

Data Cleaning and preprocessing of Titanic Dataset using Python

```
titanic_preprocessing.py > ...
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  from sklearn.preprocessing import StandardScaler
6
7  # Load dataset
8  df = pd.read_csv("titanic.csv") # Make sure titanic.csv is in the same folder
9
10 # Show basic info
11 print(df.head())
12 print(df.info())
13
14 # Handle missing values
15 df['Age'].fillna(df['Age'].median(), inplace=True)
16 df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)
17 df.drop('Cabin', axis=1, inplace=True)
18
19 # Encode categorical variables
20 df = pd.get_dummies(df, columns=['Sex', 'Embarked'], drop_first=True)
21
22 # Normalize numerical columns
23 scaler = StandardScaler()
24 df[['Age', 'Fare']] = scaler.fit_transform(df[['Age', 'Fare']])
25
26 # Boxplot for outliers
27 sns.boxplot(x=df['Fare'])
28 plt.title("Fare Outliers")
29 plt.show()
30
31 # Remove Fare outliers using IQR
32 Q1 = df['Fare'].quantile(0.25)
33 Q3 = df['Fare'].quantile(0.75)
34 IQR = Q3 - Q1
35 df = df[~((df['Fare'] < (Q1 - 1.5 * IQR)) | (df['Fare'] > (Q3 + 1.5 * IQR)))]
36
37 # Save cleaned data
38 df.to_csv("titanic_cleaned.csv", index=False)
39 print("Data preprocessing complete.")
```

