

# Introduction

RMS Titanic was a British passenger liner operated by the White Star Line that sank in the North Atlantic Ocean on 15 April 1912, after striking an iceberg during her maiden voyage from Southampton to New York City. Of the estimated 2,224 passengers and crew aboard, more than 1,500 died, making the sinking at the time one of the deadliest of a single ship and the deadliest peacetime sinking of a supertanker or cruise ship to date. RMS Titanic was the largest ship afloat at the time she entered service and was the second of three Olympic-class ocean liners operated by the White Star Line. It was built by the Harland and Wolff shipyard in Belfast. Thomas Andrews, chief naval architect of the shipyard at the time, died in the disaster. (<https://en.wikipedia.org/wiki/Titanic>)

In this file, I will look deeply into the dataset and create a Machine Learning Model to predict the survived passenger. I will try different model and compare them. Also, I will try to find which features have a more impact on our process.

## STEP 0

In this step, we will import some libraries and useful things. After that, getting the data and describe our dataset and Data Exploration/Analysis.

```
In [1]: #some useful library and necessary things
import pandas as pd
import numpy as np

from scipy.stats import pearsonr

import matplotlib.pyplot as plt
import seaborn as sns

from IPython.display import display

In [2]: #read data from .csv file
passenger_train = pd.read_csv('train.csv')
passenger_test = pd.read_csv('test.csv')

In [3]: passenger_train.info()

Out[4]:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   PassengerId            891 non-null    int64
 1   Survived               891 non-null    int64
 2   Pclass                 891 non-null    int64
 3   Name                   891 non-null    object
 4   Sex                    891 non-null    object
 5   Age                    714 non-null    float64
 6   SibSp                  891 non-null    int64
 7   Parch                  891 non-null    int64
 8   Ticket                 891 non-null    object
 9   Fare                   891 non-null    float64
10   Cabin                  204 non-null    object
11   Embarked               889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB

The training data set have 891 examples and 11 features + one target values ('Survived'). Features have a different type, 5 of them are integers, 5 are objects and 2 are floats. I can give a short description of the features:
```

```
In [4]: #PassengerId:Unique Id of a passenger
#Pclass: Ticket class
#Name: Passenger Name
#Sex: Sex
#Age: Age in years
#SibSp: # of siblings / spouses aboard the Titanic
#Parch: # of parents / children aboard the Titanic
#Ticket: Ticket Number
#Fare: Passenger fare
#Cabin: Cabin number
#Embarked: Port of Embarkation

passenger_train.describe()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308642	29.691118	0.523008	0.381594	32.204208
std	257.353842	0.486592	0.836071	14.526497	1.102743	0.806057	49.693429
min	1.000000	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	668.500000	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

This graph shows us our training data set the survival rate is 38.3%. Also, we can see the age range is 0.4 and 80. I will look later for missing values more deeply but we can see that also feature age has missing values.

```
ax.legend()
ax.set_title('Male')

## You can write these notes into also in histplot function
#ax=plt.hist(women[women['Survived']==1],Age.dropna(),bins=18, label = not_survived,ax=axes[0],kde = False,
#ax=plt.hist(women[women['Survived']==0],Age.dropna(),bins=18, label = not_survived,ax=axes[0],kde = False,
#ax.legend()
#ax=plt.hist(male[male['Survived']==1],Age.dropna(),bins=18, label = not_survived,ax=axes[1],kde = False,color=
#ax=plt.hist(male[male['Survived']==0],Age.dropna(),bins=18, label = not_survived,ax=axes[1],kde = False,
#ax.legend()
#ax.set_title('Male')
```

As you can see above in our data set we need to convert a lot of features into numeric ones later on for the machine learning algorithm to process them. Also, our numerical values have a different range, we need to set them as similar and close one each other. And the next step, we'll look at our dataset which has missing values.

```
In [6]: total=passenger_train.isnull().sum().sort_values(ascending=False)
percent = passenger_train.isnull().sum().div(passenger_train.isnull().count()*100)
percent_2 = round(percent, 2)
missing_data = pd.concat([total, percent_2], axis=1, keys=['Total', '%'])
missing_data.head(5)
```

	Total	%
Cabin	687	77.1
Age	177	19.9
Embarked	2	0.2
Fare	0	0.0
Ticket	0	0.0

The main problems are Cabin and age sections. Embarked has just 2 missing data and we can fill it with ease.

```
In [7]: passenger_train.cabin.values

Out[7]: array(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
        'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'], dtype=object)

We have 11 features and one target feature as 'Survived'. In the next steps, I'll look at which features have most correlated with 'Survived'. First of all, we'll start with 'Age' and 'Sex'.
```

## Step 1

In this step, as I said before we will look at each feature how they related to the Survived and analysing them on some graph.

### Step1.1 Age and Sex

```
In [8]: survived = 'survived'
not_survived = 'not survived'

fig, axes = plt.subplots(nrows=1,ncols=2,figsize=(15,5))
women = passenger_train[passenger_train['Sex']=='female']
men = passenger_train[passenger_train['Sex']=='male']
axsns.distplot(women[women['Survived']==1],Age.dropna(),bins=18, label = survived,ax=axes[0],kde = False)
axsns.distplot(women[women['Survived']==0],Age.dropna(),bins=18, label = not_survived,ax=axes[0],kde = False,
ax.legend()
ax.set_title('Female')
axsns.distplot(men[men['Survived']==1],Age.dropna(),bins=18, label = survived,ax=axes[1],kde = False)
axsns.distplot(men[men['Survived']==0],Age.dropna(),bins=18, label = not_survived,ax=axes[1],kde = False)
ax.legend()
ax.set_title('Male')

## You can write these names also in histplot function.
#axsns.distplot(women[women['Survived']==1],Age.dropna(),bins=18, label = survived,ax=axes[0],kde = False,color='b')
#axsns.distplot(women[women['Survived']==0],Age.dropna(),bins=18, label = not_survived,ax=axes[0],kde = False,color='r')
#axsns.distplot(men[men['Survived']==1],Age.dropna(),bins=18, label = survived,ax=axes[1],kde = False,color='b')
#axsns.distplot(men[men['Survived']==0],Age.dropna(),bins=18, label = not_survived,ax=axes[1],kde = False,color='r')
#axsns.legend()
```

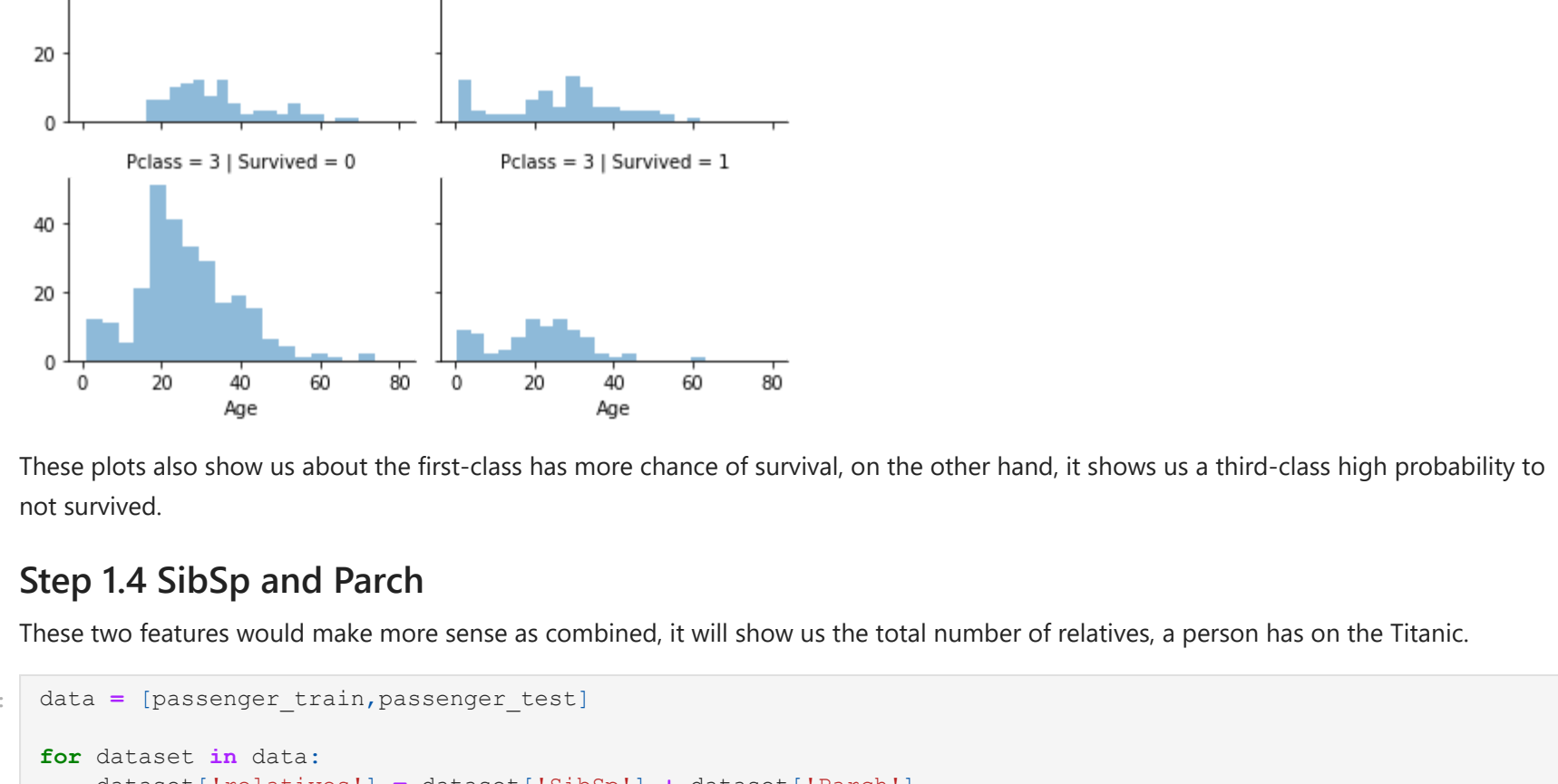
For women passengers, the survival chances are higher between 14 and 40 years old. For men, this is changing. It's also around the age of 18 and 18 but on the other side it's not the same for men.

### Step 1.2 Embarked,Pclass and Sex

```
In [9]: fig = sns.FacetGrid(passenger_train,row='Embarked',height=4,aspect=1.5)
fig.map(sns.pointplot,'Pclass','Survived','Sex',palette=None,order=None,hue_order=None)
fig.add_legend()

S = Southampton, UK
C = Cherbourg, France
Q = Queenstown, Ireland
print(passenger_train['Embarked'].value_counts())

S    644
C    168
Q     77
Name: Embarked, dtype: object
```



As we can see above graphs, embarked is correlated with survival depending on gender.

Women on port S and on port Q have a higher chance of survival. On port C, there is less chance of survival. Men on port C have a higher probability to survive otherwise on port Q and S they have less chance of survival.

### Step 1.3 Pclass

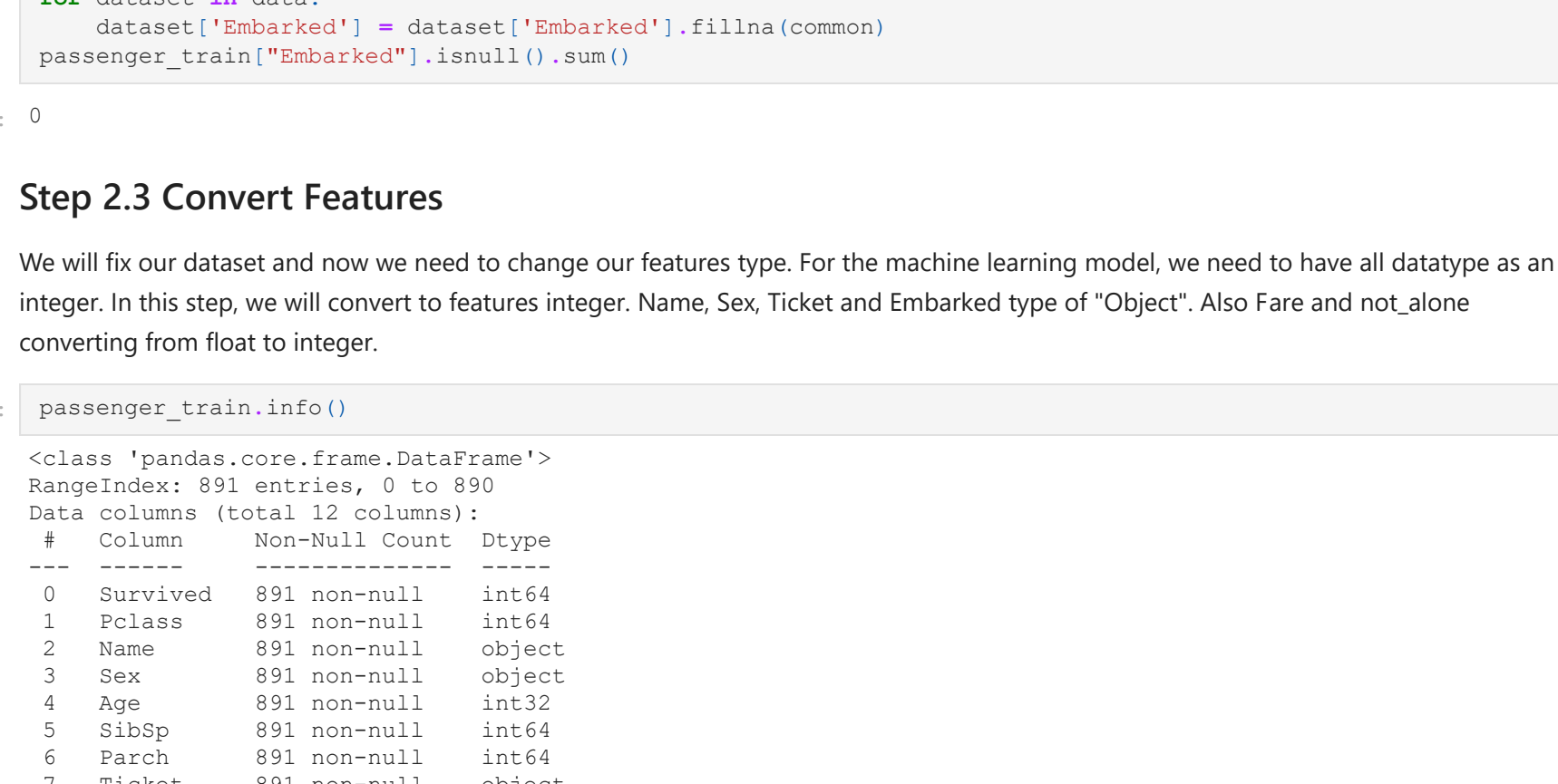
Pclass also seems to be correlated with survival that's why we will plot it.

```
In [10]: sns.barplot(x='Pclass',y='Survived',data=passenger_train)

Out[10]: <AxesSubplot: xlabel='Pclass', ylabel='Survived'>
```

As expected the first class have a higher survival chance than the other classes. We'll plot another plot for looking more deeply and understand each class effect on survival.

```
In [11]: c1=sns.FacetGrid(passenger_train,col='Survived',row='Pclass',height=2,aspect=1.5)
c1.map(sns.pointplot,'Age','Survived','Sex',alpha=0.5,bins=18)
c1.add_legend()
```



These plots also show us about the first-class has more chance of survival, on the other hand, it shows us a third-class high probability to not survived.

### Step1.4 SibSp and Parch

These two features would make more sense as combined, it will show us the total number of relatives, a person has on the Titanic.

```
In [12]: data = [passenger_train,passenger_test]

for dataset in data:
    dataset['relatives'] = dataset['SibSp'] + dataset['Parch']

    dataset.loc[dataset['relatives'] > 0, 'not_alone'] = 1
    dataset.loc[dataset['relatives'] == 0, 'not_alone'] = 0
    #dataset['not_alone'] = dataset['not_alone'].astype(int)
    passenger_train['not_alone'] = dataset['not_alone'].count()
    # Sifir yalniz insanlari posterir

In [13]: axes = sns.factorplot(x='relatives',y='Survived',data=passenger_train,height=4.5,aspect=1.5)
# axes = sns.factorplot(x='relatives',y='Survived',data=passenger_train,height=5,aspect=1.5)
```



The above plotting shows us, you had a high probability of survival with 1 to 3 relatives and if you had less than 1 or more than 3 that probability is plotting.

## Step 2 Data Preprocessing

Before the starting analysis of our missing data, we need to look at the dataset. In our dataset, we don't need to use PassengerId, Embarked(2) and Age(177). I decided to delete also Cabin. If we have missing values on features, we'll fix it.

```
In [14]: passenger_train = passenger_train.drop(['PassengerId','Cabin'],axis=1)

In [15]: data = [passenger_train,passenger_test]

for dataset in data:
    mean = passenger_train['Age'].mean()
    std = passenger_test['Age'].std()
    dataset['Age'] = dataset['Age'].fillna(mean)
    # compute random numbers between the mean, std and is null
    rand = np.random.randint(mean-std,mean+std,size=1)
    dataset.loc[dataset['Age'].isnull(),Age] = rand
    age_slice = dataset['Age'].copy()
    age_slice.loc[age_slice.isnull()] = rand
    dataset['Age'] = passenger_train['Age'].astype(int)
    passenger_train['Age'] = dataset['Age'].astype(int)
    passenger_test['Age'] = dataset['Age'].astype(int)
```

### Step2.1 Age

Just 2 missing values on this one. I will just fill these with the most common ones.

```
In [16]: passenger_train['Embarked'].describe()
```

```
count      889
unique        3
top          C
freq         644
Name: Embarked, dtype: object

common='S'
data = [passenger_train,passenger_test]
for dataset in data:
    dataset['Embarked'] = dataset['Embarked'].fillna(common)
    passenger_train['Embarked'] = dataset['Embarked'].count()
```

### Step 2.2 Embarked

Just 2 missing values on this one. I will just fill these with the most common ones.

```
In [17]: passenger_train['Embarked'].describe()
```

```
count      889
unique        3
top          C
freq         644
Name: Embarked, dtype: object

common='S'
data = [passenger_train,passenger_test]
for dataset in data:
    dataset['Embarked'] = dataset['Embarked'].fillna(common)
    passenger_train['Embarked'] = dataset['Embarked'].count()
```

### Step 2.3 Convert Features

We will fix our dataset and now we need to change our features type. For the machine learning model, we need to have all datatype as an integer. In this step, we will convert to features integer. Name, Sex, Ticket and Embarked type of 'Object'. Also Fare and not\_alone converting from float to integer.

```
In [18]: passenger_train.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Survived               891 non-null    int64
 1   Pclass                 891 non-null    int64
 2   Name                   891 non-null    object
 3   Sex                    891 non-null    object
 4   Age                    891 non-null    int64
 5   SibSp                  891 non-null    int64
 6   Parch                  891 non-null    int64
 7   Ticket                 891 non-null    object
 8   Fare                   891 non-null    float64
 9   Embarked               891 non-null    object
10   relatives              891 non-null    int64
11   not_alone              891 non-null    int64
dtypes: float64(2), int64(11), object(2)
memory usage: 80.2+ KB

In [19]: # Fare and not_alone float to int
data = [passenger_train,passenger_test]

for dataset in data:
    dataset['Fare'] = dataset['Fare'].fillna(0)
    dataset['Fare'] = dataset['Fare'].astype(int)
    dataset['not_alone'] = dataset['not_alone'].fillna(0)
    dataset['not_alone'] = dataset['not_alone'].astype(int)
```

### Step 2.3.1 Name

We'll change to Name as a Title. So that we can build a new feature.

```
In [20]: data = [passenger_train,passenger_test]

for dataset in data:
    dataset['Title'] = dataset.Name.str.extract(pat='([A-Z-a-z]+)',expand=False)
    pd.crosstab(passenger_train['Title'],passenger_train['Sex'])
```

Sex	Mr	Mrs	Ms	Rev	Sir	Other
male	517	0	0	0	1	0
female	0	125	1	0	0	0

```
In [21]: # replace titles with a more common title or as Rare
dataset['Title'] = dataset['Title'].replace(['Capt','Col','Countess','Don','Dona','Dr','Jonkheer','Lady','Major','Rev','Sir'],'Rare')
dataset['Title'] = dataset['Title'].replace(['Wile','Miss'],'Rare')
dataset['Title'] = dataset['Title'].replace(['Ms','Mrs'],'Rare')
dataset['Title'] = dataset['Title'].replace(['Mme'],'Rare')
passenger_train[['Title','Survived']].groupby('Title').mean()
```

Survived	Mr	Mrs	Ms	Rev	Sir	Other
mean	0.575000	0.702703	0.156673	0.793651	0.347826	0.347826

```
In [22]: titles = ["Mr": 1, "Mrs": 2, "Master": 3, "Rare": 4]
for dataset in data:
    dataset['Title'] = dataset['Title'].replace('Mr',1)
    dataset['Title'] = dataset['Title'].replace('Mrs',2)
    dataset['Title'] = dataset['Title'].replace('Master',3)
    dataset['Title'] = dataset['Title'].replace('Rare',4)
    # convert titles into numbers
    passenger_train = passenger_train.drop(['Name'],axis=1)
    passenger_test = passenger_test.drop(['Name'],axis=1)
```

```
In [23]: passenger_train.head()

Out[23]:
Survived  Pclass  Sex  Age  SibSp  Parch  Ticket  Fare  Embarked  relatives  not_alone  Title
0         0      3   male  22   1     0    A/5 21171  7.25  S         1         0         0
1         1      1   female  38   1     0    PC 17599  71.28  C         1         0         2
2         1      3   female  26   0     0    STON/O2 3101282  7.92  S         0         0         1
3         1      1   female  35   1     0    113803  53.10  S         1         0         2
4         0      3   male  35   0     0    373450  8.05  S         0         0         1
```

### Step 2.3.2 Sex

Now in this step, we will convert the sex feature as an integer. 0 for male and 1 for female.

```
In [24]: genders = {"male":0,"female":1}
data = [passenger_train,passenger_test]

for dataset in data:
    dataset['Sex'] = dataset['Sex'].map(genders)
```

### Step 2.3.3 Ticket

```
In [25]: passenger_train['Ticket'].describe()

Out[25]:
count      891
unique     681
top        1601
freq        7
Name: Ticket, dtype: object

681 values are unique on this dataset, that's why it would be difficult to classify them. We will drop it.

passenger_train = passenger_train.drop(['Ticket'],axis=1)
passenger_test = passenger_test.drop(['Ticket'],axis=1)
```

### Step 2.3.4 Embarked

Convert 'Embarked' into numeric values.

```
In [27]: ports = {'S':0,'C':1,'Q':2}
data = [passenger_train,passenger_test]
for dataset in data:
    dataset['Embarked'] = dataset['Embarked'].map(ports)

In [28]: passenger_train.head()

Out[28]:
Survived  Pclass  Sex  Age  SibSp  Parch  Fare  Embarked  relatives  not_alone  Title
0         0      3   male  22   1     0         7         0         1         0         0
1         1      1   female  38   1     0        71         1         1         0         2
2         1      3   female  26   0     0        7         0         0         1         1
3         1      1   female  35   1     0        53         1         0         2
4         0      3   male  35   0     0         8         0         0         1         0
```

## Step 3 - Setting Categories

Now our all features numeric values but some of them have huge range difference. As feature age value range is from 0.4 to 80. For a machine learning model, we need to have a similar range of difference between all features. In this step, we will update our features(Age and Fare) as like that

### Step3.1 Age

According to distribution we will grouping age feature as a different group. (0-10 = 0) & (11-18 = 1) & (19-22 = 2) & (23-27 = 3) & (28-33 = 4) & (34-40 = 5) & (>71-155 = 1) & (>155 - 315 = 2) & (>315-995 = 3) & (>995-2505 = 4) & (>2505... = 5)

```
In [29]: data = [passenger_train,passenger_test]

for dataset in data:
    dataset['Age'] = dataset['Age'].astype(int)
    dataset.loc[dataset['Age'] <= 10, 'Age'] = 0
    dataset.loc[dataset['Age'] > 10 & dataset['Age'] <= 18, 'Age'] = 1
    dataset.loc[dataset['Age'] > 18 & dataset['Age'] <= 22, 'Age'] = 2
    dataset.loc[dataset['Age'] > 22 & dataset['Age'] <= 27, 'Age'] = 3
    dataset.loc[dataset['Age'] > 27 & dataset['Age'] <= 33, 'Age'] = 4
    dataset.loc[dataset['Age'] > 33 & dataset['Age'] <= 40, 'Age'] = 5
    dataset.loc[dataset['Age'] > 40, 'Age'] = 6
    passenger_train['Age'] = dataset['Age'].value_counts()
```

```
In [29]: # Age
count      891
unique     161
top        155
freq        2
Name: Age, dtype: int64

Step3.1 Fare
Also in this feature, we have a very different range. Titanic has offered a different type of ticket and according to them, we have ticket class and fare. That's why we need to convert our fare features like age. According to distribution we will grouping age feature as a different group (0-7.915 = 0) & (>7.915-155 = 1) & (>155 - 315 = 2) & (>315-995 = 3) & (>995-2505 = 4) & (>2505... = 5)
```

```
In [30]: data = [passenger_train,passenger_test]

for dataset in data:
    dataset.loc[dataset['Fare'] <= 7.91, 'Fare'] = 0
    dataset.loc[dataset['Fare'] > 7.91 & dataset['Fare'] <= 15, 'Fare'] = 1
    dataset.loc[dataset['Fare'] > 15 & dataset['Fare'] <= 31, 'Fare'] = 2
    dataset.loc[dataset['Fare'] > 31 & dataset['Fare'] <= 99, 'Fare'] = 3
    dataset.loc[dataset['Fare'] > 99 & dataset['Fare'] <= 250, 'Fare'] = 4
    dataset.loc[dataset['Fare'] > 250, 'Fare'] = 5
    dataset['Fare'] = dataset['Fare'].astype(int)
```

```
In [31]: passenger_train['Fare'].describe()

Out[31]:
count      891.000000
mean       1.49270500
std        0.00000000
min         0.00000000
25%         0.00000000
50%         0.00000000
75%         0.00000000
max         5.00000000
Name: Fare, dtype: float64

In [32]: #Afterwards starting machine learning models, let's take a last look at our trainin dataset
passenger_train.describe()
```

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	relatives	not_alone	Title
count	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000
mean	0.383838	2.308642	0.324113	35.25253	0.523008	0.381594	1.492705	0.361392	0.904602	0.602694	0.728395
std	0.486592	0.836071	0.477990	1.851262	1.102743	0.806057	1.250933	0.635673	1.613459	0.489615	1.003039
min	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.000000	2.000000	0.000000	2.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	0.000000	3.000000	0.000000	4.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	1.000000	3.000000	1.000000	5.000000	1.000000	0.000000	2.000000	1.000000	1.000000	1.000000	0.000000
max	1.000000	3.000000	1.000000	6.000000	8.000000	6.000000	5.000000	2.000000	10.000000	1.000000	4.000000

We have 11 features as you can see above. All of them are integer and have a close range.

## Step 4 Machine Learning Models

In this part, we will train 3 models and compare their results. Firstly, we need to drop survived features on our training dataset. Because the dataset does not provide labels for their testing set, we used the predictions on the training set to compare the algorithms with each other. Also, drop the passenger Id and cabin on test data.

```
In [33]: # Algorithms
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import SGDClassifier

In [34]: X_train = passenger_train.drop('Survived',axis=1)
Y_train = passenger_train['Survived']
scores = cross_val_score(X_train,Y_train,RandomForestClassifier(),axis=1).copy()
```

```
In [35]: X_train

Out[35]:
Pclass  Sex  Age  SibSp  Parch  Fare  Embarked  relatives  not_alone  Title
0      3   0   2     1     0     0         7         0         1         0
1      1   1   5     1     0     3         1         1         0         2
2      3   1   3     0     0     0         1         0         0         1
3      1   1   5     1     0     3         1         1         0         2
4      3   0   5     0     0     1         0         0         0         1
...
886     2     0     3     0     0     1         0         0         1         4
887     1     1     2     0     0     2         0         0         1         1
888     3     1     3     1     2     2         0         3         0         1
889     1     0     3     0     0     2         1         0         0         1
890     3     0     4     0     0     0         2         0         1         0
891 rows x 10 columns
```

```
In [36
```



```
In [42]: ## Stochastic Gradient Descent
sgd = SGDClassifier()
scores = cross_val_score(sgd, X_train, Y_train, cv=10, scoring = "accuracy")
print("Scores:", scores)
print("Mean:", scores.mean())
print("Standard Deviation:", scores.std())
Score_rand = scores.mean()

Scores: [0.71111111 0.82022472 0.38202247 0.74157303 0.68539326 0.73033708
0.7752809 0.65168539 0.79775281 0.79775281]
Mean: 0.70931358302224
Standard Deviation: 0.120305914875269
```

```
In [43]: ## Random Forest K-Fold Cross Validation

rf = RandomForestClassifier(n_estimators=100)
scores = cross_val_score(rf, X_train, Y_train, cv=10, scoring = "accuracy")
print("Scores:", scores)
print("Mean:", scores.mean())
print("Standard Deviation:", scores.std())
Score_rand = scores.mean()

Scores: [0.73333333 0.83146067 0.76404494 0.83146067 0.87640449 0.86516854
0.84249663 0.78651685 0.82022472 0.83146067]
Mean: 0.8182771535505024
Standard Deviation: 0.04226631171418877
```

```
In [44]: # We can show all K-fold results on a table, we can compare easily our result
results = pd.DataFrame({'Model': ('Logistic Regression','Stochastic Gradient Descent','Random Forest'),
                        'Score': (Score_logr,Score_sgd,Score_rand)})

resultdf = results.sort_values(by='Score',ascending = False)

resultdf = resultdf.set_index('Score')
resultdf.head()
```

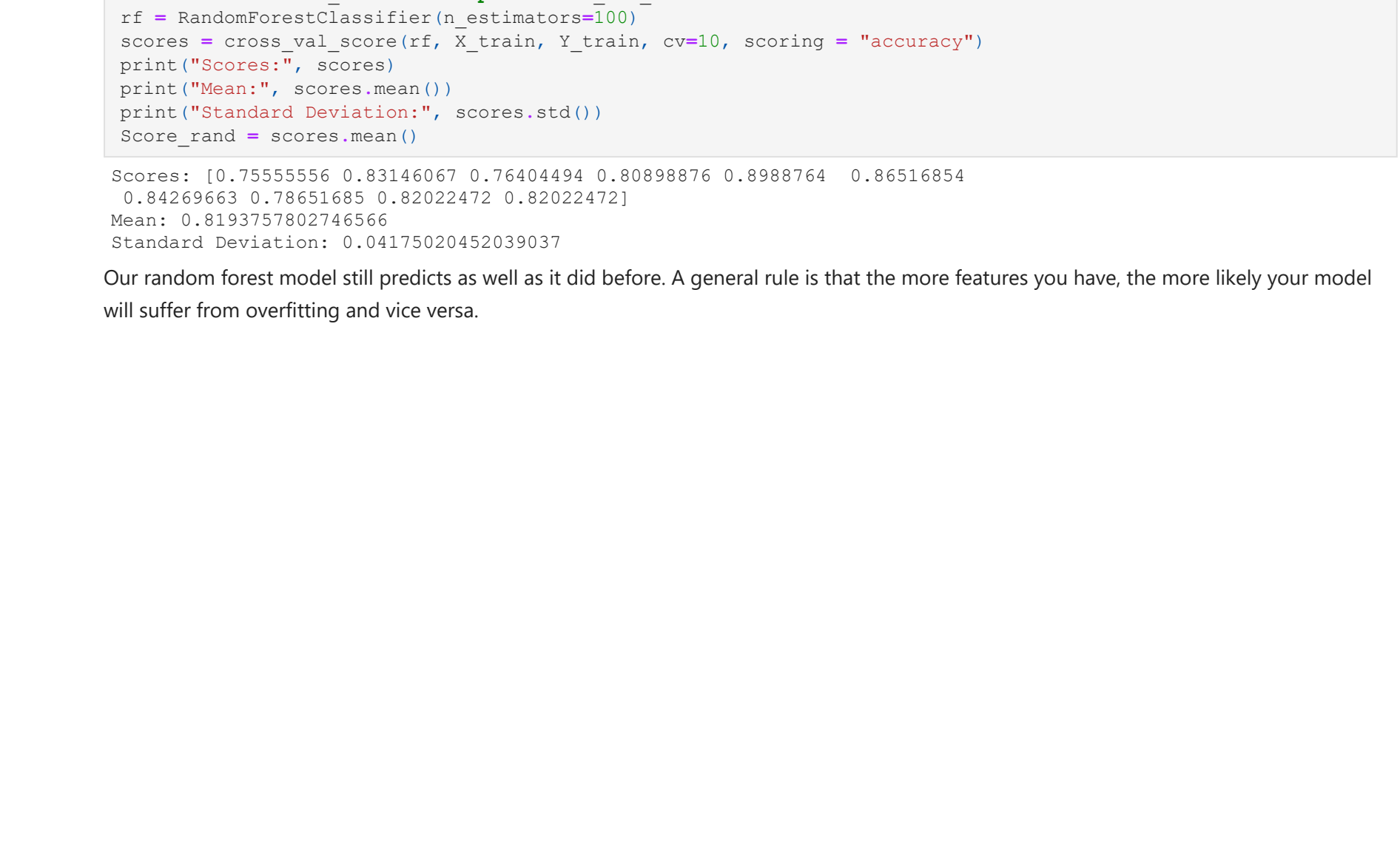
	Model
Score	
<b>0.818277</b>	Random Forest
<b>0.806966</b>	Logistic Regression
<b>0.709313</b>	Stochastic Gradient Descent

These results looks much more better and realistic than firs models. Still random forest is the best one and we can continue with that. Our model has an average accuracy of 82% with a standard deviation of 4%. The standard deviation shows us, how precise the estimates are. This means our model can differ +4%.

In the next steps, we will try to understand which feature has more important for our models. For this step, we will use the Random Forest model result because we reach the best solution with it and easy to use model.

## Feature Importance

```
In [45]: importances = pd.DataFrame({'feature':X_train.columns,'importance':np.round(randomF.feature_importances_,4)})
importances = importances.sort_values('Importance',ascending=False).set_index('feature')
importances.head(10)
```



```
In [46]: importances.plot.bar()
plt.title('Importance of Features')
plt.xlabel('Features')
plt.ylabel('Value')
```

Out [46]: Text (0, 0.5, 'Value')

As we expect our result has the most correlation between "Title-Sex-Age-Pclass and fare". On the other hand, Parch and not\_alone are not too important for our random forest prediction model. That's why we can drop these two features in the dataset and train again.

```
In [47]: passenger_train = passenger_train.drop(['Parch','not_alone','SibSp'],axis=1)
passenger_test = passenger_test.drop(['Parch','not_alone','SibSp'],axis=1)
```

```
In [48]: ## Random Forest
from sklearn.model_selection import cross_val_score
rf = RandomForestClassifier(n_estimators=100)
scores = cross_val_score(rf, X_train, Y_train, cv=10, scoring = "accuracy")
print("Scores:", scores)
print("Mean:", scores.mean())
print("Standard Deviation:", scores.std())
Score_rand = scores.mean()

Scores: [0.75555556 0.83146067 0.76404494 0.80898876 0.8988764 0.86516854
0.84249663 0.78651685 0.82022472 0.82022472]
Mean: 0.8193757802746566
Standard Deviation: 0.04175020452039037
```

Our random forest model still predicts as well as it did before. A general rule is that the more features you have, the more likely your model will suffer from overfitting and vice versa.