

```
Task-01: Titanic Classification (Data Science Intern)
```

## Import Libraries

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

The Dataset

```
In [5]: train = pd.read_csv('titanic_classification.csv')
```

```
In [6]: train.head()
```

```
Out[6]:
```

	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	Heikkinen, 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

## Exploratory Data Analysis

### Missing Data

we can use seaborn to create a simple heatmap to see where we are missing data!

```
In [7]: train.isnull()
```

```
Out[7]:
```

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

891 rows x 12 columns

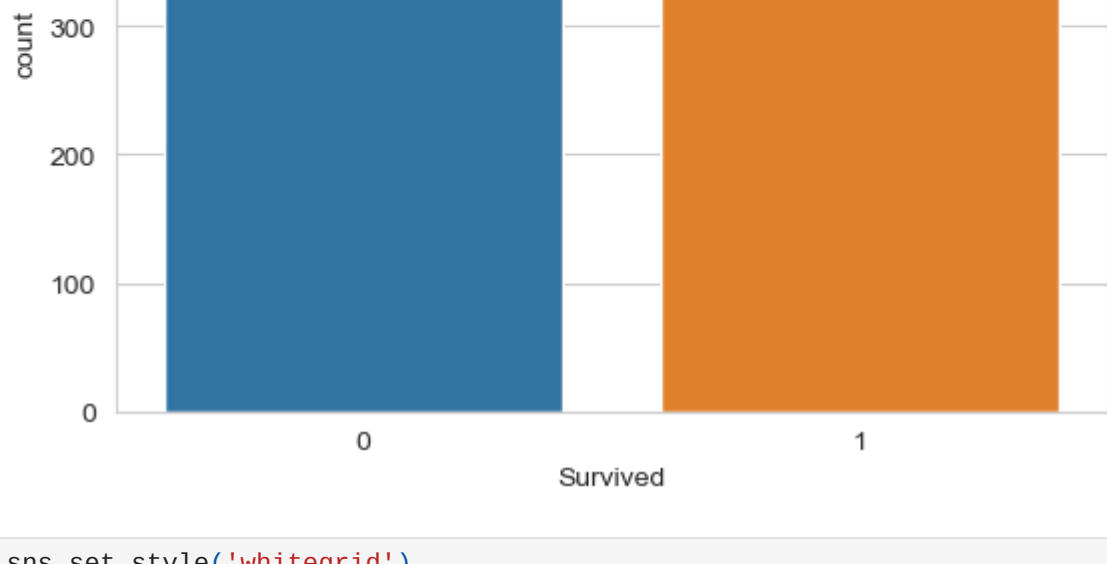
```
In [10]: sns.heatmap(train.isnull(),yticklabels=False,cbar=False,cmap='viridis')
```

```
Out[10]: >
```



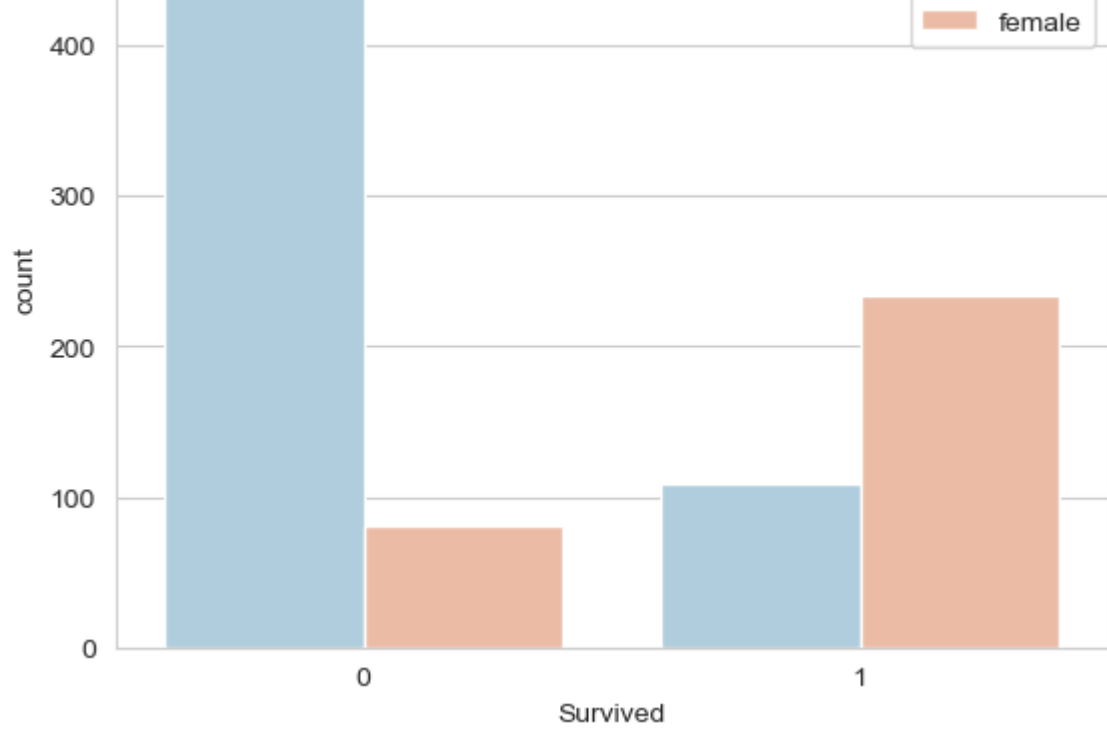
```
In [11]: sns.set_style('whitegrid')
sns.countplot(x='Survived',data=train)
```

```
Out[11]: <Axes: xlabel='Survived', ylabel='count'>
```



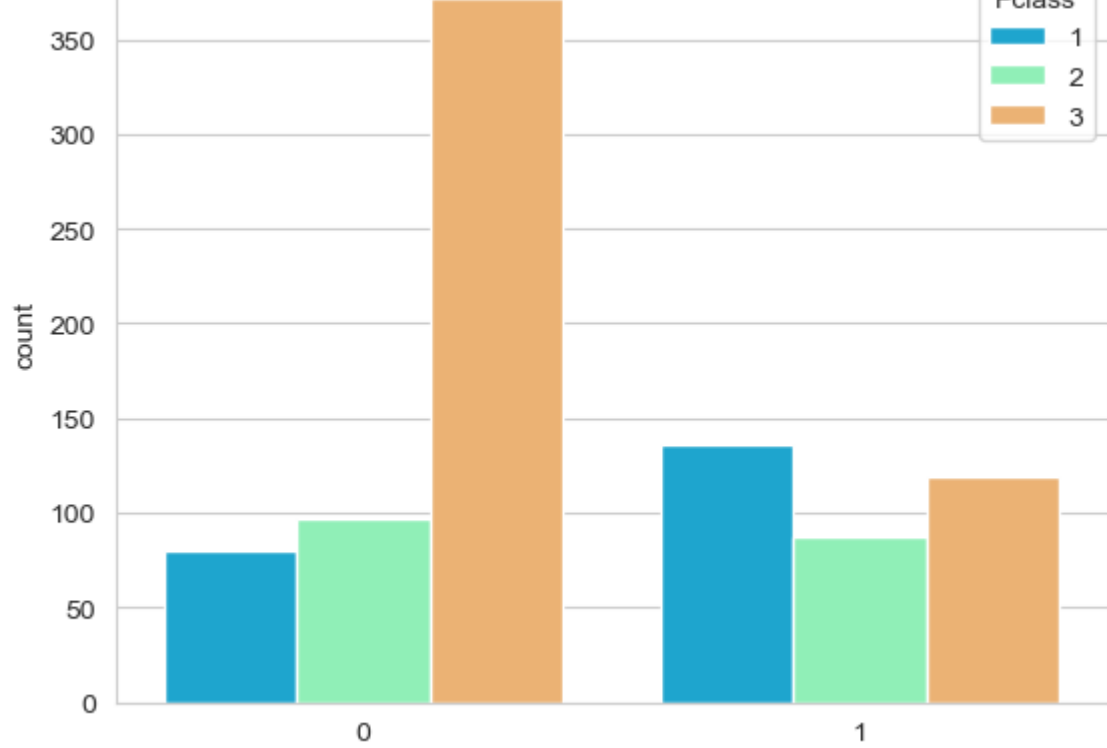
```
In [12]: sns.set_style('whitegrid')
sns.countplot(x='Survived',hue='Sex',data=train,palette='RdBu_r')
```

```
Out[12]: <Axes: xlabel='Survived', ylabel='count'>
```



```
In [13]: sns.set_style('whitegrid')
sns.countplot(x='Survived',hue='Pclass',data=train,palette='rainbow')
```

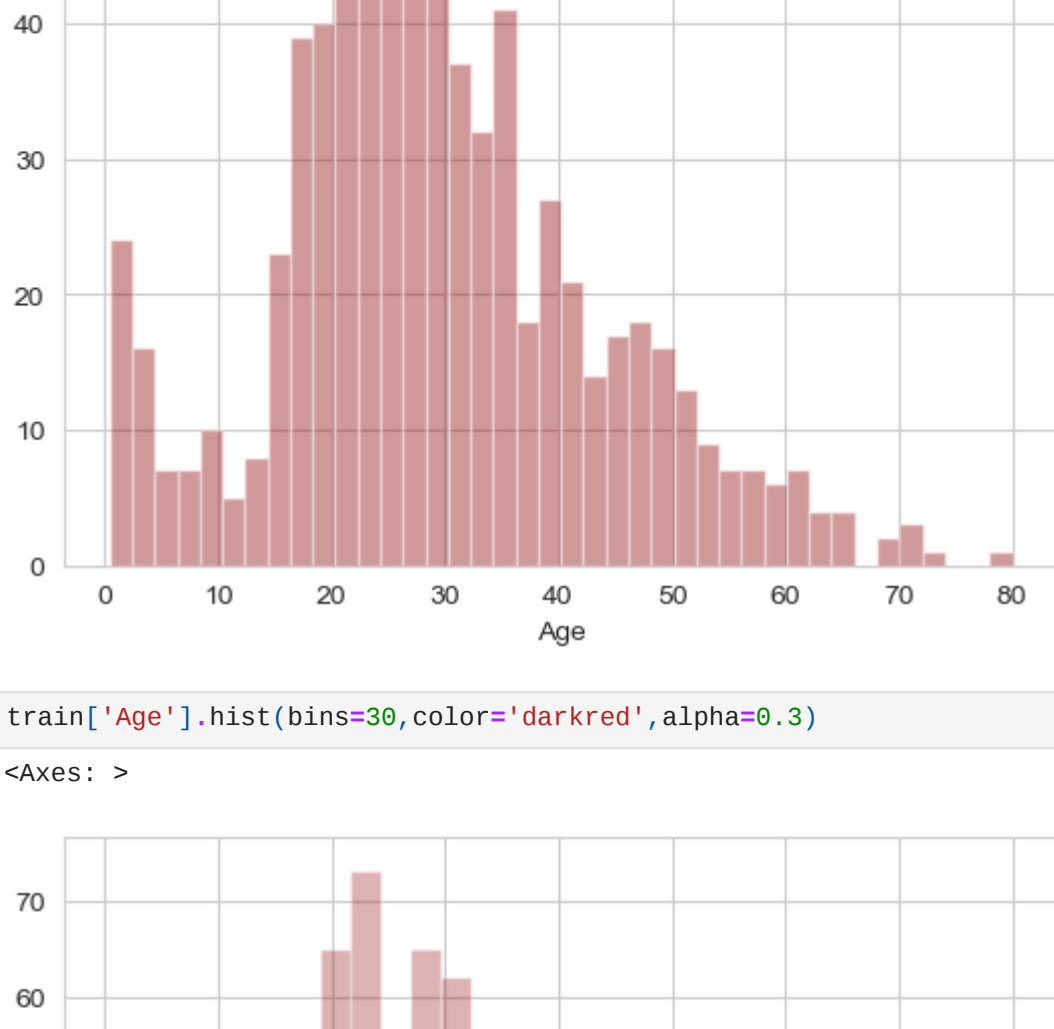
```
Out[13]: <Axes: xlabel='Survived', ylabel='count'>
```



```
In [14]: sns.distplot(train['Age'].dropna(),kde=False,color='darkred',bins=40)
```

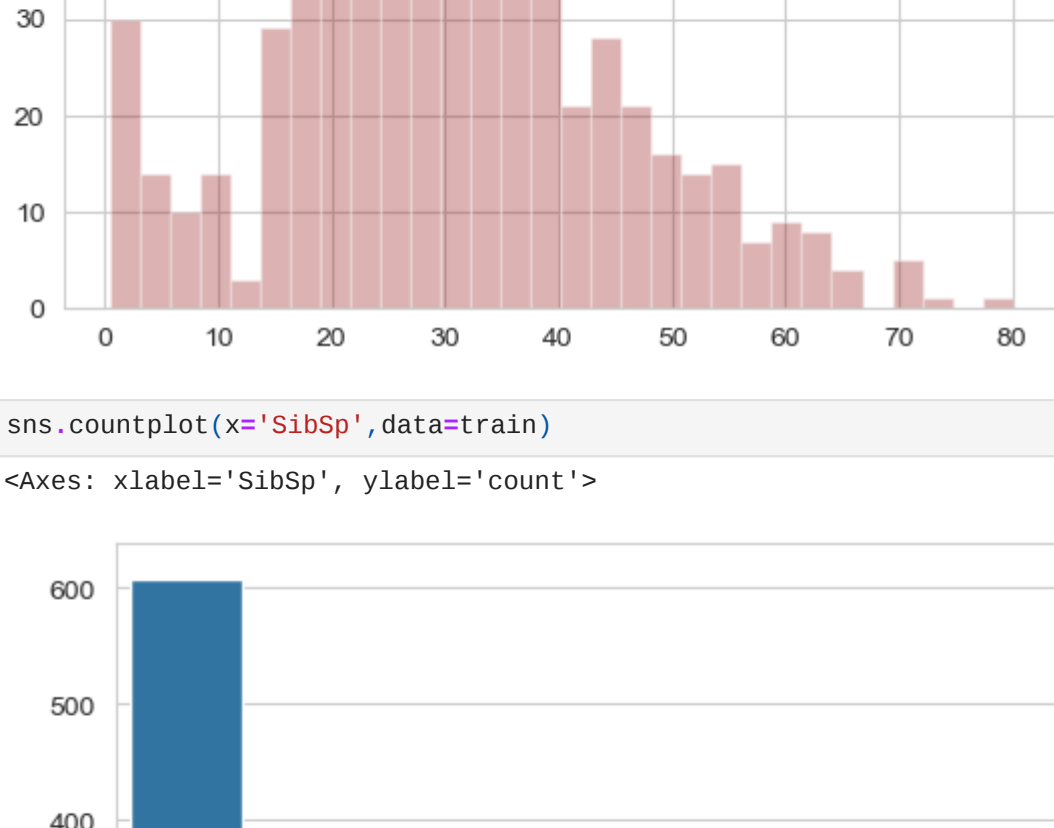
C:\Users\INDIA\AppData\Local\Temp\ipykernel\_17136\2002818437.py:1: UserWarning:  
"distplot" is a deprecated function and will be removed in seaborn v0.14.0.  
Please adapt your code to use either "displot" (a figure-level function with similar flexibility) or "histplot" (an axes-level function for histograms).  
For a guide to updating your code to use the new functions, please see  
<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
Out[14]: <Axes: xlabel='Age'>
```



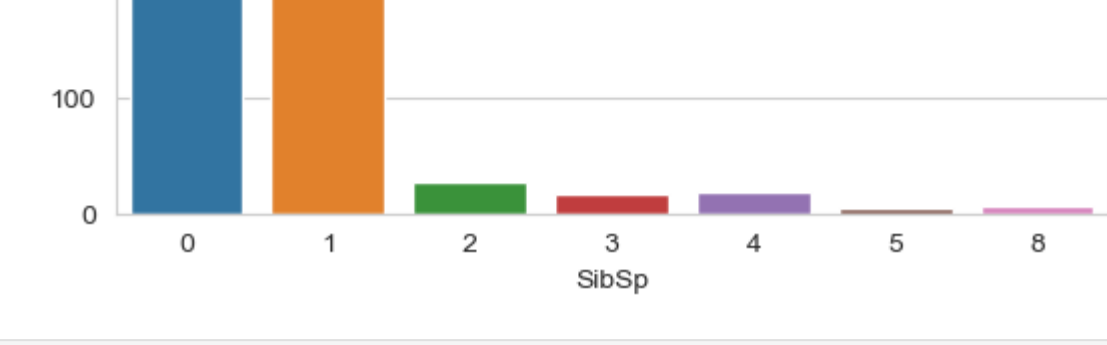
```
In [15]: train['Age'].hist(bins=30,color='darkred',alpha=0.3)
```

```
Out[15]: <Axes: >
```



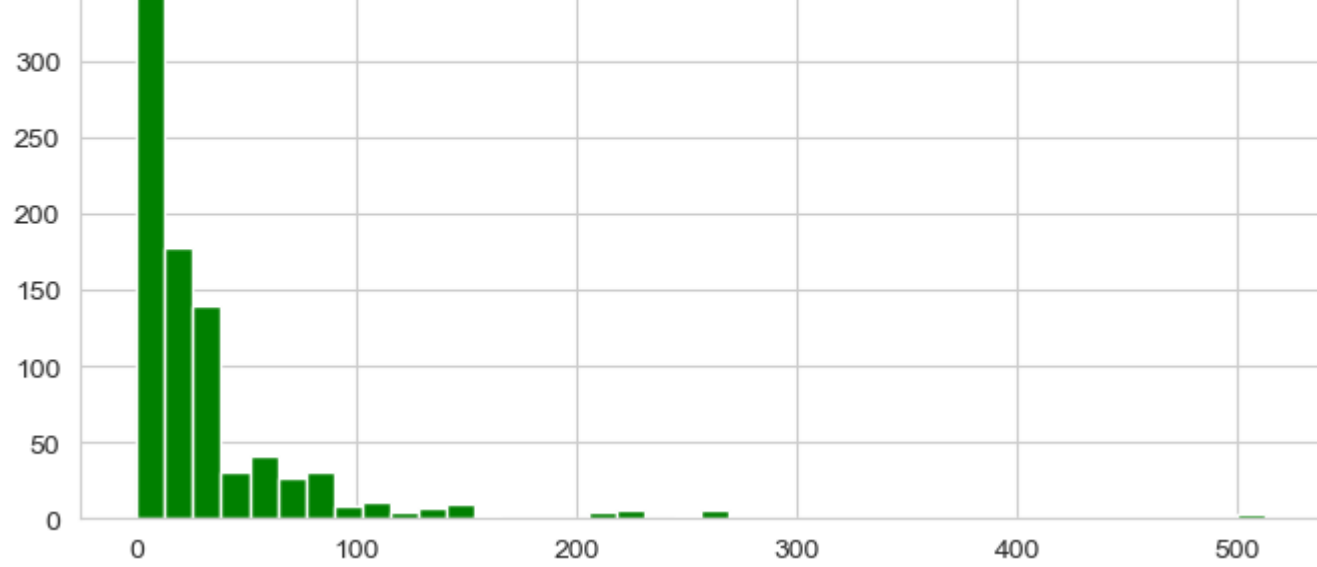
```
In [16]: sns.countplot(x='SibSp',data=train)
```

```
Out[16]: <Axes: xlabel='SibSp', ylabel='count'>
```



```
In [17]: train['Fare'].hist(color='green',bins=40,figsize=(8,4))
```

```
Out[17]: <Axes: >
```



Cufflinks for plots Let's take a quick moment to show an example of cufflinks!

```
In [ ]: import cufflinks as cf
cf.go_offline()
```

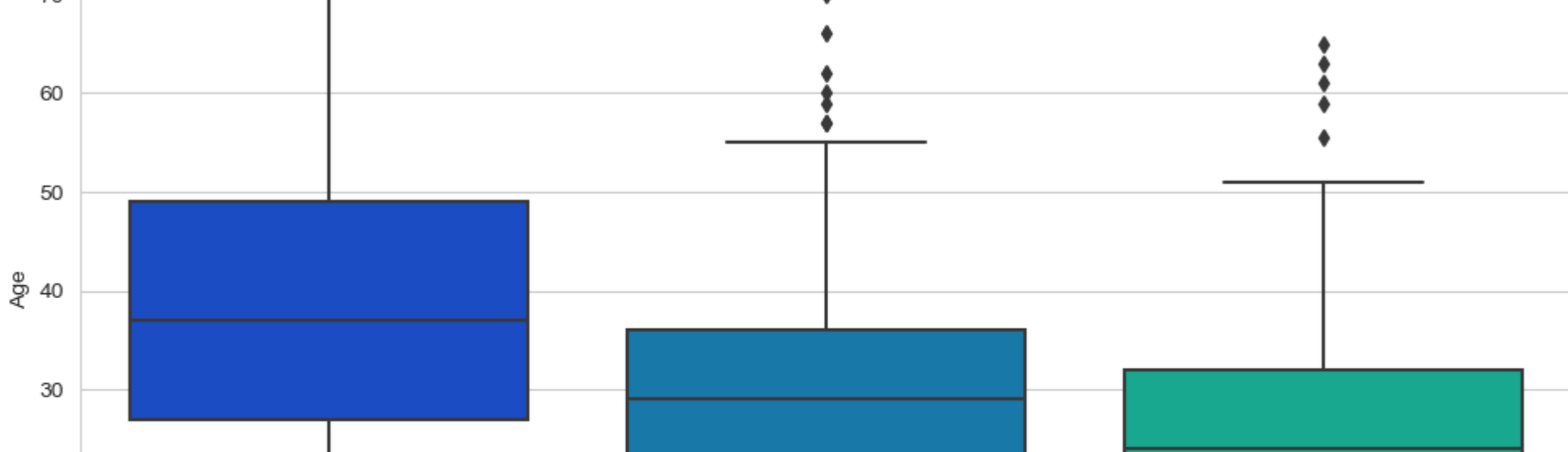
```
In [ ]: train['Fare'].iplot(kind='hist',bins=30,color='green')
```

## Data Cleaning

We want to fill in missing age data instead of just dropping the missing age data rows. One way to do this is by filling in the mean age of all the passengers (imputation). However we can be smarter about this and check the average age by passenger class. For example:

```
In [22]: plt.figure(figsize=(12, 7))
sns.boxplot(x='Pclass',y='Age',data=train,palette='winter')
```

```
Out[22]: <Axes: xlabel='Pclass', ylabel='Age'>
```



#We can see the wealthier passengers in the higher classes tend to be older, which makes sense. We'll use these average age values to impute based on Pclass for

```
In [23]: def impute_age(cols):
Age = cols[0]
Pclass = cols[1]

if pd.isnull(Age):

    if Pclass == 1:
        return 37

    elif Pclass == 2:
        return 29

    else:
        return 24

    else:
        return Age
```

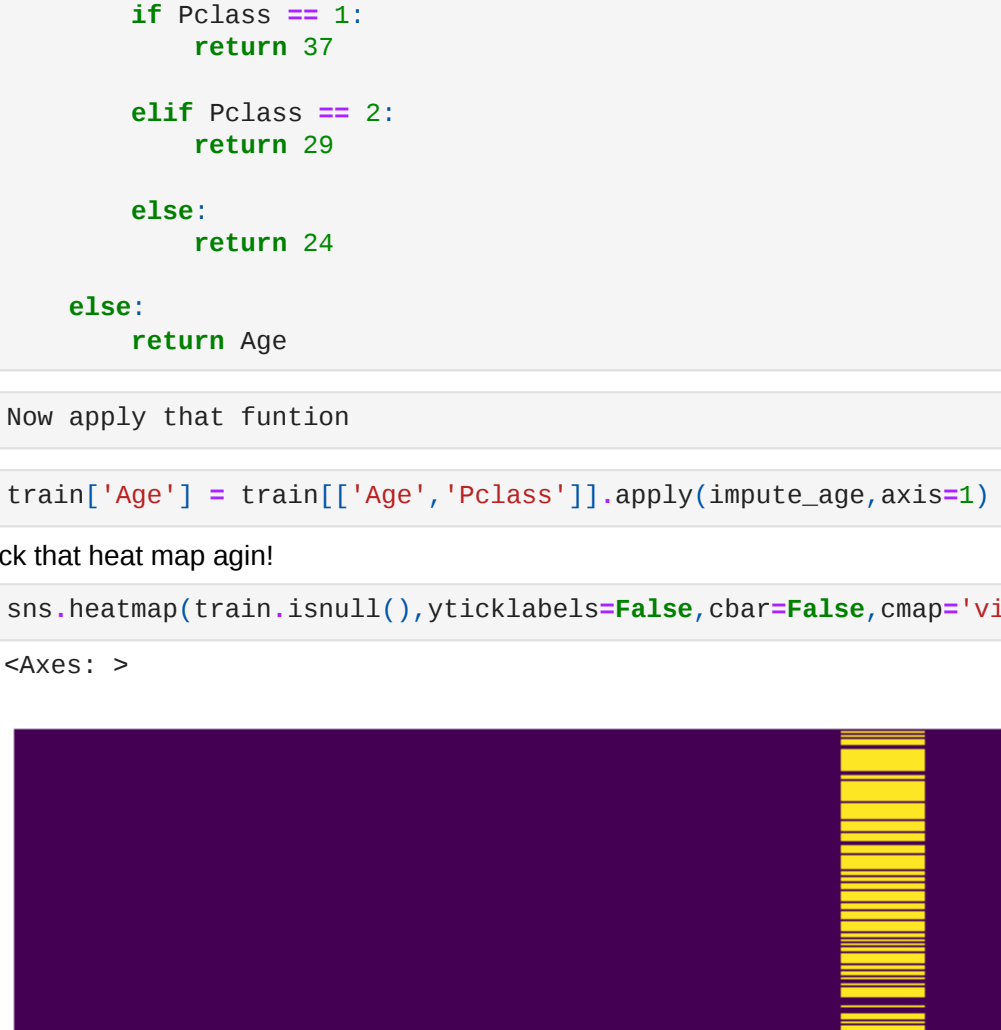
In [ ]: Now apply that funtion

```
In [25]: train['Age'] = train[['Age','Pclass']].apply(impute_age,axis=1)
```

Now lets Check that heat map again!

```
In [27]: sns.heatmap(train.isnull(),yticklabels=False,cbar=False,cmap='viridis')
```

```
Out[27]: <Axes: >
```



## ahead and drop the Cabin column and the row in Embarked that is NaN

```
In [ ]: train.drop('Cabin',axis=1,inplace=True)
```

```
In [29]: train.head()
```

```
Out[29]:
```

	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	Heikkinen, 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 [ ]: tain.dropna(inplace=True)
```

## Converting Categorical Features

Convert categorical features to dummy variables using pandas! Otherwise our machine learning algorithm won't be able to directly take in those features as inputs.

```
In [30]: train.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                   891 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 [31]: pd.get_dummies(train['Embarked'],drop_first=True).head()
```

```
Out[31]:
```

	Q	S
0	0	1
1	0	0
2	0	1
3	0	1
4	0	1

```
In [ ]: embark = pd.get_dummies(train['Sex'],drop_first=True)
sexb = pd.get_dummies(train['Embarked'],drop_first=True)
```

```
In [ ]: train.drop(['Sex','Embarked','Name','Ticket'],axis=1,inplace=True)
```

```
In [32]: train.head()
```

```
Out[32]:
```

	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	Heikkinen, 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 [ ]: train = pd.concat([train,sexb,embark],axis=1)
```

```
In [61]: train.head()
```

```
Out[61]:
```

	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	Heikkinen, 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

#Building A Logistic Regression

```
In [ ]: Train Test Split
```

```
In [62]: train.drop('Survived',axis=1).head()
```

```
Out[62]:
```

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

```
In [63]: train['Survived'].head()
```

```
Out[63]:
```

	Survived
0	0
1	1
2	1
3	1
4	0

Name: Survived, dtype: int64

```
In [ ]: from sklearn.model_selection import train_test_split
```

```
In [ ]: x_train,x_test,y_train,y_test = train_test_split(train_split(train.drop('Survived',axis=1),train['Survived']), test_size=0.30,random_state=101)
```

Training and Predicting

```
In [ ]: from sklearn.linear_model import LogisticRegression
```

```
In [ ]: predictions = logmodel.predict(X_test)
```

```
In [ ]: from sklearn.metrics import confusion_matrix
```

```
In [ ]: accuracy=confusion_matrix(y_test,predictions)
```

```
In [ ]: from sklearn.metrics import accuracy_score
```

```
In [ ]: accuracy=accuracy_score(y_test,predictions)
```

## Evaluation

```
In [ ]: from sklearn.metrics import classification_report
```

```
In [ ]: print(classification_report(y_test,predictions))
```

precision recall f1-score support

	0	0.81	0.93	0.86	163
1	1	0.85	0.65	0.74	104

avg / total 0.82 0.82 0.81 267

```
In [ ]: """End Of The Code"""
```

```
In [ ]: *****THANK YOU!*****
```