```
1 from google.colab import files
2 import pandas as pd
3
4 # Upload the file
5 uploaded = files.upload()
```

Choose Files  Titanic-Dataset.csv
  • Titanic-Dataset.csv(text/csv) - 61194 bytes, last modified: 12/24/2021 - 100% done
  Saving Titanic-Dataset.csv to Titanic-Dataset.csv

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score, classification_report
```

```
1 titanic_df = pd.read_csv('Titanic-Dataset.csv')
2
3 # Display the first few rows of the dataset
4 titanic_df.head()
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7. |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs | female | 38.0 | 1 | 0 | PC 17599 | 71. |

Next steps:   Generate code with  titanic_df    ⊙ View recommended plots    New interactive sheet

```
1 titanic_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
```

```
1 # Drop columns that won't be used in the model
2 titanic_df.drop(['PassengerId', 'Name', 'Ticket', 'Cabin'], axis=1, inplace=True)
3
4 # Handle missing values
5 titanic_df['Age'].fillna(titanic_df['Age'].median(), inplace=True)
6 titanic_df['Embarked'].fillna(titanic_df['Embarked'].mode()[0], inplace=True)
7
8 # Convert categorical variables to numeric
9 titanic_df = pd.get_dummies(titanic_df, columns=['Sex', 'Embarked'], drop_first=Tru
10
11 # Display the processed dataframe
12 print(titanic_df.head())
```

```
   Survived  Pclass   Age  SibSp  Parch     Fare  Sex_male  Embarked_Q  \
0         0       3  22.0      1      0   7.2500      True       False
1         1       1  38.0      1      0  71.2833     False       False
2         1       3  26.0      0      0   7.9250     False       False
3         1       1  35.0      1      0  53.1000     False       False
4         0       3  35.0      0      0   8.0500      True       False

   Embarked_S
0        True
1       False
2        True
3        True
4        True
```

```
1 # Define features (X) and target (y)
2 X = titanic_df.drop('Survived', axis=1)
3 y = titanic_df['Survived']
4
5 # Split the data into training and testing sets
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_sta
7
8 # Standardize the features
9 scaler = StandardScaler()
10 X_train = scaler.fit_transform(X_train)
11 X_test = scaler.transform(X_test)
```

```
1 # Initialize the model
2 model = LogisticRegression()
3
4 # Train the model
5 model.fit(X_train, y_train)
```

```
▾ LogisticRegression
LogisticRegression()
```

```
1 # Predict on the test set
2 y_pred = model.predict(X_test)
3
4 # Evaluate the model
5 accuracy = accuracy_score(y_test, y_pred)
6 report = classification_report(y_test, y_pred)
7
8 print(f"Accuracy: {accuracy}")
9 print("Classification Report:\n", report)
```

```
Accuracy: 0.8100558659217877
Classification Report:
              precision    recall  f1-score   support

           0       0.83      0.86      0.84       105
           1       0.79      0.74      0.76        74

    accuracy                           0.81       179
   macro avg       0.81      0.80      0.80       179
weighted avg       0.81      0.81      0.81       179
```
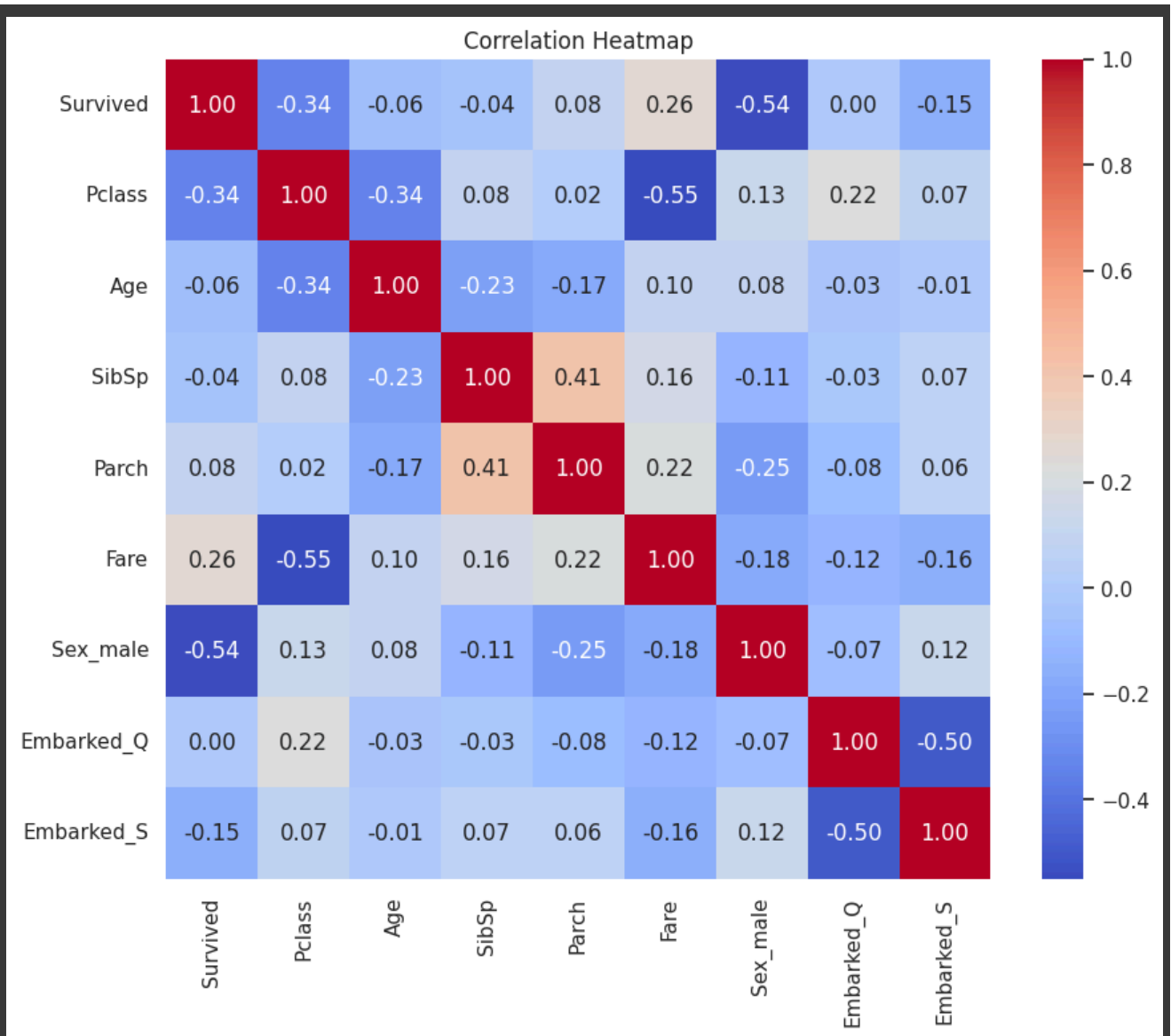
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 from sklearn.metrics import roc_curve, roc_auc_score, confusion_matrix
4
5 # Set the style for seaborn
6 sns.set(style="whitegrid")
```

```
1 # Plotting the correlation heatmap
2 plt.figure(figsize=(10, 8))
3 sns.heatmap(titanic_df.corr(), annot=True, fmt=".2f", cmap='coolwarm')
4 plt.title("Correlation Heatmap")
5 plt.show()
```
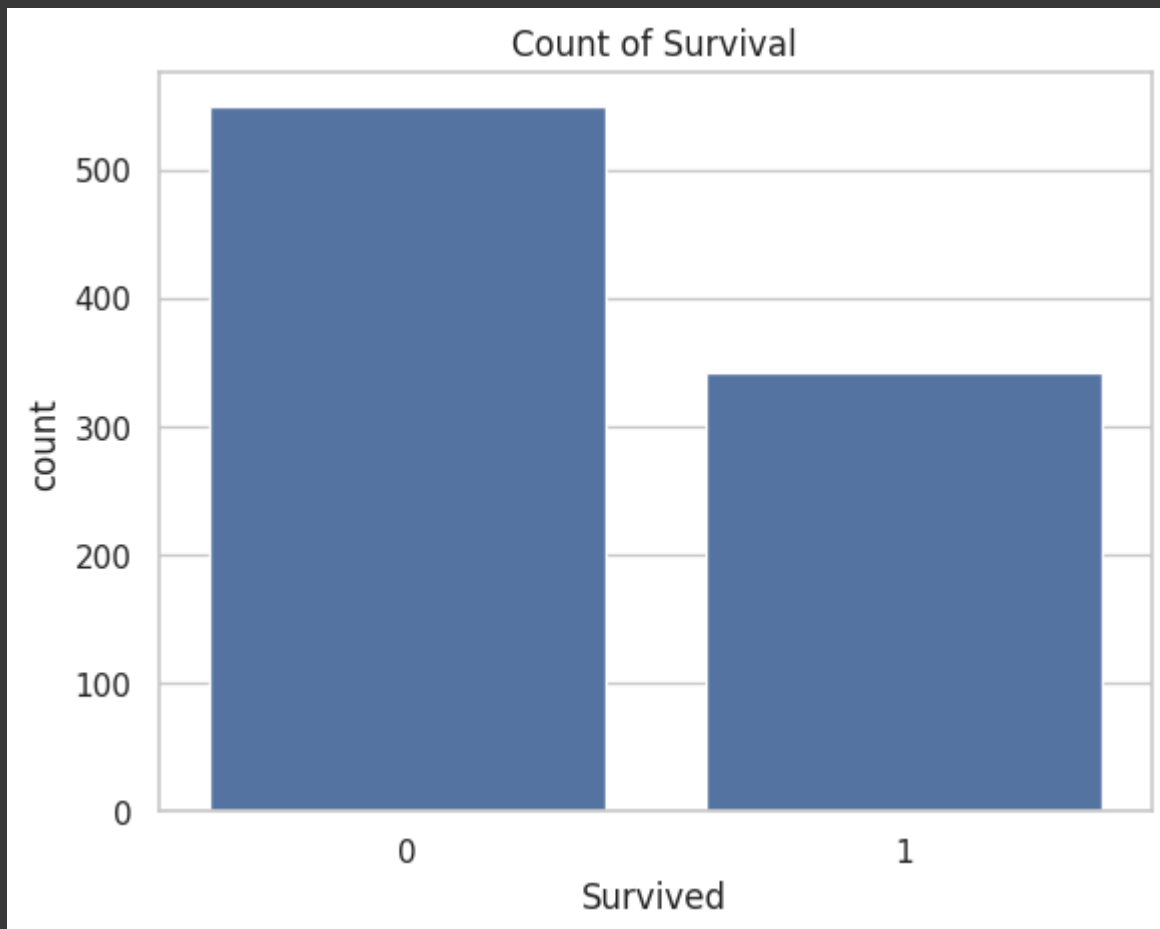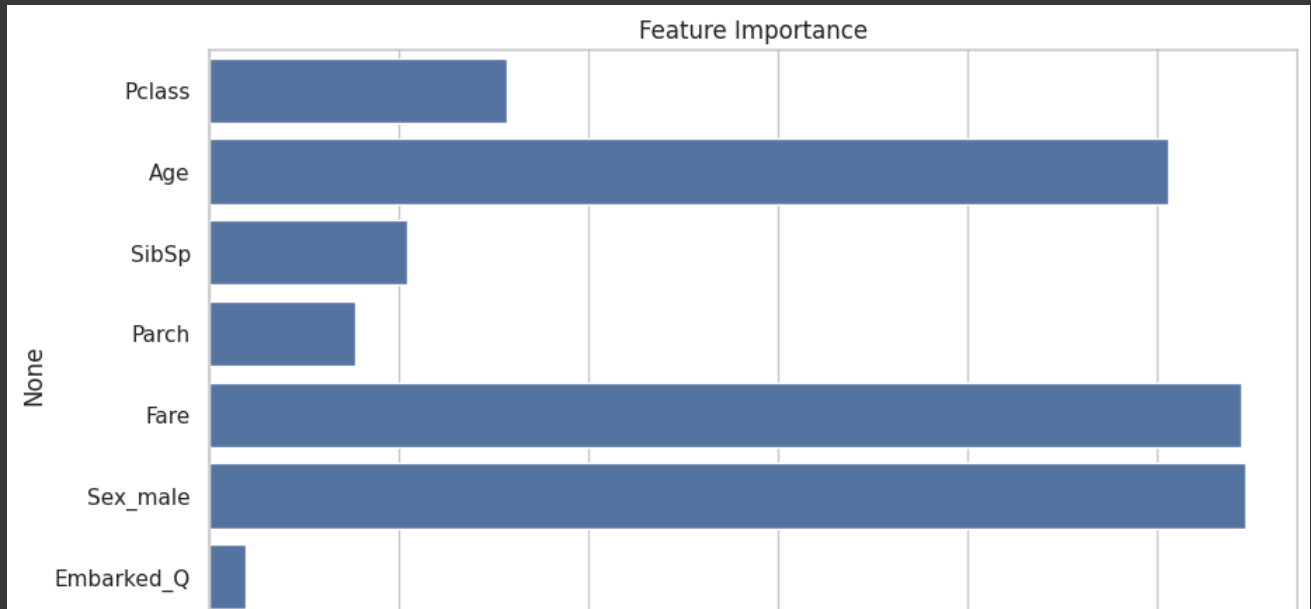
Correlation Heatmap

```
1 # Plotting the survival count
2 sns.countplot(x='Survived', data=titanic_df)
3 plt.title("Count of Survival")
4 plt.show()
```

Count of Survival

```
 1 from sklearn.ensemble import RandomForestClassifier
 2
 3 # Train a Random Forest model
 4 rf_model = RandomForestClassifier(random_state=42)
 5 rf_model.fit(X_train, y_train)
 6
 7 # Get feature importances
 8 importances = rf_model.feature_importances_
 9 feature_names = X.columns
10
11 # Plot feature importances
12 plt.figure(figsize=(10, 6))
13 sns.barplot(x=importances, y=feature_names)
14 plt.title("Feature Importance")
15 plt.show()
```

Feature Importance

```
1  # Calculate ROC curve and AUC
2  y_pred_prob = model.predict_proba(X_test)[:, 1]
3  fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob)
4  roc_auc = roc_auc_score(y_test, y_pred_prob)
5
6  # Plot ROC curve
7  plt.figure(figsize=(8, 6))
8  plt.plot(fpr, tpr, label=f"ROC Curve (area = {roc_auc:.2f})")
9  plt.plot([0, 1], [0, 1], 'k--')
10 plt.xlabel("False Positive Rate")
11 plt.ylabel("True Positive Rate")
12 plt.title("ROC Curve")
13 plt.legend(loc="lower right")
14 plt.show()
```



ROC Curve