

# PRODIGY INFOTECH

## TASK-2

Perform data cleaning and exploratory data analysis(EDA) on a dataset of your choice, such as the Titanic dataset from kaggle. Explore the relationships between variables and identify patterns and trends in the data.

```
In [44]: #Import required libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [45]: #Load the dataset of titanic
Data=pd.read_csv("C:/Users/steph/OneDrive/Documents/titanic.csv")
print(Data)
```

```
PassengerId  Survived  Pclass \
0            1         0      3
1            2         1      1
2            3         1      3
3            4         1      1
4            5         0      3
..          ...
886          887        0      2
887          888        1      1
888          889        0      3
889          890        1      1
890          891        0      3

Name      Sex   Age  SibSp \
0    Braund, Mr. Owen Harris   male  22.0     1
1  Cumings, Mrs. John Bradley (Florence Briggs Th... female  38.0     1
2           Heikkinen, Miss. Laina  female  26.0     0
3    Futrelle, Mrs. Jacques Heath (Lily May Peel) female  35.0     1
4           Allen, Mr. William Henry   male  35.0     0
..          ...
886          Montvila, Rev. Juozas   male  27.0     0
887          Graham, Miss. Margaret Edith female  19.0     0
888          Johnston, Miss. Catherine Helen "Carrie" female   NaN     1
889          Behr, Mr. Karl Howell   male  26.0     0
890          Dooley, Mr. Patrick   male  32.0     0

Parch      Ticket     Fare Cabin Embarked
0          0       A/5 21171  7.2500   NaN      S
1          0       PC 17599  71.2833  C85      C
2          0  STON/O2. 3101282  7.9250   NaN      S
3          0       113803  53.1000  C123      S
4          0       373450  8.0500   NaN      S
..          ...
886          0       211536 13.0000   NaN      S
887          0       112053 30.0000  B42      S
888          2      W./C. 6607 23.4500   NaN      S
889          0       111369 30.0000  C148      C
890          0       370376  7.7500   NaN      Q

[891 rows x 12 columns]
```

```
In [46]: # Create a dataframe for the data
df=pd.DataFrame(Data)
df
```

Out[46]:		PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cat
	<b>0</b>	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	N
	<b>1</b>	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	female	38.0	1	0	PC 17599	71.2833	C
	<b>2</b>	3	1	3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	26.0	0	0	STON/O2. 3101282	7.9250	N
	<b>3</b>	4	1	1	Allen, Mr. William Henry	male	35.0	1	0	113803	53.1000	C1
	<b>4</b>	5	0	3	Montvila, Rev. Juozas	male	35.0	0	0	373450	8.0500	N
	...	...	...	...	...	...	...	...	...	...	...	...
	<b>886</b>	887	0	2	Graham, Miss. Margaret Edith	female	19.0	0	0	211536	13.0000	N
	<b>887</b>	888	1	1	Johnston, Miss. Catherine Helen "Carrie"	female	Nan	1	2	W./C. 6607	23.4500	N
	<b>888</b>	889	0	3	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C1
	<b>889</b>	890	1	1	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	N
	<b>890</b>	891	0	3								

891 rows × 12 columns

In [47]: #First 5 rows of the dataset  
df.head()

Out[47]:

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

In [48]: #Last 5 rows of the dataset  
df.tail()

Out[48]:

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

In [49]: #Check for all the columns of the dataset  
df.columns

```
Out[49]: Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',  
               'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],  
               dtype='object')
```

```
In [50]: #Check for the number of rows and columns of the dataset  
df.shape
```

```
Out[50]: (891, 12)
```

```
In [51]: #Check for the information ,i.e, dtype and null value for each column  
df.info
```

```
Out[51]: <bound method DataFrame.info of  
          PassengerId  Survived  Pclass \\\n          0            1        0      3  
          1            2        1      1  
          2            3        1      3  
          3            4        1      1  
          4            5        0      3  
          ..          ...      ...  ...  
          886           887      0      2  
          887           888      1      1  
          888           889      0      3  
          889           890      1      1  
          890           891      0      3  
  
                           Name    Sex   Age  SibSp \\\n          0      Braund, Mr. Owen Harris    male  22.0     1  
          1  Cumings, Mrs. John Bradley (Florence Briggs Th... female  38.0     1  
          2      Heikkinen, Miss. Laina    female  26.0     0  
          3      Futrelle, Mrs. Jacques Heath (Lily May Peel) female  35.0     1  
          4      Allen, Mr. William Henry    male  35.0     0  
          ..          ...      ...  ...  
          886      Montvila, Rev. Juozas    male  27.0     0  
          887      Graham, Miss. Margaret Edith female  19.0     0  
          888      Johnston, Miss. Catherine Helen "Carrie" female  NaN     1  
          889      Behr, Mr. Karl Howell    male  26.0     0  
          890      Dooley, Mr. Patrick    male  32.0     0  
  
          Parch      Ticket      Fare Cabin Embarked  
          0          0       A/5 21171    7.2500   NaN      S  
          1          0       PC 17599   71.2833   C85      C  
          2          0      STON/O2. 3101282   7.9250   NaN      S  
          3          0       113803  53.1000  C123      S  
          4          0       373450  8.0500   NaN      S  
          ..          ...      ...  ...  ...  
          886          0       211536  13.0000   NaN      S  
          887          0       112053  30.0000  B42      S  
          888          2      W./C. 6607  23.4500   NaN      S  
          889          0       111369  30.0000  C148      C  
          890          0       370376  7.7500   NaN      Q  
  
[891 rows x 12 columns]>
```

```
In [52]: #Check for Statistical Analysis  
df.describe
```

```
Out[52]: <bound method NDFrame.describe of
   PassengerId  Survived  Pclass \
0              1         0      3
1              2         1      1
2              3         1      3
3              4         1      1
4              5         0      3
..             ...       ...
886            887       0      2
887            888       1      1
888            889       0      3
889            890       1      1
890            891       0      3

                                                Name     Sex   Age  SibSp \
0           Braund, Mr. Owen Harris    male  22.0     1
1  Cumings, Mrs. John Bradley (Florence Briggs Th... female  38.0     1
2           Heikkinen, Miss. Laina  female  26.0     0
3  Futrelle, Mrs. Jacques Heath (Lily May Peel) female  35.0     1
4           Allen, Mr. William Henry    male  35.0     0
..             ...
886          Montvila, Rev. Juozas    male  27.0     0
887          Graham, Miss. Margaret Edith female  19.0     0
888          Johnston, Miss. Catherine Helen "Carrie" female  NaN     1
889            Behr, Mr. Karl Howell    male  26.0     0
890          Dooley, Mr. Patrick    male  32.0     0

   Parch      Ticket     Fare Cabin Embarked
0      0        A/5 21171  7.2500   NaN      S
1      0         PC 17599  71.2833  C85      C
2      0  STON/O2. 3101282  7.9250   NaN      S
3      0        113803  53.1000  C123      S
4      0        373450  8.0500   NaN      S
..     ...
886      0        211536 13.0000   NaN      S
887      0        112053 30.0000  B42      S
888      2        W./C. 6607 23.4500   NaN      S
889      0        111369 30.0000  C148      C
890      0        370376  7.7500   NaN      Q
```

[891 rows x 12 columns]>

```
In [53]: #Check for the null values
print(df.isnull().sum())
```

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

# DATA CLEANING

```
In [54]: # Impute missing values with mean for numerical variables and mode for categorical
df['Age'].fillna(df['Age'].mean(), inplace=True)
df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)

#Drop irrelevant columns
df.drop(['PassengerId', 'Name', 'Ticket', 'Cabin'], axis=1, inplace=True)

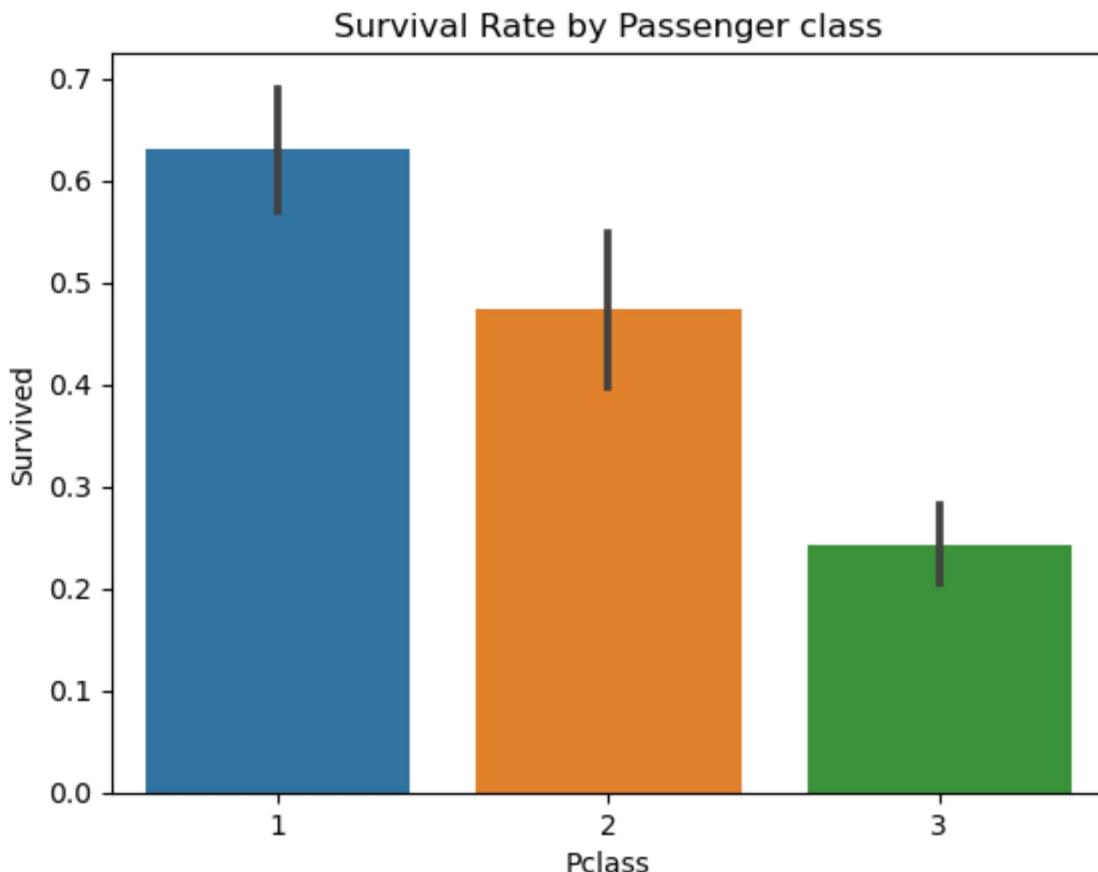
#Check for the missing values
print(df.isnull().sum())
```

```
Survived      0
Pclass        0
Sex           0
Age           0
SibSp         0
Parch         0
Fare          0
Embarked      0
dtype: int64
```

# EXPLORATORY DATA ANALYSIS

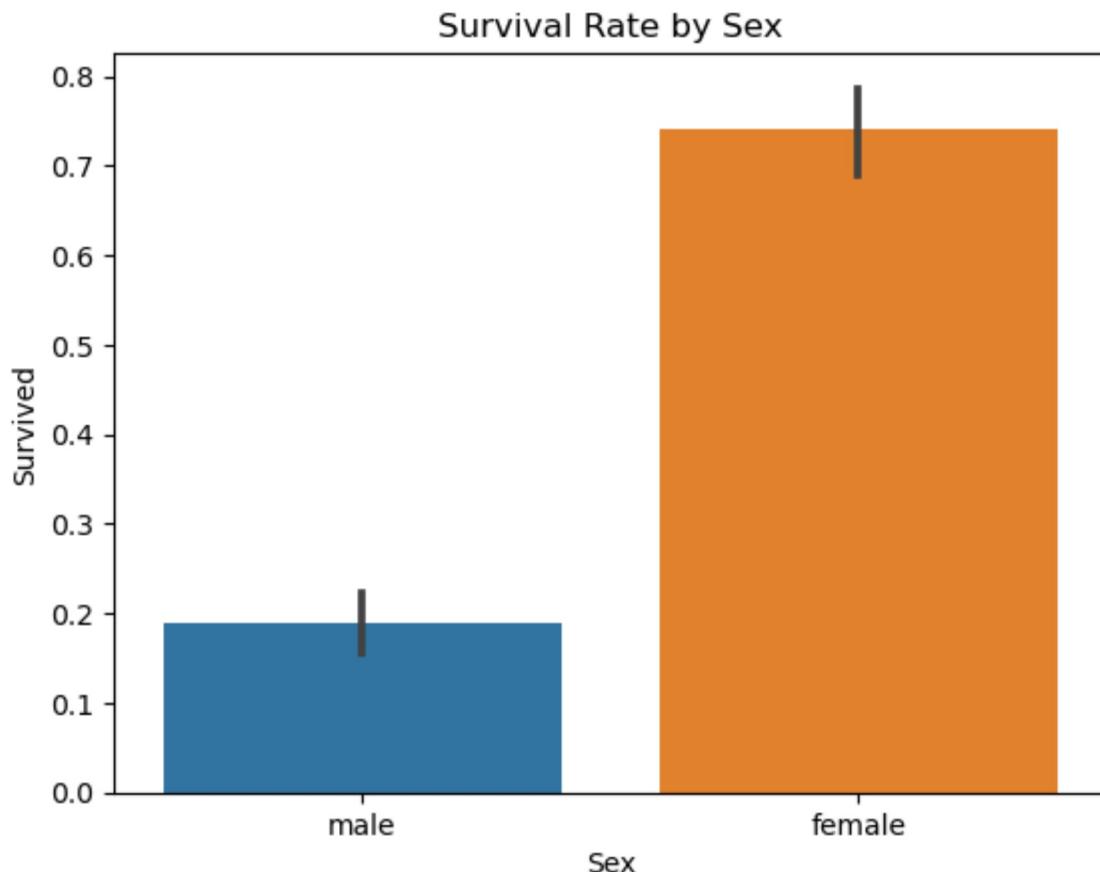
## 1. Survival rate by Passenger class

```
In [65]: sns.barplot(x='Pclass', y='Survived', data=df)
plt.title('Survival Rate by Passenger class')
plt.show()
```



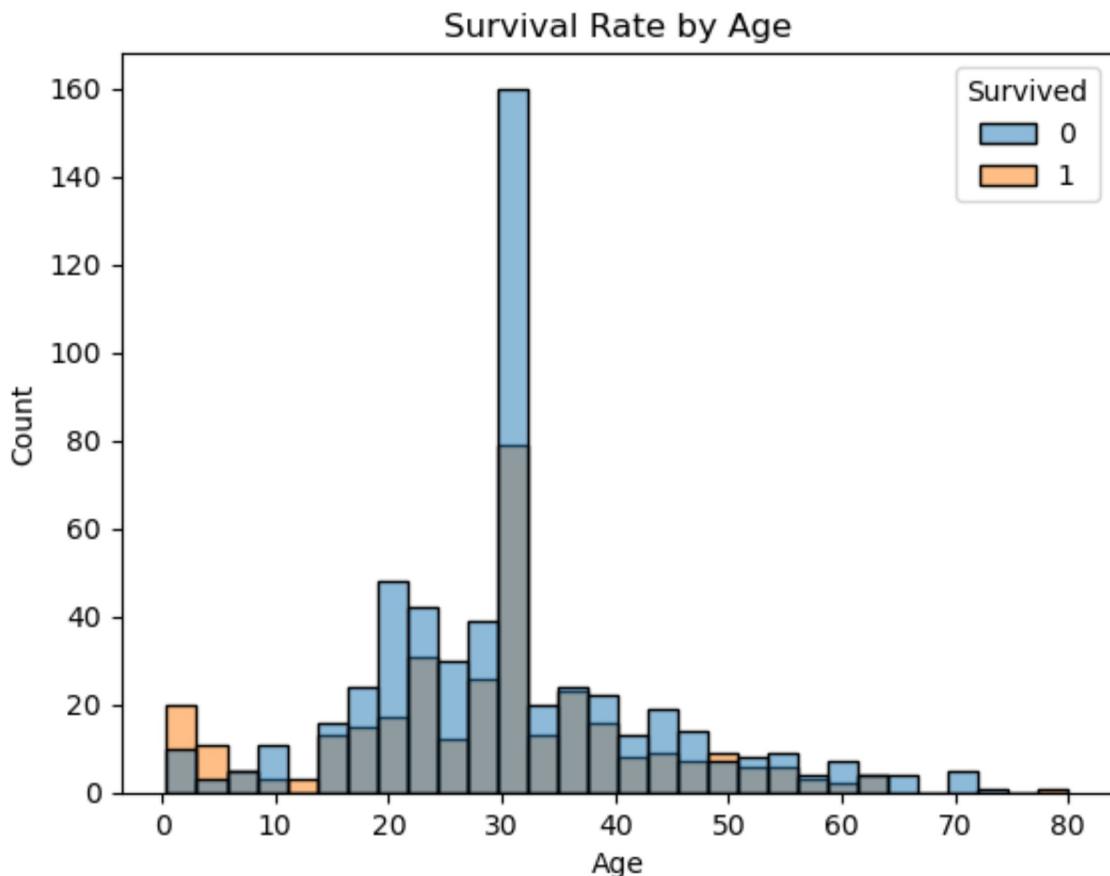
## 2. Survival rate by Sex

```
In [57]: sns.barplot(x='Sex', y='Survived', data=df)
plt.title('Survival Rate by Sex')
plt.show()
```



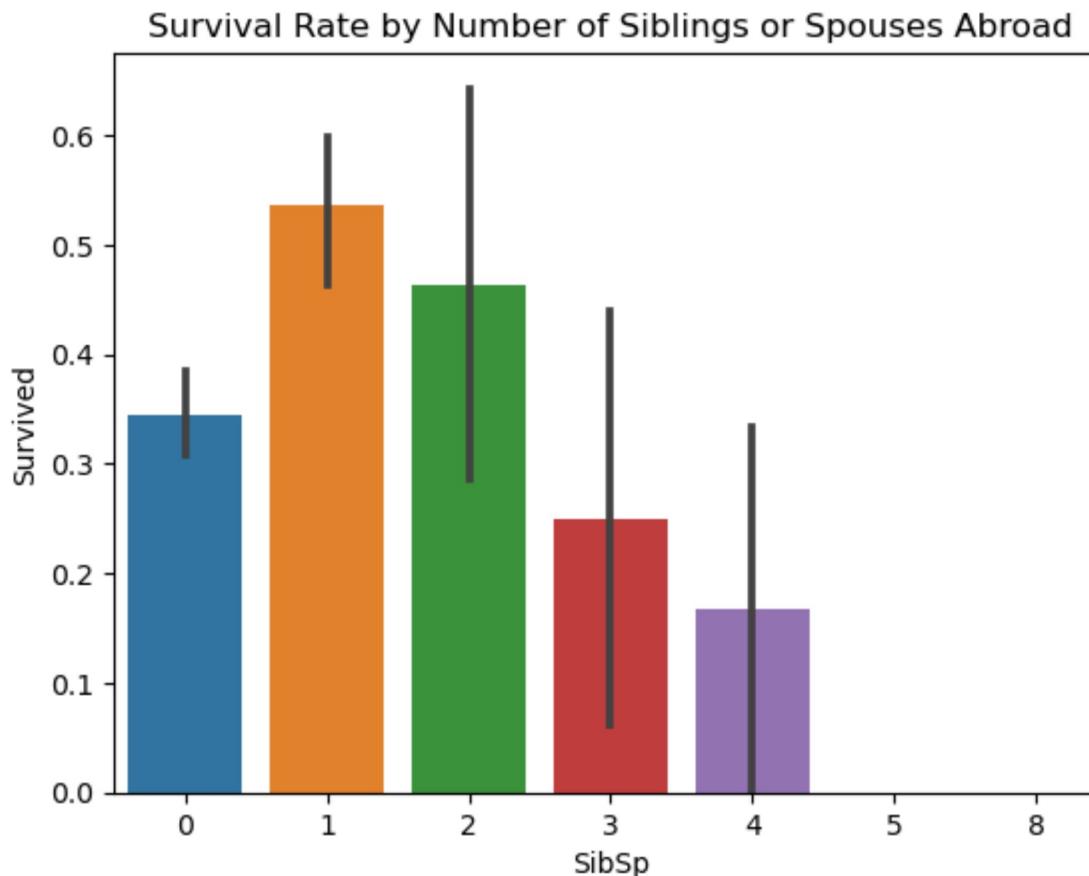
### 3. Survival rate by Age

```
In [74]: sns.histplot(x='Age', hue='Survived', data=df)
plt.title('Survival Rate by Age')
plt.show()
```



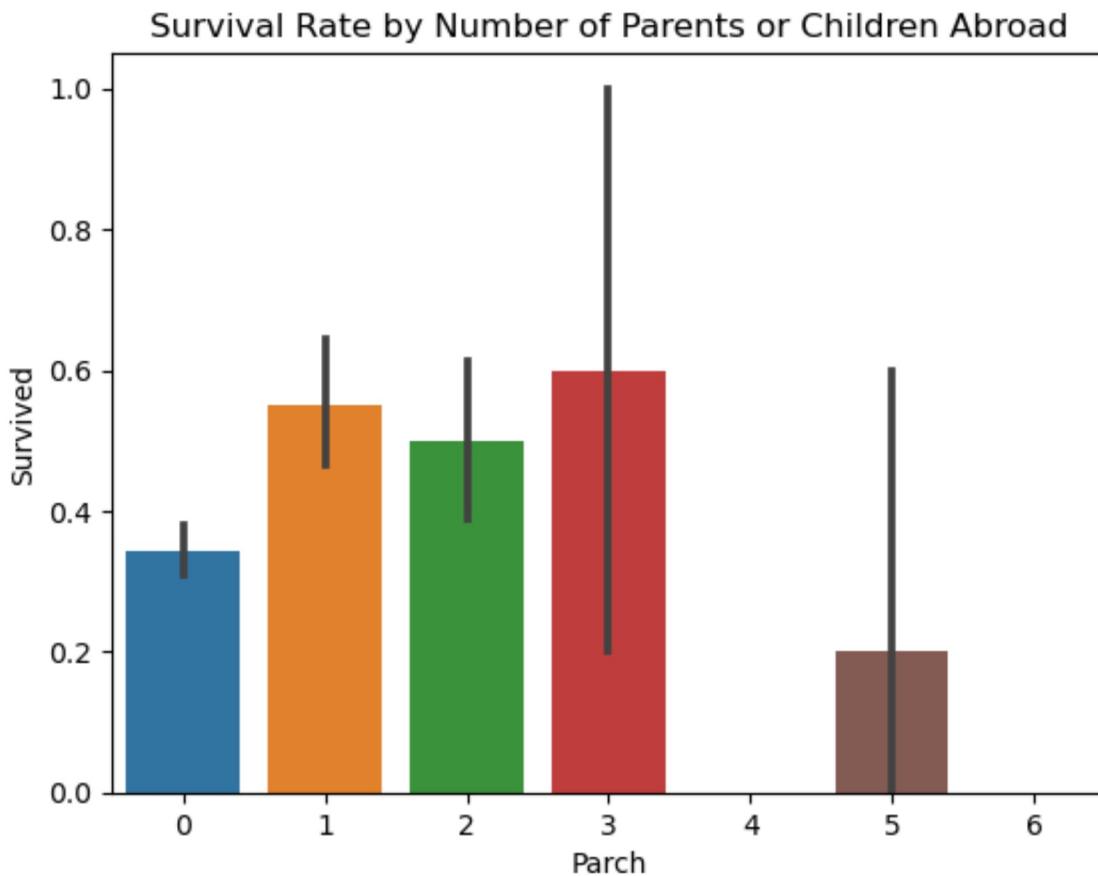
#### 4. Survival rate by Number of Siblings or Spouses Abroad

```
In [67]: sns.barplot(x='SibSp', y='Survived', data=df)
plt.title('Survival Rate by Number of Siblings or Spouses Abroad')
plt.show()
```



## 5. Survival rate by Number of Parents or Children Abroad

```
In [75]: sns.barplot(x='Parch', y='Survived', data=df)
plt.title('Survival Rate by Number of Parents or Children Abroad')
plt.show()
```



## 6. Survival rate by Embarked

```
In [64]: sns.barplot(x='Embarked', y='Survived', data=df)
plt.title('Survival Rate by Embarked')
plt.show()
```

