
        train_ticket_survival_rate_NA.append(0)
```

```

    else:
        train_ticket_survival_rate.append(mean_survival_rate)
        train_ticket_survival_rate_NA.append(0)

    for i in range(len(df_test)):
        if df_test['Ticket'].iloc[i] in ticket_rates:
            test_ticket_survival_rate.append(ticket_rates[df_test['Ticket'].iloc[i]])
            test_ticket_survival_rate_NA.append(1)
        else:
            test_ticket_survival_rate.append(mean_survival_rate)
            test_ticket_survival_rate_NA.append(0)

    df_train['Ticket_Survival_Rate'] = train_ticket_survival_rate
    df_train['Ticket_Survival_Rate_NA'] = train_ticket_survival_rate_NA
    df_test['Ticket_Survival_Rate'] = test_ticket_survival_rate
    df_test['Ticket_Survival_Rate_NA'] = test_ticket_survival_rate_NA

```

In [37]:

```

for df in [df_train, df_test]:
    df['Survival_Rate'] = (df['Ticket_Survival_Rate'] + df['Family_Survival_Rate']) / 2
    df['Survival_Rate_NA'] = (df['Ticket_Survival_Rate_NA'] + df['Family_Survival_Rate_NA']) / 2

```

2.5 Feature Transformation

2.5.1 Label Encoding Non-Numerical Features

`Embarked`, `Sex`, `Deck`, `Title` and `Family_Size_Grouped` are object type, and `Age` and `Fare` features are category type. They are converted to numerical type with `LabelEncoder`. `LabelEncoder` basically labels the classes from `0` to `n`. This process is necessary for models to learn from those features.

```
In [38]:
```

```
non_numeric_features = ['Embarked', 'Sex', 'Deck', 'Title', 'Family_Size_Grouped', 'Age', 'Fare']

for df in dfs:
    for feature in non_numeric_features:
        df[feature] = LabelEncoder().fit_transform(df[feature])
```

2.5.2 One-Hot Encoding the Categorical Features

The categorical features (Pclass , Sex , Deck , Embarked , Title) are converted to one-hot encoded features with OneHotEncoder . Age and Fare features are not converted because they are ordinal unlike the previous ones.

```
In [39]:
```

```
cat_features = ['Pclass', 'Sex', 'Deck', '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(-1, 1)).toarray()
        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[:6]], axis=1)
df_test = pd.concat([df_test, *encoded_features[6:]], axis=1)
```

2.6 Conclusion

`Age` and `Fare` features are binned. Binning helped dealing with outliers and it revealed some homogeneous groups in those features. `Family_Size` is created by adding `Parch` and `SibSp` features and `1`. `Ticket_Frequency` is created by counting the occurrence of `Ticket` values.

`Name` feature is very useful. First, `Title` and `Is_Married` features are created from the title prefix in the names. Second, `Family_Survival_Rate` and `Family_Survival_Rate_NA` features are created by target encoding the surname of the passengers. `Ticket_Survival_Rate` is created by target encoding the `Ticket` feature. `Survival_Rate` feature is created by averaging the `Family_Survival_Rate` and `Ticket_Survival_Rate` features.

Finally, the non-numeric type features are label encoded and categorical features are one-hot encoded. Created **5** new features (`Family_Size` , `Title` , `Is_Married` , `Survival_Rate` and `Survival_Rate_NA`) and dropped the useless features after encoding.

>Show hidden code

Out[40]:

	Age	Deck_1	Deck_2	Deck_3	Deck_4	Embarked_1	Embarked_2	Embarked_3	Family_S
0	2	0.0	0.0	0.0	1.0	0.0	0.0	1.0	0.0
1	7	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
2	4	0.0	0.0	0.0	1.0	0.0	0.0	1.0	1.0
3	7	1.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0
4	7	0.0	0.0	0.0	1.0	0.0	0.0	1.0	1.0

3. Model

```
In [41]:  
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))
```

```
X_train shape: (891, 26)  
y_train shape: (891, )  
X_test shape: (418, 26)
```

3.1 Random Forest

Created 2 `RandomForestClassifier`'s. One of them is a single model and the other is for k-fold cross validation.

The highest accuracy of the `single_best_model` is **0.82775** in public leaderboard. However, it doesn't perform better in k-fold cross validation. It is a good model to start experimenting and hyperparameter tuning.

The highest accuracy of `leaderboard_model` is **0.83732** in public leaderboard with 5-fold cross validation. This model is created for leaderboard score and it is tuned to overfit slightly. It is designed to overfit because the estimated probabilities of `X_test` in every fold are going to be divided by **N** (fold count). If this model is used as a single model, it would struggle to predict lots of samples correctly.

Which model should I use?

- `leaderboard_model` overfits to test set so it's not suggested to use models like this in real life projects.
- `single_best_model` is a good model to start experimenting and learning about decision trees.

In [42]:

```
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)
```

`StratifiedKFold` is used for stratifying the target variable. The folds are made by preserving the percentage of samples for each class in target variable (`Survived`).


```

N = 5
oob = 0
probs = pd.DataFrame(np.zeros((len(X_test), N * 2)), columns=['Fold_{}{}_Prob_{}'.format(i, j) for i in range(1, N + 1) for j in range(2)])
importances = pd.DataFrame(np.zeros((X_train.shape[1], N)), columns=['Fold_{}'.format(i) for i in range(1, N + 1)], index=df_all.columns)
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 {}\n'.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])[:, 1])
    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])[:, 1])
    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_sc

```

```
ore_))

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

Fold 1

```
[Parallel(n_jobs=-1)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=-1)]: Done  42 tasks      | elapsed:   0.1s  
[Parallel(n_jobs=-1)]: Done 192 tasks      | elapsed:   0.3s  
[Parallel(n_jobs=-1)]: Done 442 tasks      | elapsed:   0.8s  
[Parallel(n_jobs=-1)]: Done 792 tasks      | elapsed:   1.4s  
[Parallel(n_jobs=-1)]: Done 1242 tasks     | elapsed:   2.2s  
[Parallel(n_jobs=-1)]: Done 1750 out of 1750 | elapsed:   3.0s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:   0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:   0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:   0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:   0.4s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:   0.7s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:   0.9s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:   0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:   0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:   0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:   0.3s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:   0.5s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:   0.6s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:   0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:   0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:   0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:   0.4s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:   0.6s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:   0.8s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:   0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:   0.1s
```

```
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:    0.2s
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:    0.4s
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:    0.6s
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:    0.8s finished
```

Fold 1 OOB Score: 0.8455056179775281

Fold 2

```
[Parallel(n_jobs=-1)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=-1)]: Done  42 tasks      | elapsed:    0.1s  
[Parallel(n_jobs=-1)]: Done 192 tasks      | elapsed:    0.3s  
[Parallel(n_jobs=-1)]: Done 442 tasks      | elapsed:    0.8s  
[Parallel(n_jobs=-1)]: Done 792 tasks      | elapsed:    1.4s  
[Parallel(n_jobs=-1)]: Done 1242 tasks     | elapsed:    2.1s  
[Parallel(n_jobs=-1)]: Done 1750 out of 1750 | elapsed:    3.0s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:    0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:    0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:    0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:    0.4s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:    0.6s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:    0.9s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:    0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:    0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:    0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:    0.3s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:    0.4s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:    0.6s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:    0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:    0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:    0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:    0.4s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:    0.6s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:    0.8s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:    0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:    0.1s
```

```
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:    0.2s
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:    0.4s
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:    0.6s
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:    0.8s finished
```

Fold 2 OOB Score: 0.8469101123595506

Fold 3

```
[Parallel(n_jobs=-1)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=-1)]: Done  42 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=-1)]: Done 192 tasks      | elapsed:  0.3s  
[Parallel(n_jobs=-1)]: Done 442 tasks      | elapsed:  0.8s  
[Parallel(n_jobs=-1)]: Done 792 tasks      | elapsed:  1.4s  
[Parallel(n_jobs=-1)]: Done 1242 tasks     | elapsed:  2.2s  
[Parallel(n_jobs=-1)]: Done 1750 out of 1750 | elapsed:  3.0s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:  0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:  0.4s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:  0.6s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:  0.9s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:  0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:  0.3s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:  0.4s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:  0.6s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:  0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:  0.4s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:  0.6s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:  0.8s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s
```

```
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:    0.2s
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:    0.4s
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:    0.6s
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:    0.8s finished
```

Fold 3 OOB Score: 0.8345021037868162

Fold 4

```
[Parallel(n_jobs=-1)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=-1)]: Done  42 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=-1)]: Done 192 tasks      | elapsed:  0.3s  
[Parallel(n_jobs=-1)]: Done 442 tasks      | elapsed:  0.8s  
[Parallel(n_jobs=-1)]: Done 792 tasks      | elapsed:  1.3s  
[Parallel(n_jobs=-1)]: Done 1242 tasks     | elapsed:  2.1s  
[Parallel(n_jobs=-1)]: Done 1750 out of 1750 | elapsed:  3.0s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:  0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:  0.4s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:  0.6s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:  0.9s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:  0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:  0.3s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:  0.4s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:  0.6s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:  0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:  0.4s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:  0.6s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:  0.8s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s
```

```
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:    0.2s
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:    0.4s
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:    0.6s
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:    0.8s finished
```

Fold 4 OOB Score: 0.8387096774193549

Fold 5

```
[Parallel(n_jobs=-1)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=-1)]: Done  42 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=-1)]: Done 192 tasks      | elapsed:  0.3s  
[Parallel(n_jobs=-1)]: Done 442 tasks      | elapsed:  0.8s  
[Parallel(n_jobs=-1)]: Done 792 tasks      | elapsed:  1.4s  
[Parallel(n_jobs=-1)]: Done 1242 tasks     | elapsed:  2.2s  
[Parallel(n_jobs=-1)]: Done 1750 out of 1750 | elapsed:  3.0s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:  0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:  0.4s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:  0.7s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:  0.9s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:  0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:  0.3s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:  0.4s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:  0.6s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s  
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:  0.2s  
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:  0.4s  
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:  0.6s  
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:  0.8s finished  
  
[Parallel(n_jobs=4)]: Using backend ThreadingBackend with 4 concurrent workers.  
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  0.0s  
[Parallel(n_jobs=4)]: Done 192 tasks      | elapsed:  0.1s
```

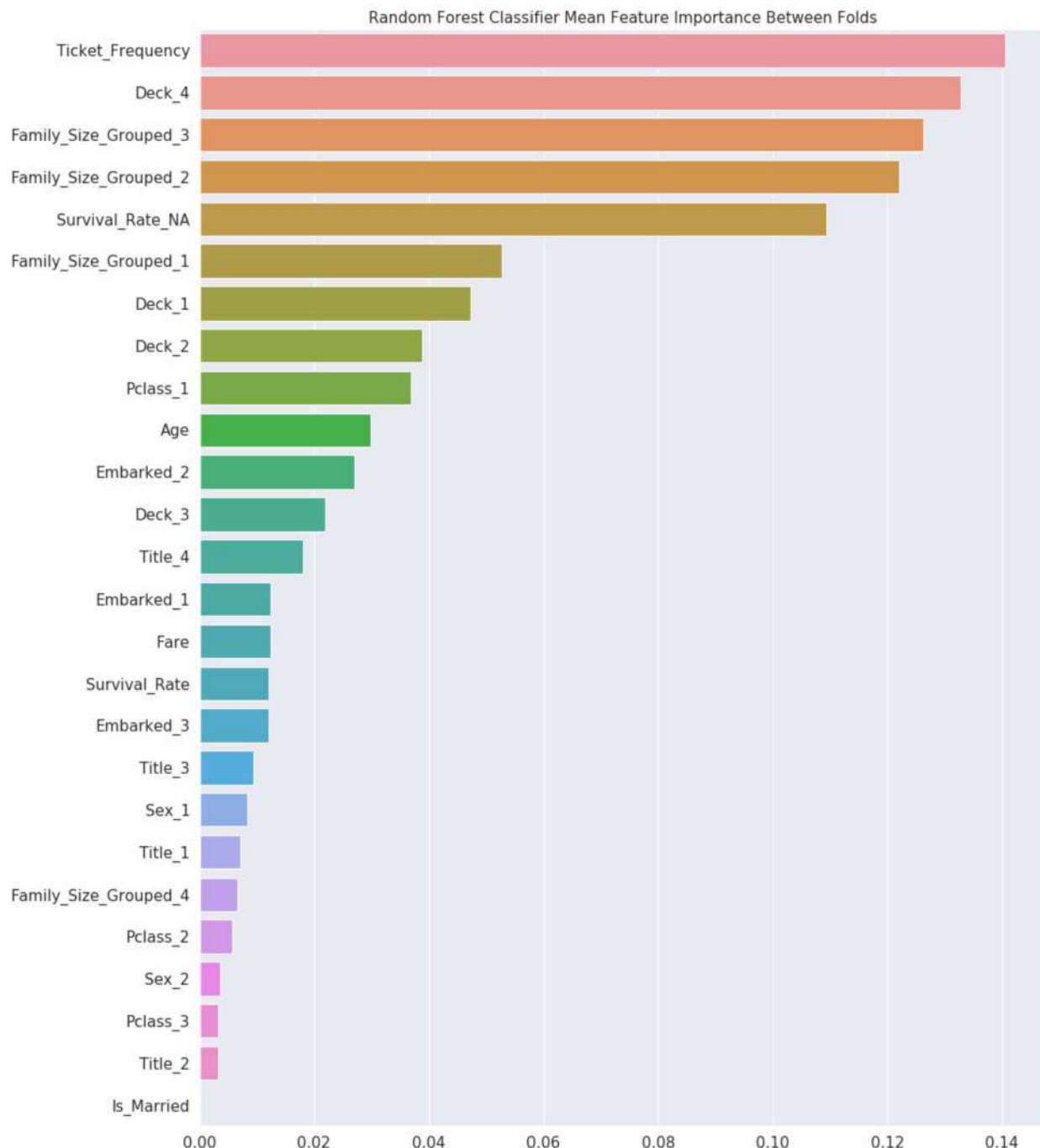
```
[Parallel(n_jobs=4)]: Done 442 tasks      | elapsed:    0.2s
[Parallel(n_jobs=4)]: Done 792 tasks      | elapsed:    0.4s
[Parallel(n_jobs=4)]: Done 1242 tasks     | elapsed:    0.6s
[Parallel(n_jobs=4)]: Done 1750 out of 1750 | elapsed:    0.8s finished
```

Fold 5 OOB Score: 0.8529411764705882

Average OOB Score: 0.8437137376027675

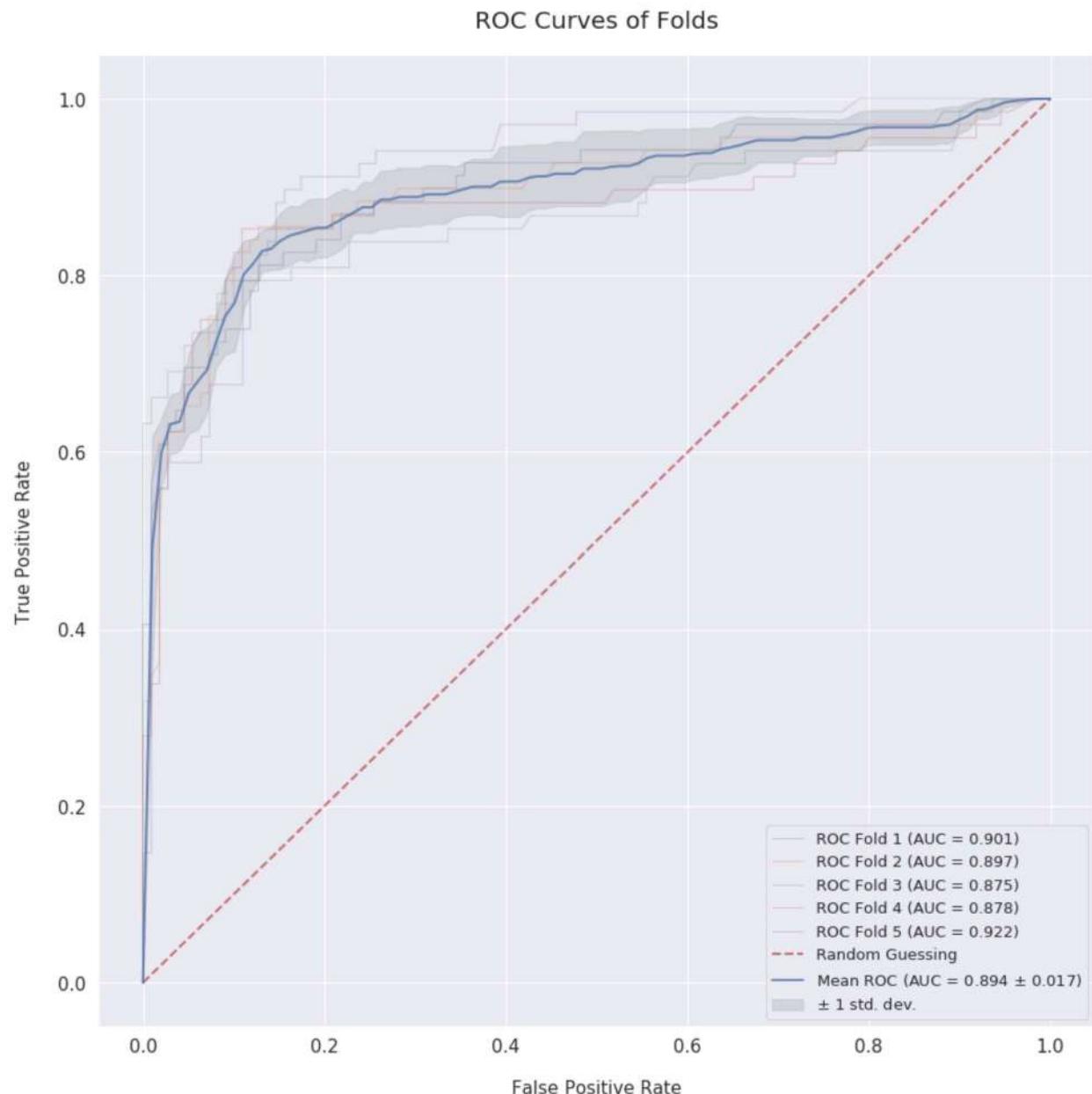
3.2 Feature Importance

>Show hidden code



3.3 ROC Curve

❖ Show hidden code



3.4 Submission

In [46]:

```
class_survived = [col for col in probs.columns if col.endswith('Prob_1')]
probs['1'] = probs[class_survived].sum(axis=1) / N
probs['0'] = probs.drop(columns=class_survived).sum(axis=1) / N
probs['pred'] = 0
pos = probs[probs['1'] >= 0.5].index
probs.loc[pos, 'pred'] = 1

y_pred = probs['pred'].astype(int)

submission_df = pd.DataFrame(columns=['PassengerId', 'Survived'])
submission_df['PassengerId'] = df_test['PassengerId']
submission_df['Survived'] = y_pred.values
submission_df.to_csv('submissions.csv', header=True, index=False)
submission_df.head(10)
```

Out[46]:

	PassengerId	Survived
891	892	0
892	893	1
893	894	0
894	895	0
895	896	1
896	897	0
897	898	1
898	899	0
899	900	1
900	901	0