# Data Cleaning & Analysis

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

In [2]: ▶|
```python
titanic=pd.read_csv(r'C:\Users\yogay\OneDrive\Desktop\Yogita_Yadav\Data Science\1st\Titanic dataset analysis\DATASET\train.csv',header = 0, dtype={'Age': np.float64
```

In [3]: ▶|
```python
titanic.tail()
```

Out[3]:

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **886** | 887 | 0 | 2 | Montvila, Rev. Juozas | male | 27.0 | 0 | 0 | 211536 | 13.00 | NaN | S |
| **887** | 888 | 1 | 1 | Graham, Miss. Margaret Edith | female | 19.0 | 0 | 0 | 112053 | 30.00 | B42 | S |
| **888** | 889 | 0 | 3 | Johnston, Miss. Catherine Helen "Carrie" | female | NaN | 1 | 2 | W./C. 6607 | 23.45 | NaN | S |
| **889** | 890 | 1 | 1 | Behr, Mr. Karl Howell | male | 26.0 | 0 | 0 | 111369 | 30.00 | C148 | C |
| **890** | 891 | 0 | 3 | Dooley, Mr. Patrick | male | 32.0 | 0 | 0 | 370376 | 7.75 | NaN | Q |

In [4]: ▶|
```python
titanic.describe()
```

Out[4]:

| | 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 [5]: ▶|
```python
titanic.shape
```

Out[5]: (891, 12)

In [8]: ▶
```python
del titanic["Name"]
titanic.head()
```

Out[8]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| **1** | 2 | 1 | 1 | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| **2** | 3 | 1 | 3 | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| **3** | 4 | 1 | 1 | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| **4** | 5 | 0 | 3 | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |

In [9]: ▶
```python
del titanic["Ticket"]
titanic.head()
```

Out[9]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | male | 22.0 | 1 | 0 | 7.2500 | NaN | S |
| **1** | 2 | 1 | 1 | female | 38.0 | 1 | 0 | 71.2833 | C85 | C |
| **2** | 3 | 1 | 3 | female | 26.0 | 0 | 0 | 7.9250 | NaN | S |
| **3** | 4 | 1 | 1 | female | 35.0 | 1 | 0 | 53.1000 | C123 | S |
| **4** | 5 | 0 | 3 | male | 35.0 | 0 | 0 | 8.0500 | NaN | S |

In [10]: ▶
```python
del titanic["Fare"]
titanic.head()
```

Out[10]:

| | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | male | 22.0 | 1 | 0 | NaN | S |
| **1** | 2 | 1 | 1 | female | 38.0 | 1 | 0 | C85 | C |
| **2** | 3 | 1 | 3 | female | 26.0 | 0 | 0 | NaN | S |
| **3** | 4 | 1 | 1 | female | 35.0 | 1 | 0 | C123 | S |
| **4** | 5 | 0 | 3 | male | 35.0 | 0 | 0 | NaN | S |

In [11]: ▶
```python
del titanic["Cabin"]
titanic.head()
```

Out[11]:

|   | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Embarked |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | male | 22.0 | 1 | 0 | S |
| **1** | 2 | 1 | 1 | female | 38.0 | 1 | 0 | C |
| **2** | 3 | 1 | 3 | female | 26.0 | 0 | 0 | S |
| **3** | 4 | 1 | 1 | female | 35.0 | 1 | 0 | S |
| **4** | 5 | 0 | 3 | male | 35.0 | 0 | 0 | S |

In [12]: ▶
```python
# Changing Value for "Male, Female" string values to numeric values , male=1 and female=2
def getNumber(str):
    if str=="male":
        return 1
    else:
        return 2
titanic["Gender"]=titanic["Sex"].apply(getNumber)
#We have created a new column called "Gender" and
#filling it with values 1,2 based on the values of sex column
titanic.head()
```

Out[12]:

|   | PassengerId | Survived | Pclass | Sex | Age | SibSp | Parch | Embarked | Gender |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | male | 22.0 | 1 | 0 | S | 1 |
| **1** | 2 | 1 | 1 | female | 38.0 | 1 | 0 | C | 2 |
| **2** | 3 | 1 | 3 | female | 26.0 | 0 | 0 | S | 2 |
| **3** | 4 | 1 | 1 | female | 35.0 | 1 | 0 | S | 2 |
| **4** | 5 | 0 | 3 | male | 35.0 | 0 | 0 | S | 1 |

In [13]: ▶
```python
del titanic["Sex"]
titanic.head()
```

Out[13]:

|   | PassengerId | Survived | Pclass | Age | SibSp | Parch | Embarked | Gender |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | 22.0 | 1 | 0 | S | 1 |
| **1** | 2 | 1 | 1 | 38.0 | 1 | 0 | C | 2 |
| **2** | 3 | 1 | 3 | 26.0 | 0 | 0 | S | 2 |
| **3** | 4 | 1 | 1 | 35.0 | 1 | 0 | S | 2 |
| **4** | 5 | 0 | 3 | 35.0 | 0 | 0 | S | 1 |

```
In [14]:  ▶  titanic.isnull().sum()
```

```
Out[14]:  PassengerId      0
          Survived         0
          Pclass           0
          Age            177
          SibSp            0
          Parch            0
          Embarked         2
          Gender           0
          dtype: int64
```

```
In [15]:  ▶  meanS= titanic[titanic.Survived==1].Age.mean()
             meanS
```

```
Out[15]:  28.343689655172415
```

```
In [16]:  ▶  titanic["age"]=np.where(pd.isnull(titanic.Age) & titanic["Survived"]==1 ,meanS, titanic["Age"])
             titanic.head()
```

Out[16]:

| | PassengerId | Survived | Pclass | Age | SibSp | Parch | Embarked | Gender | age |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 22.0 | 1 | 0 | S | 1 | 22.0 |
| 1 | 2 | 1 | 1 | 38.0 | 1 | 0 | C | 2 | 38.0 |
| 2 | 3 | 1 | 3 | 26.0 | 0 | 0 | S | 2 | 26.0 |
| 3 | 4 | 1 | 1 | 35.0 | 1 | 0 | S | 2 | 35.0 |
| 4 | 5 | 0 | 3 | 35.0 | 0 | 0 | S | 1 | 35.0 |

```
In [17]:  ▶  titanic.isnull().sum()
```

```
Out[17]:  PassengerId      0
          Survived         0
          Pclass           0
          Age            177
          SibSp            0
          Parch            0
          Embarked         2
          Gender           0
          age            125
          dtype: int64
```

```
In [18]:  ▶  meanNS=titanic[titanic.Survived==0].Age.mean()
             meanNS
```

```
Out[18]:  30.62617924528302
```

In [19]: ▶|
```python
titanic.age.fillna(meanNS,inplace=True)
titanic.head()
```

Out[19]:

| | PassengerId | Survived | Pclass | Age | SibSp | Parch | Embarked | Gender | age |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | 22.0 | 1 | 0 | S | 1 | 22.0 |
| **1** | 2 | 1 | 1 | 38.0 | 1 | 0 | C | 2 | 38.0 |
| **2** | 3 | 1 | 3 | 26.0 | 0 | 0 | S | 2 | 26.0 |
| **3** | 4 | 1 | 1 | 35.0 | 1 | 0 | S | 2 | 35.0 |
| **4** | 5 | 0 | 3 | 35.0 | 0 | 0 | S | 1 | 35.0 |

In [20]: ▶|
```python
titanic.isnull().sum()
```

Out[20]:
```
PassengerId      0
Survived         0
Pclass           0
Age            177
SibSp            0
Parch            0
Embarked         2
Gender           0
age              0
dtype: int64
```

In [21]: ▶|
```python
del titanic['Age']
titanic.head()
```

Out[21]:

| | PassengerId | Survived | Pclass | SibSp | Parch | Embarked | Gender | age |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | 1 | 0 | S | 1 | 22.0 |
| **1** | 2 | 1 | 1 | 1 | 0 | C | 2 | 38.0 |
| **2** | 3 | 1 | 3 | 0 | 0 | S | 2 | 26.0 |
| **3** | 4 | 1 | 1 | 1 | 0 | S | 2 | 35.0 |
| **4** | 5 | 0 | 3 | 0 | 0 | S | 1 | 35.0 |

In [24]: ▶|
```python
import warnings
warnings.filterwarnings('ignore')
```

In [25]: ▶|
```python
survivedQ = titanic[titanic.Embarked == 'Q'][titanic.Survived == 1].shape[0]
survivedC = titanic[titanic.Embarked == 'C'][titanic.Survived == 1].shape[0]
survivedS = titanic[titanic.Embarked == 'S'][titanic.Survived == 1].shape[0]
print(survivedQ)
print(survivedC)
print(survivedS)
```

```
30
93
217
```

In [26]: ▶|
```python
survivedQ = titanic[titanic.Embarked == 'Q'][titanic.Survived == 0].shape[0]
survivedC = titanic[titanic.Embarked == 'C'][titanic.Survived == 0].shape[0]
survivedS = titanic[titanic.Embarked == 'S'][titanic.Survived == 0].shape[0]
print(survivedQ)
print(survivedC)
print(survivedS)
```

```
47
75
427
```

In [27]: ▶|
```python
titanic.dropna(inplace=True)
titanic.head()
```

Out[27]:

| | PassengerId | Survived | Pclass | SibSp | Parch | Embarked | Gender | age |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | 1 | 0 | S | 1 | 22.0 |
| 1 | 2 | 1 | 1 | 1 | 0 | C | 2 | 38.0 |
| 2 | 3 | 1 | 3 | 0 | 0 | S | 2 | 26.0 |
| 3 | 4 | 1 | 1 | 1 | 0 | S | 2 | 35.0 |
| 4 | 5 | 0 | 3 | 0 | 0 | S | 1 | 35.0 |

In [28]: ▶|
```python
titanic.isnull().sum()
```

Out[28]:
```
PassengerId    0
Survived       0
Pclass         0
SibSp          0
Parch          0
Embarked       0
Gender         0
age            0
dtype: int64
```

In [29]:  ▶| 
```python
titanic.rename(columns={'age':'Age'}, inplace=True)
titanic.head()
```

Out[29]:

|   | PassengerId | Survived | Pclass | SibSp | Parch | Embarked | Gender | Age |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | 1 | 0 | S | 1 | 22.0 |
| **1** | 2 | 1 | 1 | 1 | 0 | C | 2 | 38.0 |
| **2** | 3 | 1 | 3 | 0 | 0 | S | 2 | 26.0 |
| **3** | 4 | 1 | 1 | 1 | 0 | S | 2 | 35.0 |
| **4** | 5 | 0 | 3 | 0 | 0 | S | 1 | 35.0 |

In [30]:  ▶| 
```python
titanic.rename(columns={'Gender':'Sex'}, inplace=True)
titanic.head()
```

Out[30]:

|   | PassengerId | Survived | Pclass | SibSp | Parch | Embarked | Sex | Age |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | 1 | 0 | S | 1 | 22.0 |
| **1** | 2 | 1 | 1 | 1 | 0 | C | 2 | 38.0 |
| **2** | 3 | 1 | 3 | 0 | 0 | S | 2 | 26.0 |
| **3** | 4 | 1 | 1 | 1 | 0 | S | 2 | 35.0 |
| **4** | 5 | 0 | 3 | 0 | 0 | S | 1 | 35.0 |

In [31]:  ▶| 
```python
def getEmb(str):
    if str=="S":
        return 1
    elif str=='Q':
        return 2
    else:
        return 3
titanic["Embark"]=titanic["Embarked"].apply(getEmb)
titanic.head()
```

Out[31]:

|   | PassengerId | Survived | Pclass | SibSp | Parch | Embarked | Sex | Age | Embark |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | 1 | 0 | S | 1 | 22.0 | 1 |
| **1** | 2 | 1 | 1 | 1 | 0 | C | 2 | 38.0 | 3 |
| **2** | 3 | 1 | 3 | 0 | 0 | S | 2 | 26.0 | 1 |
| **3** | 4 | 1 | 1 | 1 | 0 | S | 2 | 35.0 | 1 |
| **4** | 5 | 0 | 3 | 0 | 0 | S | 1 | 35.0 | 1 |

In [32]: ▶| ```python
del titanic['Embarked']
titanic.rename(columns={'Embark':'Embarked'}, inplace=True)
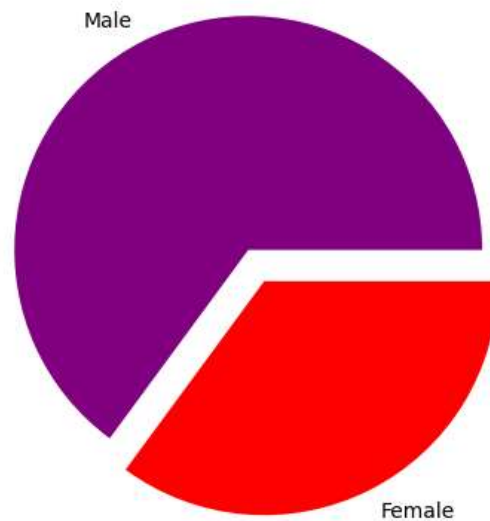titanic.head()
```

Out[32]:

| | PassengerId | Survived | Pclass | SibSp | Parch | Sex | Age | Embarked |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | 1 | 0 | 1 | 22.0 | 1 |
| **1** | 2 | 1 | 1 | 1 | 0 | 2 | 38.0 | 3 |
| **2** | 3 | 1 | 3 | 0 | 0 | 2 | 26.0 | 1 |
| **3** | 4 | 1 | 1 | 1 | 0 | 2 | 35.0 | 1 |
| **4** | 5 | 0 | 3 | 0 | 0 | 1 | 35.0 | 1 |

In [33]:

```python
#Drawing a pie chart for number of males and females aboard
import matplotlib.pyplot as plt
from matplotlib import style

males = (titanic['Sex'] == 1).sum()
#Summing up all the values of column gender with a
#condition for male and similary for females
females = (titanic['Sex'] == 2).sum()
print(males)
print(females)
p = [males, females]
plt.pie(p,      #giving array
        labels = ['Male', 'Female'], #Correspndingly giving labels
        colors = ['purple', 'red'],   # Corresponding colors
        explode = (0.15, 0),     #How much the gap should me there between the pies
        startangle = 0)  #what start angle should be given
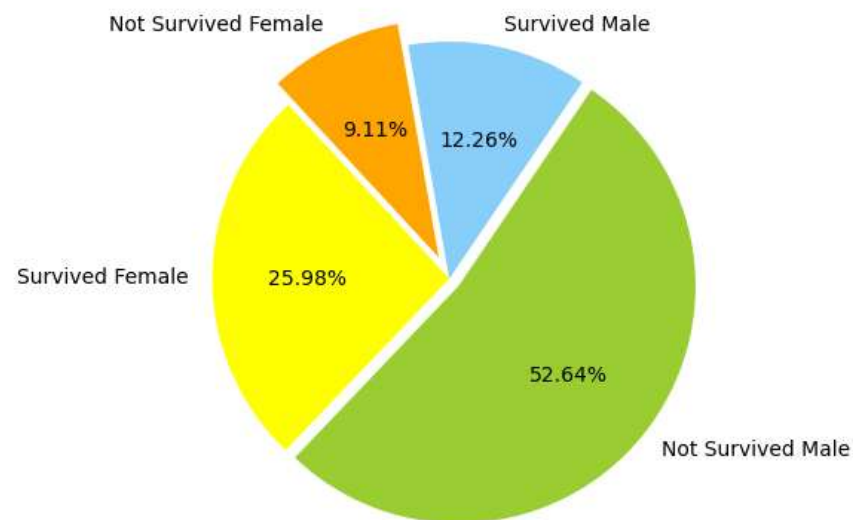plt.axis('equal')
plt.show()
```

577
312

In [34]: ▶|
```python
MaleS=titanic[titanic.Sex==1][titanic.Survived==1].shape[0]
print(MaleS)
MaleN=titanic[titanic.Sex==1][titanic.Survived==0].shape[0]
print(MaleN)
FemaleS=titanic[titanic.Sex==2][titanic.Survived==1].shape[0]
print(FemaleS)
FemaleN=titanic[titanic.Sex==2][titanic.Survived==0].shape[0]
print(FemaleN)
```

```
109
468
231
81
```

In [35]: ▶|
```python
chart=[MaleS,MaleN,FemaleS,FemaleN]
colors=['lightskyblue','yellowgreen','Yellow','Orange']
labels=["Survived Male","Not Survived Male","Survived Female","Not Survived Female"]
explode=[0,0.05,0,0.1]
plt.pie(chart,labels=labels,colors=colors,explode=explode,startangle=100,counterclock=False,autopct="%.2f%%")
plt.axis("equal")
plt.show()
```



In [ ]: ▶|