

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
import warnings
warnings.filterwarnings("ignore")
%matplotlib inline
```

```
In [3]: #Load the dataset
#dataset used https://www.kaggle.com/competitions/titanic/data
data = pd.read_csv('train.csv')
```

In [4]: data

Out[4]:

	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	7.2500
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	1	0	PC 17599	71.2833
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500


891 rows × 12 columns



In [5]: data.head()

Out[5]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	I
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	I
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	I



In [6]: data.head(10)

Out[6]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	I
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	I
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	I
5	6	0	3	Moran, Mr. James	male	NaN	0	0	330877	8.4583	I
6	7	0	1	McCarthy, Mr. Timothy J	male	54.0	0	0	17463	51.8625	
7	8	0	3	Palsson, Master. Gosta Leonard	male	2.0	3	1	349909	21.0750	I
8	9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27.0	0	2	347742	11.1333	I
9	10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14.0	1	0	237736	30.0708	I



```
In [7]: data.tail()
```

Out[7]:

	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	NaN
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00	B42
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.45	NaN
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00	C14E
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75	NaN

In [8]: data.tail(10)

Out[8]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
881	882	0	3	Markun, Mr. Johann	male	33.0	0	0	349257	7.8958
882	883	0	3	Dahlberg, Miss. Gerda Ulrika	female	22.0	0	0	7552	10.5167
883	884	0	2	Banfield, Mr. Frederick James	male	28.0	0	0	C.A./SOTON 34068	10.5000
884	885	0	3	Sutehall, Mr. Henry Jr	male	25.0	0	0	SOTON/OQ 392076	7.0500
885	886	0	3	Rice, Mrs. William (Margaret Norton)	female	39.0	0	5	382652	29.1250
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500



In [9]: *# Data Preprocessing*

```
#display information about data set  
data.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 [10]: data.columns

Out[10]: Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
 'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],
 dtype='object')

In [11]: data.describe()

Out[11]:

	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
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

```
In [12]: data.describe(include='all')
```

```
Out[12]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch
count	891.000000	891.000000	891.000000	891	891	714.000000	891.000000	891.000000
unique	NaN	NaN	NaN	891	2	NaN	NaN	NaN
top	NaN	NaN	NaN	Braund, Mr. Owen Harris	male	NaN	NaN	NaN
freq	NaN	NaN	NaN	1	577	NaN	NaN	NaN
mean	446.000000	0.383838	2.308642	NaN	NaN	29.699118	0.523008	0.381594
std	257.353842	0.486592	0.836071	NaN	NaN	14.526497	1.102743	0.806057
min	1.000000	0.000000	1.000000	NaN	NaN	0.420000	0.000000	0.000000
25%	223.500000	0.000000	2.000000	NaN	NaN	20.125000	0.000000	0.000000
50%	446.000000	0.000000	3.000000	NaN	NaN	28.000000	0.000000	0.000000
75%	668.500000	1.000000	3.000000	NaN	NaN	38.000000	1.000000	0.000000
max	891.000000	1.000000	3.000000	NaN	NaN	80.000000	8.000000	6.000000

```
In [13]: data.shape
```

```
Out[13]: (891, 12)
```

```
In [14]: data.dtypes
```

```
Out[14]: PassengerId      int64
Survived      int64
Pclass        int64
Name          object
Sex           object
Age          float64
SibSp         int64
Parch         int64
Ticket        object
Fare          float64
Cabin         object
Embarked      object
dtype: object
```

```
In [15]: data.index
```

```
Out[15]: RangeIndex(start=0, stop=891, step=1)
```




```
In [16]: # Check The Missing Value in data using pandas isnull()
data.isnull()
```

```
Out[16]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	False	False	False	False	False	False	False	False	False	False	True	
1	False	False	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	False	True
3	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	False	True
...
886	False	False	False	False	False	False	False	False	False	False	False	True
887	False	False	False	False	False	False	False	False	False	False	False	False
888	False	False	False	False	False	True	False	False	False	False	False	True
889	False	False	False	False	False	False	False	False	False	False	False	False
890	False	False	False	False	False	False	False	False	False	False	False	True

891 rows × 12 columns



```
In [17]: data.isnull().any()
```

```
Out[17]: 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
```

```
In [18]: data.isnull().sum()
```

```
Out[18]: 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 [19]: data['Age'].fillna(data['Age'].mean(),inplace=True)
```

```
In [20]: data.isnull().sum()
```

```
Out[20]: 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
```

Visualization

```
In [21]: data['Name']
```

```
Out[21]: 0      Braund, Mr. Owen Harris
1  Cumings, Mrs. John Bradley (Florence Briggs Th...
2      Heikkinen, Miss. Laina
3  Futrelle, Mrs. Jacques Heath (Lily May Peel)
4      Allen, Mr. William Henry
...
886  Montvila, Rev. Juozas
887  Graham, Miss. Margaret Edith
888  Johnston, Miss. Catherine Helen "Carrie"
889  Behr, Mr. Karl Howell
890  Dooley, Mr. Patrick
Name: Name, Length: 891, dtype: object
```

```
In [22]: data['Sex'].value_counts()
```

```
Out[22]: male      577  
female    314  
Name: Sex, dtype: int64
```

```
In [23]: data['Ticket'].value_counts()
```

```
Out[23]: 347082      7  
CA. 2343      7  
1601          7  
3101295       6  
CA 2144        6  
..           ..  
9234          1  
19988         1  
2693          1  
PC 17612       1  
370376         1  
Name: Ticket, Length: 681, dtype: int64
```

```
In [24]: data['Cabin'].value_counts()
```

```
Out[24]: B96 B98      4  
G6          4  
C23 C25 C27    4  
C22 C26       3  
F33          3  
..           ..  
E34          1  
C7           1  
C54          1  
E36          1  
C148         1  
Name: Cabin, Length: 147, dtype: int64
```

```
In [25]: data['Embarked'].value_counts()
```

```
Out[25]: S      644  
C      168  
Q       77  
Name: Embarked, dtype: int64
```

```
In [26]: def fun1(value):  
         if (value == "male"):  
             return 1  
         else:  
             return 0
```

```
In [27]: def fun2(value):  
        if (value == 'S'):  
            return 0  
        elif (value == 'C'):  
            return 1  
        elif (value == 'Q'):  
            return 2  
        else:  
            return 0
```

```
In [28]: data["Sex"] = data["Sex"].apply(fun1)
```

```
In [29]: data["Embarked"] = data["Embarked"].apply(fun2)
```

```
In [30]: data.isnull().sum()
```

```
Out[30]: 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
```

```
In [31]: data.columns
```

```
Out[31]: Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',  
               'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],  
              dtype='object')
```

```
In [32]: data = data.drop("Cabin", axis=1)
```

```
In [33]: data.columns
```

```
Out[33]: Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',  
               'Parch', 'Ticket', 'Fare', 'Embarked'],  
              dtype='object')
```

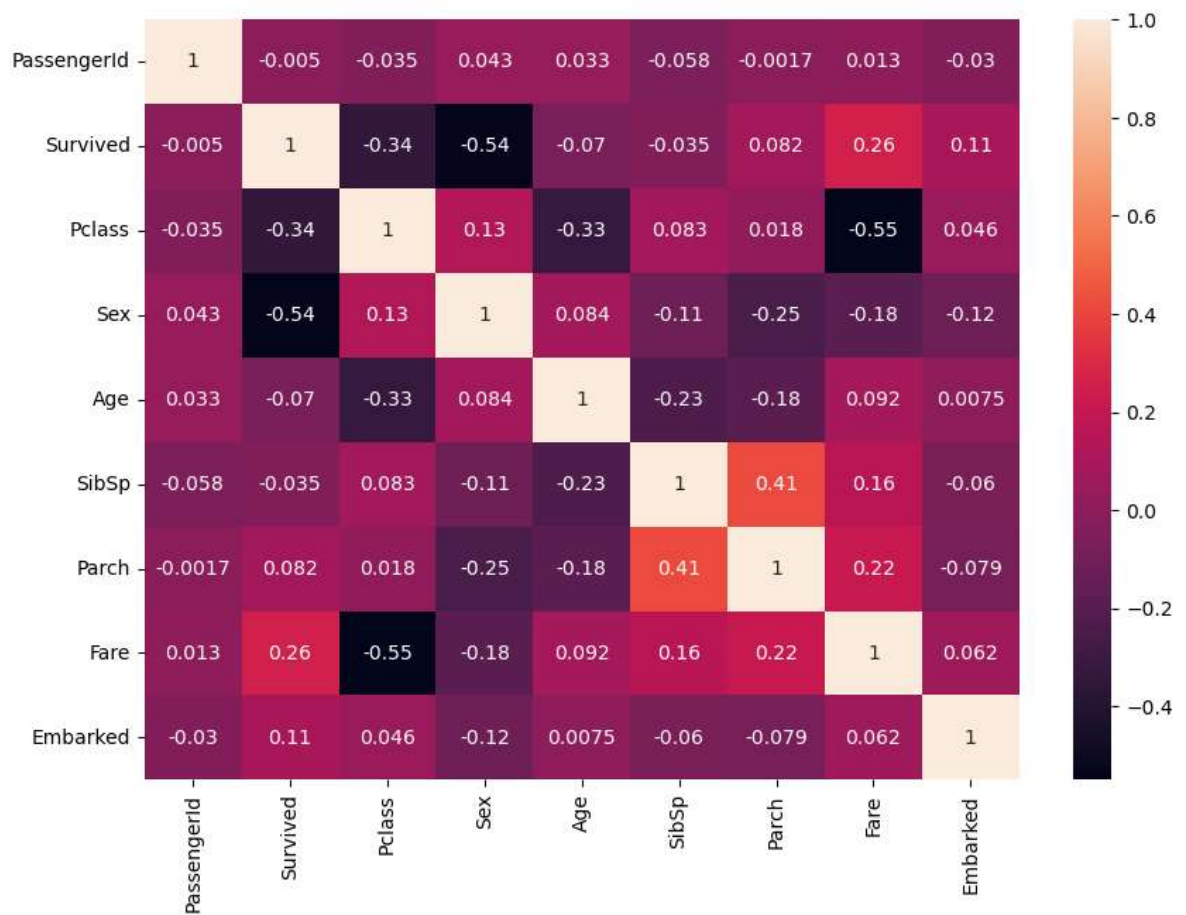
```
In [34]: data.isnull().sum()
```

```
Out[34]: PassengerId    0
Survived    0
Pclass      0
Name        0
Sex         0
Age         0
SibSp       0
Parch       0
Ticket      0
Fare        0
Embarked    0
dtype: int64
```

```
In [35]: data.shape
```

```
Out[35]: (891, 11)
```

```
In [37]: #Age has a lot of null values and is one of the attributes we need to use.
plt.figure(figsize=(10,7))
sns.heatmap(data.corr(), annot=True)
plt.show()
```



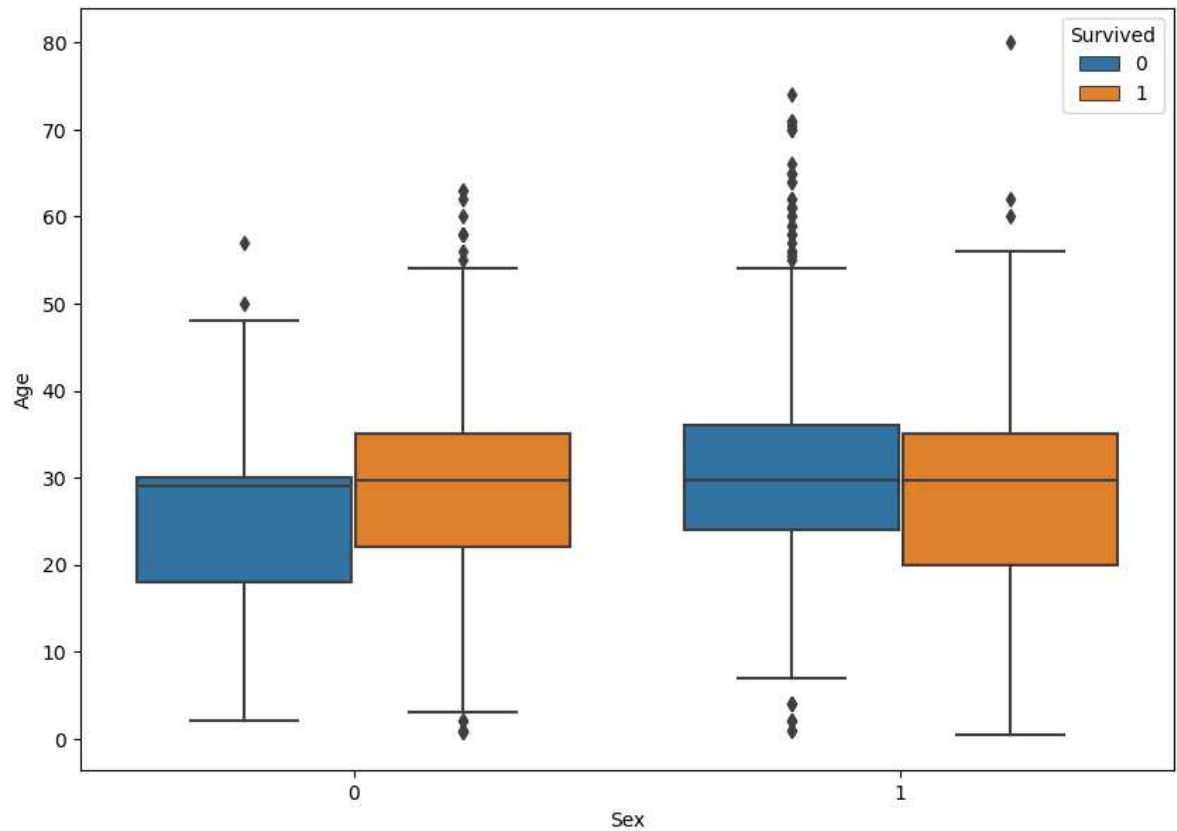
```
In [38]: #From the above corealtion matrix we can see that the attribute 'Age' is not h  
#This means we can randomly fill in the missing data for 'Age' within the vali
```

```
In [40]: px.box(data["Sex"], data["Age"], color=data["Survived"])
```

```
In [43]: # Assuming you have a DataFrame named `data` with columns "Sex", "Age", and "Survived"

# Create a boxplot with "Sex" on the x-axis, "Age" on the y-axis, and "Survived" as the hue
plt.figure(figsize=(10, 7))
box = sns.boxplot(x="Sex", y="Age", hue="Survived", data=data)

# Display the plot
plt.show()
```



In []: