

Task - 2[Train]

Importing necessary libraries

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

Reading the dataset

```
In [2]: df=pd.read_csv("train.csv")
```

```
In [3]: df
```

Out[3]:	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
...
886	887	0	2	Montvila, 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	Johnston, Miss. Catherine Helen "Carrie"	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 × 12 columns

```
In [4]: df.head()
```

Out[4]:	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

Some information about the dataset

```
In [5]: df.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 [6]: df.describe()
```

	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

Checking for the duplicate values

```
In [7]: df.duplicated().sum()
```

```
Out[7]: np.int64(0)
```

```
In [8]: numeric_columns = df.select_dtypes(include=['number']).columns
df.groupby('Survived')[numeric_columns].mean()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
Survived							
0	447.016393	0.0	2.531876	30.626179	0.553734	0.329690	22.117887
1	444.368421	1.0	1.950292	28.343690	0.473684	0.464912	48.395408

```
In [9]: df['Age'].fillna(df['Age'].mean(),inplace=True)
df['Embarked'].fillna(df['Embarked'].mode()[0],inplace=True)
```

```
In [10]: df.isnull().sum()
```

```
Out[10]: PassengerId    0
Survived              0
Pclass               0
Name                 0
Sex                  0
Age                  0
SibSp                0
Parch                0
Ticket              0
Fare                 0
Cabin               687
Embarked            0
dtype: int64
```

```
In [11]: df['Cabin'].fillna('Unknown', inplace=True)
```

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

```
Out[12]: PassengerId    0
         Survived      0
         Pclass      0
         Name        0
         Sex         0
         Age         0
         SibSp       0
         Parch       0
         Ticket      0
         Fare        0
         Cabin       0
         Embarked    0
         dtype: int64
```

```
In [13]: df['Survived'].value_counts()
```

```
Out[13]: Survived
0      549
1      342
Name: count, dtype: int64
```

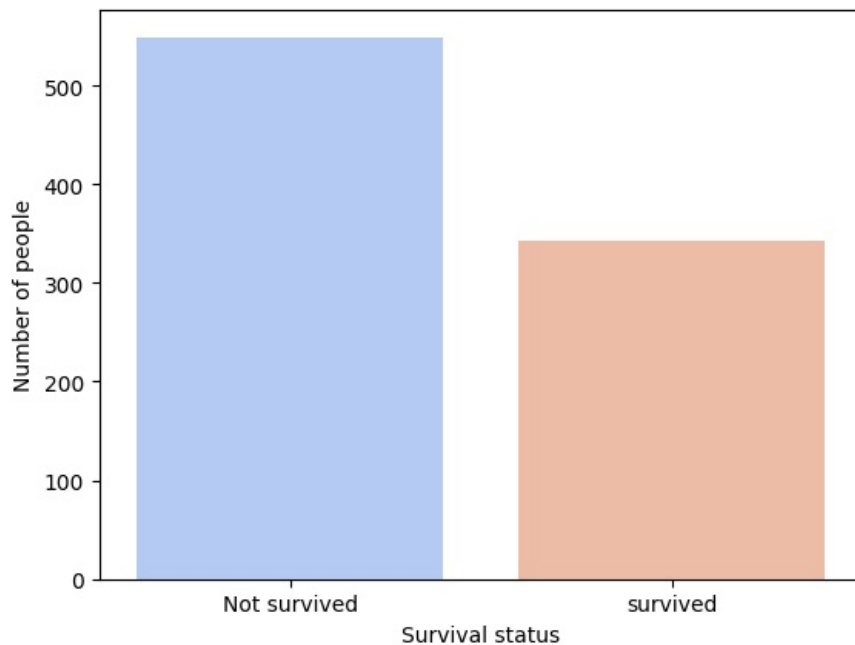
Visualisation

```
In [14]: sns.countplot(x='Survived',data=df,palette='coolwarm',)
plt.xlabel("Survival status")
plt.ylabel("Number of people")
plt.xticks(ticks=[0,1],labels=['Not survived','survived'])
plt.show()
```

C:\Users\Komal\AppData\Local\Temp\ipykernel_28112\1486452071.py:1: FutureWarning:

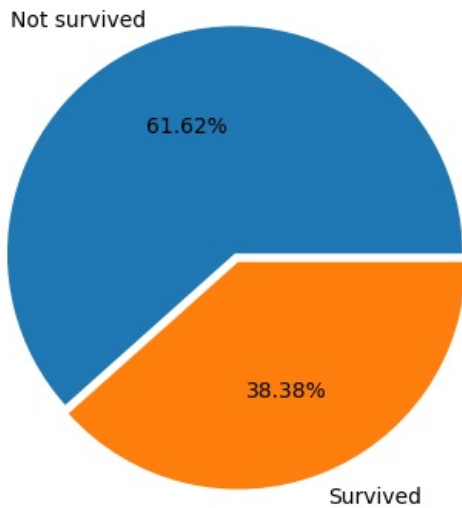
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.countplot(x='Survived',data=df,palette='coolwarm',)
```



```
In [15]: plt.pie(df['Survived'].value_counts(),explode=[0,0.04],autopct="%1.2f%%",labels=['Not survived','Survived'])
plt.title("Survival of people")
plt.show()
```

Survival of people

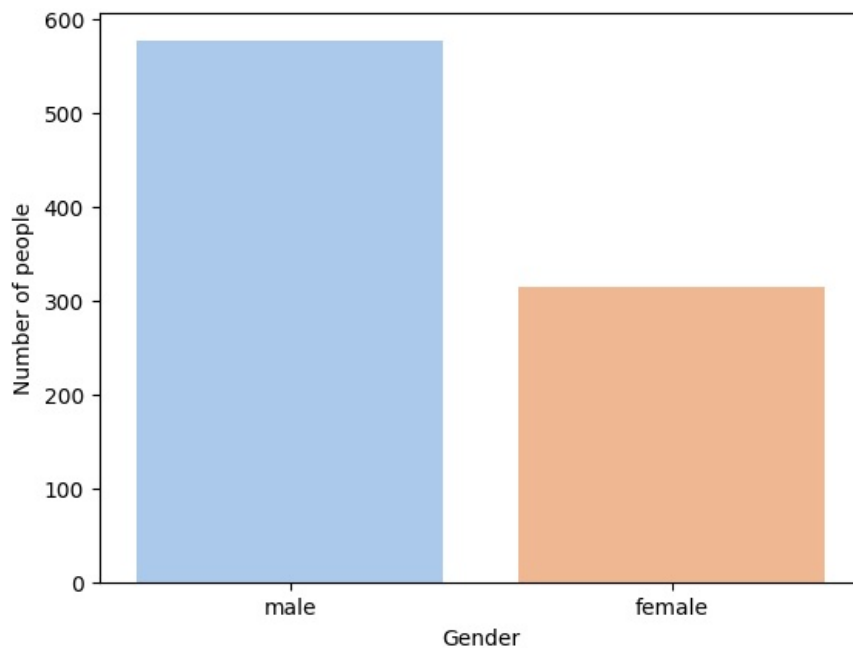


```
In [16]: sns.countplot(x='Sex',data=df,palette='pastel',)  
plt.xlabel("Gender")  
plt.ylabel("Number of people")  
plt.show()
```

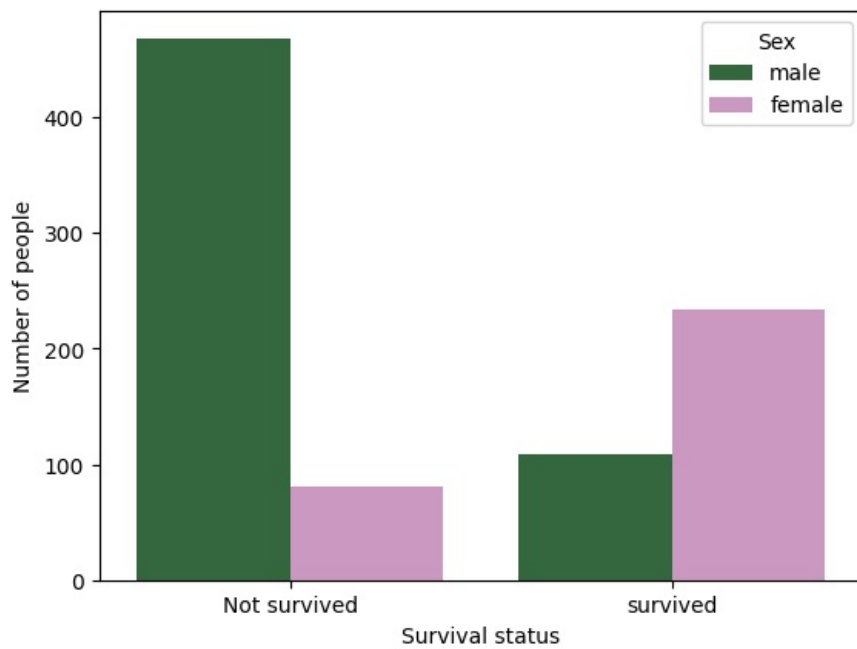
C:\Users\Komal\AppData\Local\Temp\ipykernel_28112\1148348333.py:1: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

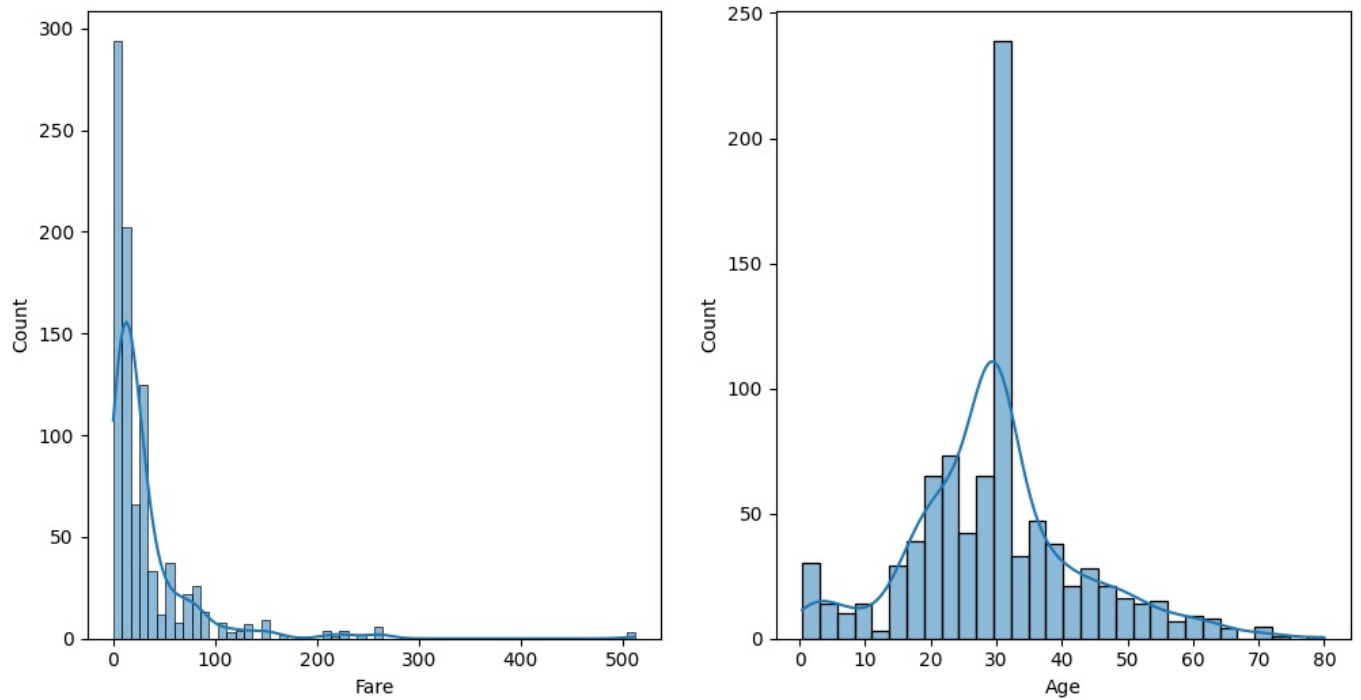
```
sns.countplot(x='Sex',data=df,palette='pastel',)
```



```
In [17]: sns.countplot(x='Survived',hue='Sex',data=df,palette='cubehelix',)  
plt.xlabel("Survival status")  
plt.ylabel("Number of people")  
plt.xticks(ticks=[0,1],labels=['Not survived','survived'])  
plt.show()
```



```
In [18]: fig, axes = plt.subplots(1, 2, figsize=(12, 6))
sns.histplot(df['Fare'], kde=True, ax=axes[0])
sns.histplot(df['Age'].dropna(), kde=True, ax=axes[1])
plt.show()
```



```
In [ ]:
```