

```
In [1]: 1 import pandas as pd
        2 import warnings
        3 warnings.filterwarnings("ignore")
        4 import numpy as np
```

```
In [2]: 1 data=pd.read_csv("/home/placement/Downloads/Titanic Dataset.csv")
```

```
In [3]: 1 data.describe()
```

```
Out[3]:
```

	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	29.699118	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

In [4]:

1 data.info()

```
<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
```

In [5]:

1 data.head()

Out[5]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

```
In [6]: 1 data.isna().sum()
```

```
Out[6]: PassengerId      0  
Survived      0  
Pclass        0  
Name          0  
Sex           0  
Age          177  
SibSp         0  
Parch         0  
Ticket        0  
Fare          0  
Cabin        687  
Embarked      2  
dtype: int64
```

```
In [7]: 1 data=data.drop(['PassengerId', 'Name', 'Ticket', 'SibSp', 'Parch', 'Cabin'],axis=1)
```

In [8]:

```
1 data
```

Out[8]:

	Survived	Pclass	Sex	Age	Fare	Embarked
0	0	3	male	22.0	7.2500	S
1	1	1	female	38.0	71.2833	C
2	1	3	female	26.0	7.9250	S
3	1	1	female	35.0	53.1000	S
4	0	3	male	35.0	8.0500	S
...
886	0	2	male	27.0	13.0000	S
887	1	1	female	19.0	30.0000	S
888	0	3	female	NaN	23.4500	S
889	1	1	male	26.0	30.0000	C
890	0	3	male	32.0	7.7500	Q

891 rows × 6 columns

In [9]:

```
1 data['Sex']=data['Sex'].map({'male':1,'female':0})
```

```
In [10]: 1 data
```

```
Out[10]:
```

	Survived	Pclass	Sex	Age	Fare	Embarked
0	0	3	1	22.0	7.2500	S
1	1	1	0	38.0	71.2833	C
2	1	3	0	26.0	7.9250	S
3	1	1	0	35.0	53.1000	S
4	0	3	1	35.0	8.0500	S
...
886	0	2	1	27.0	13.0000	S
887	1	1	0	19.0	30.0000	S
888	0	3	0	NaN	23.4500	S
889	1	1	1	26.0	30.0000	C
890	0	3	1	32.0	7.7500	Q

891 rows × 6 columns

```
In [11]: 1 data=data.fillna(data.median())
```

```
In [12]: 1 data
```

```
Out[12]:
```

	Survived	Pclass	Sex	Age	Fare	Embarked
0	0	3	1	22.0	7.2500	S
1	1	1	0	38.0	71.2833	C
2	1	3	0	26.0	7.9250	S
3	1	1	0	35.0	53.1000	S
4	0	3	1	35.0	8.0500	S
...
886	0	2	1	27.0	13.0000	S
887	1	1	0	19.0	30.0000	S
888	0	3	0	28.0	23.4500	S
889	1	1	1	26.0	30.0000	C
890	0	3	1	32.0	7.7500	Q

891 rows × 6 columns

```
In [13]: 1 data['Pclass'].unique()
```

```
Out[13]: array([3, 1, 2])
```

```
In [14]: 1 data.fillna(35,inplace=True)
```

In [15]:

```
1 data
```

Out[15]:

	Survived	Pclass	Sex	Age	Fare	Embarked
0	0	3	1	22.0	7.2500	S
1	1	1	0	38.0	71.2833	C
2	1	3	0	26.0	7.9250	S
3	1	1	0	35.0	53.1000	S
4	0	3	1	35.0	8.0500	S
...
886	0	2	1	27.0	13.0000	S
887	1	1	0	19.0	30.0000	S
888	0	3	0	28.0	23.4500	S
889	1	1	1	26.0	30.0000	C
890	0	3	1	32.0	7.7500	Q

891 rows × 6 columns

In [16]:

```
1 data.isna().sum()
```

Out[16]:

Survived	0
Pclass	0
Sex	0
Age	0
Fare	0
Embarked	0
dtype:	int64

```
In [17]: 1 data.dtypes
```

```
Out[17]: Survived      int64  
Pclass        int64  
Sex           int64  
Age          float64  
Fare         float64  
Embarked      object  
dtype: object
```

```
In [18]: 1 y=data['Survived']  
        2 x=data.drop('Survived',axis=1)
```

```
In [19]: 1 y
```

```
Out[19]: 0      0  
        1      1  
        2      1  
        3      1  
        4      0  
        ..  
      886      0  
      887      1  
      888      0  
      889      1  
      890      0  
Name: Survived, Length: 891, dtype: int64
```


In [20]:

```
1 x
```

Out[20]:

	Pclass	Sex	Age	Fare	Embarked
0	3	1	22.0	7.2500	S
1	1	0	38.0	71.2833	C
2	3	0	26.0	7.9250	S
3	1	0	35.0	53.1000	S
4	3	1	35.0	8.0500	S
...
886	2	1	27.0	13.0000	S
887	1	0	19.0	30.0000	S
888	3	0	28.0	23.4500	S
889	1	1	26.0	30.0000	C
890	3	1	32.0	7.7500	Q

891 rows × 5 columns

In [21]:

```
1 x=pd.get_dummies(x)
```

In [22]:

```
1 x
```

Out[22]:

	Pclass	Sex	Age	Fare	Embarked_35	Embarked_C	Embarked_Q	Embarked_S
0	3	1	22.0	7.2500	0	0	0	1
1	1	0	38.0	71.2833	0	1	0	0
2	3	0	26.0	7.9250	0	0	0	1
3	1	0	35.0	53.1000	0	0	0	1
4	3	1	35.0	8.0500	0	0	0	1
...
886	2	1	27.0	13.0000	0	0	0	1
887	1	0	19.0	30.0000	0	0	0	1
888	3	0	28.0	23.4500	0	0	0	1
889	1	1	26.0	30.0000	0	1	0	0
890	3	1	32.0	7.7500	0	0	1	0

891 rows × 8 columns

In [23]:

```
1 from sklearn.model_selection import train_test_split
2 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.33,random_state=42)
```

In [24]:

```
1 x_train.isna().sum()
```

Out[24]:

```
Pclass      0
Sex          0
Age          0
Fare         0
Embarked_35  0
Embarked_C   0
Embarked_Q   0
Embarked_S   0
dtype: int64
```

Random Forest

```
In [25]: 1 #importing Randaom Forest Classifier from sklearn.ensemble
          2 %time
          3 from sklearn.model_selection import GridSearchCV #GridSearchCV is for parameter tuning
          4 from sklearn.ensemble import RandomForestClassifier
          5 cls=RandomForestClassifier()
          6 n_estimators=[25,50,75,100,125,150,175,200] #number of decision trees in the forest, default = 100
          7 criterion=['gini','entropy'] #criteria for choosing nodes default = 'gini'
          8 max_depth=[3,5,10] #maximum number of nodes in a tree default = None (it will go till all possible nodes)
          9 parameters={'n_estimators': n_estimators, 'criterion': criterion, 'max_depth': max_depth} #this will undergo
          10 RFC_cls = GridSearchCV(cls, parameters)
          11 RFC_cls.fit(x_train,y_train)
```

CPU times: user 7 µs, sys: 0 ns, total: 7 µs
Wall time: 14.5 µs

```
Out[25]: GridSearchCV(estimator=RandomForestClassifier(),
                      param_grid={'criterion': ['gini', 'entropy'],
                                   'max_depth': [3, 5, 10],
                                   'n_estimators': [25, 50, 75, 100, 125, 150, 175, 200]})
```

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

```
In [26]: 1 RFC_cls.best_params_
```

```
Out[26]: {'criterion': 'gini', 'max_depth': 10, 'n_estimators': 200}
```

```
In [27]: 1 cls=RandomForestClassifier(n_estimators=175,criterion='entropy',max_depth=5)
```

```
In [28]: 1 cls.fit(x_train,y_train)
```

```
Out[28]: RandomForestClassifier(criterion='entropy', max_depth=5, n_estimators=175)
```

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

```
In [29]: 1 rfy_pred=cls.predict(x_test)
```

```
In [30]: 1 rfy_pred
```

```
Out[30]: array([0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1,
                0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,
                1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0,
                0, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1,
                0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0,
                0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0,
                1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0,
                0, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 0,
                0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,
                1, 0, 1, 0, 0, 0, 1, 1, 0])
```

```
In [31]: 1 from sklearn.metrics import confusion_matrix
        2 confusion_matrix(y_test,rfy_pred)
```

```
Out[31]: array([[156,  19],
                [ 39,  81]])
```

```
In [32]: 1 from sklearn.metrics import accuracy_score
        2 accuracy_score(y_test,rfy_pred)
```

```
Out[32]: 0.8033898305084746
```

Logistic Regression

```
In [33]: 1 from sklearn.linear_model import LogisticRegression
          2 classifier=LogisticRegression()
          3 classifier.fit(x_train,y_train)
```

Out[33]: LogisticRegression()
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [34]: 1 x_test
```

Out[34]:

	Pclass	Sex	Age	Fare	Embarked_35	Embarked_C	Embarked_Q	Embarked_S
709	3	1	28.0	15.2458	0	1	0	0
439	2	1	31.0	10.5000	0	0	0	1
840	3	1	20.0	7.9250	0	0	0	1
720	2	0	6.0	33.0000	0	0	0	1
39	3	0	14.0	11.2417	0	1	0	0
...
715	3	1	19.0	7.6500	0	0	0	1
525	3	1	40.5	7.7500	0	0	1	0
381	3	0	1.0	15.7417	0	1	0	0
140	3	0	28.0	15.2458	0	1	0	0
173	3	1	21.0	7.9250	0	0	0	1

295 rows × 8 columns

```
In [35]: 1 y_pred=classifier.predict(x_test)
```

```
In [36]: 1 y_pred
```

```
Out[36]: array([0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0,
                1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0,
                1, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1,
                0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1,
                0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1,
                1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0,
                0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1,
                0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0,
                0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0,
                1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0,
                0, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1,
                0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0,
                0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0,
                1, 0, 0, 0, 0, 0, 1, 1, 0])
```

```
In [37]: 1 from sklearn.metrics import confusion_matrix
          2 confusion_matrix(y_test,y_pred)
```

```
Out[37]: array([[149, 26],
                [ 31, 89]])
```

```
In [38]: 1 from sklearn.metrics import accuracy_score
          2 accuracy_score(y_test,y_pred)
```

```
Out[38]: 0.8067796610169492
```

```
In [ ]: 1
```