# MINI PROJECT : TITANIC SURVIVAL PREDICTION

**AIM:** To analyze the accuracy of survival chances during the Titanic disaster using machine learning techniques.

**Description:** This project contains test data and traning data. predictive model to classify passengers as survivors or non-survivors. By using random forest classifier algorithm we can predict the accuracy of survival chances.

## Algorithm Used:

**Random forest classifier algorithm:** Using the Random Forest classification algorithm, we aim to build a predictive model that will allow us to estimate the likelihood of survival for each individual aboard the Titanic.

**Work flow:**

1) Data selection and Data preprocessing
2) Optimizing data for model training
3) Data transformation
4) Model training
5) Predict accuracy
6) Data Visualization

## Data set:

#Train dataset:

| | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | PassengerI | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
| 2 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22 | 1 | 0 | A/5 21171 | 7.25 | | S |
| 3 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Fl( female | female | 38 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 4 | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26 | 0 | 0 | STON/O2. | 7.925 | | S |
| 5 | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (L | female | 35 | 1 | 0 | 113803 | 53.1 | C123 | S |
| 6 | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35 | 0 | 0 | 373450 | 8.05 | | S |
| 7 | 6 | 0 | 3 | Moran, Mr. James | male | | 0 | 0 | 330877 | 8.4583 | | Q |
| 8 | 7 | 0 | 1 | McCarthy, Mr. Timothy J | male | 54 | 0 | 0 | 17463 | 51.8625 | E46 | S |
| 9 | 8 | 0 | 3 | Palsson, Master. Gosta Leonar | male | 2 | 3 | 1 | 349909 | 21.075 | | S |
| 10 | 9 | 1 | 3 | Johnson, Mrs. Oscar W (Elisabe | female | 27 | 0 | 2 | 347742 | 11.1333 | | S |
| 11 | 10 | 1 | 2 | Nasser, Mrs. Nicholas (Adele A | female | 14 | 1 | 0 | 237736 | 30.0708 | | C |
| 12 | 11 | 1 | 3 | Sandstrom, Miss. Marguerite R | female | 4 | 1 | 1 | PP 9549 | 16.7 | G6 | S |
| 13 | 12 | 1 | 1 | Bonnell, Miss. Elizabeth | female | 58 | 0 | 0 | 113783 | 26.55 | C103 | S |
| 14 | 13 | 0 | 3 | Saundercock, Mr. William Henr | male | 20 | 0 | 0 | A/5. 2151 | 8.05 | | S |
| 15 | 14 | 0 | 3 | Andersson, Mr. Anders Johan | male | 39 | 1 | 5 | 347082 | 31.275 | | S |

**Source Code:**

```
import warnings

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

import seaborn as sns

plt.style.use('fivethirtyeight')

%matplotlib inline

warnings.filterwarnings('ignore')

train = pd.read_csv('train.csv')

test = pd.read_csv('test.csv')

# To know number of columns and rows

train.shape

# (891, 12)

train = train.drop(['Cabin'], axis=1)

test = test.drop(['Cabin'], axis=1)

train = train.drop(['Ticket'], axis=1)

test = test.drop(['Ticket'], axis=1)

# replacing the missing values in

# the Embarked feature with S

train = train.fillna({"Embarked": "S"})

# sort the ages into logical categories

train["Age"] = train["Age"].fillna(-0.5)

test["Age"] = test["Age"].fillna(-0.5)

bins = [-1, 0, 5, 12, 18, 24, 35, 60, np.inf]

labels = ['Unknown', 'Baby', 'Child', 'Teenager',

        'Student', 'Young Adult', 'Adult', 'Senior']

train['AgeGroup'] = pd.cut(train["Age"], bins, labels=labels)

test['AgeGroup'] = pd.cut(test["Age"], bins, labels=labels)
```

```python
# create a combined group of both datasets
combine = [train, test]
# extract a title for each Name in the
# train and test datasets
for dataset in combine:
    dataset['Title'] = dataset.Name.str.extract(' ([A-Za-z]+)\.', expand=False)
pd.crosstab(train['Title'], train['Sex'])
# replace various titles with more common names
for dataset in combine:
    dataset['Title'] = dataset['Title'].replace(['Lady', 'Capt', 'Col',
                                                  'Don', 'Dr', 'Major',
                                                  'Rev', 'Jonkheer', 'Dona'],
                                                 'Rare')
    dataset['Title'] = dataset['Title'].replace(
        ['Countess', 'Lady', 'Sir'], 'Royal')
    dataset['Title'] = dataset['Title'].replace('Mlle', 'Miss')
    dataset['Title'] = dataset['Title'].replace('Ms', 'Miss')
    dataset['Title'] = dataset['Title'].replace('Mme', 'Mrs')


train[['Title', 'Survived']].groupby(['Title'], as_index=False).mean()


# map each of the title groups to a numerical value
title_mapping = {"Mr": 1, "Miss": 2, "Mrs": 3,
                 "Master": 4, "Royal": 5, "Rare": 6}
for dataset in combine:
    dataset['Title'] = dataset['Title'].map(title_mapping)
    dataset['Title'] = dataset['Title'].fillna(0)
mr_age = train[train["Title"] == 1]["AgeGroup"].mode()  # Young Adult
miss_age = train[train["Title"] == 2]["AgeGroup"].mode()  # Student
```

```python
mrs_age = train[train["Title"] == 3]["AgeGroup"].mode()  # Adult
master_age = train[train["Title"] == 4]["AgeGroup"].mode()  # Baby
royal_age = train[train["Title"] == 5]["AgeGroup"].mode()  # Adult
rare_age = train[train["Title"] == 6]["AgeGroup"].mode()  # Adult
age_title_mapping = {1: "Young Adult", 2: "Student",
                3: "Adult", 4: "Baby", 5: "Adult", 6: "Adult"}
for x in range(len(train["AgeGroup"])):
    if train["AgeGroup"][x] == "Unknown":
        train["AgeGroup"][x] = age_title_mapping[train["Title"][x]]
for x in range(len(test["AgeGroup"])):
    if test["AgeGroup"][x] == "Unknown":
        test["AgeGroup"][x] = age_title_mapping[test["Title"][x]]
# map each Age value to a numerical value
age_mapping = {'Baby': 1, 'Child': 2, 'Teenager': 3,
            'Student': 4, 'Young Adult': 5, 'Adult': 6,
            'Senior': 7}
train['AgeGroup'] = train['AgeGroup'].map(age_mapping)
test['AgeGroup'] = test['AgeGroup'].map(age_mapping)
train.head()
# dropping the Age feature for now, might change
train = train.drop(['Age'], axis=1)
test = test.drop(['Age'], axis=1)
train = train.drop(['Name'], axis=1)
test = test.drop(['Name'], axis=1)
sex_mapping = {"male": 0, "female": 1}
train['Sex'] = train['Sex'].map(sex_mapping)
test['Sex'] = test['Sex'].map(sex_mapping)
embarked_mapping = {"S": 1, "C": 2, "Q": 3}
train['Embarked'] = train['Embarked'].map(embarked_mapping)
```

```python
test['Embarked'] = test['Embarked'].map(embarked_mapping)
for x in range(len(test["Fare"])):
    if pd.isnull(test["Fare"][x]):
        pclass = test["Pclass"][x]  # Pclass = 3
        test["Fare"][x] = round(
            train[train["Pclass"] == pclass]["Fare"].mean(), 4)
# map Fare values into groups of
# numerical values
train['FareBand'] = pd.qcut(train['Fare'], 4,
                 labels=[1, 2, 3, 4])
test['FareBand'] = pd.qcut(test['Fare'], 4,
                 labels=[1, 2, 3, 4])
# drop Fare values
train = train.drop(['Fare'], axis=1)
test = test.drop(['Fare'], axis=1)
from sklearn.model_selection import train_test_split
# Drop the Survived and PassengerId
# column from the trainset
predictors = train.drop(['Survived', 'PassengerId'], axis=1)
target = train["Survived"]
x_train, x_val, y_train, y_val = train_test_split(
    predictors, target, test_size=0.2, random_state=0)
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
randomforest = RandomForestClassifier()
# Fit the training data along with its output
randomforest.fit(x_train, y_train)
y_pred = randomforest.predict(x_val)
# Find the accuracy score of the model
```

acc_randomforest = round(accuracy_score(y_pred, y_val) * 100, 2)

print(acc_randomforest)

f, ax = plt.subplots(1, 2, figsize=(12, 4))

train['Survived'].value_counts().plot.pie(

      explode=[0, 0.1], autopct='%1.1f%%', ax=ax[0], shadow=False)

ax[0].set_title('Survivors (1) and the dead (0)')

ax[0].set_ylabel('')

sns.countplot(x='Survived', data=train, ax=ax[1])

ax[1].set_ylabel('Quantity')

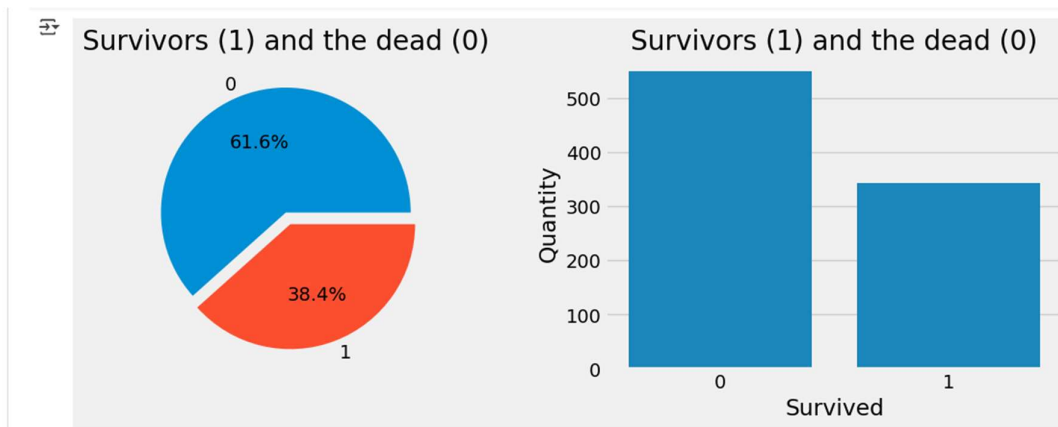ax[1].set_title('Survivors (1) and the dead (0)')

plt.show()

**Output:**



**Result:** Hence,the analyzing the accuracy and visualization of survival chances during Titanic disaster using random forest classification algorithm has been executed successfully.