

▼ KILLARI ABHILESH

21BCE9306

MORNING BATCH

ASSIGNMENT - 3

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
```

```
df=pd.read_csv("titanic.csv")
df
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Emba
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	



```
df.head()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embark
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	
				Heikkinen								

```
df.shape
```

(891, 12)

```
df.describe()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare	
count	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000	
mean	446.000000	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208	

```
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
```

```
df.isnull().any()

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
```

```
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
```

```
df.Embarked.value_counts()

S    644
C    168
Q     77
Name: Embarked, dtype: int64
```

```
df["Age"].fillna(df["Age"].mean(),inplace=True)
```

```
df["Age"]

0      22.000000
1      38.000000
2      26.000000
3      35.000000
4      35.000000
...
886    27.000000
887    19.000000
888    29.699118
889    26.000000
890    32.000000
Name: Age, Length: 891, dtype: float64
```

```
df.isnull().any()

PassengerId    False
Survived        False
Pclass          False
Name            False
Sex             False
Age             False
SibSp           False
Parch           False
Ticket          False
Fare            False
Cabin           True
Embarked        True
dtype: bool

df["Embarked"].fillna(df["Embarked"].mode()[0],inplace=True)
df["Embarked"]



0      S
1      C
2      S
3      S
4      S
..
886    S
887    S
888    S
889    C
890    Q
Name: Embarked, Length: 891, dtype: object
```

```
df.isnull().sum()

PassengerId    0
Survived        0
Pclass          0
Name            0
Sex             0
Age             0
SibSp           0
Parch           0
Ticket          0
Fare            0
Cabin          687
Embarked        0
dtype: int64
```

```
corr=df.corr()
corr

<ipython-input-119-e92159fe13f9>:1: FutureWarning: The default value of numeric_only in DataFrame.corr
corr=df.corr()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare	
PassengerId	1.000000	-0.005007	-0.035144	0.033207	-0.057527	-0.001652	0.012658	
Survived	-0.005007	1.000000	-0.338481	-0.069809	-0.035322	0.081629	0.257307	
Pclass	-0.035144	-0.338481	1.000000	-0.331339	0.083081	0.018443	-0.549500	
Age	0.033207	-0.069809	-0.331339	1.000000	-0.232625	-0.179191	0.091566	
SibSp	-0.057527	-0.035322	0.083081	-0.232625	1.000000	0.414838	0.159651	
Parch	-0.001652	0.081629	0.018443	-0.179191	0.414838	1.000000	0.216225	
Fare	0.012658	0.257307	-0.549500	0.091566	0.159651	0.216225	1.000000	

```
plt.subplots(figsize=(10,10))
sns.heatmap(corr,annot=True)
```



```
df.drop(['Cabin'],axis=1, inplace=True)
df.drop(['Ticket'],axis=1, inplace=True)
df.drop(['Name'],axis=1,inplace=True)
df
```

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	1	0	3	male	22.000000	1	0	7.2500	S
1	2	1	1	female	38.000000	1	0	71.2833	C
2	3	1	3	female	26.000000	0	0	7.9250	S
3	4	1	1	female	35.000000	1	0	53.1000	S
4	5	0	3	male	35.000000	0	0	8.0500	S
...
886	887	0	2	male	27.000000	0	0	13.0000	S
887	888	1	1	female	19.000000	0	0	30.0000	S
888	889	0	3	female	29.699118	1	2	23.4500	S
889	890	1	1	male	26.000000	0	0	30.0000	C
890	891	0	3	male	32.000000	0	0	7.7500	Q

891 rows x 9 columns

```
df.Survived.value_counts()

0    549
1    342
Name: Survived, dtype: int64

df.Sex.value_counts()

male    577
female  314
Name: Sex, dtype: int64

df.Embarked.value_counts()
```

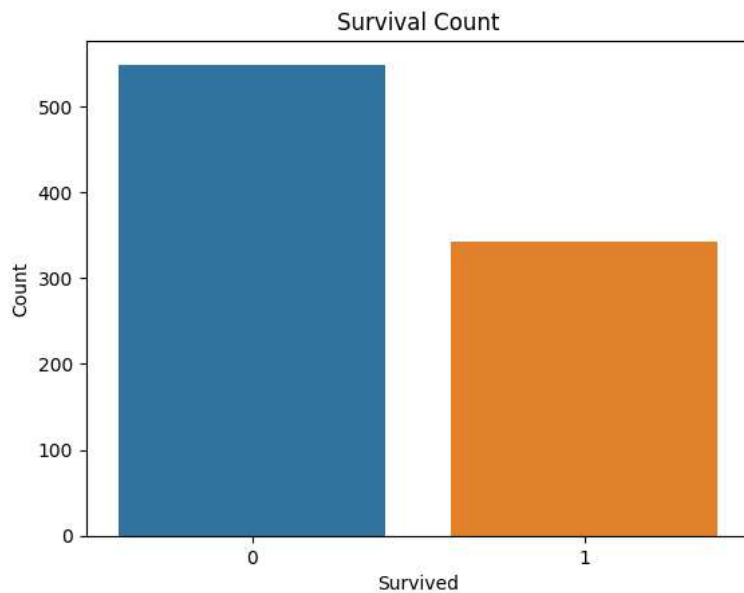
```
S    646
C    168
Q     77
Name: Embarked, dtype: int64
```

▼ DATA VISUALIZATION

```
print(df.isnull().sum())
```

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

```
# Visualize the distribution of the 'Survived' column (0 = Not Survived, 1 = Survived)
sns.countplot(data=df, x='Survived')
plt.title('Survival Count')
plt.xlabel('Survived')
plt.ylabel('Count')
plt.show()
```

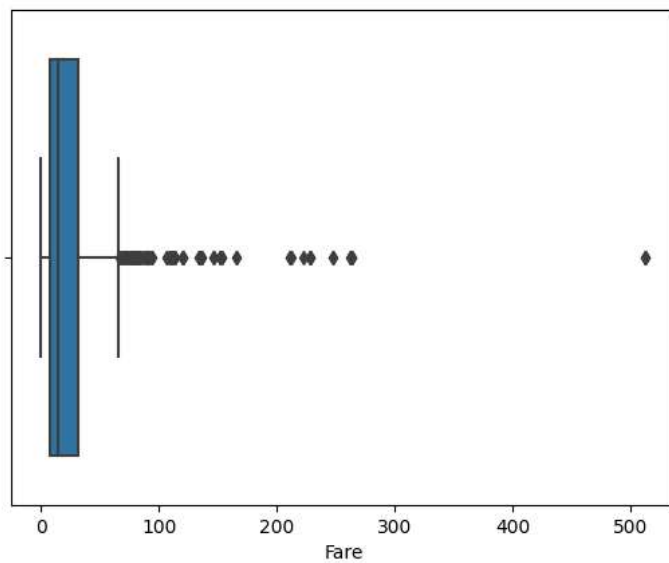


```
sns.histplot(data=df, x='Age', bins=20, kde=True)
plt.title('Age Distribution')
plt.xlabel('Age')
plt.ylabel('Count')
plt.show()
```

Age Distribution

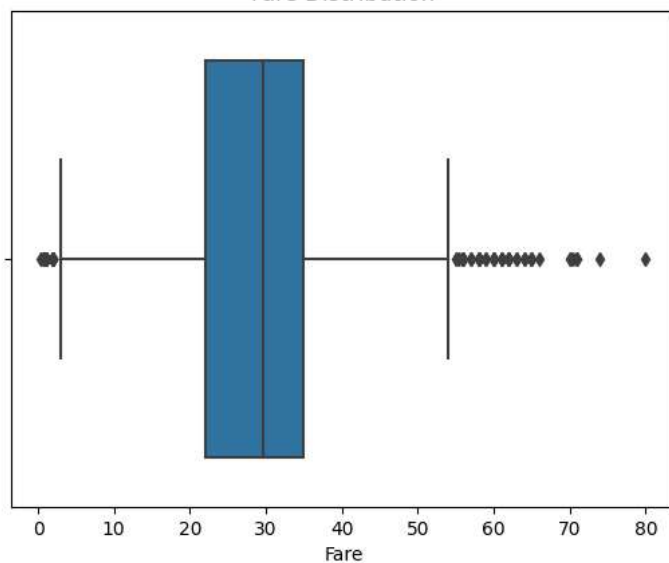
```
sns.boxplot(data=df, x='Fare')
plt.title('Fare Distribution')
plt.xlabel('Fare')
plt.show()
```

Fare Distribution



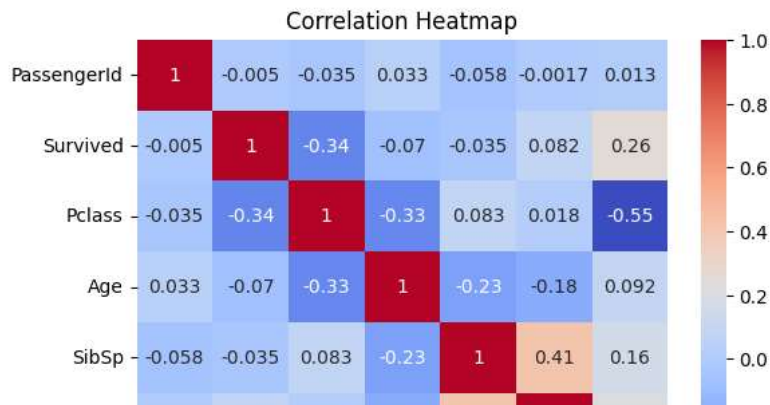
```
sns.boxplot(data=df, x='Age')
plt.title('Fare Distribution')
plt.xlabel('Fare')
plt.show()
```

Fare Distribution



```
corr_matrix = df.corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
plt.title('Correlation Heatmap')
plt.show()
```

```
<ipython-input-133-8dcdb071fff3>:1: FutureWarning: The default value of numeric_only in DataFrame.corr
corr_matrix = df.corr()
```



▼ Handling Outliers

Fare -0.013 0.26 -0.55 0.092 0.16 0.22 1

```
z_scores = np.abs(stats.zscore(df['Age']))
max_threshold=3
outliers = df['Age'][z_scores > max_threshold]

# Print and visualize the outliers
print("Outliers detected using Z-Score:")
print(outliers)
```

```
Outliers detected using Z-Score:
96    71.0
116    70.5
493    71.0
630    80.0
672    70.0
745    70.0
851    74.0
Name: Age, dtype: float64
```

```
z_scores = np.abs(stats.zscore(df['Fare']))
max_threshold=3
outliers = df['Fare'][z_scores > max_threshold]

# Print and visualize the outliers
print("Outliers detected using Z-Score:")
print(outliers)
```

```
Outliers detected 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
```

```
column_name="Fare"
Q1 = df[column_name].quantile(0.25)
Q3 = df[column_name].quantile(0.75)

# Calculate the IQR
IQR = Q3 - Q1

# Define the lower and upper bounds for outliers
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Filter rows with values outside the IQR bounds
```

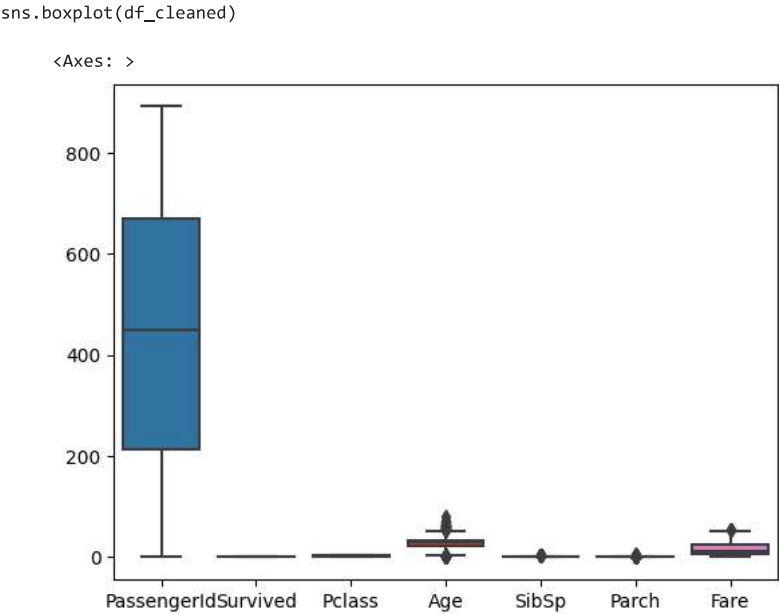
```
df_cleaned = df[(df[column_name] > lower_bound) & (df[column_name] < upper_bound)]

# Display the original and cleaned DataFrame sizes
print(f"Original DataFrame size: {df.shape}")
print(f"CleaneD DataFrame size: {df_cleaned.shape}")
df_cleaned
```

Original DataFrame size: (775, 9)
Cleaned DataFrame size: (750, 9)

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	1	0	3	male	22.000000	1	0	7.2500	S
2	3	1	3	female	26.000000	0	0	7.9250	S
3	4	1	1	female	35.000000	1	0	53.1000	S
4	5	0	3	male	35.000000	0	0	8.0500	S
5	6	0	3	male	29.699118	0	0	8.4583	Q
...
886	887	0	2	male	27.000000	0	0	13.0000	S
887	888	1	1	female	19.000000	0	0	30.0000	S
888	889	0	3	female	29.699118	1	2	23.4500	S
889	890	1	1	male	26.000000	0	0	30.0000	C
890	891	0	3	male	32.000000	0	0	7.7500	Q

750 rows × 9 columns



```
df=df_cleaned

x=df.drop('Survived', axis=1)
y=df['Survived']

x.head()
```

	PassengerId	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	1	3	male	22.000000	1	0	7.2500	S
2	3	3	female	26.000000	0	0	7.9250	S
3	4	1	female	35.000000	1	0	53.1000	S
4	5	3	male	35.000000	0	0	8.0500	S
5	6	3	male	29.699118	0	0	8.4583	Q

```
y.head()
```



```
0    0
2    1
3    1
4    0
5    0
Name: Survived, dtype: int64
```

▼ Encoding

```
le = LabelEncoder()
x['Sex'] = le.fit_transform(x['Sex'])
```

```
x.head()
```

	PassengerId	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	1	3	1	22.000000	1	0	7.2500	S
2	3	3	0	26.000000	0	0	7.9250	S
3	4	1	0	35.000000	1	0	53.1000	S
4	5	3	1	35.000000	0	0	8.0500	S
5	6	3	1	29.699118	0	0	8.4583	Q

```
x = pd.get_dummies(x,columns=['Embarked'])
```

```
x.head()
```

	PassengerId	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked_C	Embarked_Q	Embarked_S
0	1	3	1	22.000000	1	0	7.2500	0	0	1
2	3	3	0	26.000000	0	0	7.9250	0	0	1
3	4	1	0	35.000000	1	0	53.1000	0	0	1
4	5	3	1	35.000000	0	0	8.0500	0	0	1
5	6	3	1	29.699118	0	0	8.4583	0	1	0

▼ FEATURE SCALING

```
scale = StandardScaler()
x[['Age', 'Fare']] = scale.fit_transform(x[['Age', 'Fare']])
```

```
x.head()
```

	PassengerId	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked_C	Embarked_Q	Embarked_S
0	1	3	1	-0.556219	1	0	-0.779117	0	0	1
2	3	3	0	-0.243027	0	0	-0.729373	0	0	1
3	4	1	0	0.461654	1	0	2.599828	0	0	1
4	5	3	1	0.461654	0	0	-0.720161	0	0	1
5	6	3	1	0.046606	0	0	-0.690071	0	1	0

```
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=42)
```

```
print(x_train.shape)
print(x_test.shape)
print(y_train.shape)
print(y_test.shape)
```

```
(542, 10)
(233, 10)
(542,)
(233,)
```

