

```
In [67]: 1 # Importing necessary Libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.preprocessing import LabelEncoder, StandardScaler
7 from sklearn.model_selection import train_test_split
```

```
In [68]: 1 # Importing the dataset.
2 dataset=pd.read_csv("Titanic-Dataset.csv")
3 dataset.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 [69]: 1 # Checking for Null Values.
2 dataset.isnull().any()
```

```
Out[69]: 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 [70]: 1 dataset.isnull().sum()
```

```
Out[70]: 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 [71]: 1 # Handling null values
2 # Null values are present in 3 columns - Age, Cabin and Embarked
3
4 # The 'Age' column contains some missing values, replacing those with median
5 dataset['Age'] = dataset['Age'].replace(np.NaN,dataset['Age'].median())
6
7 # As there are too many null values in the 'Cabin' column, removing the column
8 dataset = dataset.drop(['Cabin'], axis=1)
9
10 # As there are very few null values in 'Embarked' column, removing the column
11 dataset.dropna(subset=['Embarked'],how='any',inplace=True)
```

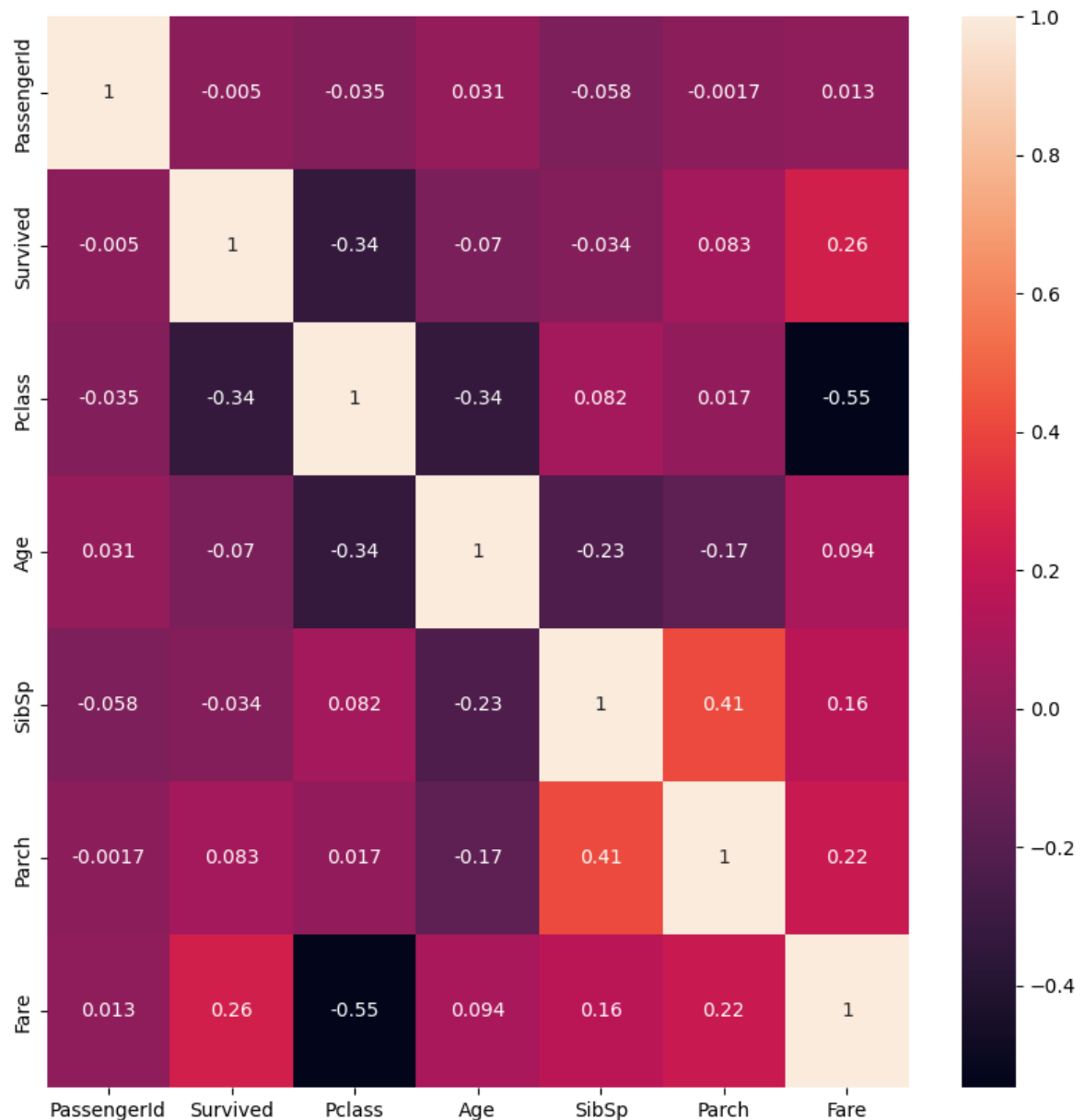
```
In [72]: 1 dataset.isnull().sum()
```

```
Out[72]: PassengerId      0
Survived      0
Pclass        0
Name          0
Sex           0
Age           0
SibSp         0
Parch         0
Ticket        0
Fare          0
Embarked      0
dtype: int64
```

```
In [73]: 1 # Data Visualization.
2 # Heatmap
3 corr=dataset.corr()
4 plt.subplots(figsize=(10,10))
5 sns.heatmap(corr,annot=True)
```

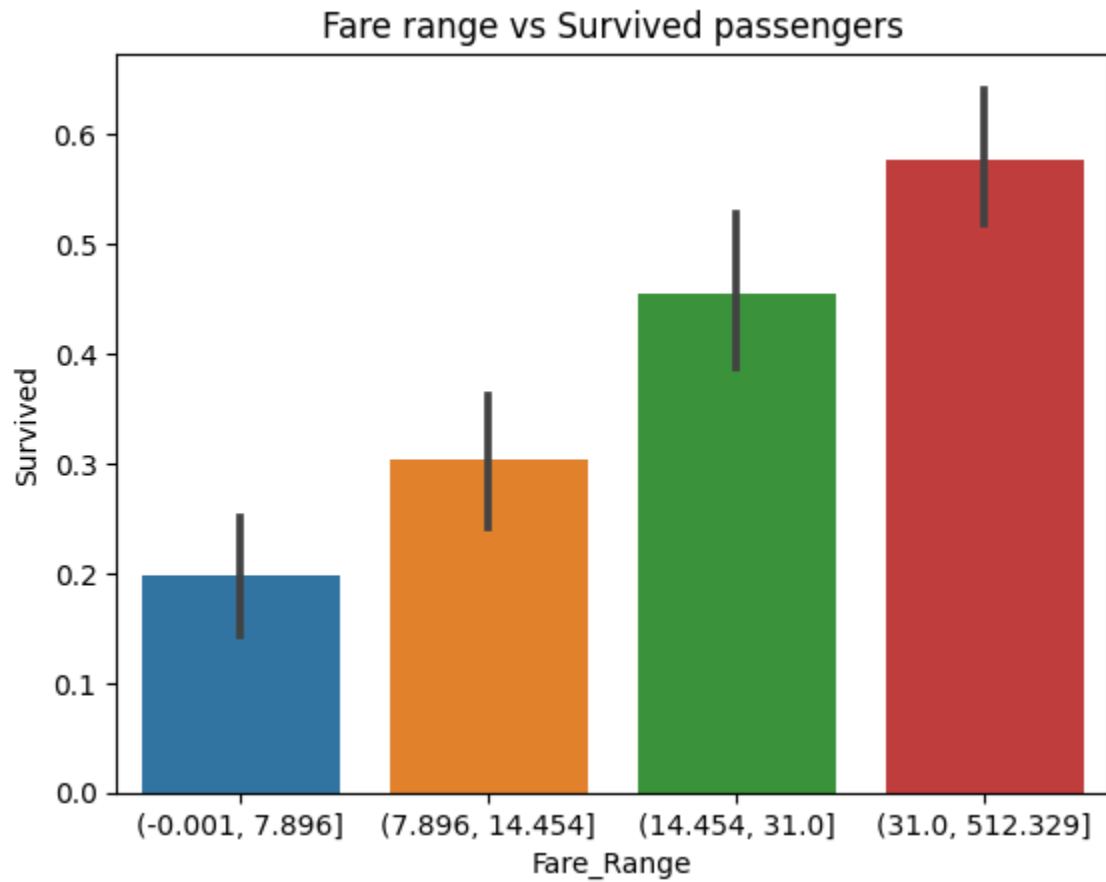
<ipython-input-73-af9811d18692>:3: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
corr=dataset.corr()

Out[73]: <Axes: >

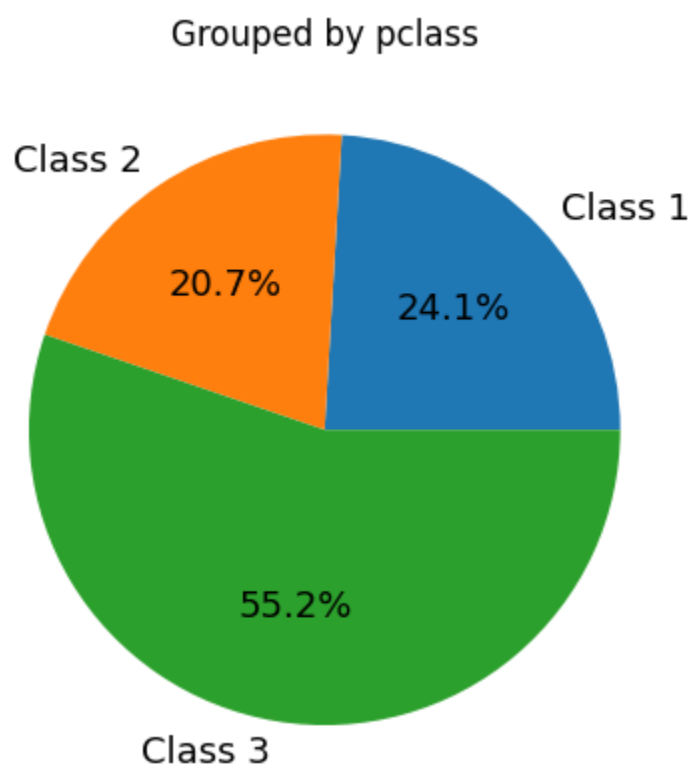


```
In [74]: 1 #Barplot
2 dataset['Fare_Range'] = pd.qcut(dataset['Fare'], 4)
3 plt.title('Fare range vs Survived passengers')
4 sns.barplot(x='Fare_Range', y='Survived', data=dataset)
```

```
Out[74]: <Axes: title={'center': 'Fare range vs Survived passengers'}, xlabel='Fare_Range', ylabel='Survived'>
```



```
In [75]: 1 #Piechart
2 pclass_count = dataset.groupby('Pclass')['Pclass'].count()
3 plt.title('Grouped by pclass')
4 plt.pie(pclass_count.values, labels=['Class 1', 'Class 2', 'Class 3'], autopct='%1.1f%%')
5 plt.show()
```



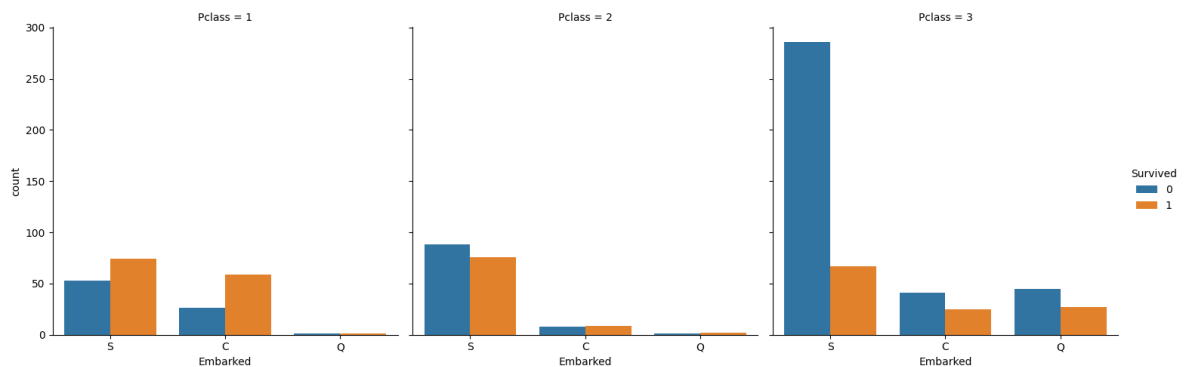
```
In [76]: 1 # Violinplot
2 plt.title('Age vs Survived')
3 sns.violinplot(x = "Sex", y = "Age", hue = "Survived", data = dataset, split :
```

Out[76]: <Axes: title={'center': 'Age vs Survived'}, xlabel='Sex', ylabel='Age'>



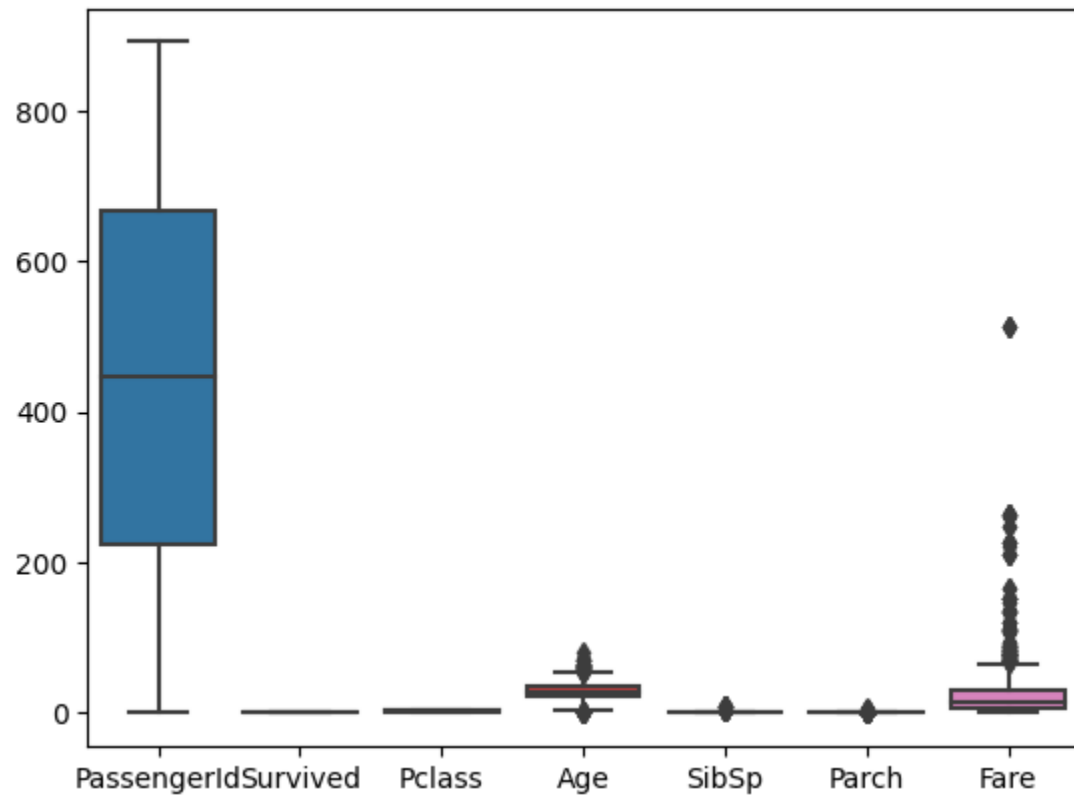
```
In [77]: 1 # Countplot
2 sns.catplot(x = 'Embarked', hue = 'Survived', kind = 'count', col = 'Pclass',
```

Out[77]: <seaborn.axisgrid.FacetGrid at 0x79e5639a7490>



```
In [78]: 1 # Outlier Detection  
2 sns.boxplot(dataset)
```

Out[78]: <Axes: >

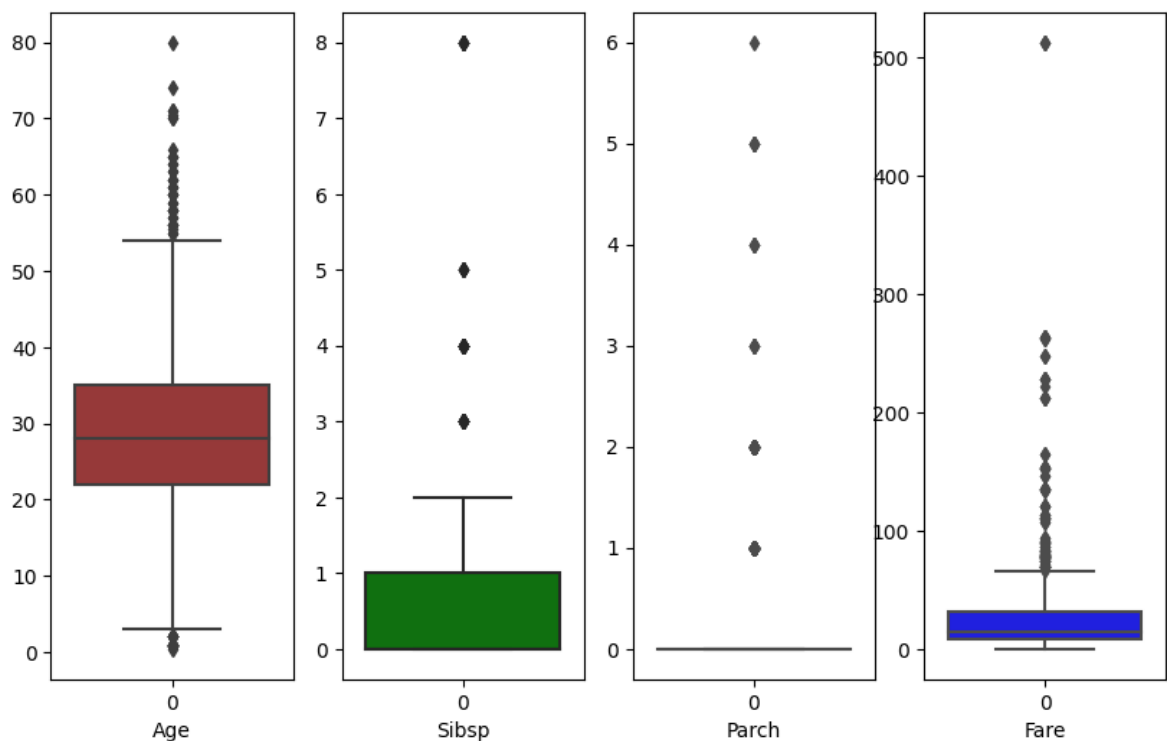


```

In [79]: 1 # Outliers are present in Age, SibSp, Parch, Fare classes
2 fig, ax = plt.subplots(1, 4, figsize=(10, 6))
3
4 sns.boxplot(data=dataset['Age'], ax=ax[0], color='brown')
5 ax[0].set_xlabel('Age')
6
7 sns.boxplot(data=dataset['SibSp'], ax=ax[1], color='green')
8 ax[1].set_xlabel('Sibsp')
9
10 sns.boxplot(data=dataset['Parch'], ax=ax[2], color='yellow')
11 ax[2].set_xlabel('Parch')
12
13 sns.boxplot(data=dataset['Fare'], ax=ax[3], color='blue')
14 ax[3].set_xlabel('Fare')

```

Out[79]: Text(0.5, 0, 'Fare')



Splitting Dependent and Independent variables

```

In [80]: 1 # Independent variables - Name, SibSp, Parch, Ticket
2 x = dataset.drop(['Name', 'SibSp', 'Parch', 'Ticket'], axis=1)
3 y = dataset['Survived']

```

```

In [81]: 1 # Perform Encoding
2 # Performing Label encoding for Sex and Embarked columns
3 encoder = LabelEncoder()
4 x['Sex'] = encoder.fit_transform(x['Sex'])
5 x['Embarked'] = encoder.fit_transform(x['Embarked'])

```



```
In [82]: 1 x.head() # Values in Sex and Embarked columns into numerical values
```

```
Out[82]:
```

	PassengerId	Survived	Pclass	Sex	Age	Fare	Embarked	Fare_Range
0	1	0	3	1	22.0	7.2500	2	(-0.001, 7.896]
1	2	1	1	0	38.0	71.2833	0	(31.0, 512.329]
2	3	1	3	0	26.0	7.9250	2	(7.896, 14.454]
3	4	1	1	0	35.0	53.1000	2	(31.0, 512.329]
4	5	0	3	1	35.0	8.0500	2	(7.896, 14.454]

```
In [87]: 1 x=x.drop(['Fare_Range'],axis=1)
2 # Feature Scaling
3 scaler = StandardScaler()
4 x_scaled = scaler.fit_transform(x)
```

```
In [88]: 1 x_scaled
```

```
Out[88]: array([[ -1.73250451, -0.78696114,  0.82520863, ..., -0.56367407,
        -0.50023975,  0.58683958],
        [ -1.72861124,  1.27071078, -1.57221121, ...,  0.66921696,
         0.78894661, -1.93955453],
        [ -1.72471797,  1.27071078,  0.82520863, ..., -0.25545131,
        -0.48664993,  0.58683958],
        ...,
        [  1.72471797, -0.78696114,  0.82520863, ..., -0.10133993,
        -0.17408416,  0.58683958],
        [  1.72861124,  1.27071078, -1.57221121, ..., -0.25545131,
        -0.0422126 , -1.93955453],
        [  1.73250451, -0.78696114,  0.82520863, ...,  0.20688282,
        -0.49017322, -0.67635748]])
```

```
In [90]: 1 # Splitting Data into Train and Test
2 x_train,x_test,y_train,y_test = train_test_split(x_scaled,y,test_size=0.3)
```

```
In [91]: 1 print("Shape of x_train:",x_train.shape)
2 print("Shape of x_test:",x_test.shape)
3 print("Shape of y_train:",y_train.shape)
4 print("Shape of y_test:",y_test.shape)
```

Shape of x_train: (622, 7)

Shape of x_test: (267, 7)

Shape of y_train: (622,)

Shape of y_test: (267,)