

```
In [1]: import numpy as np
import pandas as pd
from sklearn import preprocessing
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style='white')
sns.set(style='whitegrid', color_codes=True)
import warnings
warnings.simplefilter(action='ignore')
```

```
In [2]: train_df=pd.read_csv(r"C:\Users\Lenovo\OneDrive\Desktop\Data Sets\train.gender_submission.csv")
train_df
```

Out[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	Cumings, 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
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

891 rows × 12 columns

```
In [3]: test_df=pd.read_csv(r"C:\Users\Lenovo\OneDrive\Desktop\Data Sets\train.gender_submission.csv")
test_df
```

Out[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	Cumings, 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
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

891 rows × 12 columns

In [4]:

train_df.head()

Out[4]:

	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	Cumings, 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 [6]:

train_df.shape

Out[6]:

(891, 12)

```
In [8]: train_df.describe
```

```
Out[8]: <bound method NDFrame.describe of      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]>
```

```
In [10]: train_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 [12]: test_df.describe

Out[12]: <bound method NDFrame.describe of 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]>

In [13]: test_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
```

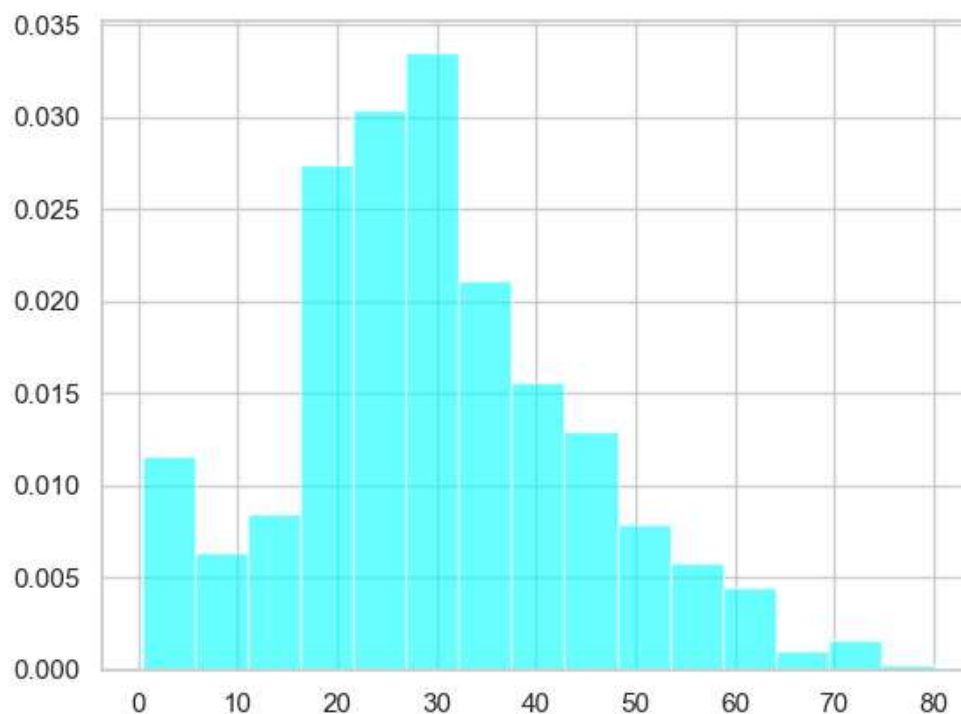
```
In [14]: train_df.isnull().sum()
```

```
Out[14]: 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
```

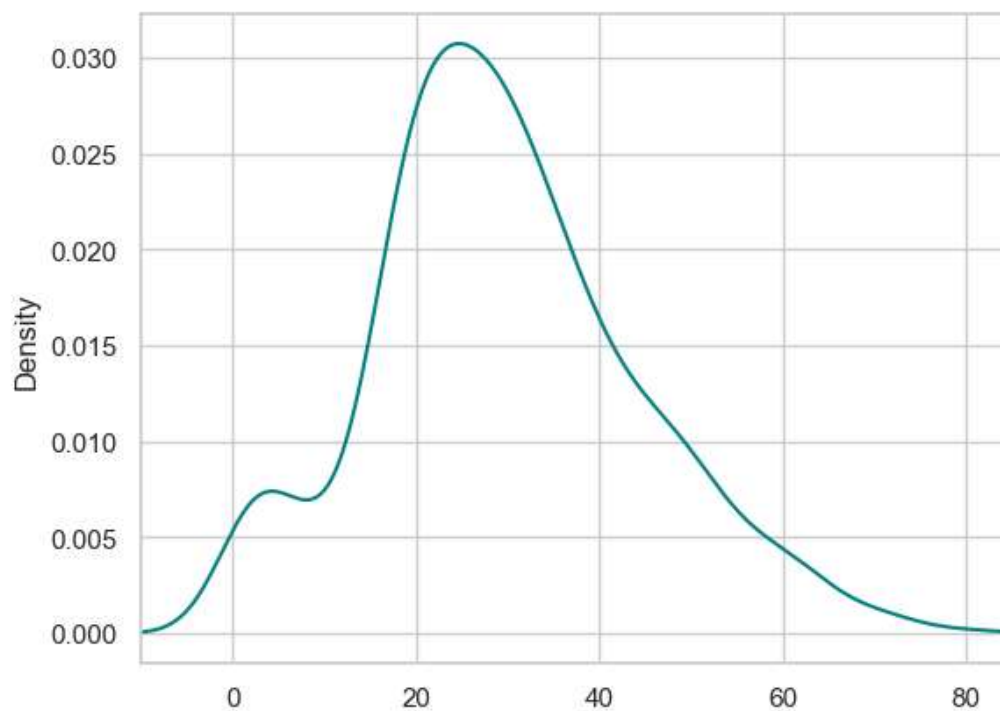
```
In [15]: test_df.isnull().sum()
```

```
Out[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
```

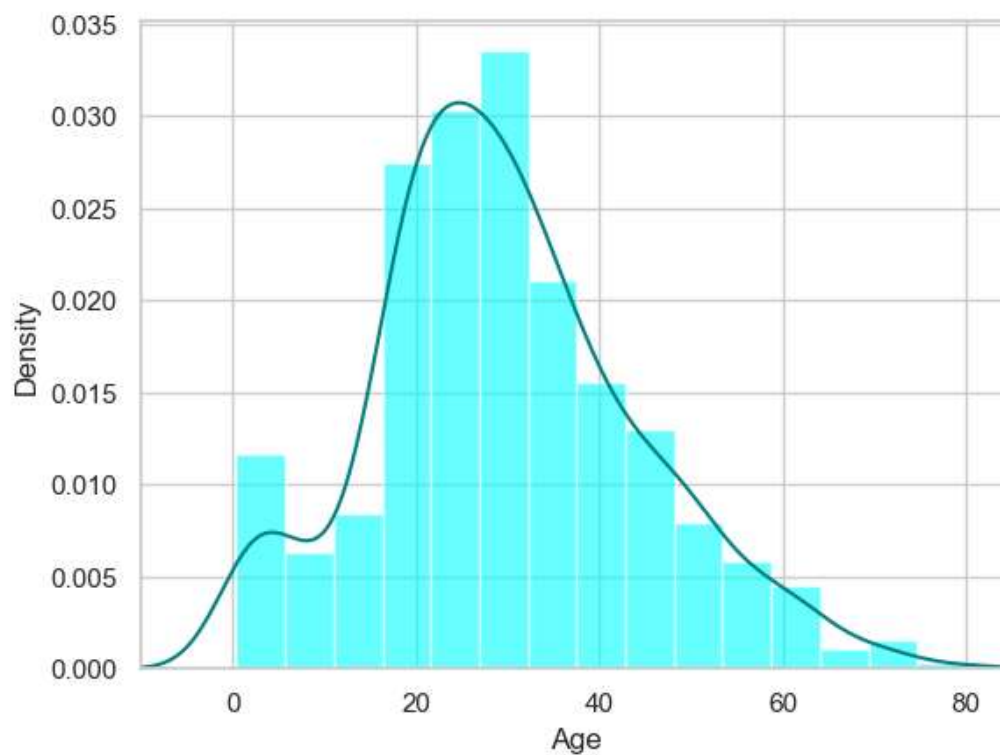
```
In [16]: ax=train_df["Age"].hist(bins=15,density=True,stacked=True,color='cyan',alpha=0.6)
```



```
In [17]: train_df["Age"].plot(kind='density',color='teal')
ax.set(xlabel='Age')
plt.xlim(-10,85)
plt.show()
```



```
In [18]: ax=train_df["Age"].hist(bins=15,density=True,stacked=True,color='cyan',alpha=0.6)
train_df["Age"].plot(kind='density',color='teal')
ax.set(xlabel='Age')
plt.xlim(-10,85)
plt.show()
```



```
In [19]: print(train_df['Age'].mean(skipna=True))  
         print(train_df['Age'].median(skipna=True))
```

```
29.69911764705882  
28.0
```

```
In [20]: print((train_df['Cabin'].isnull().sum()/train_df.shape[0])*100)
```

```
77.10437710437711
```

```
In [21]: print((train_df['Embarked'].isnull().sum()/train_df.shape[0])*100)
```

```
0.22446689113355783
```

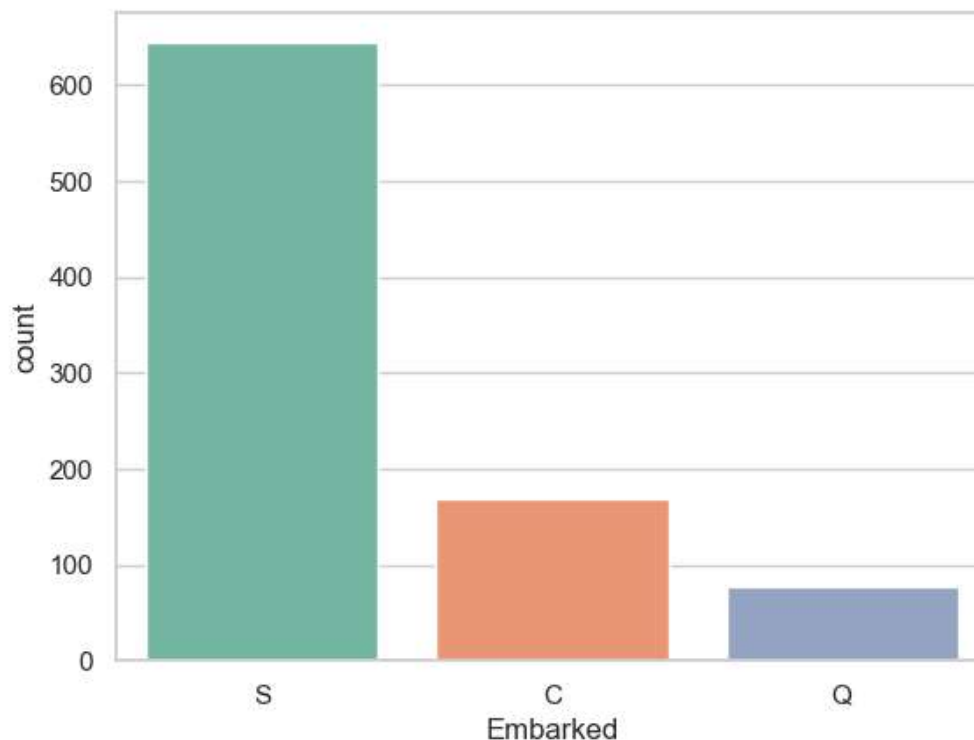
```
In [26]: print("Board passengers grouped by part of embartion(C=cherbourg,Q=Queenstown,S=Southmapton):"
```

```
Board passengers grouped by part of embartion(C=cherbourg,Q=Queenstown,S=Southmapton):
```

```
In [27]: print(train_df['Embarked'].value_counts())
```

```
Embarked  
S      644  
C      168  
Q       77  
Name: count, dtype: int64
```

```
In [28]: sns.countplot(x='Embarked',data=train_df,palette='Set2')  
         plt.show()
```



```
In [29]: print(train_df['Embarked'].value_counts().idxmax())
```

```
S
```



```
In [30]: train_data=train_df.copy()
train_data['Age'].fillna(train_df['Age'].median(skipna=True),inplace=True)
train_data['Embarked'].fillna(train_df['Embarked'].value_counts().idxmax(),inplace=True)
```

```
In [31]: train_data.drop('Cabin',axis=1,inplace=True)
```

```
In [32]: train_data.isnull().sum()
```

```
Out[32]: PassengerId    0
Survived              0
Pclass               0
Name                 0
Sex                  0
Age                  0
SibSp                0
Parch                0
Ticket               0
Fare                 0
Embarked             0
dtype: int64
```

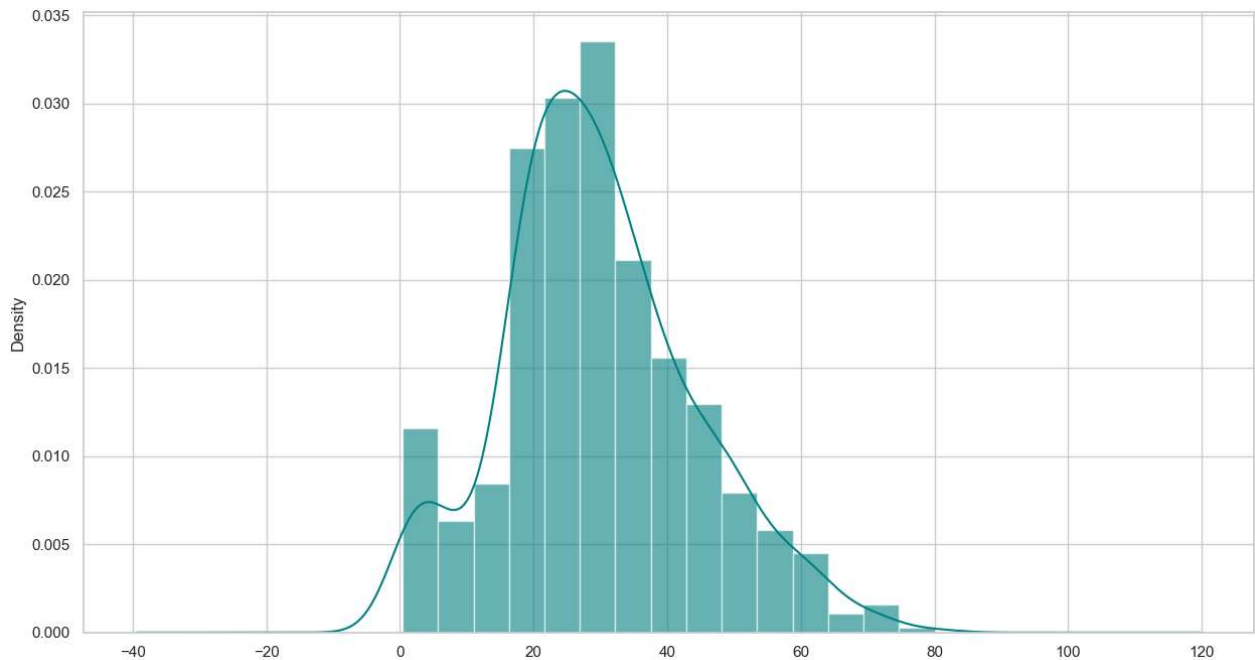
```
In [33]: train_data.head()
```

```
Out[33]:
```

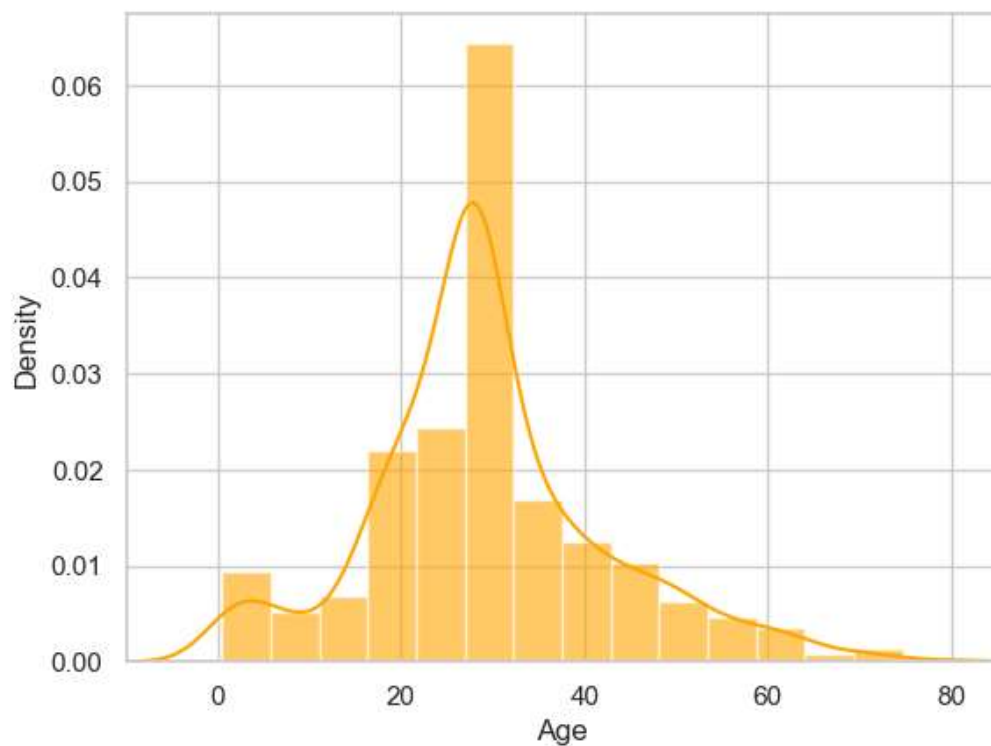
	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250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

```
In [34]: plt.figure(figsize=[15,8])  
ax=train_df['Age'].hist(bins=15,density=True,stacked=True,color='teal',alpha=0.6)  
train_df['Age'].plot(kind='density',color='teal')
```

Out[34]: <Axes: ylabel='Density'>



```
In [35]: ax=train_data['Age'].hist(bins=15,density=True,stacked=True,color='orange',alpha=0.6)  
train_data['Age'].plot(kind='density',color='orange')  
ax.set(xlabel='Age')  
plt.xlim(-10,85)  
plt.show()
```



```
In [36]: #create catagorical variable for travelling alone
train_data['TravelAlone']=np.where((train_data['SibSp']+train_data['Parch'])>0,0,1)
```

```
In [37]: train_data.drop("SibSp",axis=1,inplace=True)
train_data.drop("Parch",axis=1,inplace=True)
```

```
In [38]: #ctreate catagorical variables and drop some variables
training=pd.get_dummies(train_data,columns=["Pclass","Embarked","Sex"])
training.drop("Sex_female",axis=1,inplace=True)
training.drop("PassengerId",axis=1,inplace=True)
training.drop("Name",axis=1,inplace=True)
training.drop("Ticket",axis=1,inplace=True)
final_train=training
final_train.head()
```

Out[38]:

	Survived	Age	Fare	TravelAlone	Pclass_1	Pclass_2	Pclass_3	Embarked_C	Embarked_Q	Embarked_S	Sex
0	0	22.0	7.2500	0	False	False	True	False	False	True	
1	1	38.0	71.2833	0	True	False	False	True	False	False	
2	1	26.0	7.9250	1	False	False	True	False	False	True	
3	1	35.0	53.1000	0	True	False	False	False	False	True	
4	0	35.0	8.0500	1	False	False	True	False	False	True	

```
In [39]: test_df.isnull().sum()
```

```
Out[39]: 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 [40]: test_data=test_df.copy()
test_data['Age'].fillna(test_df['Age'].median(skipna=True),inplace=True)
test_data['Embarked'].fillna(test_df['Embarked'].value_counts().idxmax(),inplace=True)
test_data.drop('Cabin',axis=1,inplace=True)
```

```
In [41]: test_data['TravelAlone']=np.where((test_data['SibSp']+test_data['Parch'])>0,0,1)
test_data.drop("SibSp",axis=1,inplace=True)
test_data.drop("Parch",axis=1,inplace=True)
```

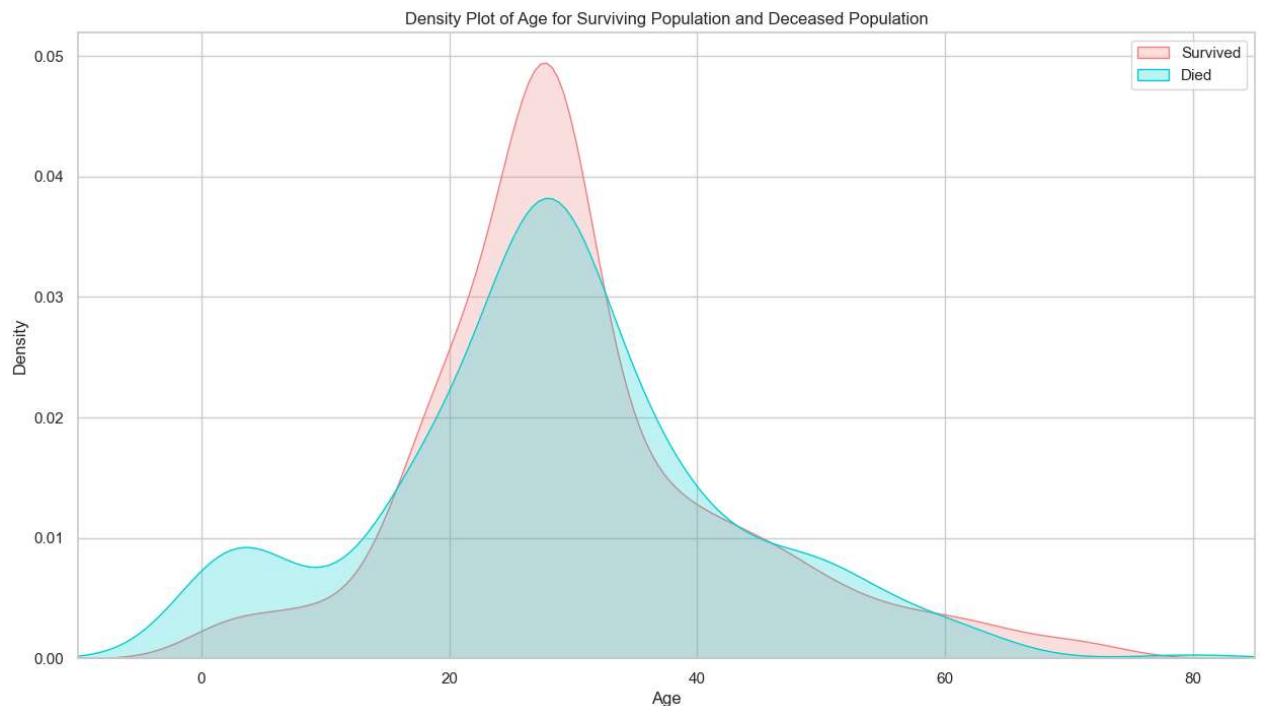
```
In [42]: testing=pd.get_dummies(test_data,columns=["Pclass","Embarked","Sex"])
testing.drop("Sex_female",axis=1,inplace=True)
testing.drop("PassengerId",axis=1,inplace=True)
testing.drop("Name",axis=1,inplace=True)
testing.drop("Ticket",axis=1,inplace=True)
final_test=testing
final_test.head()
```

Out[42]:

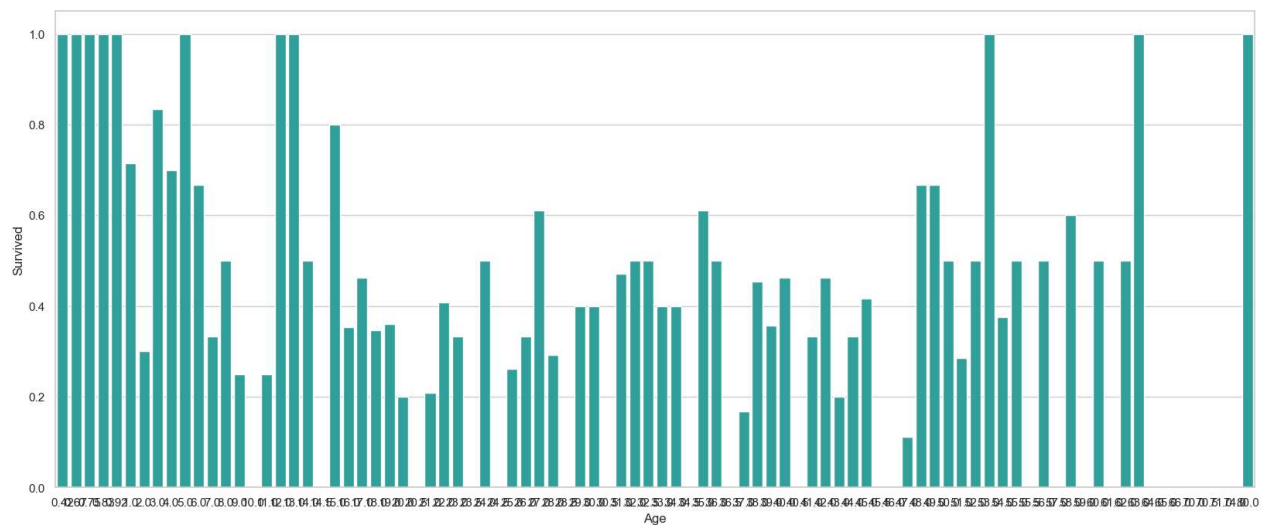
	Survived	Age	Fare	TravelAlone	Pclass_1	Pclass_2	Pclass_3	Embarked_C	Embarked_Q	Embarked_S	Sex
0	0	22.0	7.2500	0	False	False	True	False	False	True	
1	1	38.0	71.2833	0	True	False	False	True	False	False	
2	1	26.0	7.9250	1	False	False	True	False	False	True	
3	1	35.0	53.1000	0	True	False	False	False	False	True	
4	0	35.0	8.0500	1	False	False	True	False	False	True	

```
In [44]: plt.figure(figsize=(15,8))
ax = sns.kdeplot(final_train["Age"][final_train.Survived == 1], color="darkturquoise", shade=
sns.kdeplot(final_train["Age"][final_train.Survived == 0], color="lightcoral", shade=True))
plt.legend(['Survived', 'Died'])
plt.title('Density Plot of Age for Surviving Population and Deceased Population')
ax.set(xlabel='Age')
plt.xlim(-10,85)
```

Out[44]: (-10.0, 85.0)



```
In [47]: plt.figure(figsize=(20,8))
avg_survival_byage = final_train[["Age", "Survived"]].groupby(['Age'], as_index=False).mean()
g = sns.barplot(x='Age', y='Survived', data=avg_survival_byage, color="LightSeaGreen")
plt.show()
```



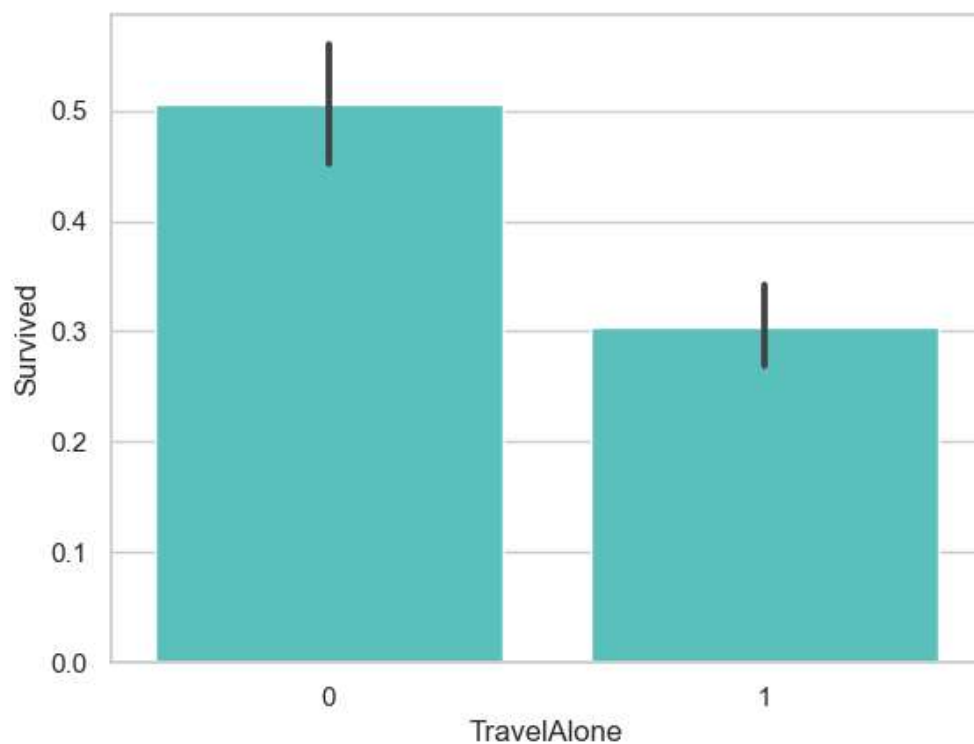
```
In [48]: final_train['IsMinor']=np.where(final_train['Age']<=16, 1, 0)
print(final_train['IsMinor'])
```

```
0      0
1      0
2      0
3      0
4      0
..
886    0
887    0
888    0
889    0
890    0
Name: IsMinor, Length: 891, dtype: int32
```

```
In [49]: final_test['IsMinor']=np.where(final_test['Age']<=16, 1, 0)
print(final_test['IsMinor'])
```

```
0      0
1      0
2      0
3      0
4      0
..
886    0
887    0
888    0
889    0
890    0
Name: IsMinor, Length: 891, dtype: int32
```

```
In [50]: sns.barplot(x='TravelAlone', y='Survived', data=final_train, color="mediumturquoise")  
plt.show()
```



```
In [51]: import seaborn as sns  
import matplotlib.pyplot as plt  
# Assuming 'train_df' is your DataFrame containing the data  
sns.barplot(x='Sex', y='Survived', data=train_df, color='aquamarine')  
plt.show()
```

