

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

UsageError: Line magic function '%' not found.

In [2]: titanic=pd.read_csv("train.csv")
titanic.head()
```

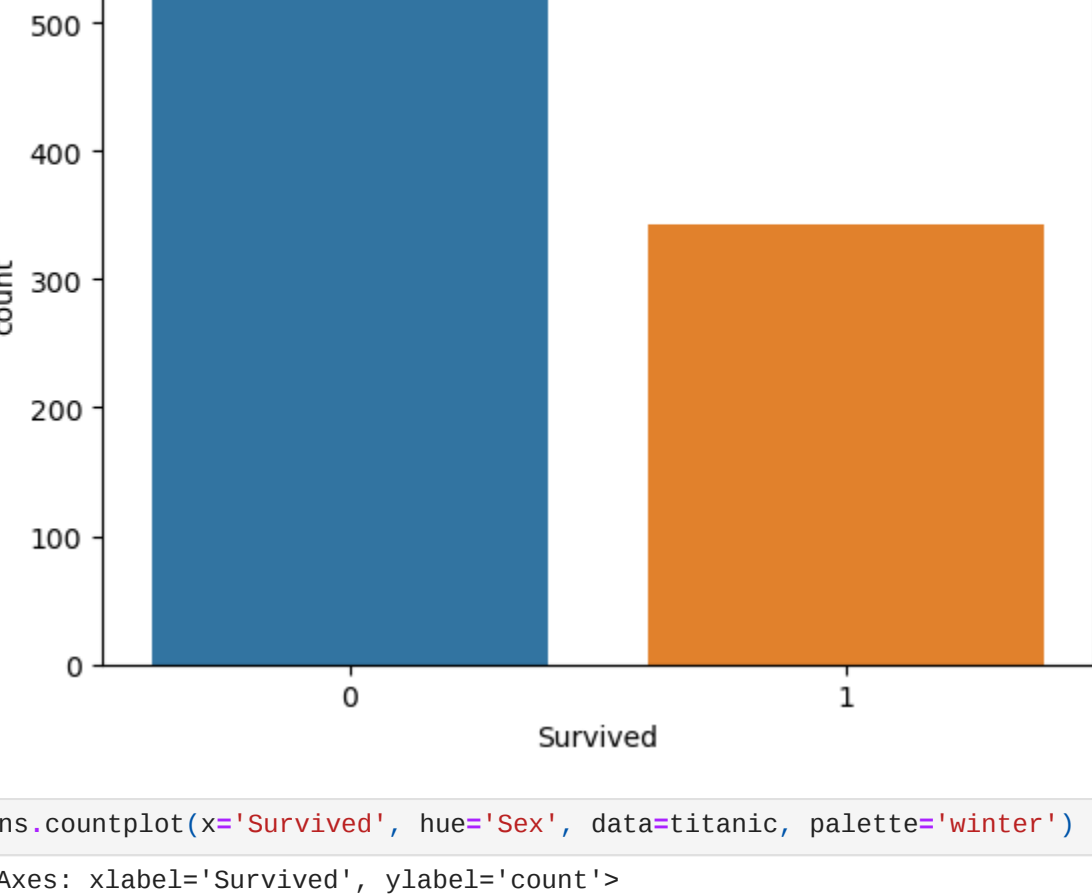
	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

```
In [3]: titanic.shape
# to obtain the no of rows and column

Out[3]: (891, 12)

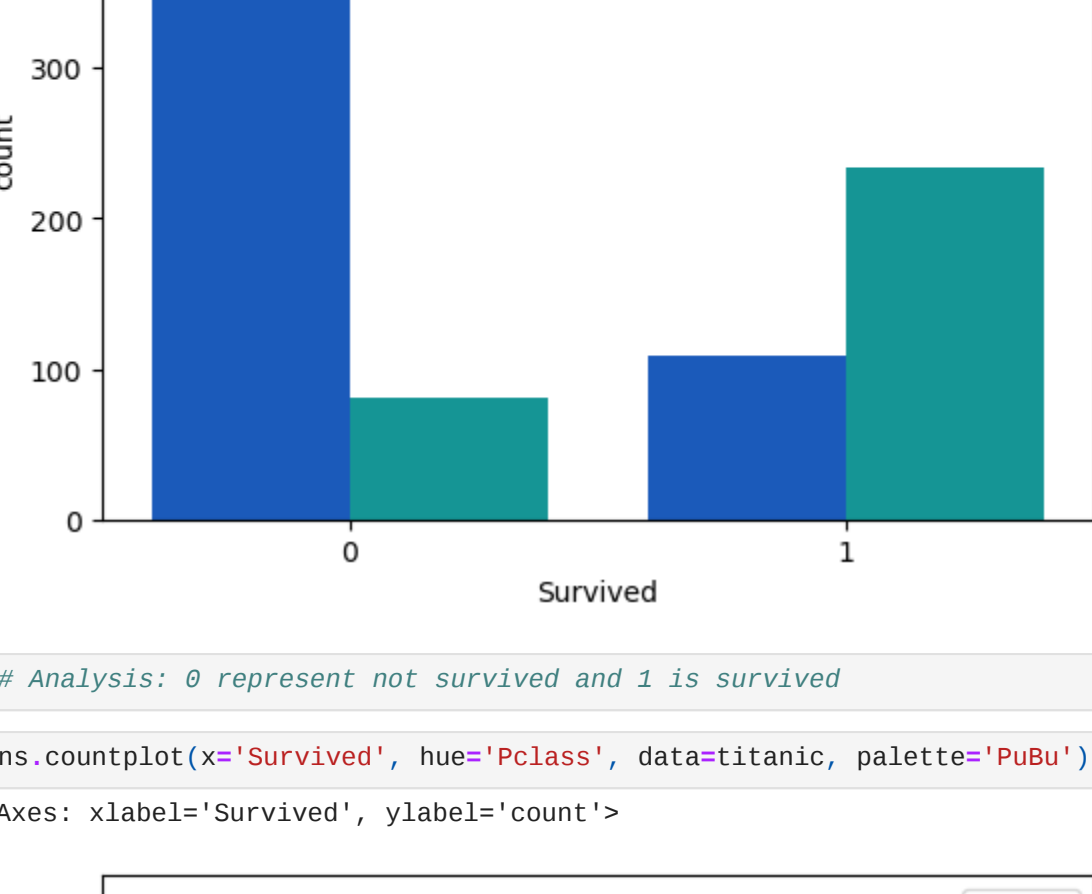
In [4]: sns.countplot(x='Survived', data=titanic)
titanic.head()
```

<Axes: xlabel='Survived', ylabel='count'>



```
In [5]: sns.countplot(x='Survived', hue='Sex', data=titanic, palette='winter')
titanic.head()
```

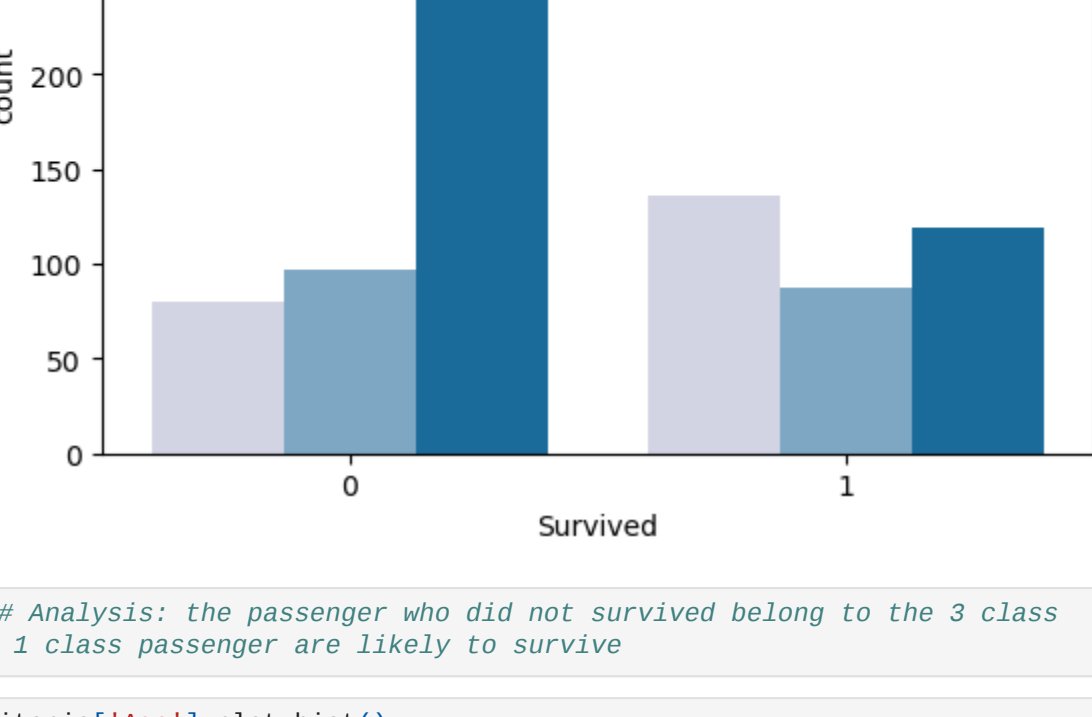
<Axes: xlabel='Survived', ylabel='count'>



```
In [ ]: ## Analysis: 0 represent not survived and 1 is survived

In [6]: sns.countplot(x='Survived', hue='Pclass', data=titanic, palette='PuBu')
titanic.head()
```

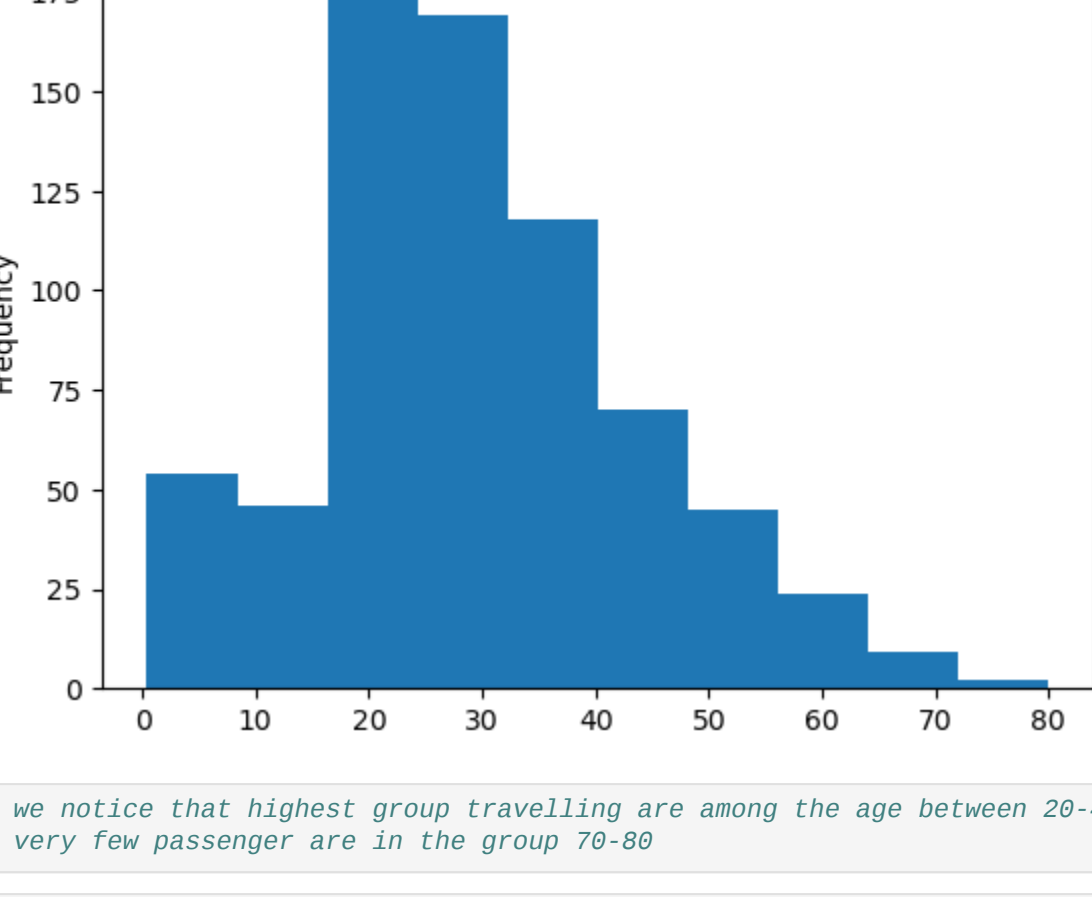
<Axes: xlabel='Survived', ylabel='count'>



```
In [ ]: ## Analysis: the passenger who did not survived belong to the 3 class
# 1 class passenger are likely to survive

In [7]: titanic['Age'].plot.hist()
```

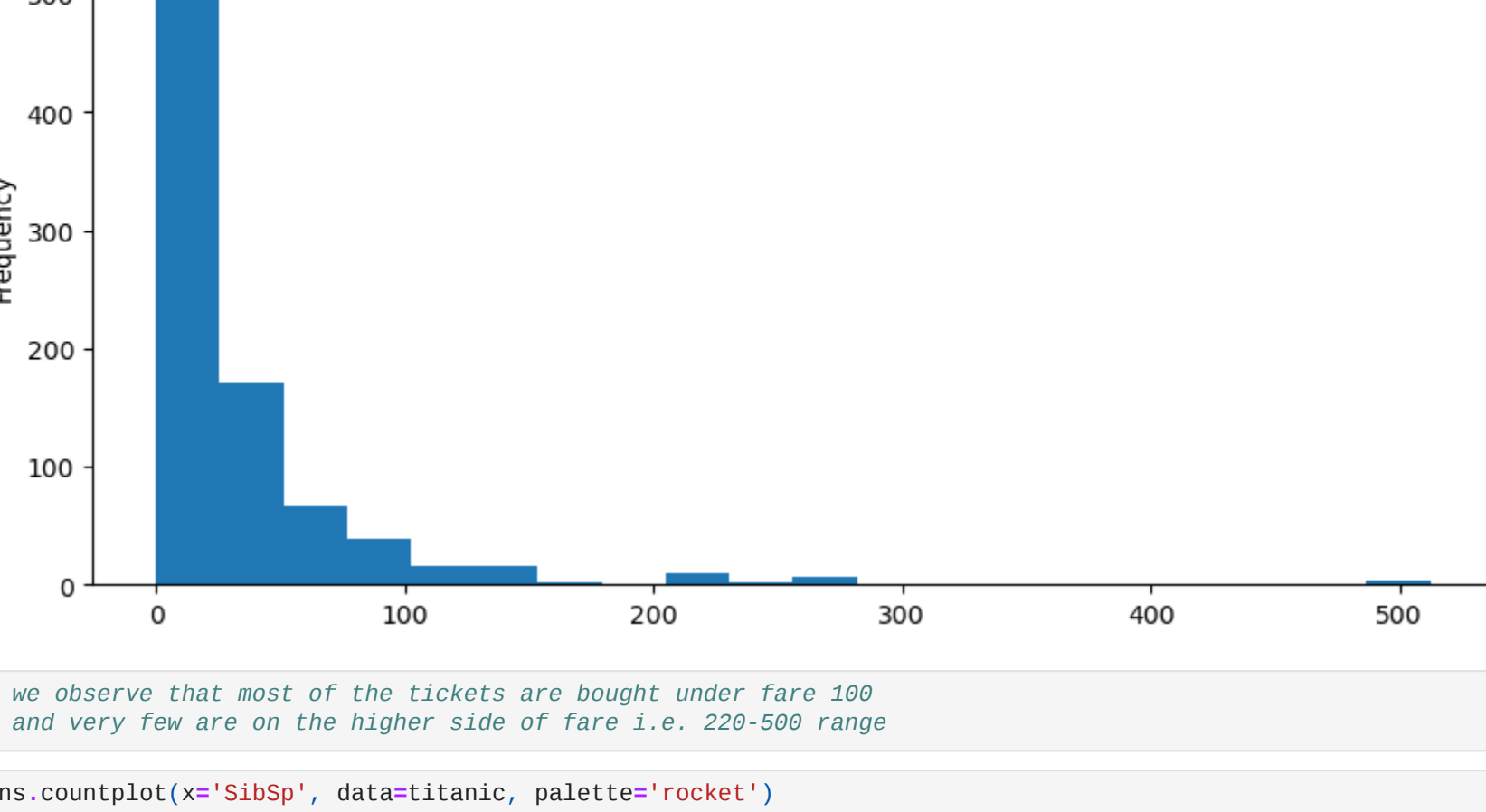
<Axes: ylabel='Frequency'>



```
In [ ]: # we notice that highest group travelling are among the age between 20-40
# very few passenger are in the group 70-80

In [8]: titanic['Fare'].plot.hist(bins=20, figsize=(10,5))
titanic.head()
```

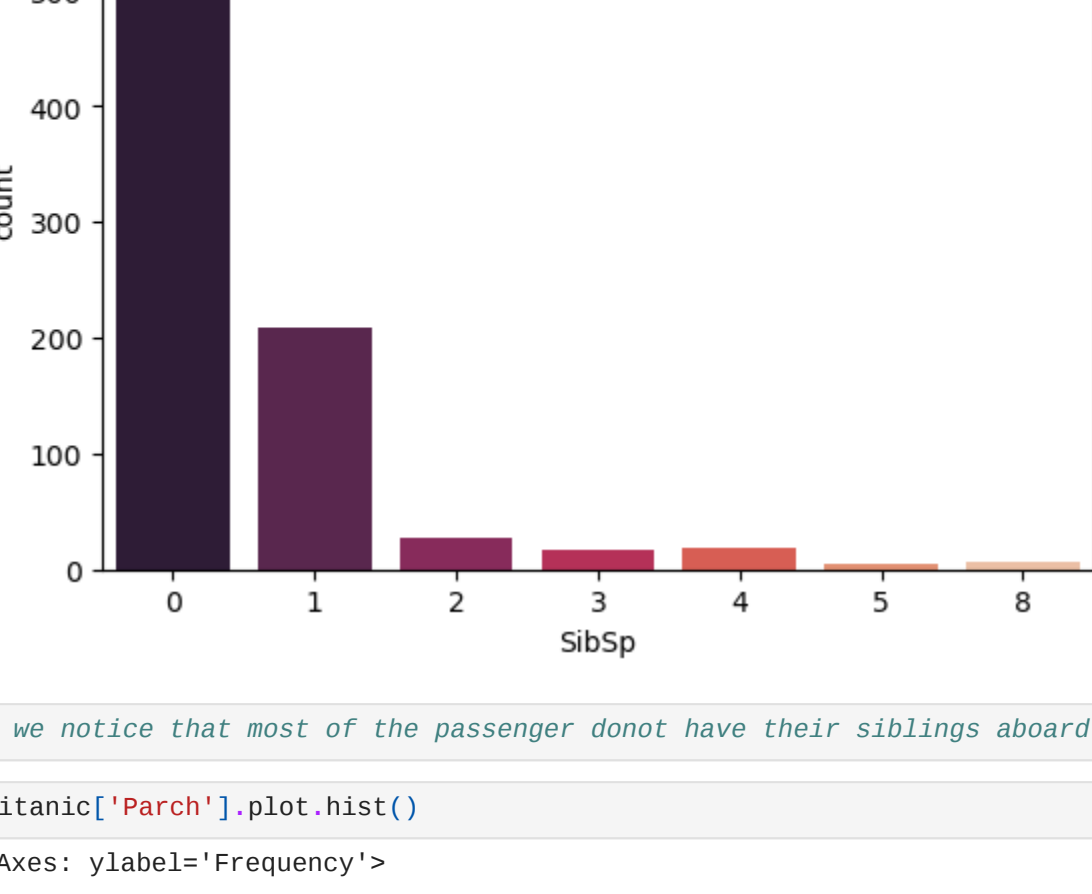
<Axes: ylabel='Frequency'>



```
In [ ]: # we observe that most of the tickets are bought under fare 100
# and very few are on the higher side of fare i.e. 220-500 range

In [9]: sns.countplot(x='SibSp', data=titanic, palette='rocket')
titanic.head()
```

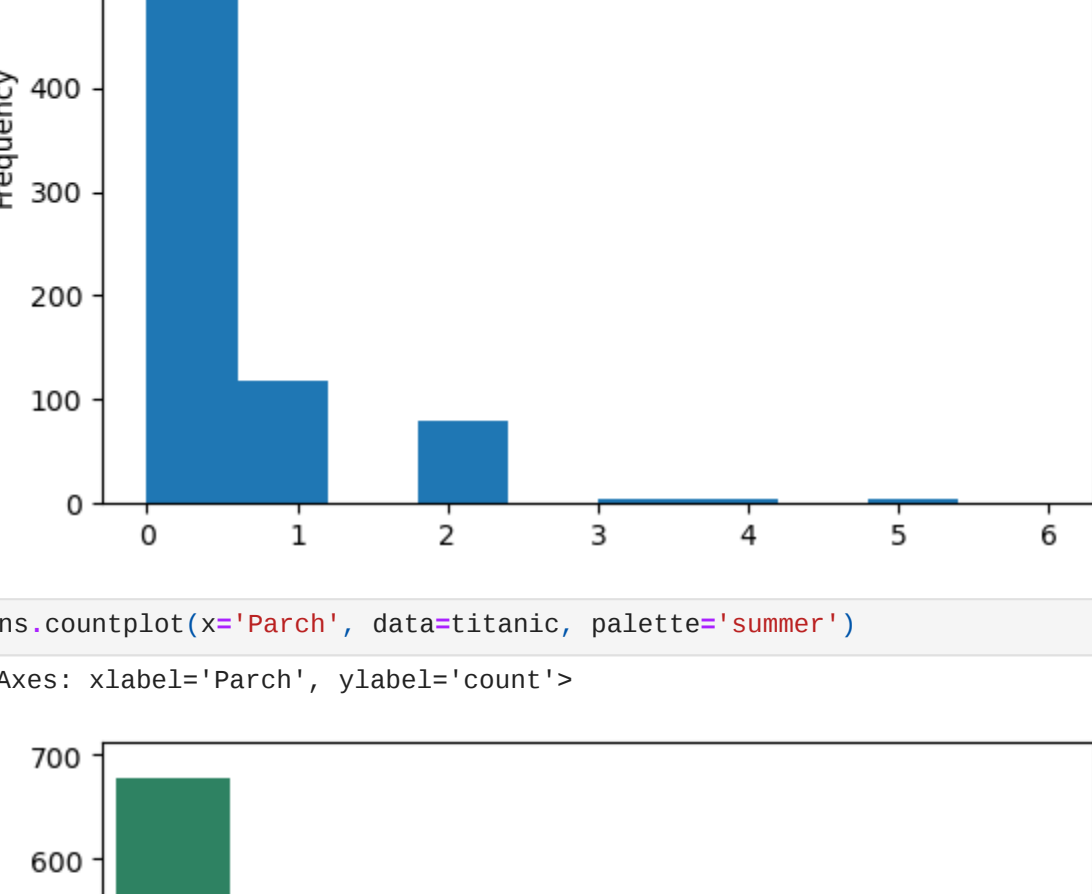
<Axes: xlabel='SibSp', ylabel='count'>



```
In [ ]: # we notice that most of the passenger donot have their siblings aboard.


In [10]: titanic['Parch'].plot.hist()
```

<Axes: ylabel='Frequency'>



```
In [11]: sns.countplot(x='Parch', data=titanic, palette='summer')
titanic.head()
```

<Axes: xlabel='Parch', ylabel='count'>



```
In [ ]: # data wrangling means cleaning the data, removing the null values
## dropping unwanted columns, adding new ones if needed.


In [12]: titanic.isnull().sum()
```

```
Out[12]: 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 [ ]: ## age and cabin has most null values and embarked too has null values
## we can plot it on heat map

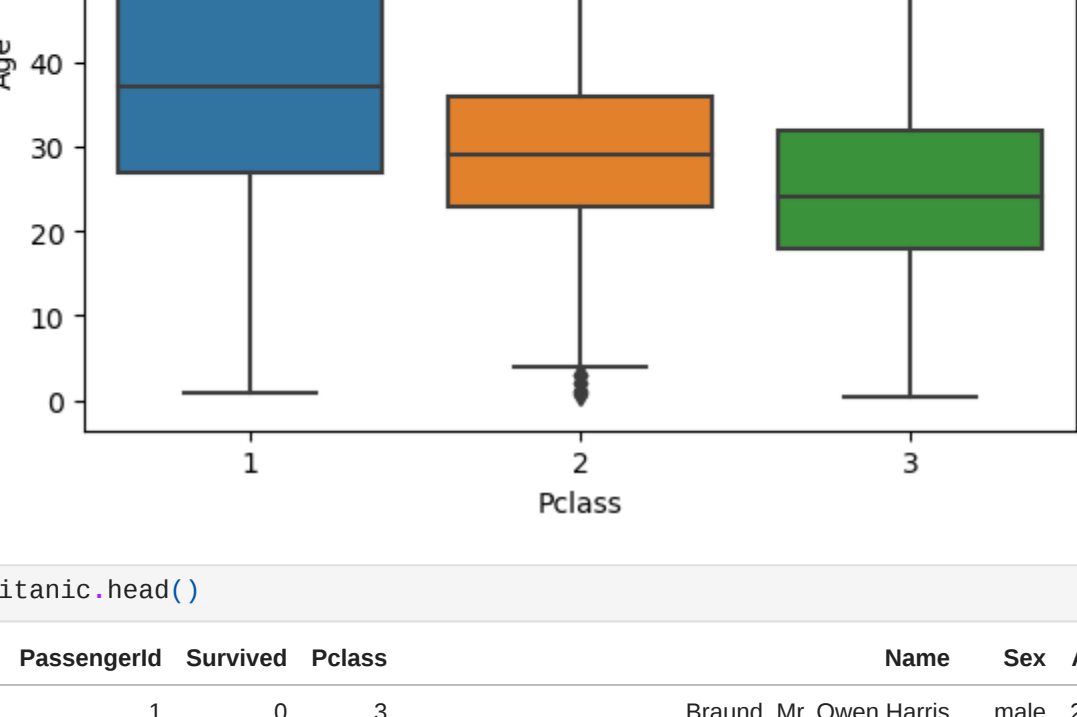
In [13]: sns.heatmap(titanic.isnull(), cmap='spring')
titanic.head()
```

<Axes: >



```
In [15]: sns.boxplot(x='Pclass', y='Age', data=titanic)
titanic.head()
```

<Axes: xlabel='Pclass', ylabel='Age'>



```
In [16]: titanic.head()
```

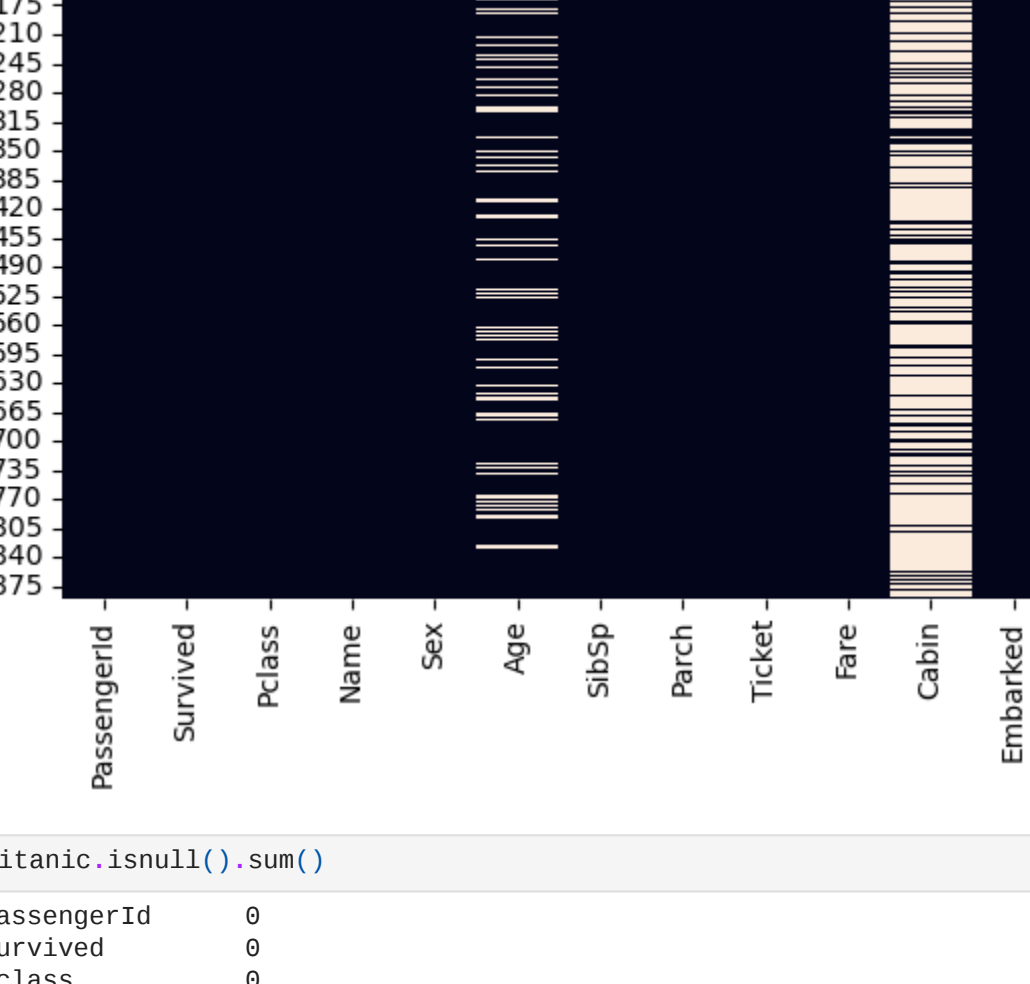
	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

```
In [17]: titanic.head(3)
```

	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	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	7.9250	NaN	S

```
In [18]: sns.heatmap(titanic.isnull(), cbar=False)
titanic.head()
```

<Axes: >



```
In [19]: titanic.isnull().sum()
```

```
Out[19]: 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 [20]: titanic.head(2)
```

	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	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C

```
In [21]: pd.get_dummies(titanic['Sex']).head()
```

	female	male
0	0	1
1	1	0
2	1	0
3	1	0
4	0	1

```
In [24]: sex=pd.get_dummies(titanic['Sex'], drop_first=True)
sex.head(3)
```

<Axes: >

	male
0	1
1	0
2	0

```
In [26]: embark=pd.get_dummies(titanic['Embarked'])
titanic.head(3)
```

<Axes: >

	C	Q	S
0	0	0	1
1	1	0	0
2	0	0	1

```
In [28]: embark=pd.get_dummies(titanic['Embarked'], drop_first=True)
titanic.head(3)
```

<Axes: >

	Q	S
0	0	1
1	0	0
2	0	1

```
In [30]: Pcl=pd.get_dummies(titanic['Pclass'], drop_first=True)
Pcl.head(3)
```

	2	3
0	0	1
1	0	0
2	0	1

```
In [ ]: # our data is converted into categorical data

In [31]: titanic=pd.concat([titanic, sex, embark,Pcl], axis=1)
titanic.head(3)
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	male	Q	S	2	3
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S	1	0	1	0	1
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C	0	0	0	0	0
2	3	1	3	Heikinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	7.9250	NaN	S	0	0	1	0	1

```
In [ ]:
In [ ]:
In [ ]:
In [ ]:
In [ ]:
```