

```
In [28]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
```

```
In [4]: data=pd.read_csv("task2.csv")
data
```

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	Cummings, 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	3	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	889	0	1	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 [5]: data.head(5)
```

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	Cummings, 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]: data.describe()
```

Out[6]:

	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 [7]: data.columns
```

```
Out[7]: Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp', 'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],
      dtype='object')
```

```
In [8]: 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 [9]: data.isnull().any()
```

```
Out[9]: PassengerId    False
Survived          False
Pclass            False
Name              False
Sex              False
Age              True
SibSp            False
Parch            False
Ticket           False
Fare             False
Cabin            True
Embarked         True
dtype: bool
```

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

```
Out[10]: 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 [11]: data['Age']=data['Age'].fillna(data['Age'].median())
data['Age']
```

```
Out[11]: 0    22.0
1    38.0
2    26.0
3    35.0
4    35.0
...
886   27.0
887   19.0
888   28.0
889   26.0
890   32.0
Name: Age, Length: 891, dtype: float64
```

```
In [14]: def fill_cabin(pclass,cabin):
if pd.isnull(cabin):
    if pclass==1:
        return'X'
    elif pclass==2:
        return'Y'
    elif pclass==3:
        return'Z'
    else:
        return cabin
data["Cabin"]=data.apply(lambda row :fill_cabin(row["Pclass"],row["Cabin"]),axis=1)
data["Cabin"]
```

```
Out[14]: 0    Z
1    C85
2    Z
3    C123
4    Z
...
886   Y
887   B42
888   Z
889   C148
890   Z
Name: Cabin, Length: 891, dtype: object
```

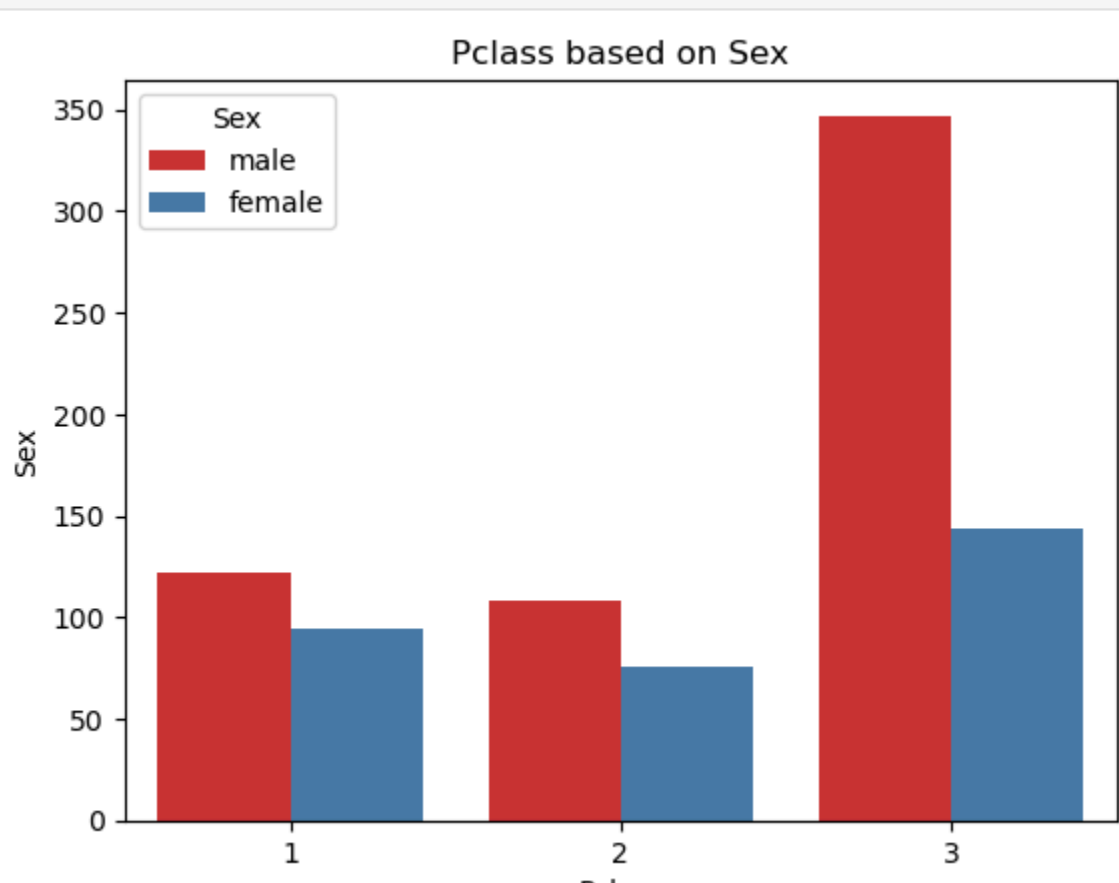
```
In [15]: data["Embarked"]=data["Embarked"].fillna(data["Embarked"].mode()[0])
data["Embarked"]
```

```
Out[15]: 0    S
1    C
2    S
3    S
4    S
...
886   S
887   S
888   S
889   C
890   Q
Name: Embarked, Length: 891, dtype: object
```

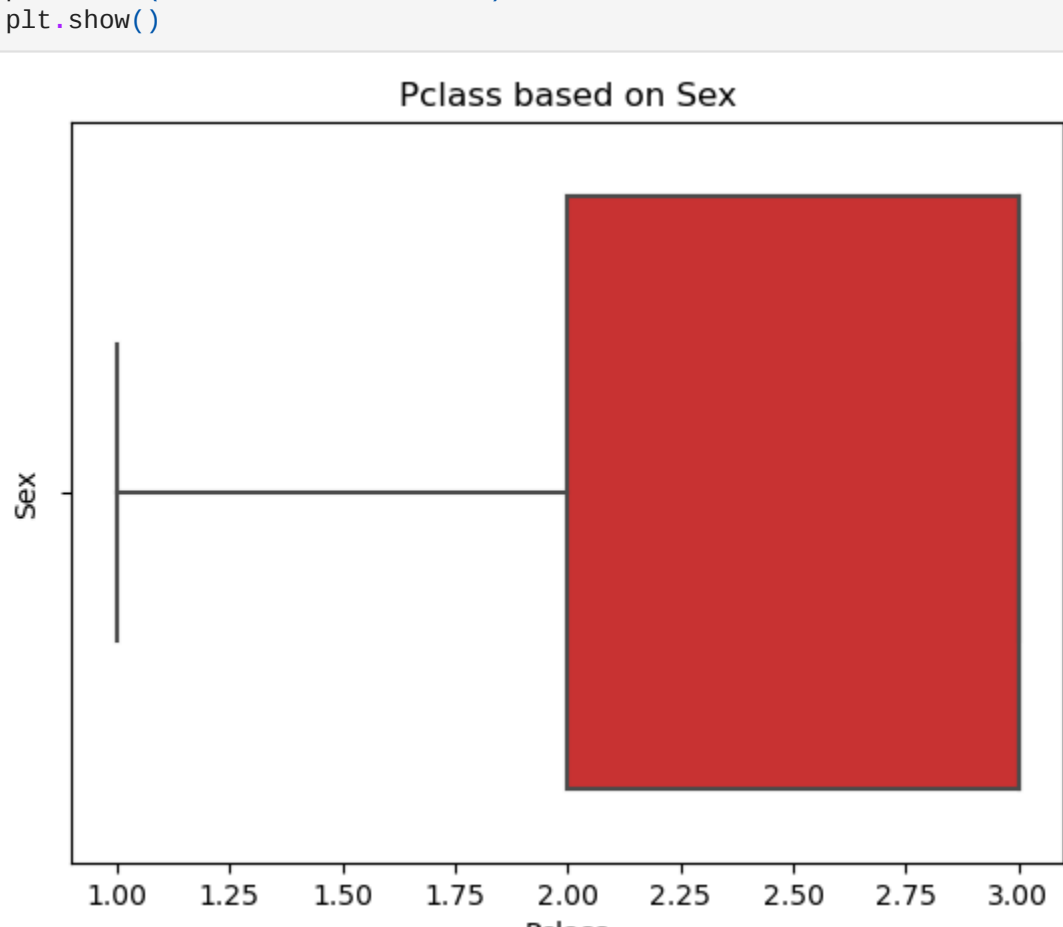
```
In [16]: data.isnull().any()
```

```
Out[16]: PassengerId    False
Survived          False
Pclass            False
Name              False
Sex              False
Age              False
SibSp            False
Parch            False
Ticket           False
Fare             False
Cabin            False
Embarked         False
dtype: bool
```

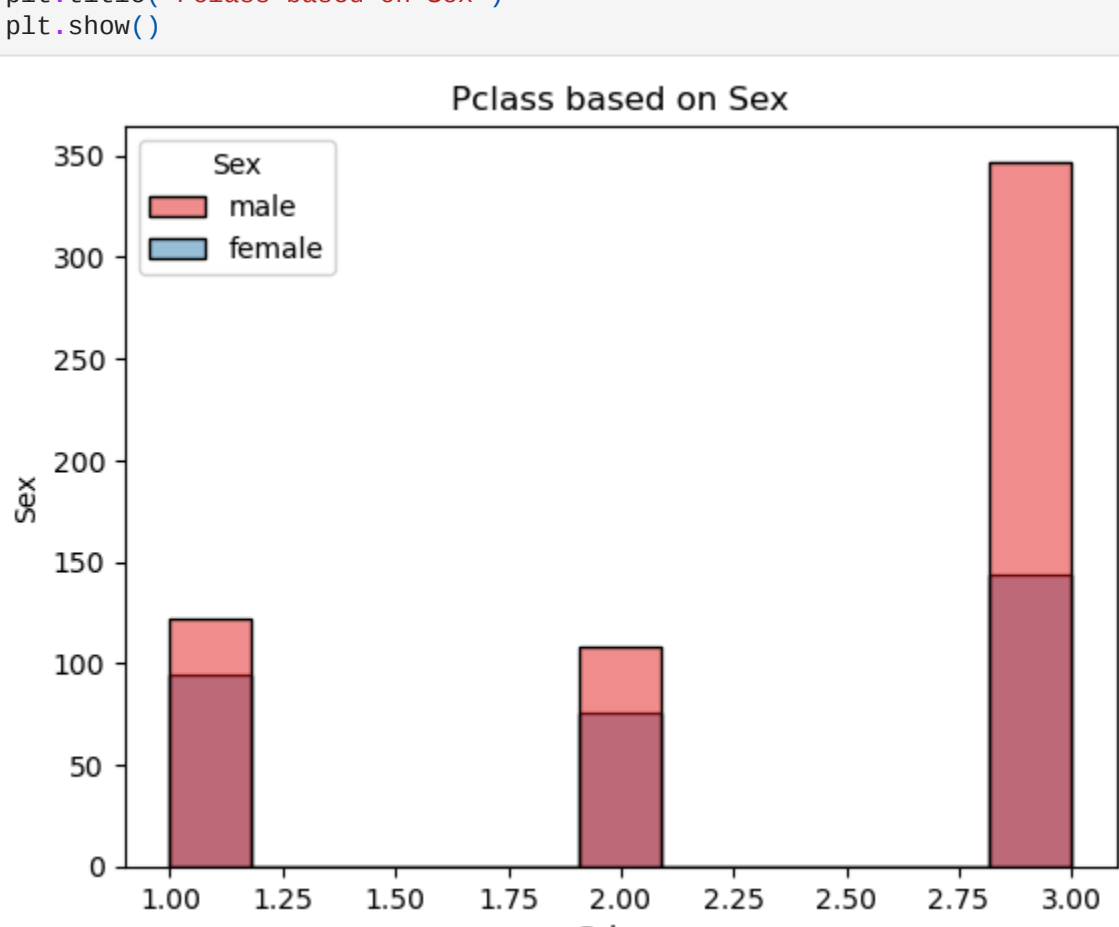
```
In [18]: sns.countplot(x="Pclass",hue="Sex",data=data,palette="Set1")
plt.xlabel("Pclass")
plt.ylabel("Sex")
plt.title("Pclass based on Sex")
plt.show()
```



```
In [19]: sns.boxplot(x="Pclass",hue="Sex",data=data,palette="Set1")
plt.xlabel("Pclass")
plt.ylabel("Sex")
plt.title("Pclass based on Sex")
plt.show()
```



```
In [20]: sns.histplot(x="Pclass",hue="Sex",data=data,palette="Set1")
plt.xlabel("pclass")
plt.ylabel("Sex")
plt.title("Pclass based on Sex")
plt.show()
```



```
In [23]: sns.heatmap(data.select_dtypes(include=np.number).corr(),annot=True,cmap="coolwarm",fmt=".1f",linewidths=0.5)
plt.title("correlation matrixx")
plt.show()
```

