

```

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")
#seaborn plots
sns.set(style="whitegrid", color_codes=True)
import warnings
warnings.simplefilter (action='ignore')

```

```
In [49]: train_df=pd.read_csv(r"C:\Users\MY HOME\Downloads\train.gender_submission.csv")
train_df
```

Out[49]:

	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 [50]: test_df=pd.read_csv(r"C:\Users\MY HOME\Downloads\test.gender_submission.csv")
test_df
```

Out[50]:

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

418 rows × 11 columns

```
In [51]: train_df.shape
```

Out[51]: (891, 12)

```
In [52]: test_df.shape
```

Out[52]: (418, 11)

In [53]: ▶ `train_df.describe`

```

Out[53]: <bound method NDFrame.describe of
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]>
```

```
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
```

```
test_df.head()
```

Out[55]:

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

```
In [56]: test_df.describe()
```

Out[56]:

	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 [57]: 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: 36.1+ KB
```

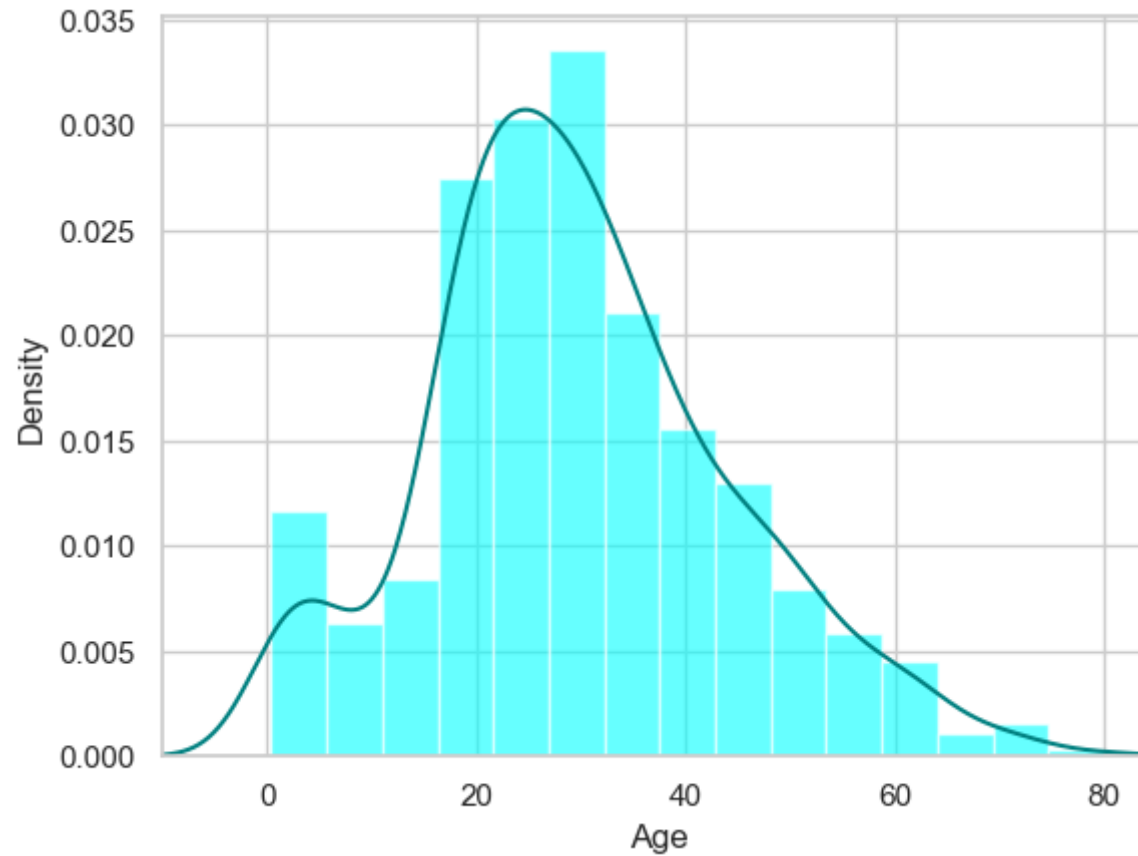
To find missing values

```
In [58]: train_df.isnull().sum()
```

```
Out[58]: PassengerId      0  
Survived                  0  
Pclass                    0  
Name                      0  
Sex                       0  
Age                        177  
SibSp                     0  
Parch                     0  
Ticket                   687  
Fare                      2  
Cabin                     91  
Embarked                  2  
dtype: int64
```



```
In [59]: ▶ 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 [60]: ▶ print(train_df["Age"].mean(skipna=True))
print(train_df["Age"].median(skipna=True))
```

29.69911764705882

28.0

77.10437710437711

0.22446689113355783

```
In [63]: ▶ print('Boarded passengers grouped by port of Embarkation(c=Cherbourg,Q=queenstown,s=Southampton):')  
print(train_df['Embarked'].value_counts())  
sns.countplot(x='Embarked',data=train_df,palette='Set2')  
plt.show()
```

Boarded passengers grouped by port of Embarkation(c=Cherbourg,Q=queenstown,s=Southampton):

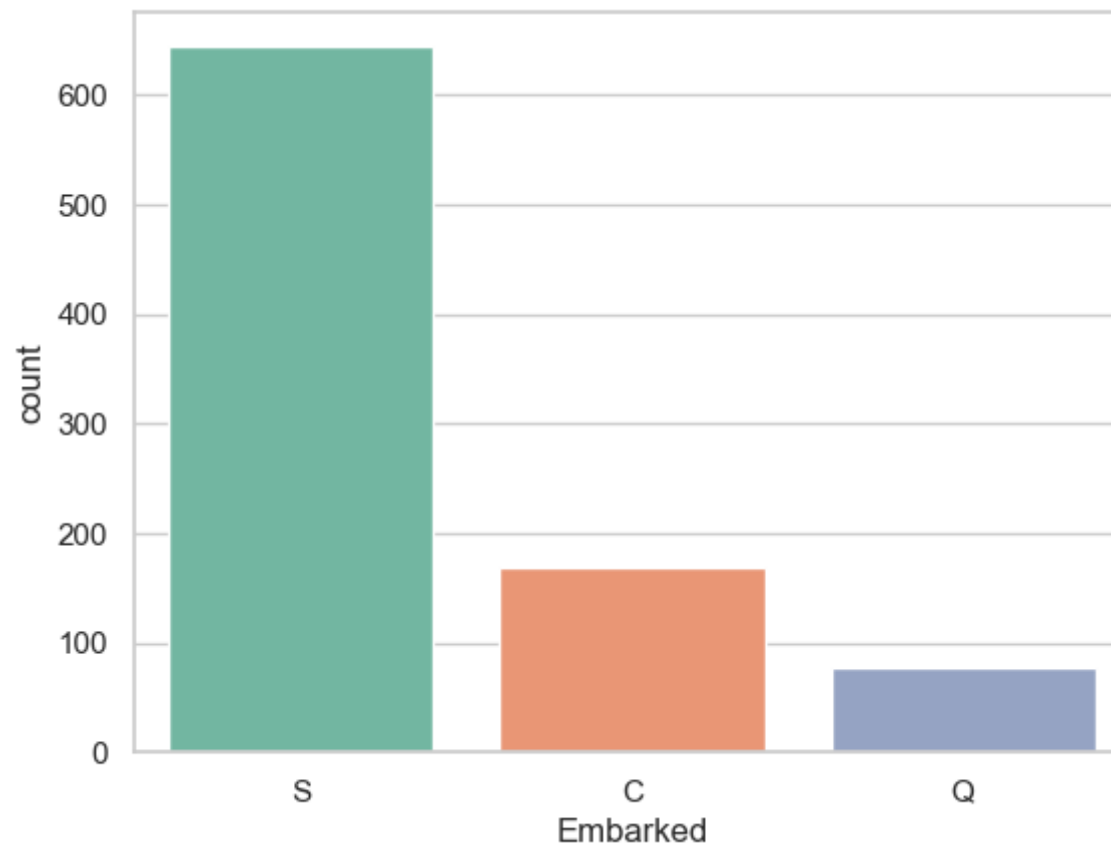
Embarked

S 644

C 168

Q 77

Name: count, dtype: int64



```
In [64]: ▶ print(train_df['Embarked'].value_counts().idxmax())
```

S

```
In [65]: ▶ 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)
train_data.drop('Cabin',axis=1,inplace=True)
```

```
In [66]: ▶ train_data.isnull().sum()
```

```
Out[66]: 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 [67]: ▶ train_data.head()
```

Out[67]:

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

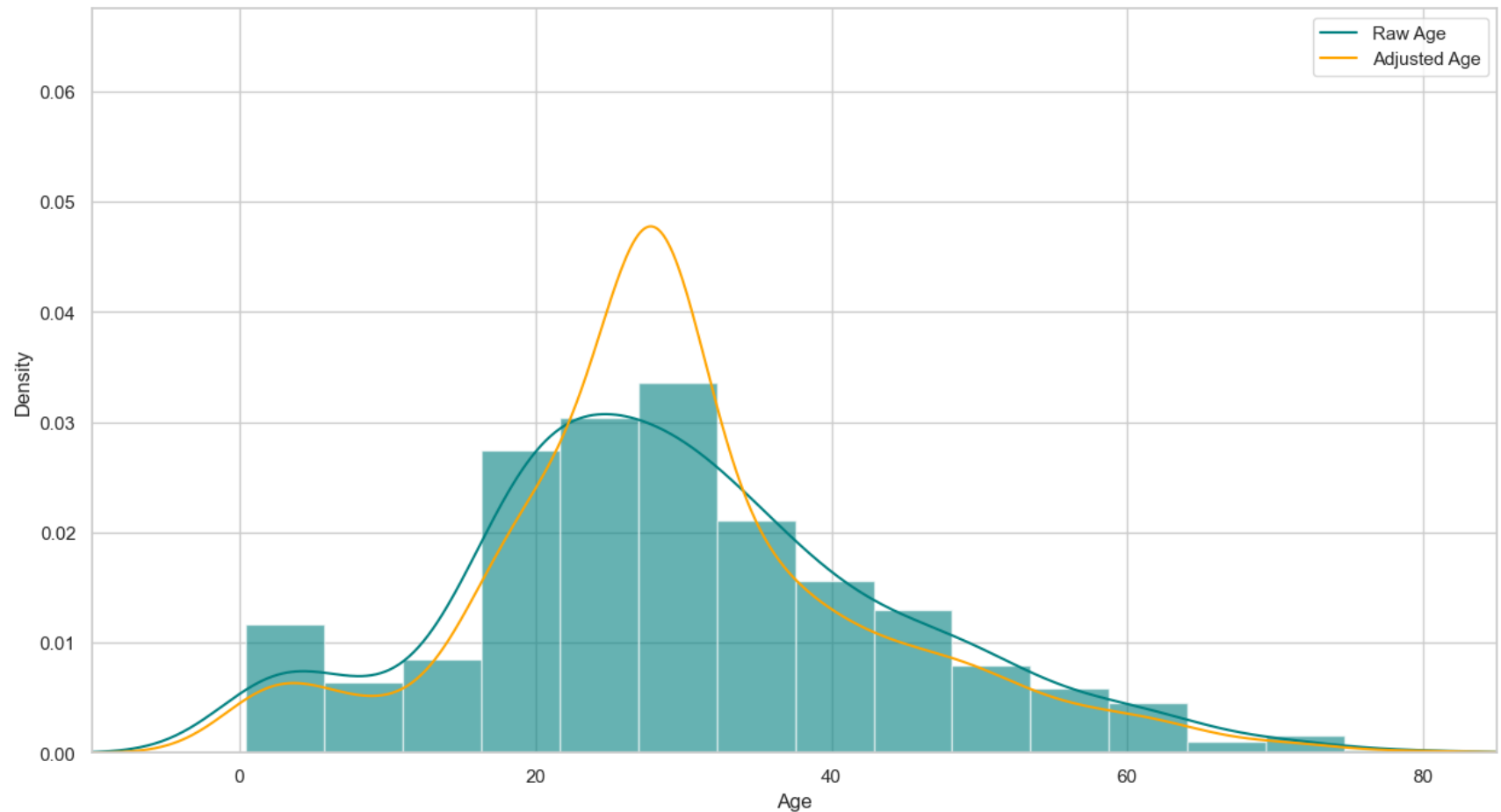
```
Out[68]: PassengerId    0
Survived                0
Pclass                  0
Name                    0
Sex                     0
Age                     0
SibSp                   0
Parch                   0
Ticket                  0
Fare                    0
Embarked                0
dtype: int64
```

```
train_data.head()
```

Out[69]:

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 [70]: ▶ 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')
ax = train_data["Age"].hist(bins=15, density=True, stacked=True, color='orange', alpha=0)
train_data["Age"].plot(kind='density', color='orange')
ax.legend(['Raw Age', 'Adjusted Age'])
ax.set(xlabel='Age')
plt.xlim(-10,85)
plt.show()
```



```
In [28]: ▶ ## Create categorical variable for traveling alone
train_data['TravelAlone']=np.where((train_data["SibSp"]+train_data["Parch"])>0,0,1)
train_data.drop('SibSp', axis=1, inplace=True)
train_data.drop('Parch', axis=1, inplace=True)
```

```
In [29]: ▶ 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[29]:

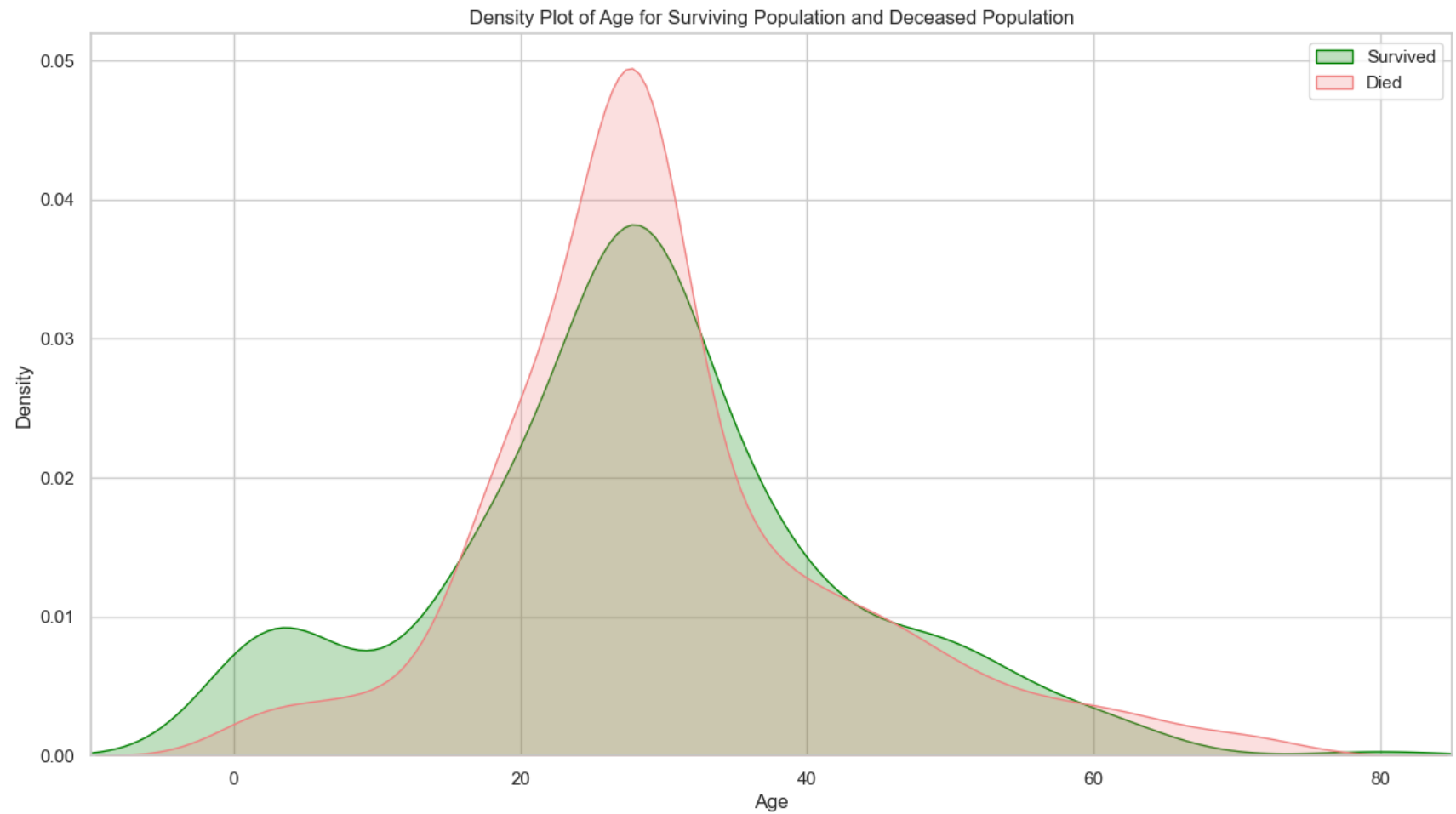
	Survived	Age	Fare	TravelAlone	Pclass_1	Pclass_2	Pclass_3	Embarked_C	Embarked_Q	Embarked_S	Sex_male
0	0	22.0	7.2500	0	False	False	True	False	False	True	True
1	1	38.0	71.2833	0	True	False	False	True	False	False	False
2	1	26.0	7.9250	1	False	False	True	False	False	True	False
3	1	35.0	53.1000	0	True	False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


```
test_df.isnull().sum()
```

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 [31]: ▶ test_data = test_df.copy()
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, 1)
test_data.drop('SibSp', axis=1, inplace=True)
test_data.drop('Parch', axis=1, inplace=True)
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()
```

```
In [34]: ▶ plt.figure(figsize=(15,8))
ax = sns.kdeplot(final_train["Age"][final_train.Survived == 1], color="green", shade=True)
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)
plt.show()
```



The bar chart displays the survival probability for each age group. The y-axis, labeled 'Survived', ranges from 0.0 to 1.0. The x-axis, labeled 'Age', shows age groups from 0.42 to 80.0. The survival probability is generally high for younger ages, drops significantly in the middle ages (around 0.42 to 0.55), and then rises again for older ages (above 0.60).

Age	Survived
0.42	1.0
0.45	1.0
0.48	1.0
0.51	1.0
0.54	1.0
0.57	1.0
0.60	0.71
0.63	0.3
0.66	0.83
0.69	0.7
0.72	1.0
0.75	0.67
0.78	0.33
0.81	0.5
0.84	0.25
0.87	0.25
0.90	1.0
0.93	1.0
0.96	0.5
0.99	0.8
1.02	0.35
1.05	0.46
1.08	0.35
1.11	0.36
1.14	0.2
1.17	0.21
1.20	0.41
1.23	0.33
1.26	0.5
1.29	0.26
1.32	0.33
1.35	0.61
1.38	0.29
1.41	0.4
1.44	0.4
1.47	0.47
1.50	0.5
1.53	0.5
1.56	0.4
1.59	0.4
1.62	0.61
1.65	0.5
1.68	0.17
1.71	0.45
1.74	0.36
1.77	0.46
1.80	0.33
1.83	0.2
1.86	0.33
1.89	0.41
1.92	0.42
1.95	0.2
1.98	0.33
2.01	0.41
2.04	0.11
2.07	0.67
2.10	0.67
2.13	0.5
2.16	0.28
2.19	0.5
2.22	1.0
2.25	0.37
2.28	0.5
2.31	0.5
2.34	0.37
2.37	0.5
2.40	0.5
2.43	0.4
2.46	0.4
2.49	0.4
2.52	0.4
2.55	0.4
2.58	0.4
2.61	0.4
2.64	0.4
2.67	0.4
2.70	0.4
2.73	0.4
2.76	0.4
2.79	0.4
2.82	0.4
2.85	0.4
2.88	0.4
2.91	0.4
2.94	0.4
2.97	0.4
3.00	0.4
3.03	0.4
3.06	0.4
3.09	0.4
3.12	0.4
3.15	0.4
3.18	0.4
3.21	0.4
3.24	0.4
3.27	0.4
3.30	0.4
3.33	0.4
3.36	0.4
3.39	0.4
3.42	0.4
3.45	0.4
3.48	0.4
3.51	0.4
3.54	0.4
3.57	0.4
3.60	0.4
3.63	0.4
3.66	0.4
3.69	0.4
3.72	0.4
3.75	0.4
3.78	0.4
3.81	0.4
3.84	0.4
3.87	0.4
3.90	0.4
3.93	0.4
3.96	0.4
3.99	0.4
4.02	0.4
4.05	0.4
4.08	0.4
4.11	0.4
4.14	0.4
4.17	0.4
4.20	0.4
4.23	0.4
4.26	0.4
4.29	0.4
4.32	0.4
4.35	0.4
4.38	0.4
4.41	0.4
4.44	0.4
4.47	0.4
4.50	0.4
4.53	0.4
4.56	0.4
4.59	0.4
4.62	0.4
4.65	0.4
4.68	0.4
4.71	0.4
4.74	0.4
4.77	0.4
4.80	0.4
4.83	0.4
4.86	0.4
4.89	0.4
4.92	0.4
4.95	0.4
4.98	0.4
5.01	0.4
5.04	0.4
5.07	0.4
5.10	0.4
5.13	0.4
5.16	0.4
5.19	0.4
5.22	0.4
5.25	0.4
5.28	0.4
5.31	0.4
5.34	0.4
5.37	0.4
5.40	0.4
5.43	0.4
5.46	0.4
5.49	0.4
5.52	0.4
5.55	0.4
5.58	0.4</

```
In [42]: 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 [43]: 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
..
413    0
414    0
415    0
416    0
417    0
Name: IsMinor, Length: 418, dtype: int32
```

A bar chart showing the survival rate for two groups: those who traveled alone (TravelAlone=0) and those who did not (TravelAlone=1). The y-axis represents the survival rate, ranging from 0.0 to 0.5. The bar for TravelAlone=0 has a height of 0.5, and the bar for TravelAlone=1 has a height of 0.3. Both bars include vertical error bars representing confidence intervals.

TravelAlone	Survived
0	0.5
1	0.3

Sex	Survived (approx.)	Lower CI (approx.)	Upper CI (approx.)
male	0.19	0.16	0.22
female	0.74	0.69	0.79

