```
In [1]: import pandas as pd
        import numpy as np
```

# read a file

```
In [26]: sai=pd.read_csv("/home/placement/Downloads/Titanic Dataset.csv")
```

```
In [3]: sai.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]: sai.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]: sa1=sai.drop(['Name','Ticket','Cabin','PassengerId','SibSp','Parch'],axis=1
```

```
In [6]: sa1
```

Out[6]:

| | 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 [7]: sa1['Sex']=sa1["Sex"].map({'male':1,'female':2})
```

```
In [8]: sa1
```

Out[8]:

| | Survived | Pclass | Sex | Age | Fare | Embarked |
|---|---|---|---|---|---|---|
| **0** | 0 | 3 | 1 | 22.0 | 7.2500 | S |
| **1** | 1 | 1 | 2 | 38.0 | 71.2833 | C |
| **2** | 1 | 3 | 2 | 26.0 | 7.9250 | S |
| **3** | 1 | 1 | 2 | 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 | 2 | 19.0 | 30.0000 | S |
| **888** | 0 | 3 | 2 | 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 [9]: sa1.isnull().sum()
```

```
Out[9]: Survived      0
        Pclass        0
        Sex           0
        Age         177
        Fare          0
        Embarked      2
        dtype: int64
```

```
In [10]: sa1['Age']=sa1['Age'].fillna(sa1['Age'].median())
         sa1
```

Out[10]:

| | Survived | Pclass | Sex | Age | Fare | Embarked |
|---|---|---|---|---|---|---|
| 0 | 0 | 3 | 1 | 22.0 | 7.2500 | S |
| 1 | 1 | 1 | 2 | 38.0 | 71.2833 | C |
| 2 | 1 | 3 | 2 | 26.0 | 7.9250 | S |
| 3 | 1 | 1 | 2 | 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 | 2 | 19.0 | 30.0000 | S |
| 888 | 0 | 3 | 2 | 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 [11]: sa1=pd.get_dummies(sa1,dtype=int)
         sa1
```

Out[11]:

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

891 rows × 8 columns

```
In [12]: sa1.isnull().sum()
```

```
Out[12]: Survived      0
         Pclass        0
         Sex           0
         Age           0
         Fare          0
         Embarked_C    0
         Embarked_Q    0
         Embarked_S    0
         dtype: int64
```

```
In [13]: y=sa1['Survived']
         x=sa1.drop('Survived',axis=1)
```

# split train and test

```
In [14]: from sklearn.model_selection import train_test_split
         x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.33,random_st
```

```
In [15]: x_train.head()
```

Out[15]:

| | Pclass | Sex | Age | Fare | Embarked_C | Embarked_Q | Embarked_S |
|---|---|---|---|---|---|---|---|
| 6 | 1 | 1 | 54.0 | 51.8625 | 0 | 0 | 1 |
| 718 | 3 | 1 | 28.0 | 15.5000 | 0 | 1 | 0 |
| 685 | 2 | 1 | 25.0 | 41.5792 | 1 | 0 | 0 |
| 73 | 3 | 1 | 26.0 | 14.4542 | 1 | 0 | 0 |
| 882 | 3 | 2 | 22.0 | 10.5167 | 0 | 0 | 1 |

```
In [16]: y_test.head()
```

```
Out[16]: 709    1
         439    0
         840    0
         720    1
         39     1
         Name: Survived, dtype: int64
```

# random forest

```
In [17]: from sklearn.model_selection import GridSearchCV #GridSearchCV is for param
         from sklearn.ensemble import RandomForestClassifier
         cls=RandomForestClassifier()
         n_estimators=[25,50,75,100,125,150,175,200] #number of decision trees in th
         criterion=['gini','entropy'] #criteria for choosing nodes default = 'gini'
         max_depth=[3,5,10] #maximum number of nodes in a tree default = None (it wi
         parameters={'n_estimators': n_estimators,'criterion':criterion,'max_depth':
         RFC_cls = GridSearchCV(cls, parameters)
         RFC_cls.fit(x_train,y_train)
```

Out[17]: 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 [18]: RFC_cls.best_params_
```

Out[18]: {'criterion': 'entropy', 'max_depth': 5, 'n_estimators': 150}

```
In [19]: cls=RandomForestClassifier(n_estimators=25,criterion='entropy',max_depth=10
```

```
In [20]: cls.fit(x_train,y_train)
```

Out[20]: RandomForestClassifier(criterion='entropy', max_depth=10, n_estimators=25)

**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 [21]: rfy_pred=cls.predict(x_test)
```

```
In [22]: rfy_pred
```

Out[22]: array([0, 0, 0, 1, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1,
                0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 0, 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, 1, 1, 1,
                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, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0, 1, 0,
                0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 1,
                0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0,
                0, 1, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0,
                1, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 1, 0,
                0, 1, 0, 1, 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, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0,
                0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0,
                1, 0, 1, 1, 0, 0, 1, 1, 0])
```

```
In [23]: from sklearn.metrics import confusion_matrix
         confusion_matrix(y_test,rfy_pred)
```

Out[23]: array([[148,  27],
                [ 35,  85]])

# EFFICENCY OF THE CONFUSION MATRIX

```
In [24]: from sklearn.metrics import accuracy_score
         accuracy_score(y_test,rfy_pred)
```

Out[24]: 0.7898305084745763