

ml-ass-2-2203a51430

March 11, 2024

```
[3]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

customer= pd.read_csv('/titanic.zip')
print(customer.describe())
```

	PassengerId	Survived	Pclass	Age	SibSp \
count	891.000000	891.000000	891.000000	714.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008
std	257.353842	0.486592	0.836071	14.526497	1.102743
min	1.000000	0.000000	1.000000	0.420000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000
50%	446.000000	0.000000	3.000000	28.000000	0.000000
75%	668.500000	1.000000	3.000000	38.000000	1.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000

	Parch	Fare
count	891.000000	891.000000
mean	0.381594	32.204208
std	0.806057	49.693429
min	0.000000	0.000000
25%	0.000000	7.910400
50%	0.000000	14.454200
75%	0.000000	31.000000
max	6.000000	512.329200

```
[4]: customer.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

```

```
[6]: print(customer.dtypes)
```

```

PassengerId    int64
Survived        int64
Pclass          int64
Name            object
Sex             object
Age             float64
SibSp           int64
Parch           int64
Ticket          object
Fare            float64
Cabin           object
Embarked        object
dtype: object

```

```
[9]: summary_stats = customer.describe(percentiles=[0.25, 0.5, 0.75, 0.9])
print("\nSummary of the dataset:")
print(summary_stats)
```

Summary of the dataset:

	PassengerId	Survived	Pclass	Age	SibSp \
count	891.000000	891.000000	891.000000	714.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008
std	257.353842	0.486592	0.836071	14.526497	1.102743
min	1.000000	0.000000	1.000000	0.420000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000
50%	446.000000	0.000000	3.000000	28.000000	0.000000
75%	668.500000	1.000000	3.000000	38.000000	1.000000
90%	802.000000	1.000000	3.000000	50.000000	1.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000

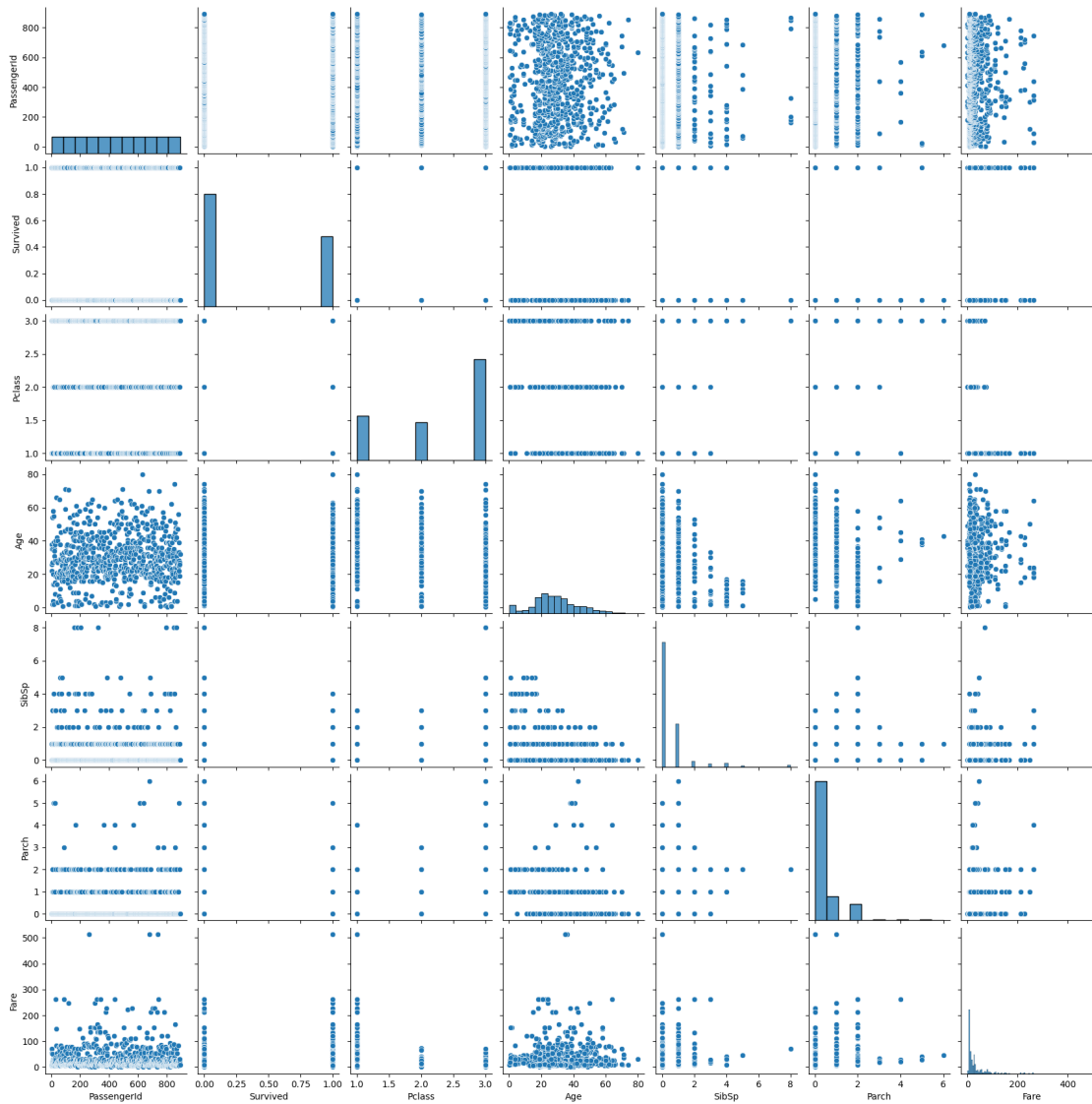
	Parch	Fare
count	891.000000	891.000000
mean	0.381594	32.204208
std	0.806057	49.693429

min	0.000000	0.000000
25%	0.000000	7.910400
50%	0.000000	14.454200
75%	0.000000	31.000000
90%	2.000000	77.958300
max	6.000000	512.329200

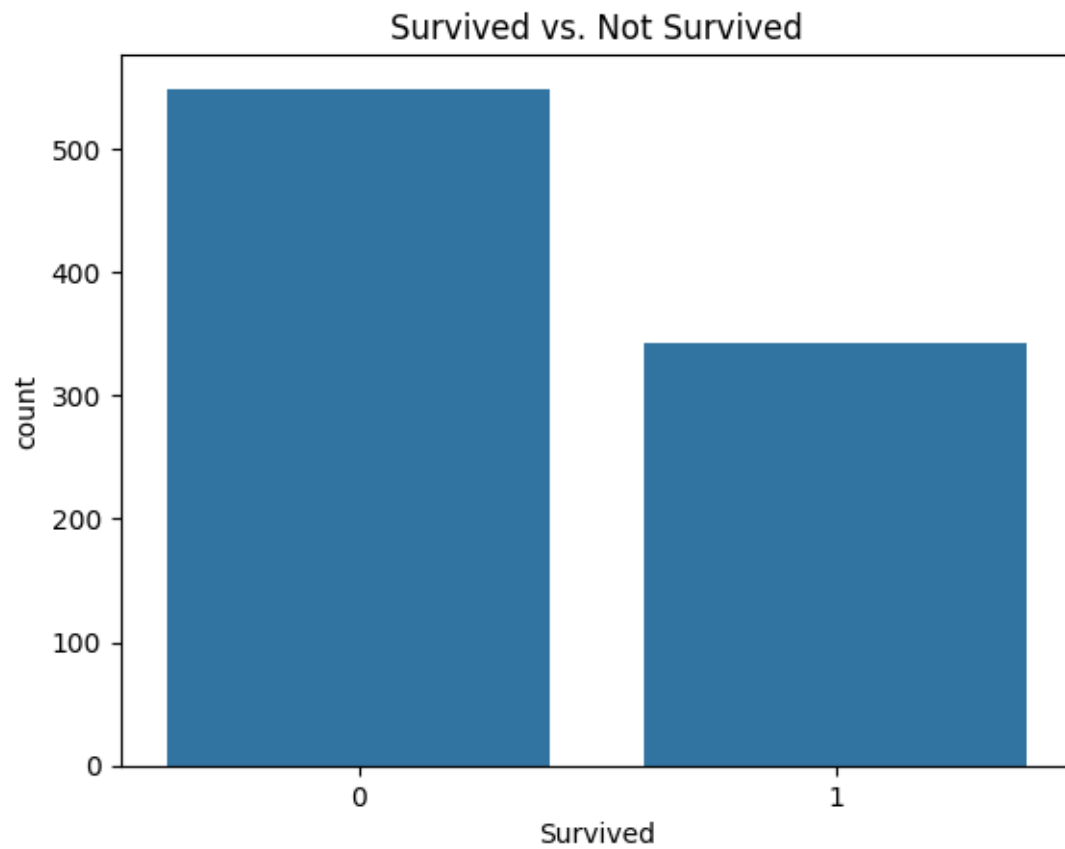
```
[10]: column=customer.columns.tolist()
      print(column)
```

```
['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp', 'Parch',
'Ticket', 'Fare', 'Cabin', 'Embarked']
```

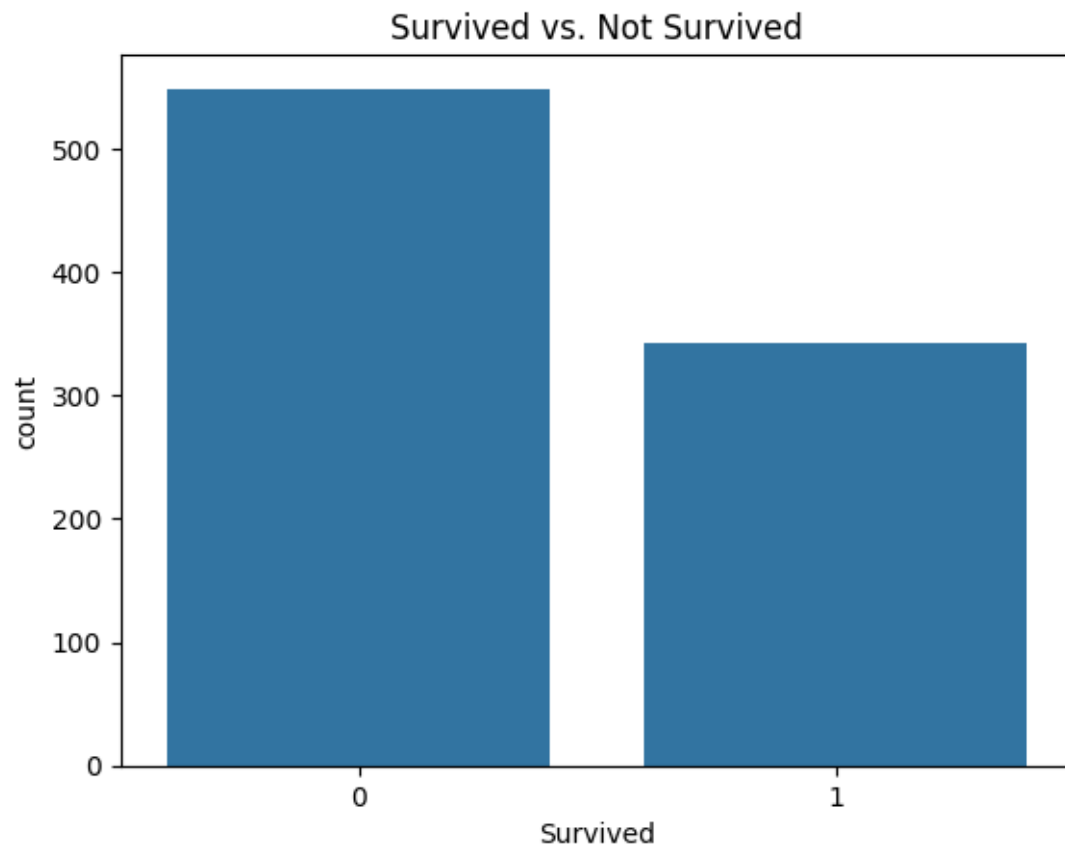
```
[11]: numeric_features = customer.select_dtypes(include=['int64', 'float64']).
      ↪columns
      sns.pairplot(customer[numeric_features])
      plt.show()
```



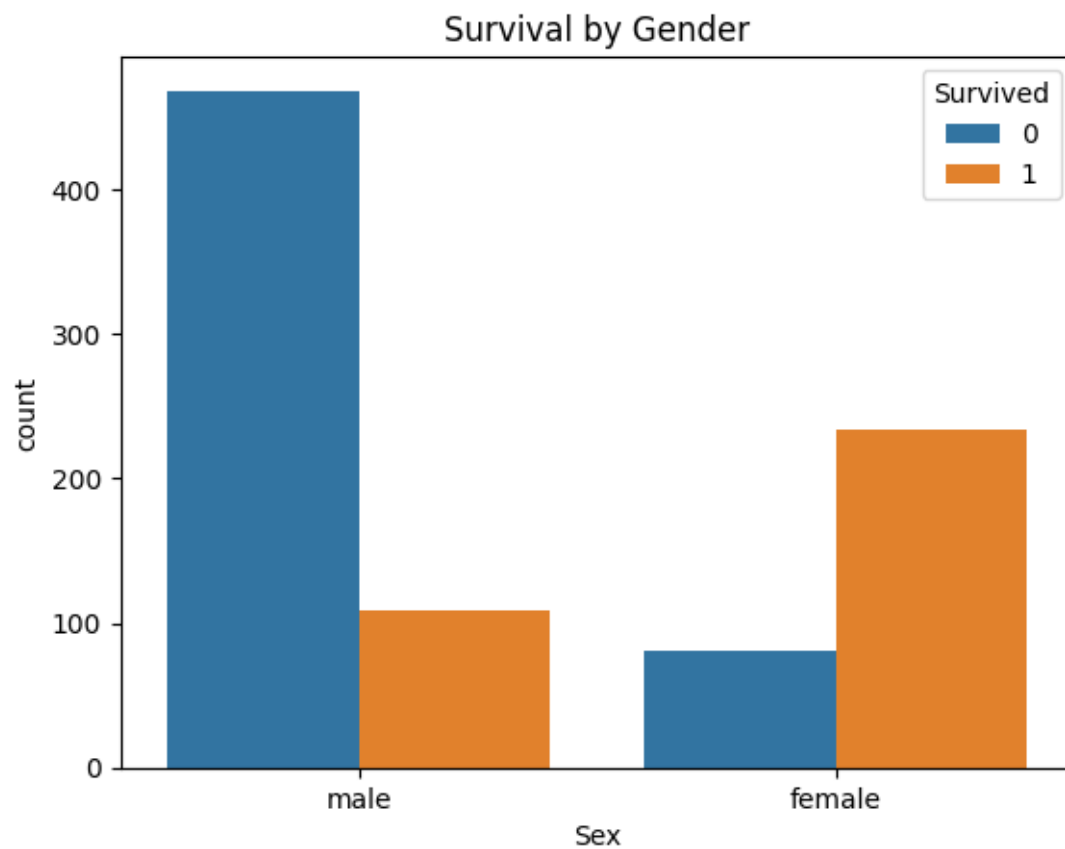
```
[15]: sns.countplot(x='Survived', data=customer)
plt.title('Survived vs. Not Survived')
plt.show()
```



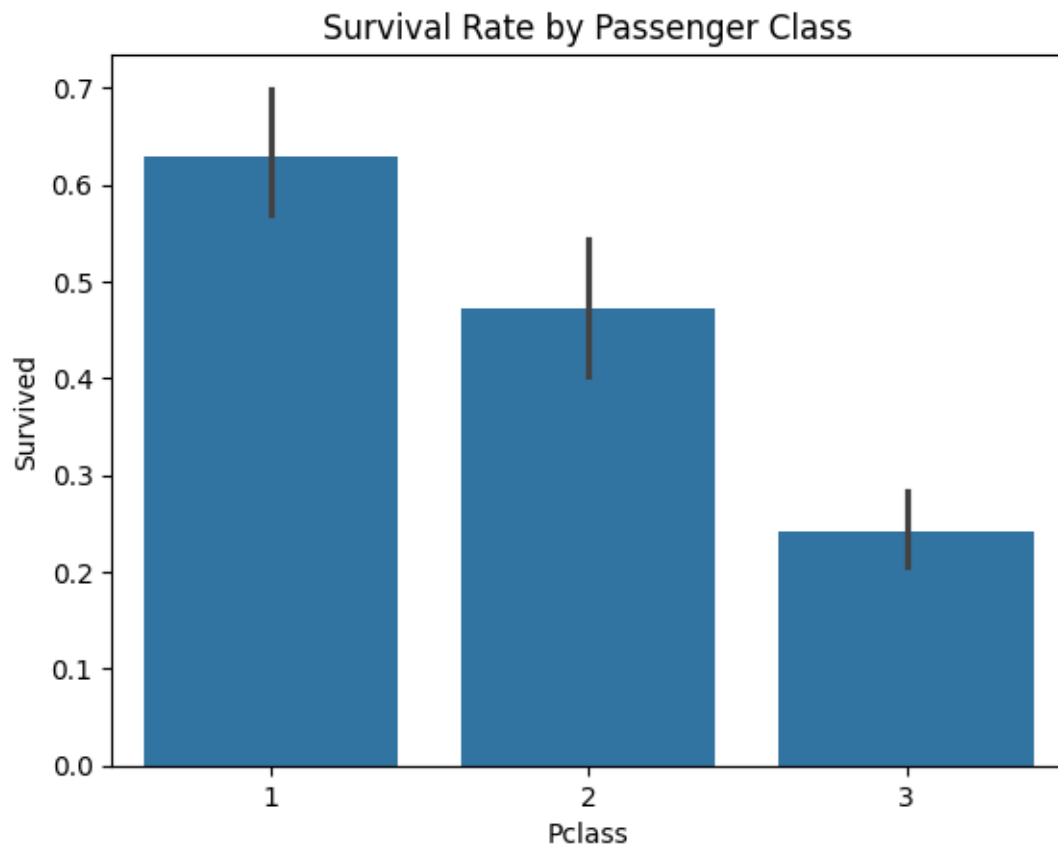
```
[16]: sns.countplot(x='Survived', data=customer)
plt.title('Survived vs. Not Survived')
plt.show()
```



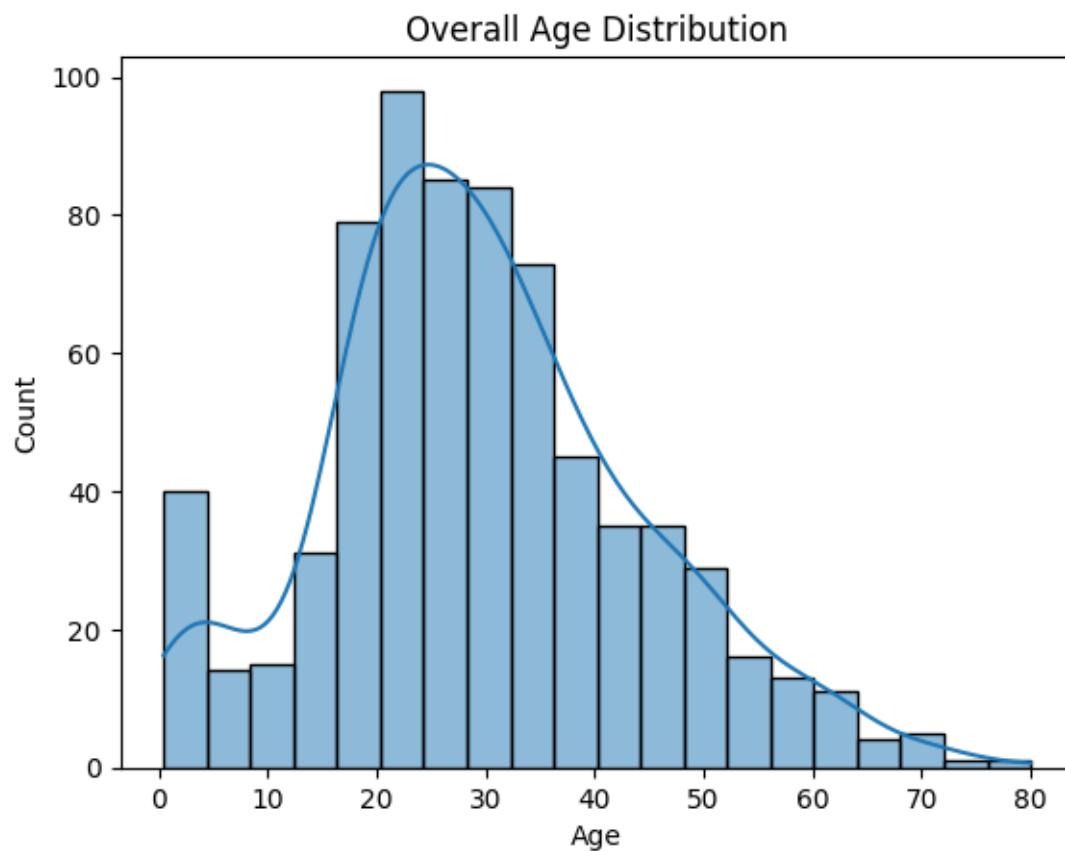
```
[17]: sns.countplot(x='Sex', data=customer, hue='Survived')  
plt.title('Survival by Gender')  
plt.show()
```



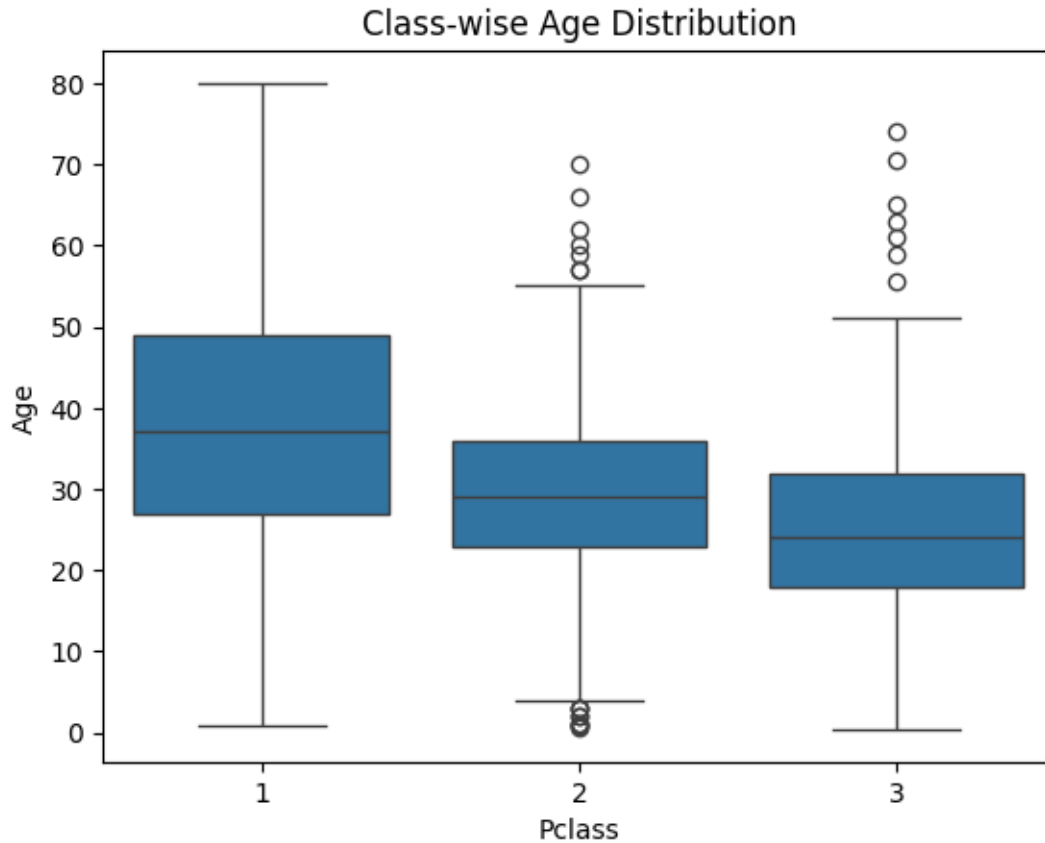
```
[18]: sns.barplot(x='Pclass', y='Survived', data=customer)
plt.title('Survival Rate by Passenger Class')
plt.show()
```



```
[19]: sns.histplot(x='Age', data=customer, kde=True)
plt.title('Overall Age Distribution')
plt.show()
```

```
[20]: sns.boxplot(x='Pclass', y='Age', data=customer)
plt.title('Class-wise Age Distribution')
plt.show()
```



```
[22]: customer['Age'].fillna(customer['Age'].mean(), inplace=True)
# Recode categorical features to a class
customer['Sex'] = customer['Sex'].map({'male': 0, 'female': 1})
customer= pd.get_dummies(customer, columns=['Embarked'], drop_first=True)
# Display the modified dataframe
print(customer.head())
```

	PassengerId	Survived	Pclass	\
0	1	0	3	
1	2	1	1	
2	3	1	3	
3	4	1	1	
4	5	0	3	

	Name	Sex	Age	SibSp	Parch	\
0	Braund, Mr. Owen Harris	0	22.0	1	0	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	1	38.0	1	0	
2	Heikkinen, Miss. Laina	1	26.0	0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	1	35.0	1	0	
4	Allen, Mr. William Henry	0	35.0	0	0	

	Ticket	Fare	Cabin	Embarked_Q	Embarked_S
0	A/5 21171	7.2500	NaN	0	1
1	PC 17599	71.2833	C85	0	0
2	STON/O2. 3101282	7.9250	NaN	0	1
3	113803	53.1000	C123	0	1
4	373450	8.0500	NaN	0	1

```
[23]: customer.drop(['Cabin', 'Ticket'], axis=1, inplace=True)
```

```
[27]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import f1_score
import matplotlib.pyplot as plt
# Assuming 'df' is your DataFrame with the provided data
# Step 1: Split the data into X (features) and Y (target)
X = customer[['Pclass', 'Age', 'SibSp', 'Parch', 'Fare']]
Y = customer['Survived']

# Split the data into training and testing sets X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42)
```

```
[31]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import f1_score
import matplotlib.pyplot as plt

X = customer[['Pclass', 'Age', 'SibSp', 'Parch', 'Fare']]
y = customer['Survived']

X = X.fillna(X.mean())

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
model = LogisticRegression()
penalty_values = [0.1, 0.5, 1, 2, 5, 10]
f1_scores = []
penalties = []
for penalty in penalty_values:

    model.set_params(C=1/penalty)
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    f1 = f1_score(y_test, y_pred)
    f1_scores.append(f1)
```

```
penalties.append(penalty)
plt.scatter(penalties, f1_scores, color='blue')
plt.title('F1 Score as a Function of Penalty')
plt.xlabel('Penalty')
plt.ylabel('F1 Score')
plt.xscale('log')
```

