

assignment-3-data-preprocessing-1

September 21, 2023

```
[1]: ##Name : kavya sree
      ##AI_ML Morning slot
```

0.0.1 Data Preprocessing.

1 1.Import the Libraries.

```
[1]: import numpy as np
      import pandas as pd
      import matplotlib.pyplot as plt
      import seaborn as sns
```

2 2.Importing the dataset.

```
[3]: df=pd.read_csv("Titanic-Dataset.csv")
```

```
[ ]: df
```

```
[ ]:
      PassengerId  Survived  Pclass  \
0               1         0       3
1               2         1       1
2               3         1       3
3               4         1       1
4               5         0       3
..            ...         ...     ...
886            887         0       2
887            888         1       1
888            889         0       3
889            890         1       1
890            891         0       3
```

```

      Name      Sex  Age  SibSp  \
0  Braund, Mr. Owen Harris    male  22.0      1
1  Cumings, Mrs. John Bradley (Florence Briggs Th...  female  38.0      1
2      Heikkinen, Miss. Laina  female  26.0      0
3  Futrelle, Mrs. Jacques Heath (Lily May Peel)  female  35.0      1
```

4		Allen, Mr. William Henry	male	35.0	0
..	
886		Montvila, Rev. Juozas	male	27.0	0
887		Graham, Miss. Margaret Edith	female	19.0	0
888		Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1
889		Behr, Mr. Karl Howell	male	26.0	0
890		Dooley, Mr. Patrick	male	32.0	0

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	NaN	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2. 3101282	7.9250	NaN	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	NaN	S
..
886	0	211536	13.0000	NaN	S
887	0	112053	30.0000	B42	S
888	2	W./C. 6607	23.4500	NaN	S
889	0	111369	30.0000	C148	C
890	0	370376	7.7500	NaN	Q

[891 rows x 12 columns]

```
[4]: df.head()
```

```
[4]: PassengerId  Survived  Pclass  \
0             1         0         3
1             2         1         1
2             3         1         3
3             4         1         1
4             5         0         3
```

	Name	Sex	Age	SibSp	\
0	Braund, Mr. Owen Harris	male	22.0	1	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	
2	Heikkinen, Miss. Laina	female	26.0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	
4	Allen, Mr. William Henry	male	35.0	0	

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	NaN	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2. 3101282	7.9250	NaN	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	NaN	S

```
[5]: df.tail()
```

```
[5]:
```

	PassengerId	Survived	Pclass	Name \
886	887	0	2	Montvila, Rev. Juozas
887	888	1	1	Graham, Miss. Margaret Edith
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"
889	890	1	1	Behr, Mr. Karl Howell
890	891	0	3	Dooley, Mr. Patrick

	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
886	male	27.0	0	0	211536	13.00	NaN	S
887	female	19.0	0	0	112053	30.00	B42	S
888	female	NaN	1	2	W./C. 6607	23.45	NaN	S
889	male	26.0	0	0	111369	30.00	C148	C
890	male	32.0	0	0	370376	7.75	NaN	Q

```
[6]: 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
```

```
[7]: df.describe()
```

```
[7]:
```

	PassengerId	Survived	Pclass	Age	SibSp \
count	891.000000	891.000000	891.000000	714.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008
std	257.353842	0.486592	0.836071	14.526497	1.102743
min	1.000000	0.000000	1.000000	0.420000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000
50%	446.000000	0.000000	3.000000	28.000000	0.000000
75%	668.500000	1.000000	3.000000	38.000000	1.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000

	Parch	Fare
count	891.000000	891.000000
mean	0.381594	32.204208
std	0.806057	49.693429
min	0.000000	0.000000
25%	0.000000	7.910400
50%	0.000000	14.454200
75%	0.000000	31.000000
max	6.000000	512.329200

```
[8]: df.shape
```

```
[8]: (891, 12)
```

```
[9]: df.corr()
```

<ipython-input-9-2f6f6606aa2c>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

```
df.corr()
```

```
[9]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	\
PassengerId	1.000000	-0.005007	-0.035144	0.036847	-0.057527	-0.001652	
Survived	-0.005007	1.000000	-0.338481	-0.077221	-0.035322	0.081629	
Pclass	-0.035144	-0.338481	1.000000	-0.369226	0.083081	0.018443	
Age	0.036847	-0.077221	-0.369226	1.000000	-0.308247	-0.189119	
SibSp	-0.057527	-0.035322	0.083081	-0.308247	1.000000	0.414838	
Parch	-0.001652	0.081629	0.018443	-0.189119	0.414838	1.000000	
Fare	0.012658	0.257307	-0.549500	0.096067	0.159651	0.216225	

	Fare
PassengerId	0.012658
Survived	0.257307
Pclass	-0.549500
Age	0.096067
SibSp	0.159651
Parch	0.216225
Fare	1.000000

```
[10]: df.corr().Fare.sort_values(ascending=False)
```

<ipython-input-10-f51f352aac84>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

```
df.corr().Fare.sort_values(ascending=False)
```

```
[10]: Fare          1.000000
      Survived      0.257307
      Parch        0.216225
      SibSp        0.159651
      Age          0.096067
      PassengerId   0.012658
      Pclass       -0.549500
      Name: Fare, dtype: float64
```

```
[11]: df.Survived.value_counts()
```

```
[11]: 0    549
      1    342
      Name: Survived, dtype: int64
```

```
[12]: df.Sex.value_counts()
```

```
[12]: male        577
      female      314
      Name: Sex, dtype: int64
```

```
[13]: df.Embarked.value_counts()
```

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

3 3.Checking for Null Values.

```
[14]: df.isnull().any()
```

```
[14]: 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
```

```
[15]: df.isnull().sum()
```

```
[15]: 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
```

```
[16]: df["Age"].mean()
```

```
[16]: 29.69911764705882
```

```
[17]: df['Age'].fillna(df['Age'].mean(),inplace=True)
```

```
[18]: df.isnull().sum()
```

```
[18]: PassengerId      0
      Survived        0
      Pclass          0
      Name            0
      Sex             0
      Age            0
      SibSp           0
      Parch           0
      Ticket          0
      Fare            0
      Cabin          687
      Embarked        2
      dtype: int64
```

```
[19]: df["Embarked"].mode()
```

```
[19]: 0    S
      Name: Embarked, dtype: object
```

```
[22]: df['Embarked'].fillna(df['Embarked'].mode(),inplace=True)
```

```
[23]: df.isnull().sum()
```

```
[23]: 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
```

```
[24]: df.drop(["Cabin"],axis=1,inplace=True)
```

```
[25]: df.drop(["Ticket"],axis=1,inplace=True)
```

```
[26]: df.drop(["Name"],axis=1,inplace=True)
```

```
[27]: df.isnull().sum()
```

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

```
[28]: df.Embarked.nunique()
```

```
[28]: 3
```

```
[29]: df.Embarked.unique()
```

```
[29]: array(['S', 'C', 'Q'], dtype=object)
```

```
[30]: df.Embarked.value_counts()
```

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

```
[31]: df.head()
```

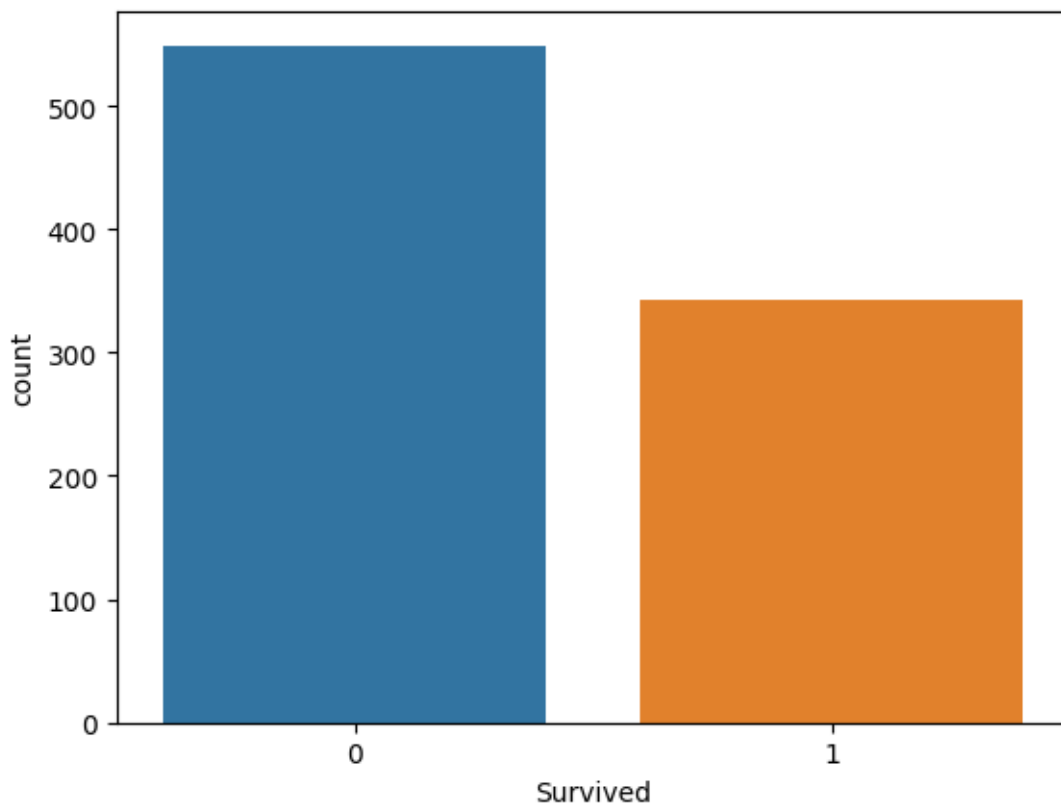
```
[31]:
```

	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	1	female	35.0	1	0	53.1000	S
4	5	0	3	male	35.0	0	0	8.0500	S

4 4.Data Visualization.

```
[32]: sns.countplot(x="Survived",data=df)
```

```
[32]: <Axes: xlabel='Survived', ylabel='count'>
```

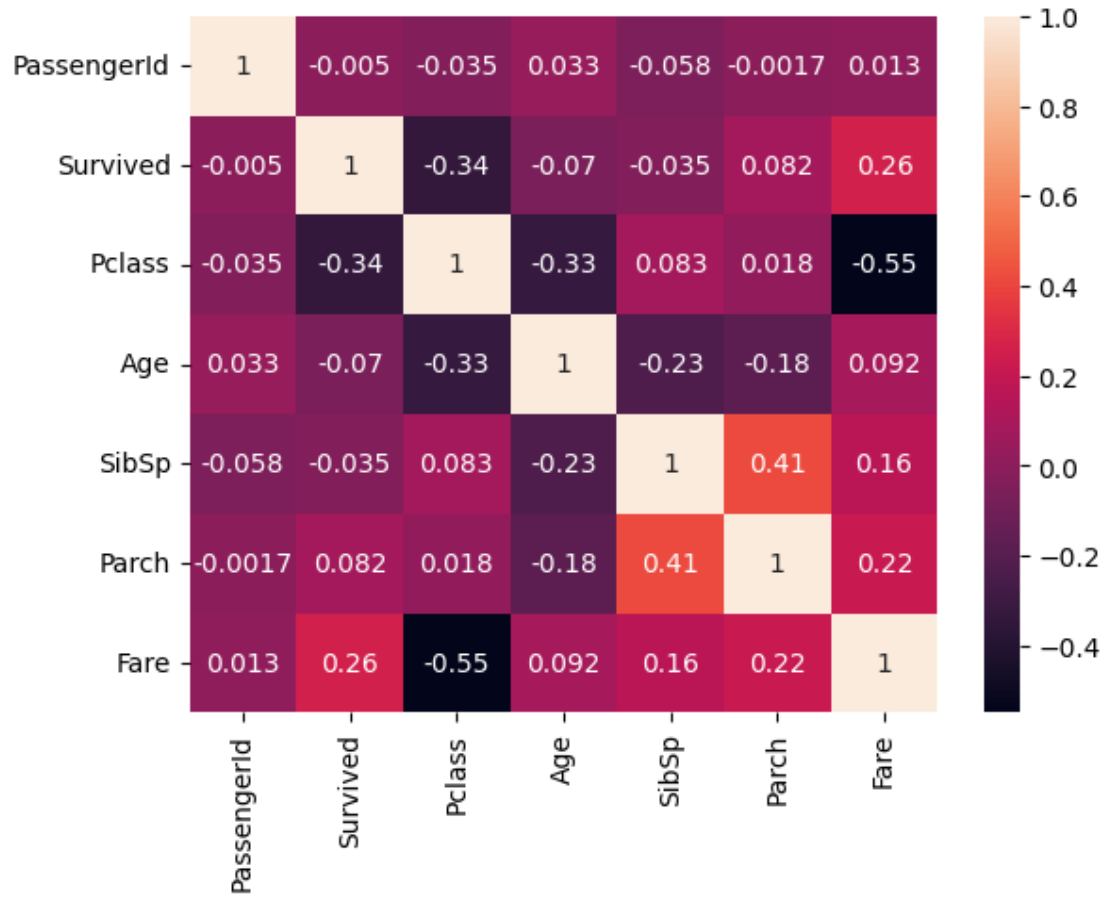


```
[33]: sns.heatmap(df.corr(),annot=True)
```

<ipython-input-33-8df7bcac526d>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only

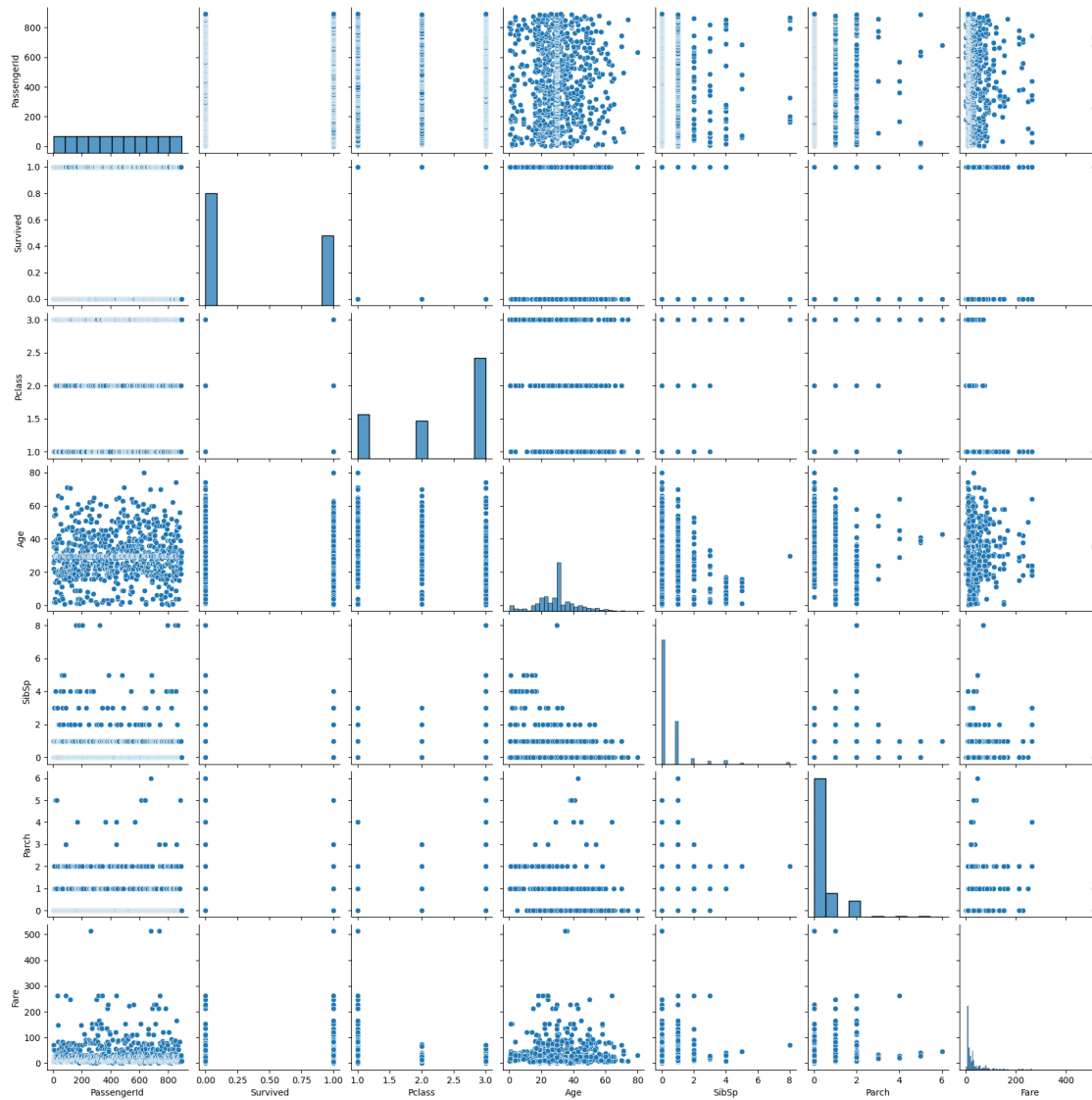

```
to silence this warning.
sns.heatmap(df.corr(),annot=True)
```

[33]: <Axes: >



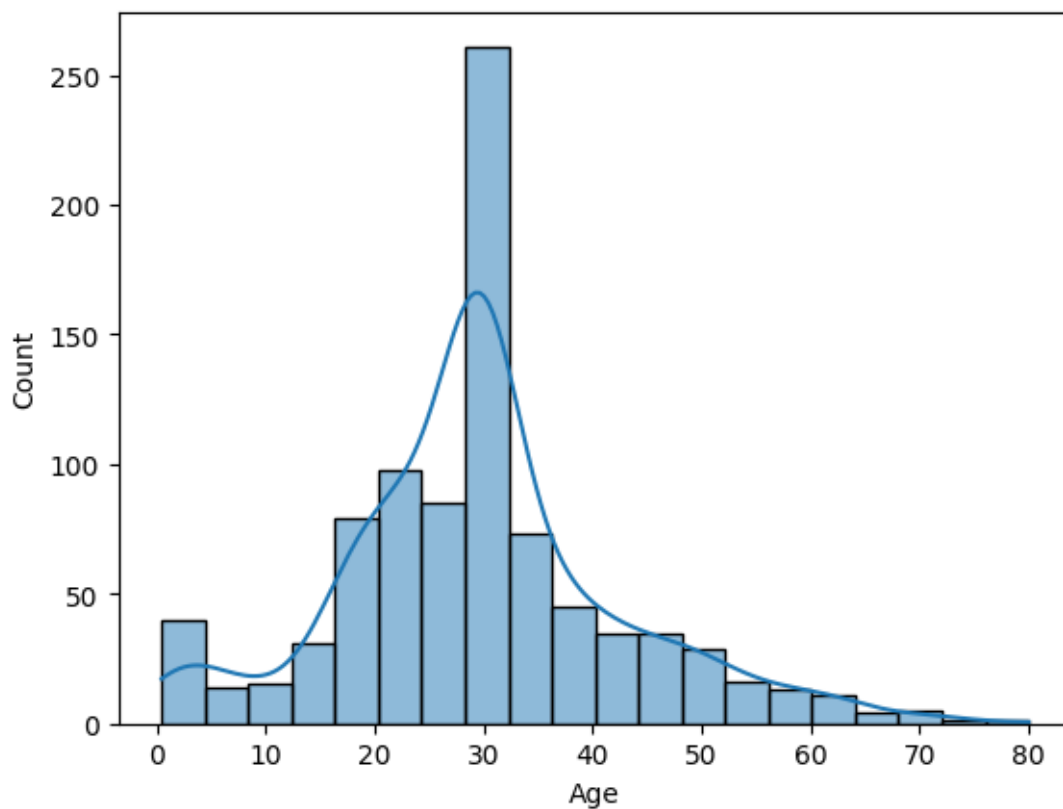
```
[34]: sns.pairplot(df)
```

[34]: <seaborn.axisgrid.PairGrid at 0x7b1e924b96f0>



```
[39]: sns.histplot(data=df,x="Age",bins=20,kde=True)
```

```
[39]: <Axes: xlabel='Age', ylabel='Count'>
```



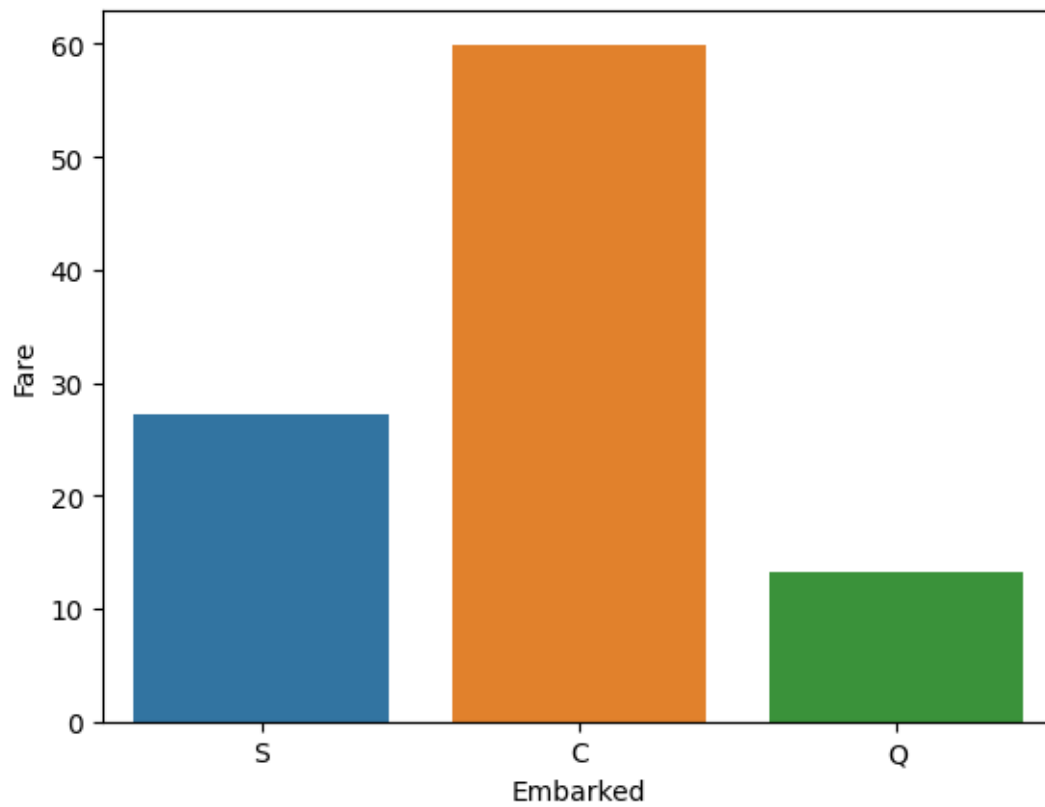
```
[36]: sns.barplot(x=df["Embarked"],y=df["Fare"],ci=None)
```

<ipython-input-36-f67c208bf54a>:1: FutureWarning:

The `ci` parameter is deprecated. Use `errorbar=None` for the same effect.

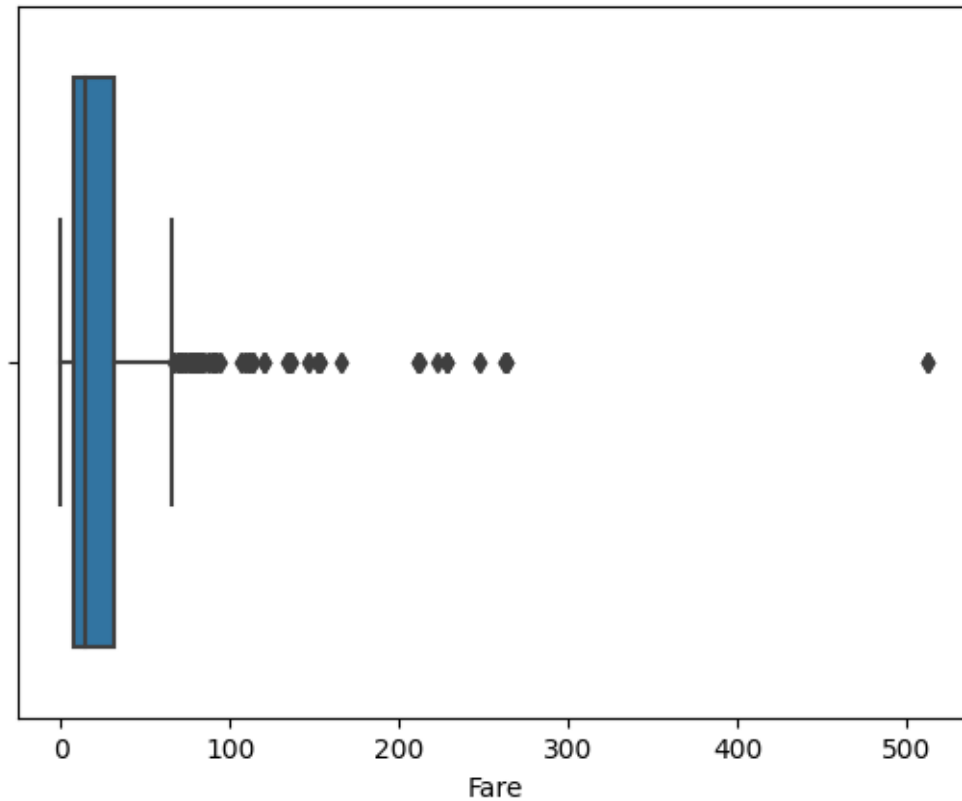
```
sns.barplot(x=df["Embarked"],y=df["Fare"],ci=None)
```

```
[36]: <Axes: xlabel='Embarked', ylabel='Fare'>
```



```
[37]: sns.boxplot(x="Fare",data=df)
```

```
[37]: <Axes: xlabel='Fare'>
```

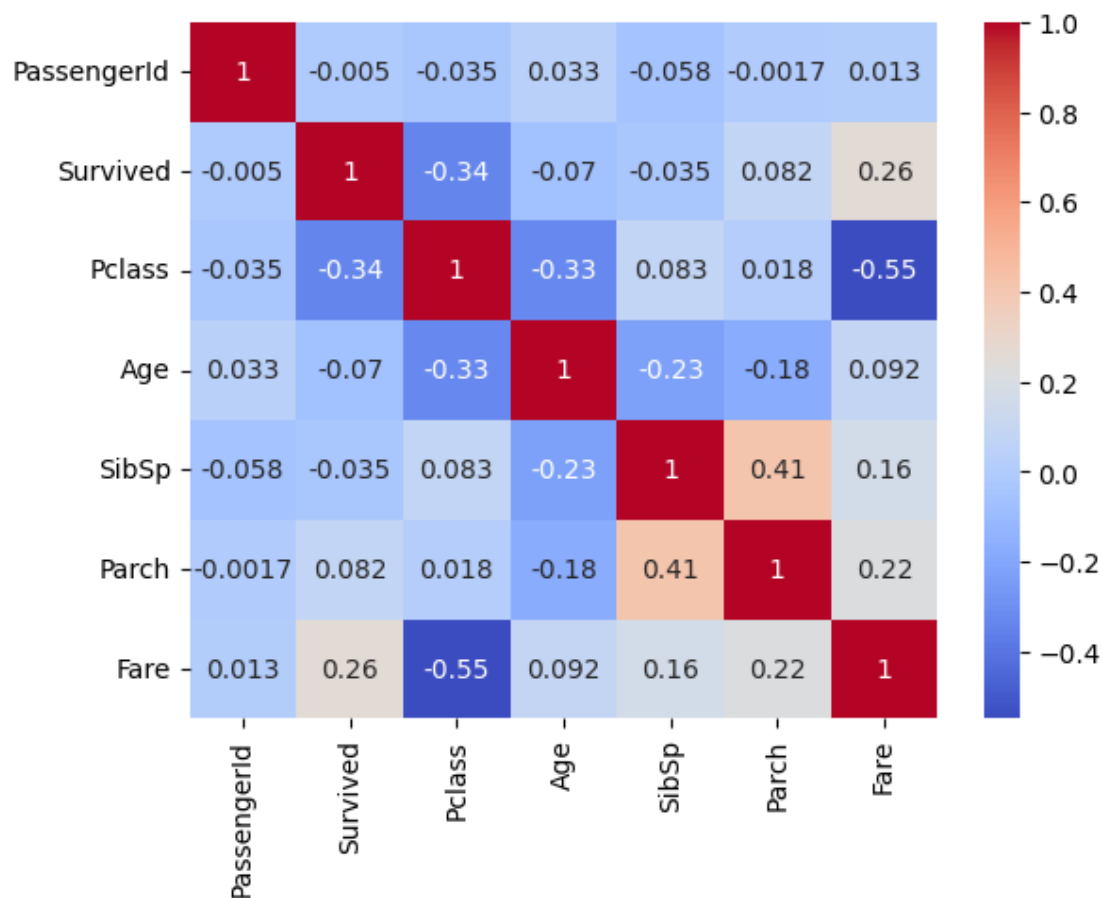


```
[38]: sns.heatmap(df.corr(),annot=True,cmap='coolwarm')
```

<ipython-input-38-407fc1d37529>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

```
sns.heatmap(df.corr(),annot=True,cmap='coolwarm')
```

```
[38]: <Axes: >
```



5 5.Outlier Detection

```
[40]: df.head()
```

```
[40]:
```

	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	1	female	35.0	1	0	53.1000	S
4	5	0	3	male	35.0	0	0	8.0500	S

```
[43]: from scipy import stats
z_scores=np.abs(stats.zscore(df["Age"]))
z_scores
```

```
[43]: 0    0.592481
      1    0.638789
      2    0.284663
```

```

3      0.407926
4      0.407926
...
886    0.207709
887    0.823344
888    0.000000
889    0.284663
890    0.177063
Name: Age, Length: 891, dtype: float64

```

```
[42]: outliers=df["Age"][z_scores>3]
```

```
[44]: outliers
```

```

[44]: 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

```

```
[45]: z_score=np.abs(stats.zscore(df["Fare"]))
      outlier=df["Fare"][z_score>3]
```

```
[46]: outlier
```

```

[46]: 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
```

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

      IQR = Q3 - Q1

      lower_bound = Q1 - 1.5 * IQR
      upper_bound = Q3 + 1.5 * IQR

      df_cleaned = df[(df["Fare"] > lower_bound) & (df["Fare"] < upper_bound)]

      print(f"Original DataFrame size: {df.shape}")
      print(f"Cleaned DataFrame size: {df_cleaned.shape}")
      df_cleaned
```

```
Original DataFrame size: (891, 9)
Cleaned DataFrame size: (775, 9)
```

```
[47]:
```

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

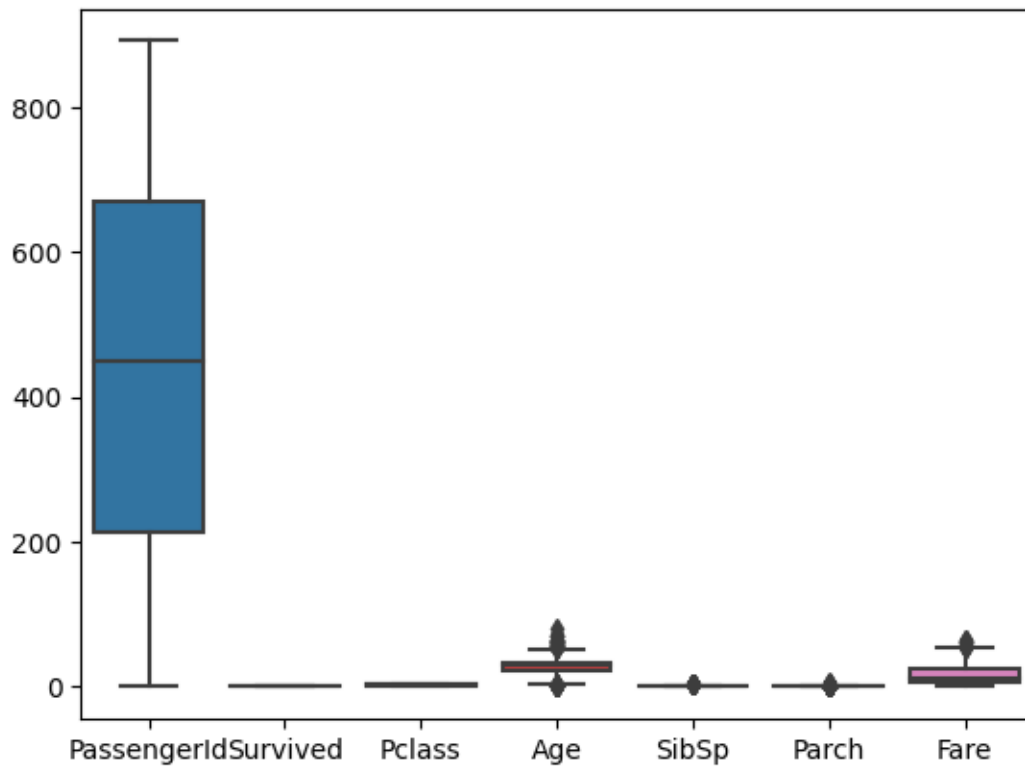
```
      Embarked
0          S
2          S
3          S
4          S
5          Q
..         ...
886         S
887         S
888         S
889         C
890         Q
```

```
[775 rows x 9 columns]
```



```
[48]: sns.boxplot(data=df_cleaned)
```

```
[48]: <Axes: >
```



6 6.Splitting Dependent and Independent variables

```
[49]: df.head()
```

```
[49]:
```

	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	1	female	35.0	1	0	53.1000	S
4	5	0	3	male	35.0	0	0	8.0500	S

```
[52]: x=df.drop(columns=["Survived"],axis=1)
```

```
[51]: x.head()
```

```
[51]:
```

	PassengerId	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
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

```
[53]: type(x)
```

```
[53]: pandas.core.frame.DataFrame
```

```
[54]: x.shape
```

```
[54]: (891, 8)
```

```
[55]: y=df["Survived"]
```

```
[56]: y.head()
```

```
[56]: 0    0
      1    1
      2    1
      3    1
      4    0
      Name: Survived, dtype: int64
```

```
[57]: type(y)
```

```
[57]: pandas.core.series.Series
```

```
[58]: y.shape
```

```
[58]: (891,)
```

7 7.Encoding

```
[59]: x.head()
```

```
[59]:   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         0        3  female  26.0     0     0    7.9250         S
3            4         1        1  female  35.0     1     0   53.1000         S
4            5         0        3   male  35.0     0     0    8.0500         S
```

```
[60]: from sklearn.preprocessing import LabelEncoder
```

```
[61]: le=LabelEncoder()
```

```
[62]: x["Embarked"]=le.fit_transform(x["Embarked"])
```

```
[63]: x.head()
```

```
[63]:
```

	PassengerId	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	1	3	male	22.0	1	0	7.2500	2
1	2	1	female	38.0	1	0	71.2833	0
2	3	3	female	26.0	0	0	7.9250	2
3	4	1	female	35.0	1	0	53.1000	2
4	5	3	male	35.0	0	0	8.0500	2

```
[64]: print(le.classes_)
```

```
['C' 'Q' 'S']
```

```
[65]: x["Sex"]=le.fit_transform(x["Sex"])
```

```
[66]: x.head()
```

```
[66]:
```

	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	0	8.0500	2

8 8.Feature Scaling.

```
[67]: from sklearn.preprocessing import StandardScaler  
sc=StandardScaler()
```

```
[68]: x[['Age', 'Fare']] = sc.fit_transform(x[['Age', 'Fare']])
```

```
[69]: x.head()
```

```
[69]:
```

	PassengerId	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	1	3	1	-0.592481	1	0	-0.502445	2
1	2	1	0	0.638789	1	0	0.786845	0
2	3	3	0	-0.284663	0	0	-0.488854	2
3	4	1	0	0.407926	1	0	0.420730	2
4	5	3	1	0.407926	0	0	-0.486337	2

9 9.Splitting Data into Train and Test.

```
[70]: from sklearn.model_selection import train_test_split
      x_train,x_test,y_train,y_test = train_test_split(x,y,test_size =0.
      ↪2,random_state =0)
```

```
[71]: print(x_train.shape,x_test.shape,y_train.shape,y_test.shape)
```

```
(712, 8) (179, 8) (712,) (179,)
```