

```

In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: titanic_data = pd.read_csv(r"C:\Users\hp\Downloads\titanic (1).csv")

In [3]: titanic_data
Out[3]:
   PassengerId  Survived  Pclass     Name  Sex  Age  SibSp  Parch    Ticket   Fare Cabin Embarked
0            0         0      3    Braund, Mr. Owen Harris   male  22.0      1      0      A/5 21171   7.2500   NaN      S
1            1         0      3  Cumings, Mrs. John Bradley (Florence Briggs Th... female  38.0      1      0      PC 17599   71.2833   C85      C
2            2         1      3    Heikinen, Miss. Laina   female  26.0      0      0  STON/O2 3101282   7.9250   NaN      S
3            3         1      1  Furelle, Mrs. Jacques Heath (Lily May Peel) female  35.0      1      0    113803  53.1000  C123      S
4            4         0      3    Allen, Mr. William Henry   male  35.0      0      0    373450   8.0500   NaN      S
...         ...         ...    ...     ...  ...  ...    ...    ...    ...     ...   ...     ...
886         887         0      2    Morvila, Rev. Juozas   male  27.0      0      0    211536  13.0000   NaN      S
887         888         1      1    Graham, Miss. Margaret Edith   female  19.0      0      0    112053  30.0000   B42      S
888         889         0      3    Johnson, Miss. Catherine Helen "Carnie" female  NaN      1      2    W/C 6607  23.4500   NaN      S
889         890         1      1    Behr, Mr. Karl Howell   male  26.0      0      0    111369  30.0000  C148      C
890         891         0      3    Dooley, Mr. Patrick   male  32.0      0      0    370376   7.7500   NaN      Q
891 rows x 12 columns

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

In [5]: titanic_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            284 non-null    object
11   Embarked          889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB

In [6]: #Handle missing values using mean mode median imputation
# Calculate Mean
mean = titanic_data["Age"].mean()
mean
Out[6]: 29.6991764785882

In [7]: #Calculate Median
median = titanic_data["Age"].median()
median
Out[7]: 29.0

In [8]: #Calculate mode
mode=titanic_data["Age"].mode()[0]
mode
Out[8]: 24.0

In [9]: titanic_data.describe()
Out[9]:
   PassengerId  Survived  Pclass     Age  SibSp  Parch     Fare
count  891.000000  891.000000  891.000000  891.000000  891.000000  891.000000
mean    440.000000  0.383838  2.309642  29.699176  0.523008  0.381594  32.204208
std     257.353842  0.486992  0.836071  14.526497  1.102743  0.809057  49.693429
min       0.000000  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%     440.000000  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  512.329000

In [10]: #Random Sample Imputation
df = pd.read_csv(r"C:\Users\hp\Downloads\titanic (1).csv", usecols=['Age', 'Fare', 'Survived'])
df.head()
Out[10]:
   Survived  Age  Fare
0          0  22.0  7.2500
1          1  38.0  71.2833
2          1  26.0  7.9250
3          1  35.0  53.1000
4          0  35.0  8.0500

In [11]: #Check null values
titanic_data.isnull().sum()
Out[11]:
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 [12]: #Where is 177 null values are present in Age feature.
#How we replace this NaN by using Random Sample Imputation
titanic_data["Age"].dropna().sample()
Out[12]:
115    21.0
Name: Age, dtype: float64

In [13]: #Where is 177 null values are present in Age feature and .sample() function return any one random value
titanic_data["Age"].dropna().sample(df["Age"].isnull().sum())
Out[13]:
899    26.0
172     1.0
230    35.0
286    21.0
588    22.0
...
90     29.0
16     2.0
382    32.0
12     26.0
726    26.0
Name: Age, Length: 177, dtype: float64

In [14]: #This function check where is null values are present and replace NaN with random sample
#We use random_state because it replace NaN with specific value only ( if we not use random_state then values change everytime when we run .)
titanic_data["Age"].dropna().sample(titanic_data["Age"].isnull().sum(), random_state=0)
Out[14]:
423    28.00
177    50.00
385     0.52
292    36.00
889    26.00
...
539    22.00
267    25.00
352    15.00
99     34.00
609    15.00
Name: Age, Length: 177, dtype: float64

In [15]: #Rows and columns
titanic_data.shape
Out[15]: (891, 12)

In [16]: #Count of males and females survived
titanic_data["Survived"].value_counts()
Out[16]:
Survived
0    549
1    342
Name: count, dtype: int64

In [17]: #Graph between survived and pclass
sns.countplot(x = titanic_data["Survived"], hue = titanic_data["Pclass"])
Out[17]:
<Axes: xlabel='Survived', ylabel='count'>


In [18]: titanic_data["Sex"]
Out[18]:
0    male
1    female
2    female
3    female
4    male
...
886   male
887   female
888   female
889   male
890   male
Name: Sex, Length: 891, dtype: object

In [19]: #Graph between survived and gender
sns.countplot(x = titanic_data["Sex"], hue = titanic_data["Survived"])
Out[19]:
<Axes: xlabel='Sex', ylabel='count'>


In [20]: from sklearn.preprocessing import LabelEncoder
labelencoder = LabelEncoder()
titanic_data["Sex"] = labelencoder.fit_transform(titanic_data["Sex"])
titanic_data.head()
Out[20]:
   PassengerId  Survived  Pclass     Name  Sex  Age  SibSp  Parch    Ticket   Fare Cabin Embarked
0            0         0      3    Braund, Mr. Owen Harris   1  22.0      1      0      A/5 21171   7.2500   NaN      S
1            1         0      3  Cumings, Mrs. John Bradley (Florence Briggs Th...   0  38.0      1      0      PC 17599   71.2833   C85      C
2            2         1      3    Heikinen, Miss. Laina   0  26.0      0      0  STON/O2 3101282   7.9250   NaN      S
3            3         1      1  Furelle, Mrs. Jacques Heath (Lily May Peel)   0  35.0      1      0    113803  53.1000  C123      S
4            4         0      3    Allen, Mr. William Henry   1  35.0      0      0    373450   8.0500   NaN      S

In [21]: titanic_data["Sex"], titanic_data["Survived"]
Out[21]:
(0      1
 1      0
 2      0
 3      0
 4      1
...
886     1
887     0
888     0
889     1
890     1
Name: Sex, Length: 891, dtype: int32,
 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 [22]: sns.countplot(x = titanic_data["Sex"], hue = titanic_data["Survived"])
Out[22]:
<Axes: xlabel='Sex', ylabel='count'>

```