	Assignment 3  Data Preprocessing  UMMALETI KUMAR  kumar.21bce9309@vitapstudent.ac.in  Data Preprocessing.
	<ol> <li>Import the Libraries.</li> <li>Importing the dataset.</li> <li>Checking for Null Values.</li> <li>Data Visualization.</li> <li>Outlier Detection</li> <li>Splitting Dependent and Independent variables</li> <li>Encoding</li> <li>Feature Scaling.</li> <li>Splitting Data into Train and Test.</li> </ol>
	1.Import the Libraries  import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns  2.Importing the dataset.
	df=pd.read_csv("Titanic-Dataset.csv")       df.head()       PassengerId Survived Pclass     Name Sex Age SibSp Parch     Ticket Fare Cabin Embarked       0     1     0     3     Braund, Mr. Owen Harris male 22.0     1     0     A/5 21171     7.2500     NaN     S       1     2     1     1     Cumings, Mrs. John Bradley (Florence Briggs Th female 38.0     1     0     PC 17599     71.2833     C85     C       2     3     1     3     Heikkinen, Miss. Laina female 26.0     0     0     STON/O2. 3101282     7.9250     NaN     S
[4]:	3
	mean         446.000000         0.383838         2.308642         29.699118         0.523008         0.381594         32.204208           std         257.353842         0.486592         0.836071         14.526497         1.102743         0.806057         49.693429           min         1.000000         0.000000         1.000000         0.420000         0.000000         0.000000         7.910400           50%         446.000000         0.000000         28.000000         0.000000         14.454200
< F D	75% 668.50000 1.00000 3.00000 1.00000 0.00000 31.00000  max 891.00000 1.00000 3.00000 80.00000 6.00000 512.329200  df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 891 entries, 0 to 890 Data columns (total 12 columns):  # Column Non-Null Count Dtype</class>
-	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 8 Ticket 891 non-null object 9 Fare 891 non-null float64
[6]: [6]:	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.corr(numeric_only=True)
	Survived       -0.005007       1.000000       -0.338481       -0.077221       -0.035322       0.081629       0.257307         Pclass       -0.035144       -0.338481       1.000000       -0.369226       0.083081       0.018443       -0.549500         Age       0.036847       -0.077221       -0.369226       1.000000       -0.189119       0.096067         SibSp       -0.057527       -0.035322       0.083081       -0.308247       1.000000       0.414838       0.159651         Parch       -0.001652       0.081629       0.018443       -0.189119       0.414838       1.000000       0.216225         Fare       0.012658       0.257307       -0.549500       0.096067       0.159651       0.216225       1.000000
	df.corr(numeric_only=True).Fare.sort_values(ascending=False)  Fare
[8]:	Name: Fare, dtype: float64  3.Checking for Null Values  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()
[9]:	PassengerId 0 Survived 0 Pclass 0 Name 0 Sex 0 Age 177 SibSp 0 Parch 0 Ticket 0 Fare 0 Cabin 687 Embarked 2
10]: 11]: 11]:	<pre>dtype: int64  df.Sex.nunique()  2  df.Sex.unique()  array(['male', 'female'], dtype=object)  df.Sex.value_counts()</pre>
12]: 13]:	Sex male 577 female 314 Name: count, dtype: int64  df.Embarked.nunique()
14]: 15]:	<pre>array(['S', 'C', 'Q', nan], dtype=object)  df.Embarked.value_counts()  Embarked S    644 C    168 Q    77 Name: count, dtype: int64</pre>
[16]:	4.Data Visualization  sns.countplot(x='Survived', data=df) plt.show()  C:\Users\kumar\AppData\Roaming\Python\Python311\site-packages\seaborn\_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a fewersion. Use isinstance(dtype, CategoricalDtype) instead if pd.api.types.is_categorical_dtype(vector): C:\Users\kumar\AppData\Roaming\Python\Python311\site-packages\seaborn\_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a fewersion. Use isinstance(dtype, CategoricalDtype) instead if pd. One is instance(dtype, CategoricalDtype) instead
	<pre>if pd.api.types.is_categorical_dtype(vector): C:\Users\kumar\AppData\Roaming\Python\Python311\site-packages\seaborn\_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a fiversion. Use isinstance(dtype, CategoricalDtype) instead     if pd.api.types.is_categorical_dtype(vector):  500 - 400 -</pre>
	100 -
	plt.scatter(df["Survived"],df["Fare"]) <matplotlib.collections.pathcollection 0x2795b158710="" at=""></matplotlib.collections.pathcollection>
	500 - 400 - 300 -
	200 - 1
18]:	sns.heatmap(df.corr(numeric_only=True), annot=True) <axes:>  PassengerId - 1 -0.005 -0.035 0.037 -0.058 -0.0017 0.013  Survived</axes:>
	Pclass
	Parch - Fare - Bassengerid0.2 0.4 0.4
20]:	<pre>import warnings # Ignore the specified warnings globally warnings.filterwarnings("ignore", category=FutureWarning) warnings.filterwarnings("ignore", category=UserWarning)  sns.pairplot(df) <seaborn.axisgrid.pairgrid 0x2795b1c1d10="" at=""></seaborn.axisgrid.pairgrid></pre>
ģ	800
	0.8
·	2.5
	5
	400
	<pre>sns.barplot(x=df["Embarked"], y=df["Fare"], ci=0)  <axes: ,="" xlabel="Embarked" ylabel="Fare">  60 -  50 -  40 -</axes:></pre>
	<u>2</u> 30 - 20 - 10 -
22]:	5. Outlier Detection  plt.figure(figsize=(8, 4)) sns.boxplot(x='Pclass', y='Fare', data=df)
	plt.show()  500 - 400 - 300 - 400 -
•	200 - 100 -
[23]:	<pre>sns.boxplot(df["SibSp"]) </pre> <pre> <pre> </pre> <pre> <pre> </pre> <pre> <p< td=""></p<></pre></pre></pre>
:	5 -
,	6.Splitting Dependent and independent Variables  df.head()
24]:	
	<pre>non_numeric_cols = ['Name', 'Ticket', 'Cabin']  # Extract numerical columns by excluding non-numeric columns numeric_cols = [col for col in df.columns if col not in non_numeric_cols]  df = df[numeric_cols]</pre> df.head()
J:	PassengerId         Survived         Pclass         Sex         Age         SibSp         Parch         Fare         Embarked           0         1         0         3         male         22.0         1         0         7.2500         S           1         2         1         1         female         38.0         1         0         71.2833         C           2         3         1         3         female         26.0         0         0         7.9250         S           3         4         1         female         35.0         1         0         53.1000         S           4         5         0         3         male         35.0         0         8.0500         S
	<pre>#indenpendent variables hould be 2 d array or dataframe X=df.drop(columns=["Survived"],axis=1) X.head()</pre>
	2 3 3 female 26.0 0 0 7.9250 S 3 4 1 female 35.0 1 0 53.1000 S 4 5 3 male 35.0 0 0 8.0500 S  X.shape (891, 8)
53]: 53]: 54]:	<pre>type(X) pandas.core.frame.DataFrame  y=df["Survived"] y.head()  0  0 1  1 2  1</pre>
55]: 55]:	3 1 4 0 Name: Survived, dtype: int64  7.Encoding  X.head()  PassengerId Pclass Sex Age SibSp Parch Fare Embarked
56]:	0       1       3       male       22.0       1       0       7.2500       S         1       2       1       female       38.0       1       0       71.2833       C         2       3       3       female       26.0       0       0       7.9250       S         3       4       1       female       35.0       1       0       53.1000       S         4       5       3       male       35.0       0       0       8.0500       S     From sklearn.preprocessing import LabelEncoder
57]:	<pre>le=LabelEncoder()  X["Sex"]=le.fit_transform(X["Sex"])  X.head()  PassengerId Pclass Sex Age SibSp Parch Fare Embarked  0</pre>
	1
60]: 60]: 61]:	<pre>mapping=dict(zip(le.classes_, range(len(le.classes_)))) mapping {'female': 0, 'male': 1}  X["Embarked"]=le.fit_transform(X["Embarked"])  X.head()</pre>
62]:	PassengerId         Pclass         Sex         Age         SibSp         Parch         Fare         Embarked           0         1         3         1         22.0         1         0         7.2500         2           1         2         1         0         38.0         1         0         71.2833         0           2         3         3         0         26.0         0         0         7.9250         2           3         4         1         0         35.0         1         0         53.1000         2           4         5         3         1         35.0         0         8.0500         2
64]: 64]:	<pre>print(le.classes_) ['C' 'Q' 'S' nan]  mapping=dict(zip(le.classes_,range(len(le.classes_)))) mapping {'C': 0, 'Q': 1, 'S': 2, nan: 3}</pre>
	8.Feature Scaling  from sklearn.preprocessing import MinMaxScaler ms=MinMaxScaler()
69]: 70]:	<pre>X_scaled = ms.fit_transform(X)  X_Scaled=pd.DataFrame(ms.fit_transform(X), columns=X.columns)  X_Scaled.head()</pre>
69]: 70]: 71]:	X_Scaled_head()  PassengerId Pclass Sex Age SibSp Parch Fare Embarked  0 0.000000 1.0 1.0 0.271174 0.125 0.0 0.014151 0.666667  1 0.001124 0.0 0.0 0.472229 0.125 0.0 0.139136 0.000000  2 0.002247 1.0 0.0 0.321438 0.000 0.0 0.015469 0.666667  3 0.003371 0.0 0.0 0.434531 0.125 0.0 0.103644 0.666667
69]: 70]: 71]: 73]:	X_Scaled=pd.DataFrame(ms.fit_transform(X),columns=X.columns)  X_Scaled.head()  Passengerid Polass Sex Age SibSp Parch Fare Embarked  0 0.000000 1.0 1.0 0.271174 0.125 0.0 0.014151 0.666667  1 0.001124 0.0 0.0 0.472229 0.125 0.0 0.139136 0.000000  2 0.002247 1.0 0.0 0.321438 0.000 0.0 0.015469 0.666667  3 0.003371 0.0 0.0 0.434531 0.125 0.0 0.103644 0.666667  4 0.004494 1.0 1.0 0.434531 0.000 0.0 0.015713 0.666667  # Splitting Data into Train and Test from sklearn.model_selection import train_test_split X_train, X_test, y_train, y_test = train_test_split(X_Scaled, y, test_size=0.2, random_state=42)  # Displaying the preprocessed data print(X_train.head()) print(X_train.head()) print(X_train.head())
70]: 71]: 71]: 73]: 73]: 73]: 73]: 73]: 73]: 73]: 73	X_Scaled=pd.DataFrame(ms.fit_transform(X),columns=X.columns)  X_Scaled.head()  PassengerId Pclass Sex Age SibSp Parch Fare Embarked  0 0.000000 1.0 1.0 0.271174 0.125 0.0 0.014151 0.666667  1 0.001124 0.0 0.0 0.472229 0.125 0.0 0.139136 0.000000  2 0.002247 1.0 0.0 0.321438 0.000 0.0 0.015469 0.666667  3 0.003371 0.0 0.0 0.434531 0.125 0.0 0.103644 0.666667  4 0.004494 1.0 1.0 0.0434531 0.000 0.0 0.015713 0.666667  # Splitting Data into Train and Test from sklearn.model_selection import train_test_split X_train, X_test, y_train, y_test = train_test_split(X_Scaled, y, test_size=0.2, random_state=42)  # Displaying the preprocessed data print(X_train.head())  print(X_test.head())