# problem statement:

```
         To predict and analyze which gender has a high chance of survival
     at the time of disaster
```

In [40]:
```python
#import libraries
import numpy as np
import pandas as pd
from sklearn import preprocessing
import matplotlib.pyplot as plt
import seaborn as sns
```

In [41]:
```python
sns.set(style="white")#white background style for seaborn plots
sns.set(style="whitegrid",color_codes=True)
```

In [42]:
```python
import warnings
warnings.simplefilter(action='ignore')
```

In [43]:
```python
train_df=pd.read_csv(r"C:\Users\my pc\Downloads\train.gender_submission.csv")
train_df
```

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | 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... | male | 35.0 | 0 | 0 | 373450 | 8.0500 |

In [44]: 
```python
test_df=pd.read_csv(r"C:\Users\my pc\Downloads\test.gender_submission.csv")
test_df
```

Out[44]:

| | PassengerId | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 892 | 3 | Kelly, Mr. James | male | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN |
| 1 | 893 | 3 | Wilkes, Mrs. James (Ellen Needs) | female | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN |
| 2 | 894 | 2 | Myles, Mr. Thomas Francis | male | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN |
| 3 | 895 | 3 | Wirz, Mr. Albert | male | 27.0 | 0 | 0 | 315154 | 8.6625 | NaN |
| 4 | 896 | 3 | Hirvonen, Mrs. Alexander (Helga E Lindqvist) | female | 22.0 | 1 | 1 | 3101298 | 12.2875 | NaN |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 413 | 1305 | 3 | Spector, Mr. Woolf | male | NaN | 0 | 0 | A.5. 3236 | 8.0500 | NaN |
| 414 | 1306 | 1 | Oliva y Ocana, Dona. Fermina | female | 39.0 | 0 | 0 | PC 17758 | 108.9000 | C105 |
| 415 | 1307 | 3 | Saether, Mr. Simon Sivertsen | male | 38.5 | 0 | 0 | SOTON/O.Q. 3101262 | 7.2500 | NaN |
| 416 | 1308 | 3 | Ware, Mr. Frederick | male | NaN | 0 | 0 | 359309 | 8.0500 | NaN |
| 417 | 1309 | 3 | Peter, Master. Michael J | male | NaN | 1 | 1 | 2668 | 22.3583 | NaN |

418 rows × 11 columns

In [45]: `train_df.head()`

Out[45]:

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Ca |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | N |
| 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 | N |
| 3 | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C1 |
| 4 | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | N |

In [46]: `train_df.shape`

Out[46]: `(891, 12)`

In [47]: `test_df.head()`

Out[47]:

| | PassengerId | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embark |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 892 | 3 | Kelly, Mr. James | male | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN | |
| 1 | 893 | 3 | Wilkes, Mrs. James (Ellen Needs) | female | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN | |
| 2 | 894 | 2 | Myles, Mr. Thomas Francis | male | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN | |
| 3 | 895 | 3 | Wirz, Mr. Albert | male | 27.0 | 0 | 0 | 315154 | 8.6625 | NaN | |
| 4 | 896 | 3 | Hirvonen, Mrs. Alexander (Helga E Lindqvist) | female | 22.0 | 1 | 1 | 3101298 | 12.2875 | NaN | |

In [48]: `test_df.shape`

Out[48]: (418, 11)

In [49]: `test_df.describe()`

Out[49]:

|        | PassengerId | Pclass     | Age        | SibSp      | Parch      | Fare       |
|--------|-------------|------------|------------|------------|------------|------------|
| count  | 418.000000  | 418.000000 | 332.000000 | 418.000000 | 418.000000 | 417.000000 |
| mean   | 1100.500000 | 2.265550   | 30.272590  | 0.447368   | 0.392344   | 35.627188  |
| std    | 120.810458  | 0.841838   | 14.181209  | 0.896760   | 0.981429   | 55.907576  |
| min    | 892.000000  | 1.000000   | 0.170000   | 0.000000   | 0.000000   | 0.000000   |
| 25%    | 996.250000  | 1.000000   | 21.000000  | 0.000000   | 0.000000   | 7.895800   |
| 50%    | 1100.500000 | 3.000000   | 27.000000  | 0.000000   | 0.000000   | 14.454200  |
| 75%    | 1204.750000 | 3.000000   | 39.000000  | 1.000000   | 0.000000   | 31.500000  |
| max    | 1309.000000 | 3.000000   | 76.000000  | 8.000000   | 9.000000   | 512.329200 |

In [50]: `train_df.describe()`

Out[50]:

|        | 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 [51]: `test_df.describe()`

Out[51]:

|        | PassengerId | Pclass     | Age        | SibSp      | Parch      | Fare       |
|--------|-------------|------------|------------|------------|------------|------------|
| count  | 418.000000  | 418.000000 | 332.000000 | 418.000000 | 418.000000 | 417.000000 |
| mean   | 1100.500000 | 2.265550   | 30.272590  | 0.447368   | 0.392344   | 35.627188  |
| std    | 120.810458  | 0.841838   | 14.181209  | 0.896760   | 0.981429   | 55.907576  |
| min    | 892.000000  | 1.000000   | 0.170000   | 0.000000   | 0.000000   | 0.000000   |
| 25%    | 996.250000  | 1.000000   | 21.000000  | 0.000000   | 0.000000   | 7.895800   |
| 50%    | 1100.500000 | 3.000000   | 27.000000  | 0.000000   | 0.000000   | 14.454200  |
| 75%    | 1204.750000 | 3.000000   | 39.000000  | 1.000000   | 0.000000   | 31.500000  |
| max    | 1309.000000 | 3.000000   | 76.000000  | 8.000000   | 9.000000   | 512.329200 |

In [52]: `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 [53]: `test_df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 418 entries, 0 to 417
Data columns (total 11 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   PassengerId  418 non-null    int64
 1   Pclass       418 non-null    int64
 2   Name         418 non-null    object
 3   Sex          418 non-null    object
 4   Age          332 non-null    float64
 5   SibSp        418 non-null    int64
 6   Parch        418 non-null    int64
 7   Ticket       418 non-null    object
 8   Fare         417 non-null    float64
 9   Cabin        91 non-null     object
 10  Embarked     418 non-null    object
dtypes: float64(2), int64(4), object(5)
memory usage: 36.0+ KB
```

# To find the missing values

In [54]: `train_df.isna().any()`

Out[54]:
```
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
```
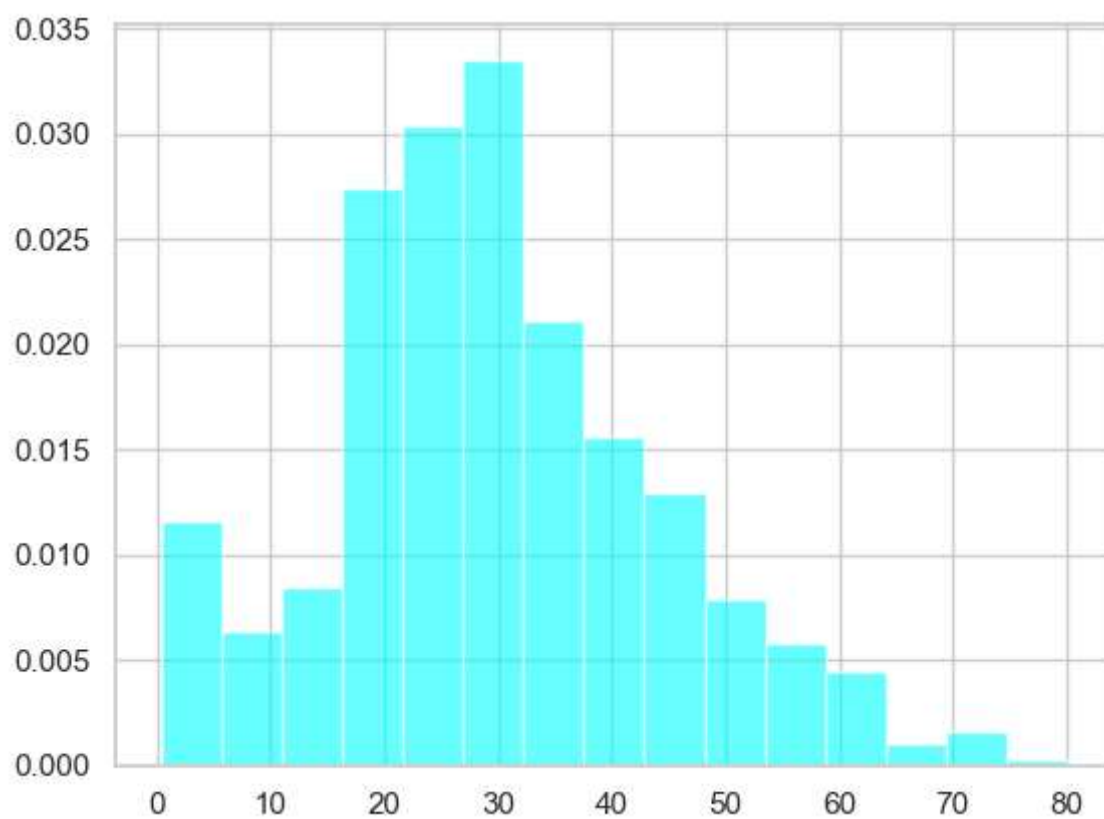
In [55]: `test_df.isna().any()`

Out[55]:
```
PassengerId    False
Pclass         False
Name           False
Sex            False
Age             True
SibSp          False
Parch          False
Ticket         False
Fare            True
Cabin           True
Embarked       False
dtype: bool
```

In [56]: `train_df.isnull().sum()`

Out[56]:
```
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 [57]: `test_df.isnull().sum()`

Out[57]:
```
PassengerId      0
Pclass           0
Name             0
Sex              0
Age             86
SibSp            0
Parch            0
Ticket           0
Fare             1
Cabin          327
Embarked         0
dtype: int64
```
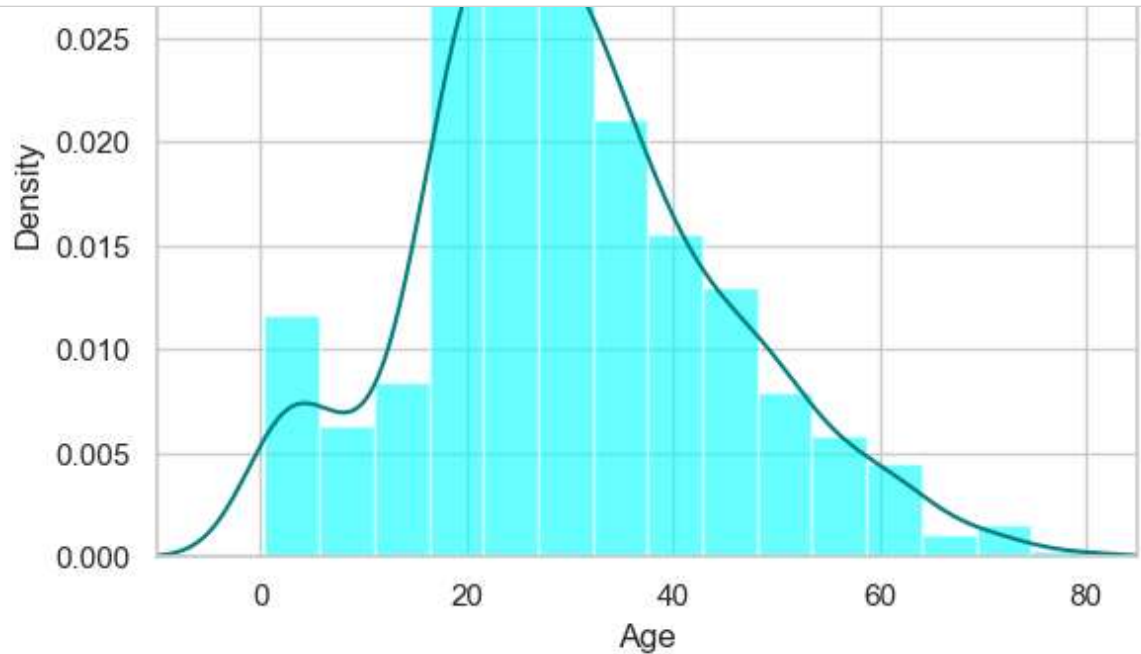
In [58]:
```
ax=train_df["Age"].hist(bins=15,density=True,stacked=True,color='cyan',alpha=0
ax
```

Out[58]: `<Axes: >`

In [59]:
```python
ax=train_df["Age"].hist(bins=15,density=True,stacked=True,color='cyan',alpha=0
train_df["Age"].plot(kind="density",color="teal")
ax.set(xlabel="Age")
plt.xlim(-10,85)
plt.show()
```



In [60]:
```python
print(train_df["Age"].mean(skipna=True))
print(train_df["Age"].median(skipna=True))
```

```
29.69911764705882
28.0
```

In [61]:
```python
print((train_df["Cabin"].isnull().sum()/train_df.shape[0])*100)
```

```
77.10437710437711
```

In [62]:
```python
print((train_df["Embarked"].isnull().sum()/train_df.shape[0])*100)
```
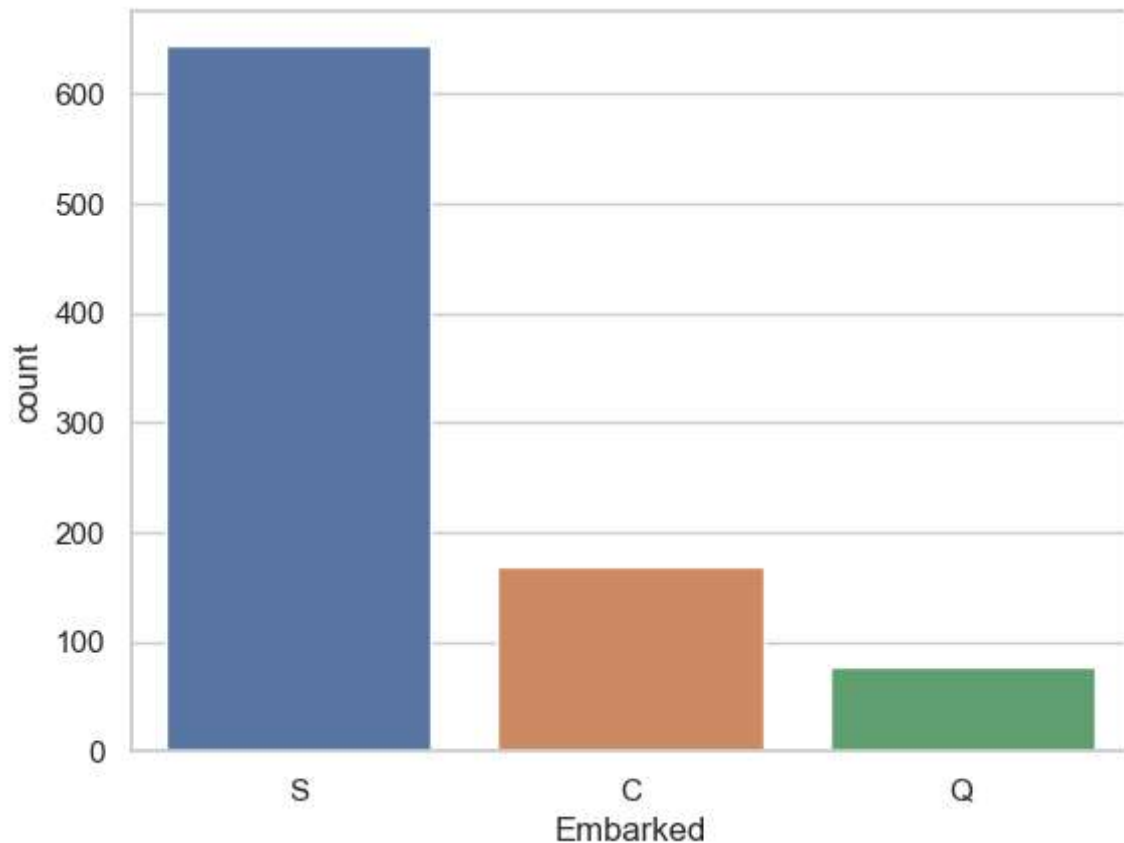
```
0.22446689113355783
```

In [63]:
```python
print("Boarded passengers grouped by port of embarkation (C-cherybourg,Q-Queen:
```

```
Boarded passengers grouped by port of embarkation (C-cherybourg,Q-Queenstown,
S=Southampton):
```

In [64]:
```python
sns.countplot(x="Embarked",data=train_df)
plt.show()
```



In [65]:
```python
print(train_df["Embarked"].value_counts().idxmax())
```

```
S
```

In [66]:
```python
train_data=train_df.copy()
```

In [67]:
```python
train_data["Age"].fillna(train_df["Age"].median(skipna=True),inplace=True)
train_data["Embarked"].fillna(train_df["Embarked"].value_counts().idxmax(),inp
train_data.drop("Cabin",axis=1,inplace=True)
```

In [68]:
```python
train_data.isnull().sum()
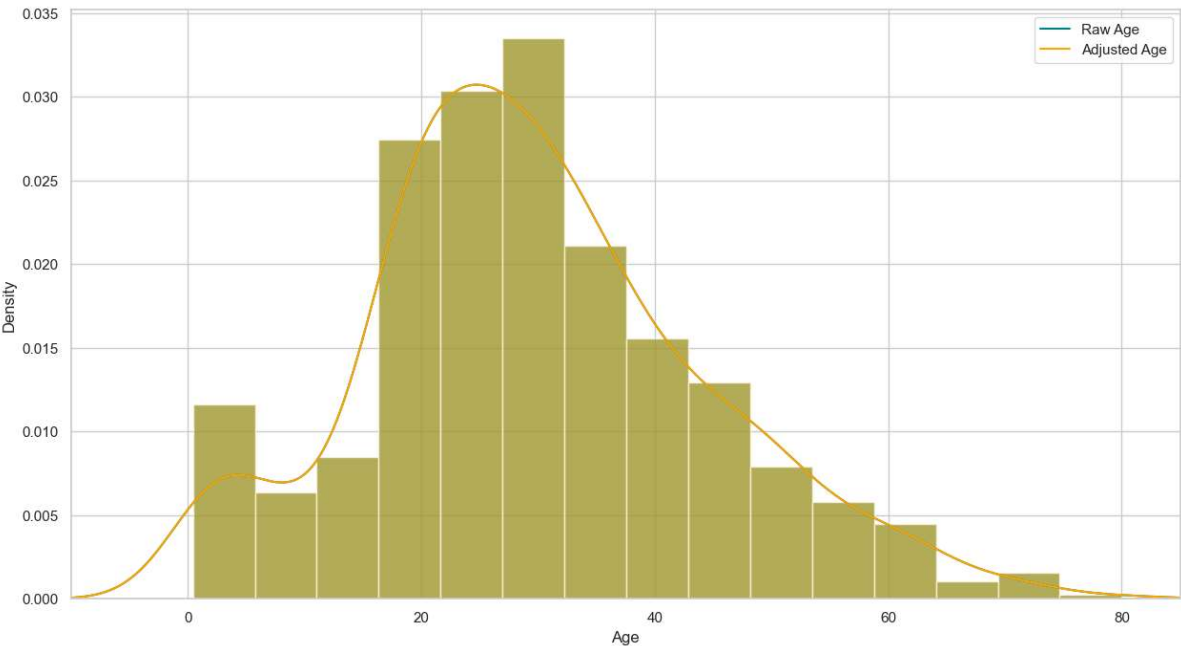```

Out[68]:
```
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 [69]: `train_data.head()`

Out[69]:

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Em |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 | |

In [70]:
```python
plt.figure(figsize=(15,8))
ax=train_df["Age"].hist(bins=15,density=True,stacked=True,color='teal',alpha=0
train_df["Age"].plot(kind="density",color="teal")
ax=train_df["Age"].hist(bins=15,density=True,stacked=True,color='orange',alpha
train_df["Age"].plot(kind="density",color="orange")
ax.legend(["Raw Age",'Adjusted Age'])
ax.set(xlabel="Age")
plt.xlim(-10,85)
plt.show()
```
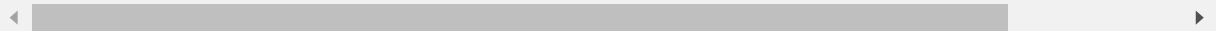
In [71]: `train_df`

Out[71]:

| | 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 [72]:
```python
#creating categorical variable for travelling alone
train_data["TravelAlone"]=np.where((train_data["SibSp"]+train_data["Parch"])>0
train_data.drop("SibSp",axis=1,inplace=True)
train_data.drop("Parch",axis=1,inplace=True)
```

In [73]: 
```python
#create categorical variables and 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[73]:

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

In [74]: 
```python
test_df.isnull().sum()
```

Out[74]: 
```
PassengerId      0
Pclass           0
Name             0
Sex              0
Age             86
SibSp            0
Parch            0
Ticket           0
Fare             1
Cabin          327
Embarked         0
dtype: int64
```

In [75]: 
```python
test_data=test_df.copy()
```

In [76]:
```python
test_data["Age"].fillna(train_df["Age"].median(skipna=True),inplace=True)
test_data["Fare"].fillna(train_df["Fare"].median(skipna=True),inplace=True)
test_data.drop("Cabin",axis=1,inplace=True)
test_data["TravelAlone"]=np.where((test_data["SibSp"]+test_data["Parch"])>0,0,
test_data.drop("SibSp",axis=1,inplace=True)
test_data.drop("Parch",axis=1,inplace=True)
testing=pd.get_dummies(train_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=training
final_test.head()
```

Out[76]:

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