

```
In [1]: import pandas as pd  
import numpy as np
```

```
In [2]: import matplotlib.pyplot as plt  
import seaborn as sns
```

```
In [6]: #READING THE DATASET  
data = pd.read_csv("C:/Users/ashwa/Downloads/archive/train.csv")
```

```
In [7]: #CONVERTING TO DATAFRAME  
df=pd.DataFrame(data)
```

```
In [8]: #return first 5 rows of the dataframe
print(df.head)
```

```
<bound method NDFrame.head of
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  Sib
Sp \
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]>
```

In [9]: *#Displays the no. of obs and features*
df.shape

Out[9]: (891, 12)

In [10]: *#displays the last 5 rows of the dataframe*
df.tail()

Out[10]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00	
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00	
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.45	
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00	
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75	

In [11]: 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
```

```
In [13]: #will return the no. of missing records in each column  
df.isnull().sum()
```

```
Out[13]: 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
```

```
In [14]: #fill missing value/null values with 0
df.fillna(0)
```

Out[14]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.25
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	1	0	PC 17599	71.28
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.92
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.10
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.05
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	0.0	1	2	W./C. 6607	23.45
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75

891 rows × 12 columns



3. Summary Statistics

```
In [15]: df.describe()
```

```
Out[15]:
```

	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	32.204200
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	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	668.500000	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

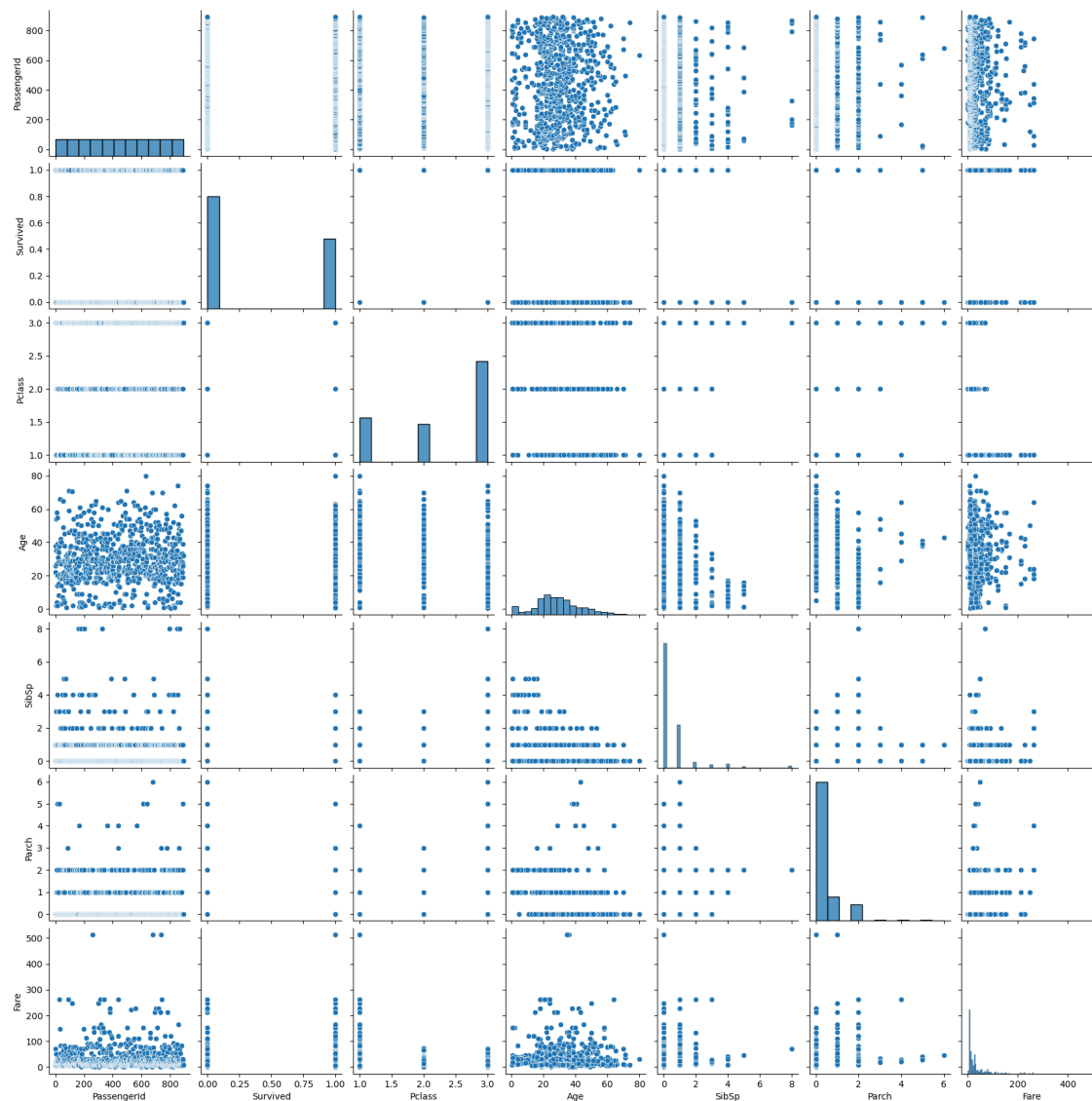
3.1. Pair Plot: Showing relationship between two categorical values

```
In [17]: # Check the column names in the DataFrame
print(data.columns)
```

```
Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
       'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],
      dtype='object')
```

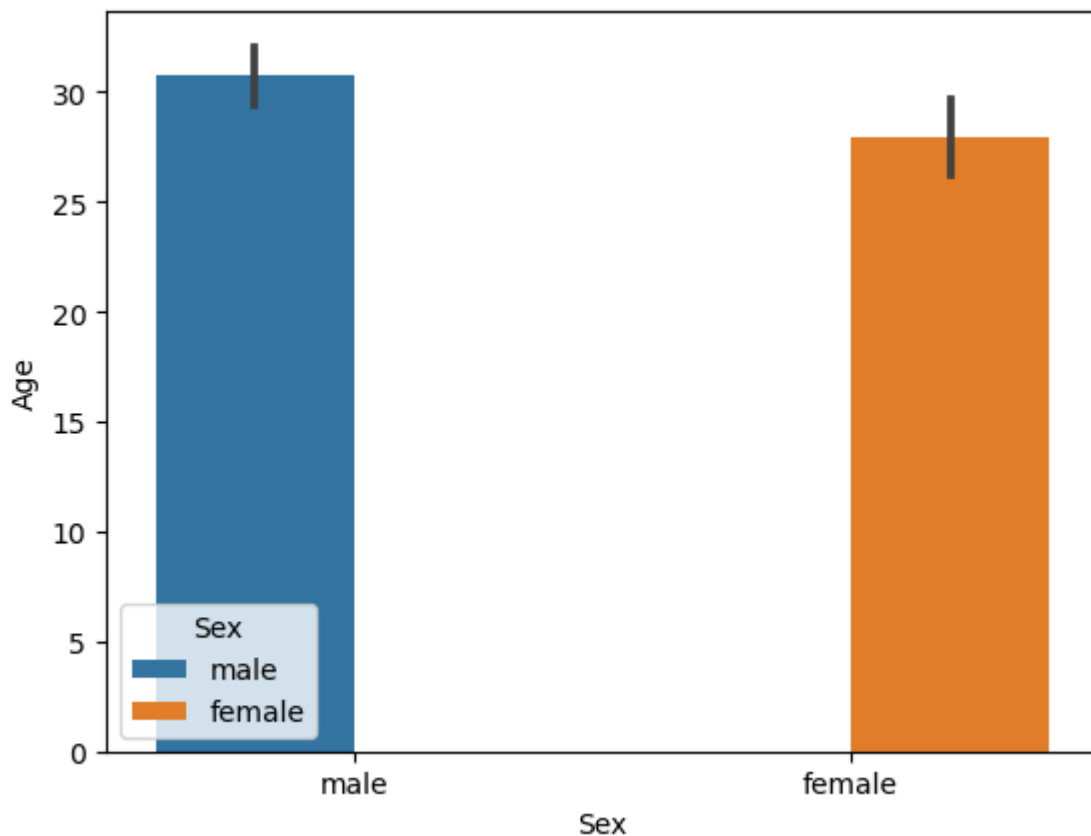
```
In [20]: # Dropping the specified columns and plotting the pairplot
plt.figure(figsize=(13, 17))
sns.pairplot(data=data.drop(['Sex', 'Embarked'], axis=1))
plt.show()
```

<Figure size 1300x1700 with 0 Axes>



3.2 Bar Plot: showing relationship between categorical variables and continuous variables

```
In [21]: sns.barplot(x='Sex',y='Age',data=df,hue='Sex')
plt.show()
```



3.3. Heatmap: showing correlation between variables

```
In [23]: dfs = df.drop(['Name','Cabin','Sex','Ticket','Embarked'],axis=1)
```

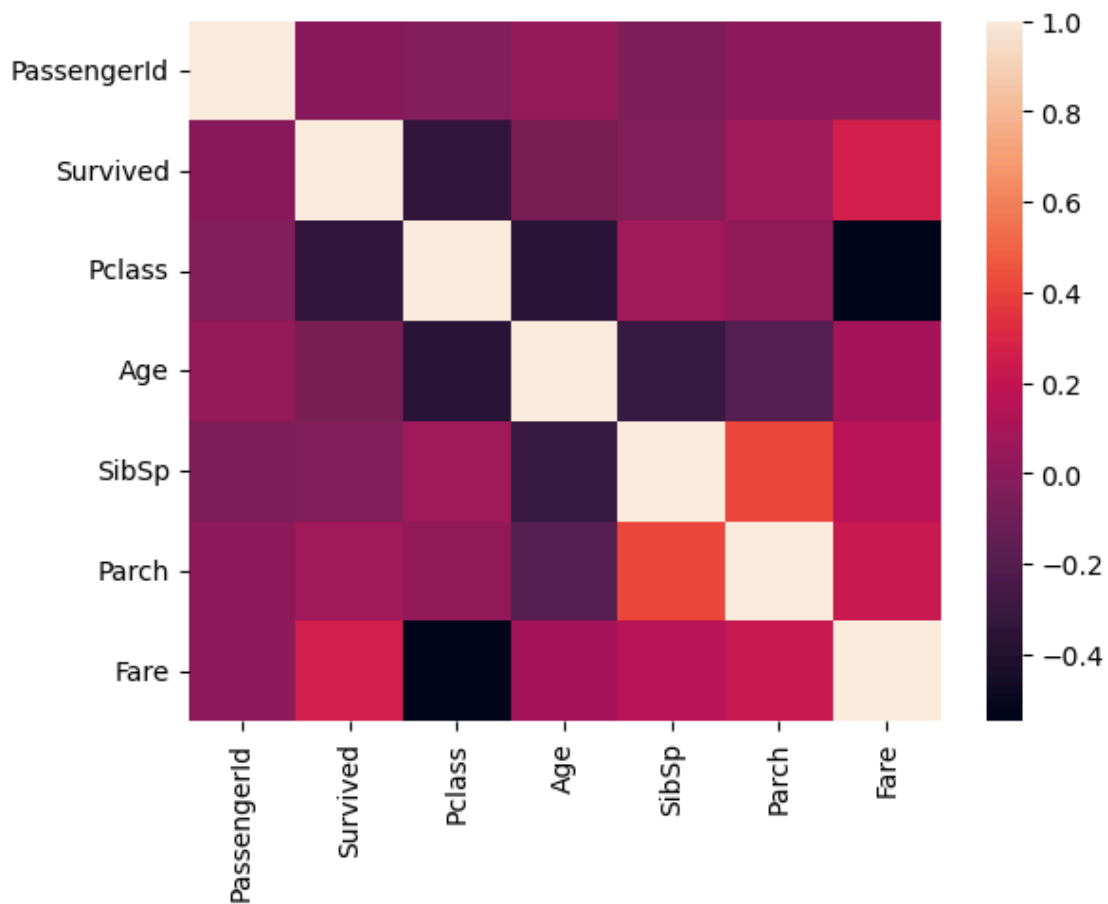
```
In [24]: dfs
```

```
Out[24]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
0	1	0	3	22.0	1	0	7.2500
1	2	1	1	38.0	1	0	71.2833
2	3	1	3	26.0	0	0	7.9250
3	4	1	1	35.0	1	0	53.1000
4	5	0	3	35.0	0	0	8.0500
...
886	887	0	2	27.0	0	0	13.0000
887	888	1	1	19.0	0	0	30.0000
888	889	0	3	NaN	1	2	23.4500
889	890	1	1	26.0	0	0	30.0000
890	891	0	3	32.0	0	0	7.7500

891 rows × 7 columns


```
In [25]: sns.heatmap(dfs.corr())  
plt.show()
```



```
In [ ]:
```