

Taanshu Babariya 03

21/01/26

Experiment-2: Implementation of an End-to-End Machine Learning Data Pipeline.

ML LAB

EXP 1-C

PID - 246001

```
In [1]: # pip install numpy
# pip install pandas
# pip install matplotlib
!pip install numpy scikit-learn
!pip install numpy seaborn
```

```
Requirement already satisfied: numpy in c:\users\hp\ml03\lib\site-packages (2.0.2)
Requirement already satisfied: scikit-learn in c:\users\hp\ml03\lib\site-packages
(1.6.1)
Requirement already satisfied: threadpoolctl>=3.1.0 in c:\users\hp\ml03\lib\site-packages
(from scikit-learn) (3.6.0)
Requirement already satisfied: scipy>=1.6.0 in c:\users\hp\ml03\lib\site-packages
(from scikit-learn) (1.13.1)
Requirement already satisfied: joblib>=1.2.0 in c:\users\hp\ml03\lib\site-packages
(from scikit-learn) (1.5.3)
```

```
WARNING: You are using pip version 20.2.3; however, version 25.3 is available.
You should consider upgrading via the 'c:\users\hp\ml03\scripts\python.exe -m pip
install --upgrade pip' command.
```

```
Requirement already satisfied: numpy in c:\users\hp\ml03\lib\site-packages (2.0.2)
Requirement already satisfied: seaborn in c:\users\hp\ml03\lib\site-packages (0.13.
2)
Requirement already satisfied: matplotlib!=3.6.1,>=3.4 in c:\users\hp\ml03\lib\site-
packages (from seaborn) (3.9.4)
Requirement already satisfied: pandas>=1.2 in c:\users\hp\ml03\lib\site-packages (fr
om seaborn) (2.3.3)
Requirement already satisfied: contourpy>=1.0.1 in c:\users\hp\ml03\lib\site-package
s (from matplotlib!=3.6.1,>=3.4->seaborn) (1.3.0)
Requirement already satisfied: cycler>=0.10 in c:\users\hp\ml03\lib\site-packages (f
rom matplotlib!=3.6.1,>=3.4->seaborn) (0.12.1)
Requirement already satisfied: importlib-resources>=3.2.0; python_version < "3.10" i
n c:\users\hp\ml03\lib\site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (6.5.2)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\hp\ml03\lib\site-pac
kages (from matplotlib!=3.6.1,>=3.4->seaborn) (2.9.0.post0)
Requirement already satisfied: packaging>=20.0 in c:\users\hp\ml03\lib\site-packages
(from matplotlib!=3.6.1,>=3.4->seaborn) (25.0)
Requirement already satisfied: pyparsing>=2.3.1 in c:\users\hp\ml03\lib\site-package
s (from matplotlib!=3.6.1,>=3.4->seaborn) (3.3.1)
Requirement already satisfied: kiwisolver>=1.3.1 in c:\users\hp\ml03\lib\site-packag
es (from matplotlib!=3.6.1,>=3.4->seaborn) (1.4.7)
Requirement already satisfied: pillow>=8 in c:\users\hp\ml03\lib\site-packages (from
matplotlib!=3.6.1,>=3.4->seaborn) (11.3.0)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\hp\ml03\lib\site-pacak
es (from matplotlib!=3.6.1,>=3.4->seaborn) (4.60.2)
Requirement already satisfied: tzdata>=2022.7 in c:\users\hp\ml03\lib\site-packages
(from pandas>=1.2->seaborn) (2025.3)
Requirement already satisfied: pytz>=2020.1 in c:\users\hp\ml03\lib\site-packages (f
rom pandas>=1.2->seaborn) (2025.2)
Requirement already satisfied: zipp>=3.1.0; python_version < "3.10" in c:\users\hp\m
l03\lib\site-packages (from importlib-resources>=3.2.0; python_version < "3.10"->mat
plotlib!=3.6.1,>=3.4->seaborn) (3.23.0)
Requirement already satisfied: six>=1.5 in c:\users\hp\ml03\lib\site-packages (from
python-dateutil>=2.7->matplotlib!=3.6.1,>=3.4->seaborn) (1.17.0)
```

WARNING: You are using pip version 20.2.3; however, version 25.3 is available.
 You should consider upgrading via the 'c:\users\hp\ml03\scripts\python.exe -m pip in
 stall --upgrade pip' command.

In [2]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

In [3]:

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import accuracy_score
```

In [4]:

```
import seaborn as sns
#Load titanic dataset
titanic_data = sns.load_dataset('titanic')
```

In [5]:

```
print(titanic_data.shape)
```

(891, 15)

```
In [6]: print(titanic_data.columns)
```

```
Index(['survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare',
       'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town',
       'alive', 'alone'],
      dtype='object')
```

```
In [7]: print(titanic_data.head())
```

```
survived  pclass    sex   age  sibsp  parch    fare embarked  class \
0         0        3  male  22.0     1      0  7.2500      S  Third
1         1        1  female  38.0     1      0  71.2833     C  First
2         1        3  female  26.0     0      0  7.9250      S  Third
3         1        1  female  35.0     1      0  53.1000     S  First
4         0        3  male  35.0     0      0  8.0500      S  Third

      who  adult_male  deck  embark_town  alive  alone
0  man      True    NaN  Southampton  no  False
1 woman    False     C  Cherbourg  yes  False
2 woman    False    NaN  Southampton  yes  True
3 woman    False     C  Southampton  yes  False
4  man      True    NaN  Southampton  no  True
```

```
In [8]: print(titanic_data.tail())
```

```
survived  pclass    sex   age  sibsp  parch    fare embarked  class \
886       0        2  male  27.0     0      0  13.00      S  Second
887       1        1  female  19.0     0      0  30.00      S  First
888       0        3  female  NaN      1      2  23.45      S  Third
889       1        1  male  26.0     0      0  30.00     C  First
890       0        3  male  32.0     0      0   7.75      Q  Third

      who  adult_male  deck  embark_town  alive  alone
886  man      True    NaN  Southampton  no  True
887 woman    False     B  Southampton  yes  True
888 woman    False    NaN  Southampton  no  False
889  man      True     C  Cherbourg  yes  True
890  man      True    NaN  Queenstown  no  True
```

```
In [9]: print(titanic_data.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 15 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   survived    891 non-null    int64  
 1   pclass      891 non-null    int64  
 2   sex         891 non-null    object  
 3   age         714 non-null    float64 
 4   sibsp       891 non-null    int64  
 5   parch       891 non-null    int64  
 6   fare         891 non-null    float64 
 7   embarked    889 non-null    object  
 8   class        891 non-null    category
 9   who          891 non-null    object  
 10  adult_male  891 non-null    bool   
 11  deck         203 non-null    category
 12  embark_town 889 non-null    object  
 13  alive        891 non-null    object  
 14  alone        891 non-null    bool   
dtypes: bool(2), category(2), float64(2), int64(4), object(5)
memory usage: 80.7+ KB
None
```

In [10]: `print(titanic_data.describe())`

	survived	pclass	age	sibsp	parch	fare
count	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208
std	0.486592	0.836071	14.526497	1.102743	0.806057	49.693429
min	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
25%	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400
50%	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

In [11]: `missing_values = titanic_data.isnull().sum()`
`print(missing_values)`

survived	0
pclass	0
sex	0
age	177
sibsp	0
parch	0
fare	0
embarked	2
class	0
who	0
adult_male	0
deck	688
embark_town	2
alive	0
alone	0
dtype: int64	

```
In [12]: new_titanic_df = titanic_data.drop(columns=['deck'])
```

```
In [13]: new_titanic_df['age'].fillna(new_titanic_df['age'].median(), inplace=True)
```

C:\Users\HP\AppData\Local\Temp\ipykernel_14364\1870140451.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
new_titanic_df['age'].fillna(new_titanic_df['age'].median(), inplace=True)
```

```
In [14]: missing_values = new_titanic_df.isnull().sum()
print(missing_values)
```

survived	0
pclass	0
sex	0
age	0
sibsp	0
parch	0
fare	0
embarked	2
class	0
who	0
adult_male	0
embark_town	2
alive	0
alone	0
dtype:	int64

```
In [15]: data = new_titanic_df
data['embark_town'].dtype
data['embark_town'].unique()
data['embark_town'].fillna(data['embark_town'].mode()[0], inplace=True)
data.isnull().sum()
```

C:\Users\HP\AppData\Local\Temp\ipykernel_14364\1968151471.py:4: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
data['embark_town'].fillna(data['embark_town'].mode()[0], inplace=True)
```

```
Out[15]: survived      0
         pclass        0
         sex          0
         age          0
         sibsp        0
         parch        0
         fare          0
         embarked      2
         class         0
         who          0
         adult_male    0
         embark_town   0
         alive         0
         alone         0
         dtype: int64
```

```
In [16]: data = new_titanic_df
data['embarked'].dtype
data['embarked'].unique()
data['embarked'].fillna(data['embarked'].mode()[0] , inplace=True)
data.isnull().sum()
```

C:\Users\HP\AppData\Local\Temp\ipykernel_14364\2633234362.py:4: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method ({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
data['embarked'].fillna(data['embarked'].mode()[0] , inplace=True)
```

```
Out[16]: survived      0
         pclass        0
         sex          0
         age          0
         sibsp        0
         parch        0
         fare          0
         embarked      0
         class         0
         who          0
         adult_male    0
         embark_town   0
         alive         0
         alone         0
         dtype: int64
```

```
In [17]: le = LabelEncoder()
data['sex'] = le.fit_transform(data['sex'])
data['embarked'] = le.fit_transform(data['embarked'])
```

```
In [19]: data.head()
```

Out[19]:

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	
0	0	3	1	22.0	1	0	7.2500		2	Third	man	True
1	1	1	0	38.0	1	0	71.2833		0	First	woman	False
2	1	3	0	26.0	0	0	7.9250		2	Third	woman	False
3	1	1	0	35.0	1	0	53.1000		2	First	woman	False
4	0	3	1	35.0	0	0	8.0500		2	Third	man	True



In [22]: `data.tail()`

Out[22]:

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_ma	
886	0	2	1	27.0	0	0	13.00		2	Second	man	Tr
887	1	1	0	19.0	0	0	30.00		2	First	woman	Fa
888	0	3	0	28.0	1	2	23.45		2	Third	woman	Fa
889	1	1	1	26.0	0	0	30.00		0	First	man	Tr
890	0	3	1	32.0	0	0	7.75		1	Third	man	Tr



In [23]: `data.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype  
 ---  --          --          --      
 0   survived    891 non-null   int64  
 1   pclass      891 non-null   int64  
 2   sex         891 non-null   int64  
 3   age         891 non-null   float64 
 4   sibsp       891 non-null   int64  
 5   parch       891 non-null   int64  
 6   fare         891 non-null   float64 
 7   embarked    891 non-null   int64  
 8   class        891 non-null   category 
 9   who          891 non-null   object  
 10  adult_male   891 non-null   bool    
 11  embark_town  891 non-null   object  
 12  alive        891 non-null   object  
 13  alone        891 non-null   bool    
dtypes: bool(2), category(1), float64(2), int64(6), object(3)
memory usage: 79.4+ KB
```

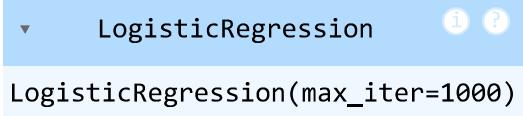
In [24]: `data = data[['pclass', 'sex', 'age', 'fare', 'embarked', 'survived']]`

```
X = data[['pclass', 'sex', 'age', 'fare', 'embarked']]
y = data['survived']
```

```
In [25]: X_train , X_test,y_train , y_test = train_test_split(X,y,test_size=0.3, random_stat
```

```
In [26]: model = LogisticRegression(max_iter = 1000)
model.fit(X_train,y_train)
```

Out[26]:



```
LogisticRegression(max_iter=1000)
```

```
In [27]: y_pred = model.predict(X_test)
```

```
In [28]: accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:",accuracy)
```

Accuracy: 0.7947761194029851

```
In [29]: new_passenger = pd.DataFrame({
    'pclass': [3],
    'sex': ['male'],
    'age' : [28],
    'fare' : [7.25],
    'embarked' : ['S']

})
```

```
In [30]: new_passenger_encoded = pd.get_dummies(new_passenger)
new_passenger_encoded = new_passenger_encoded.reindex(columns=X.columns,fill_value=
```

```
In [31]: prediction = model.predict(new_passenger_encoded)
print("Survived" if prediction[0] == 1 else "Did not survive")
```

Survived

```
In [32]: new_passenger = pd.DataFrame({
    'pclass': [1,3,2],
    'sex': ['female','male','female'],
    'age' : [38,45,14],
    'fare' : [80.0, 8.05,20.0],
    'embarked' : ['C','S','Q']

})
```

```
In [33]: new_passenger_encoded = pd.get_dummies(new_passenger)
new_passenger_encoded = new_passenger_encoded.reindex(columns=X.columns,fill_value=
```

```
In [36]: prediction = model.predict(new_passenger_encoded)
print("Survived" if prediction[0] == 1 else "Did not survive")
```

Survived

```
In [37]: predictions = model.predict(new_passenger_encoded)
for i,pred in enumerate(predictions):
```

```

    print(f"Passenger {i+1}:",
          "Survived" if pred == 1 else "Did not survive")

```

Passenger 1: Survived
Passenger 2: Survived
Passenger 3: Survived

```

In [38]: # Importing necessary libraries
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
import seaborn as sns
import pandas as pd

# Loading the Titanic dataset
titanic_df = sns.load_dataset('titanic')

# Splitting the dataset
train_data, test_data = train_test_split(titanic_df, test_size=0.2, random_state=42)

print(f"Train data shape: {train_data.shape}")
print(f"Test data shape: {test_data.shape}")

# Separating target variable
x_train = train_data.drop("survived", axis=1)
y_train = train_data["survived"]

x_test = test_data.drop("survived", axis=1)
y_test = test_data["survived"]

# Convert categorical variables into numerical form
x_train = pd.get_dummies(x_train, drop_first=True)
x_test = pd.get_dummies(x_test, drop_first=True)

# Handling missing values
x_train = x_train.fillna(0)
x_test = x_test.fillna(0)

# Align train and test columns (IMPORTANT)
x_train, x_test = x_train.align(x_test, join="left", axis=1, fill_value=0)

# Initialize Logistic Regression model
logreg = LogisticRegression(max_iter=1000)

# Train the model
logreg.fit(x_train, y_train)

# Make predictions
predictions = logreg.predict(x_test)

# Display evaluation metrics
print("Classification Report:")
print(classification_report(y_test, predictions))

print("Confusion Matrix:")
print(confusion_matrix(y_test, predictions))

```

```

print("Accuracy Score:")
print(accuracy_score(y_test, predictions))

Train data shape: (712, 15)
Test data shape: (179, 15)
Classification Report:
precision    recall   f1-score   support

          0       1.00      1.00      1.00      105
          1       1.00      1.00      1.00       74

   accuracy                           1.00      179
macro avg       1.00      1.00      1.00      179
weighted avg    1.00      1.00      1.00      179

Confusion Matrix:
[[105  0]
 [ 0  74]]
Accuracy Score:
1.0

```

PLOTTING

```

In [40]: import matplotlib.pyplot as plt
import seaborn as sns

cm = confusion_matrix(y_test, predictions)

plt.figure(figsize=(6,4))
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title("Confusion Matrix - Logistic Regression")
plt.show()

comparison_df = pd.DataFrame({
    "Actual": y_test.values,
    "Predicted": predictions
})

comparison_df.value_counts().plot(kind="bar", figsize=(6,4))
plt.title("Actual vs Predicted Survival Counts")
plt.xlabel("(Actual, Predicted)")
plt.ylabel("Count")
plt.show()
plt.figure(figsize=(6,4))
sns.countplot(x=predictions)
plt.xticks([0,1], ["Not Survived", "Survived"])
plt.title("Predicted Survival Distribution")
plt.xlabel("Prediction")
plt.ylabel("Count")
plt.show()
coefficients = pd.DataFrame({
    "Feature": x_train.columns,

```

```

    "Coefficient": logreg.coef_[0]
})

coefficients = coefficients.sort_values(by="Coefficient", ascending=False).head(10)

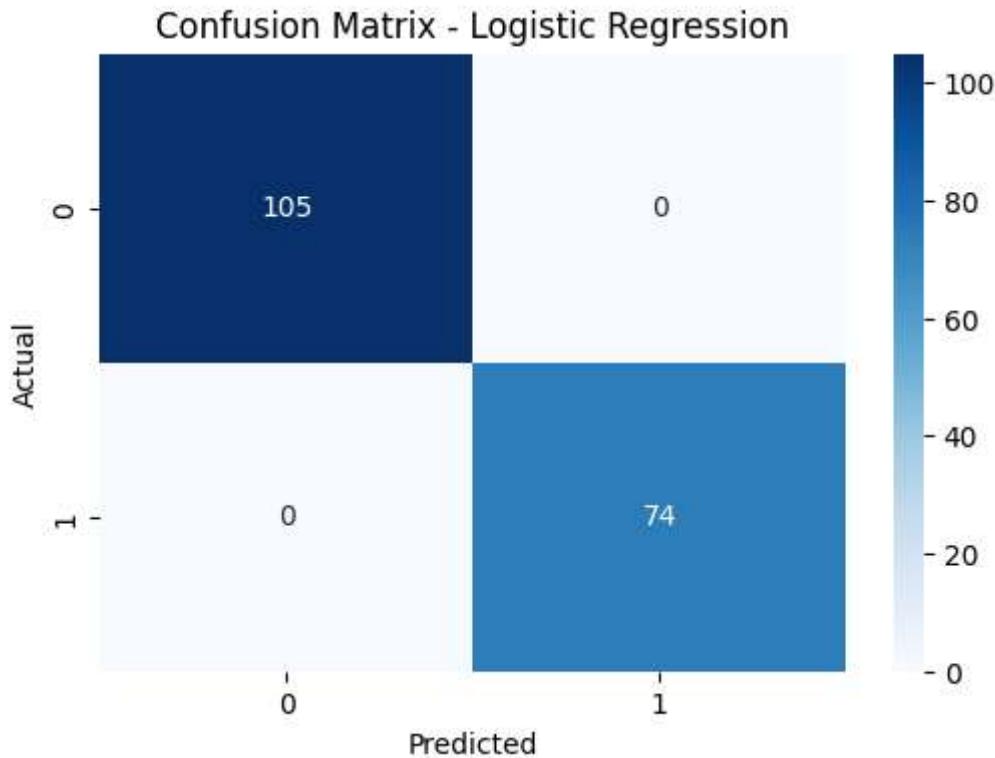
plt.figure(figsize=(8,5))
sns.barplot(x="Coefficient", y="Feature", data=coefficients)
plt.title("Top 10 Important Features (Logistic Regression)")
plt.show()

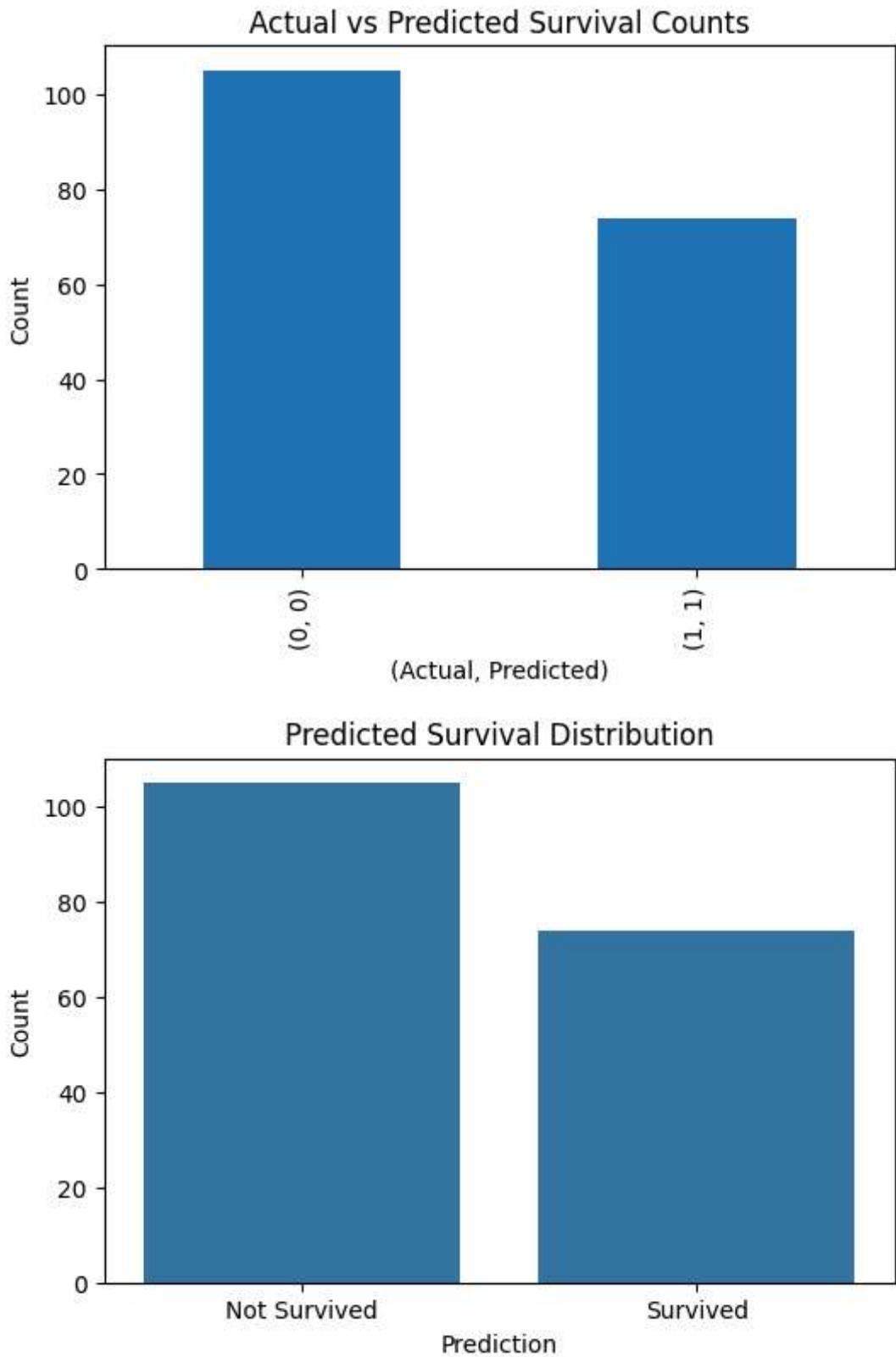
from sklearn.metrics import roc_curve, auc

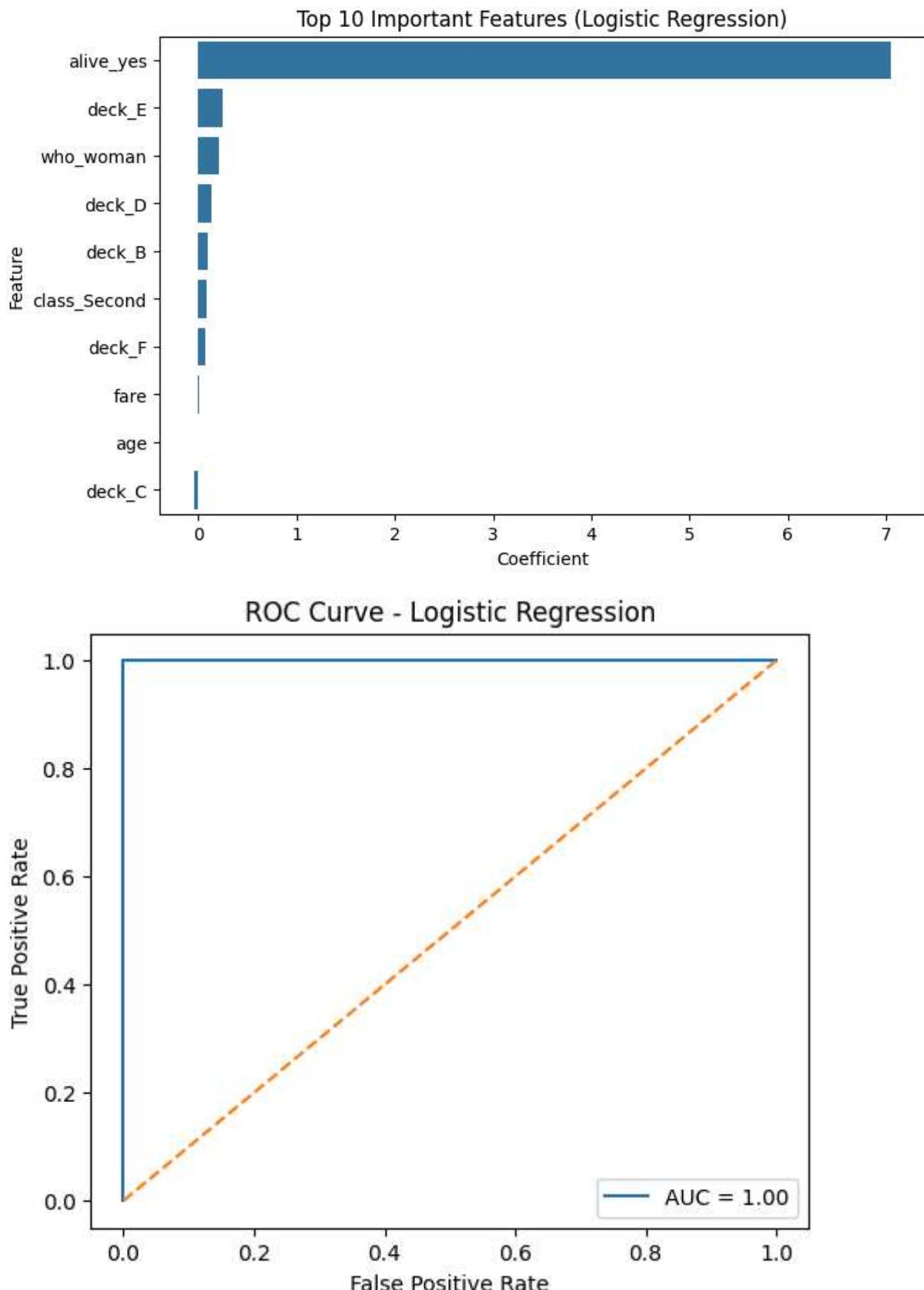
y_prob = logreg.predict_proba(x_test)[:,1]
fpr, tpr, _ = roc_curve(y_test, y_prob)
roc_auc = auc(fpr, tpr)

plt.figure(figsize=(6,5))
plt.plot(fpr, tpr, label=f"AUC = {roc_auc:.2f}")
plt.plot([0,1], [0,1], linestyle="--")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve - Logistic Regression")
plt.legend()
plt.show()

```







STANDARDIZATION

In [41]: `from sklearn.preprocessing import StandardScaler`

```

scaler = StandardScaler()

# Standardize the data
x_train_scaled = scaler.fit_transform(x_train)
x_test_scaled = scaler.transform(x_test)
logreg = LogisticRegression(max_iter=1000)
logreg.fit(x_train_scaled, y_train)

predictions = logreg.predict(x_test_scaled)

print("Accuracy:", accuracy_score(y_test, predictions))
print("Mean after standardization:", x_train_scaled.mean())
print("Std after standardization:", x_train_scaled.std())
scaler.fit_transform(x_test)

```

Accuracy: 1.0

Mean after standardization: -6.833810508195883e-18

Std after standardization: 1.0

Out[41]: array([[0.88742288, -1.34903605, 0.82036305, ..., -0.32394177,
 -1.40830868, 1.19118383],
 [-0.25537349, 0.37905601, -0.55202, ..., -0.32394177,
 0.7100716, -0.83950099],
 [0.88742288, -0.23413795, -0.55202, ..., -0.32394177,
 0.7100716, -0.83950099],
 ...,
 [0.88742288, 0.76927035, 0.82036305, ..., -0.32394177,
 0.7100716, 1.19118383],
 [-0.25537349, -0.40137266, -0.55202, ..., -0.32394177,
 0.7100716, 1.19118383],
 [0.88742288, -1.12605643, 0.82036305, ..., -0.32394177,
 0.7100716, 1.19118383]])

NORMALIZATION

In [42]:

```

from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()

# Normalize the data
x_train_norm = scaler.fit_transform(x_train)
x_test_norm = scaler.transform(x_test)
logreg = LogisticRegression(max_iter=1000)
logreg.fit(x_train_norm, y_train)

predictions = logreg.predict(x_test_norm)

print("Accuracy:", accuracy_score(y_test, predictions))
print("Min value:", x_train_norm.min())
print("Max value:", x_train_norm.max())

```

Accuracy: 1.0

Min value: 0.0

Max value: 1.0

OUTLIER DETECTION AND FILE HANDLING IN THE TITANIC DATASET

In [43]:

```

import numpy as np
data = titanic_df['fare']
mean = np.mean(data) # calculates the mean
std_dev = np.std(data) # calculates the standard deviation
Z_scores = (data - mean) / std_dev # computes the Z-scores
outliers = data[np.abs(Z_scores) > 3] # finds all the data points that are 3 standard deviations away from the mean

Q1 = titanic_df['fare'].quantile(0.25) # calculates the first quartile
Q3 = titanic_df['fare'].quantile(0.75) # calculates the third quartile
IQR = Q3 - Q1 # computes the IQR

# Below, we find all the data points that fall below the lower bound or above the upper bound
outliers = titanic_df['fare'][
    (titanic_df['fare'] < (Q1 - 1.5 * IQR)) |
    (titanic_df['fare'] > (Q3 + 1.5 * IQR))
]

mean = np.mean(titanic_df['fare']) # calculates the mean
standard_deviation = np.std(titanic_df['fare']) # calculates the standard deviation
outliers = titanic_df['fare'][np.abs(titanic_df['fare'] - mean) > 3 * standard_deviation]

```



```

import pandas as pd
import numpy as np

# Outlier detection - 'Age'
mean_age = np.mean(titanic_df['age']) # calculates the mean
std_dev_age = np.std(titanic_df['age']) # calculates the standard deviation
Z_scores_age = (titanic_df['age'] - mean_age) / std_dev_age # computes the Z-scores
outliers_age = titanic_df['age'][np.abs(Z_scores_age) > 3] # finds all the data points that are 3 standard deviations away from the mean
print("Outliers in 'Age' using Z-score: \n", outliers_age)

# Outlier detection - 'Fare'
mean_fare = np.mean(titanic_df['fare']) # calculates the mean
std_dev_fare = np.std(titanic_df['fare']) # calculates the standard deviation
Z_scores_fare = (titanic_df['fare'] - mean_fare) / std_dev_fare # computes the Z-scores
outliers_fare = titanic_df['fare'][np.abs(Z_scores_fare) > 3] # finds all the data points that are 3 standard deviations away from the mean
print("\nOutliers in 'Fare' using Z-score: \n", outliers_fare)

```



```

# Drop rows with missing 'age' values
titanic_df = titanic_df.dropna(subset=['age'])

# Calculate the upper bound for 'age'
Q1 = titanic_df['age'].quantile(0.25)
Q3 = titanic_df['age'].quantile(0.75)
IQR = Q3 - Q1

```

```

upper_bound = Q3 + 1.5 * IQR

# Cap the outliers for 'age'
titanic_df['age'] = np.where(titanic_df['age'] > upper_bound, upper_bound, titanic_

# Calculate the upper bound for 'fare'
Q1 = titanic_df['fare'].quantile(0.25)
Q3 = titanic_df['fare'].quantile(0.75)
IQR = Q3 - Q1
upper_bound = Q3 + 1.5 * IQR

# Cap the outliers for 'fare'
titanic_df['fare'] = np.where(titanic_df['fare'] > upper_bound, upper_bound, titanic_

```

Outliers in 'Age' using Z-score:

```

630    80.0
851    74.0
Name: age, dtype: float64

```

Outliers in 'Fare' using Z-score:

```

27    263.0000
88    263.0000
118   247.5208
258   512.3292
299   247.5208
311   262.3750
341   263.0000
377   211.5000
380   227.5250
438   263.0000
527   221.7792
557   227.5250
679   512.3292
689   211.3375
700   227.5250
716   227.5250
730   211.3375
737   512.3292
742   262.3750
779   211.3375
Name: fare, dtype: float64

```

PLOTTING

In [45]:

```

#Boxplot - Before Outlier Treatment
import matplotlib.pyplot as plt
import seaborn as sns

plt.figure(figsize=(12,5))

plt.subplot(1,2,1)
sns.boxplot(x=titanic_df['age'])
plt.title("Boxplot of Age")
plt.subplot(1,2,2)
sns.boxplot(x=titanic_df['fare'])

```

```

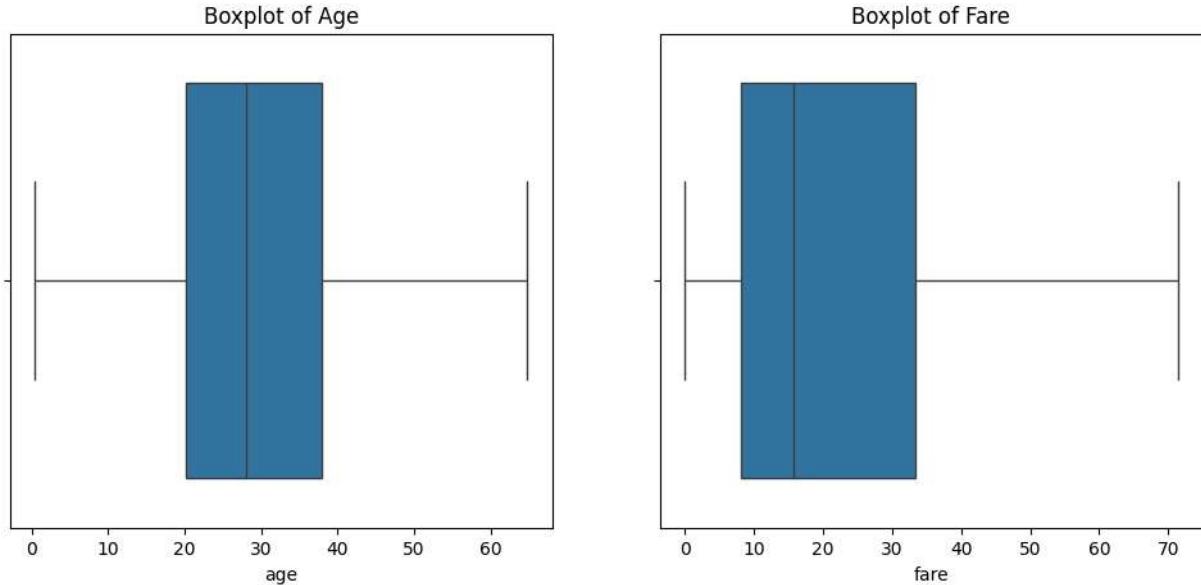
plt.title("Boxplot of Fare")
plt.show()

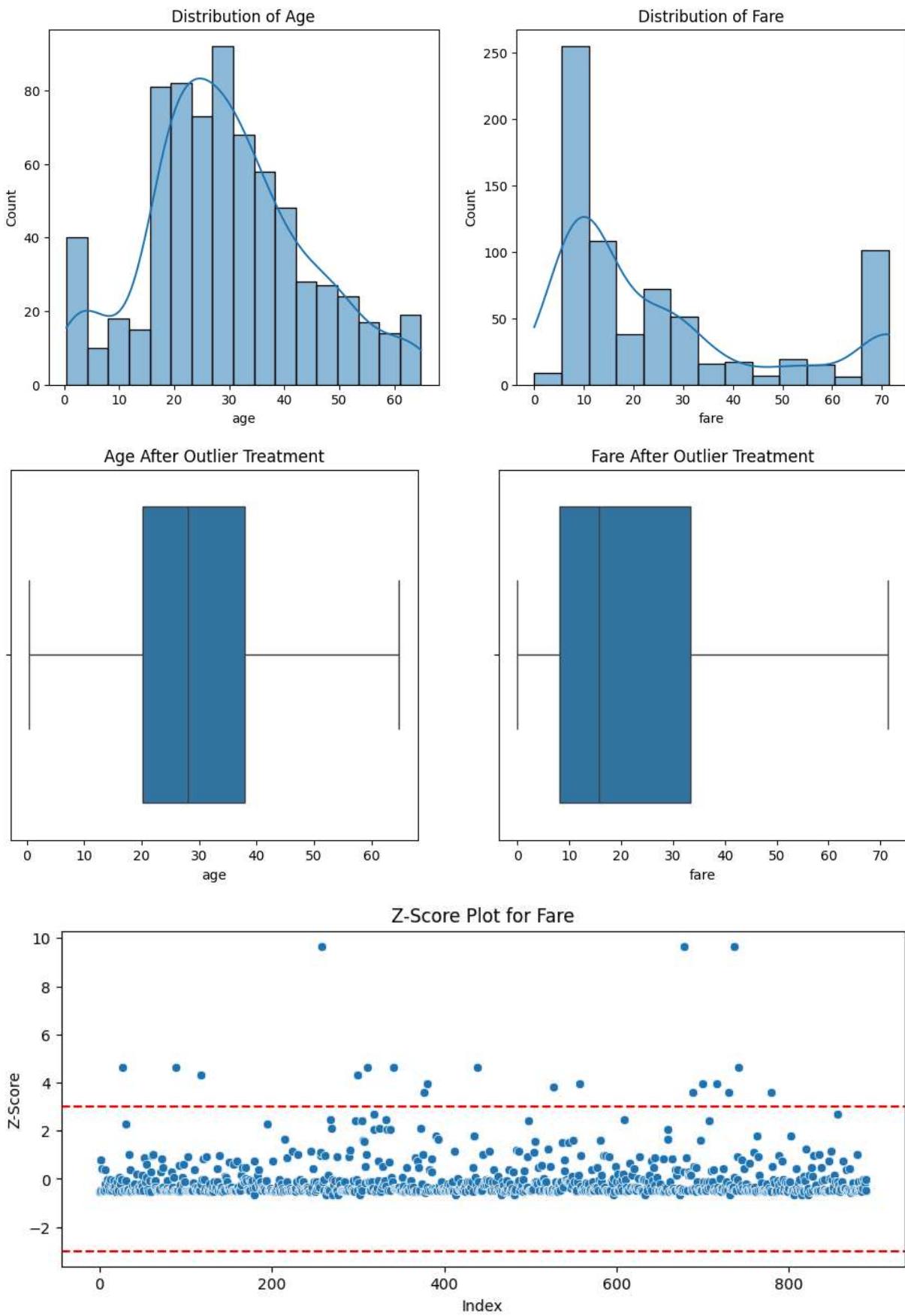
#Distribution Plot (Histogram + KDE)
plt.figure(figsize=(12,5))
plt.subplot(1,2,1)
sns.histplot(titanic_df['age'], kde=True)
plt.title("Distribution of Age")
plt.subplot(1,2,2)
sns.histplot(titanic_df['fare'], kde=True)
plt.title("Distribution of Fare")
plt.show()

#Boxplot - After Outlier Capping
fig, axes = plt.subplots(1, 2, figsize=(12,5))
sns.boxplot(x=titanic_df['age'], ax=axes[0])
axes[0].set_title("Age After Outlier Treatment")
sns.boxplot(x=titanic_df['fare'], ax=axes[1])
axes[1].set_title("Fare After Outlier Treatment")
plt.show()

#Z-Score Outliers Visualization
plt.figure(figsize=(10,4))
sns.scatterplot(x=range(len(Z_scores_fare)), y=Z_scores_fare)
plt.axhline(3, color='red', linestyle='--')
plt.axhline(-3, color='red', linestyle='--')
plt.title("Z-Score Plot for Fare")
plt.xlabel("Index")
plt.ylabel("Z-Score")
plt.show()

```





In []: