

In [1]:

```
#1 IONOSPHERE
import pandas as pd
import numpy as np
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
df=pd.read_csv(r"C:\Users\G S R KARTHIK\Downloads\archive (4)\ionosphere_data.csv")
df
```

Out[1]:

	column_a	column_b	column_c	column_d	column_e	column_f	column_g	column_h	column_i	column_j	...	column_z	column_aa	column_
0	True	False	0.99539	-0.05889	0.85243	0.02306	0.83398	-0.37708	1.00000	0.03760	...	-0.51171	0.41078	-0.461
1	True	False	1.00000	-0.18829	0.93035	-0.36156	-0.10868	-0.93597	1.00000	-0.04549	...	-0.26569	-0.20468	-0.184
2	True	False	1.00000	-0.03365	1.00000	0.00485	1.00000	-0.12062	0.88965	0.01198	...	-0.40220	0.58984	-0.221
3	True	False	1.00000	-0.45161	1.00000	1.00000	0.71216	-1.00000	0.00000	0.00000	...	0.90695	0.51613	1.000
4	True	False	1.00000	-0.02401	0.94140	0.06531	0.92106	-0.23255	0.77152	-0.16399	...	-0.65158	0.13290	-0.532
...
346	True	False	0.83508	0.08298	0.73739	-0.14706	0.84349	-0.05567	0.90441	-0.04622	...	-0.04202	0.83479	0.001
347	True	False	0.95113	0.00419	0.95183	-0.02723	0.93438	-0.01920	0.94590	0.01606	...	0.01361	0.93522	0.049
348	True	False	0.94701	-0.00034	0.93207	-0.03227	0.95177	-0.03431	0.95584	0.02446	...	0.03193	0.92489	0.025
349	True	False	0.90608	-0.01657	0.98122	-0.01989	0.95691	-0.03646	0.85746	0.00110	...	-0.02099	0.89147	-0.077
350	True	False	0.84710	0.13533	0.73638	-0.06151	0.87873	0.08260	0.88928	-0.09139	...	-0.15114	0.81147	-0.048

351 rows × 35 columns

In [2]:

```
pd.set_option('display.max_rows',1000000000)
pd.set_option('display.max_columns',1000000000)
pd.set_option('display.width',95)
print('This dataframe has %d rows and %d columns'%(df.shape))
```

This dataframe has 351 rows and 35 columns

In [3]:

```
df.head(10)
```

Out[3]:

	column_a	column_b	column_c	column_d	column_e	column_f	column_g	column_h	column_i	column_j	column_k	column_l	column_m	colu
0	True	False	0.99539	-0.05889	0.85243	0.02306	0.83398	-0.37708	1.00000	0.03760	0.85243	-0.17755	0.59755	-0.
1	True	False	1.00000	-0.18829	0.93035	-0.36156	-0.10868	-0.93597	1.00000	-0.04549	0.50874	-0.67743	0.34432	-0.
2	True	False	1.00000	-0.03365	1.00000	0.00485	1.00000	-0.12062	0.88965	0.01198	0.73082	0.05346	0.85443	0.
3	True	False	1.00000	-0.45161	1.00000	1.00000	0.71216	-1.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.
4	True	False	1.00000	-0.02401	0.94140	0.06531	0.92106	-0.23255	0.77152	-0.16399	0.52798	-0.20275	0.56409	-0.
5	True	False	0.02337	-0.00592	-0.09924	-0.11949	-0.00763	-0.11824	0.14706	0.06637	0.03786	-0.06302	0.00000	0.
6	True	False	0.97588	-0.10602	0.94601	-0.20800	0.92806	-0.28350	0.85996	-0.27342	0.79766	-0.47929	0.78225	-0.
7	False	False	0.00000	0.00000	0.00000	0.00000	1.00000	-1.00000	0.00000	0.00000	-1.00000	-1.00000	0.00000	0.
8	True	False	0.96355	-0.07198	1.00000	-0.14333	1.00000	-0.21313	1.00000	-0.36174	0.92570	-0.43569	0.94510	-0.
9	True	False	-0.01864	-0.08459	0.00000	0.00000	0.00000	0.00000	0.11470	-0.26810	-0.45663	-0.38172	0.00000	0.

In [4]:

```
features_matrix=df.iloc[:,0:34]
```

In [5]:

```
target_vector=df.iloc[:,-1]
```

In [6]:

```
print('The features matrix has %d rows and %d columns'%(features_matrix.shape))
```

The features matrix has 351 rows and 34 columns

In [7]:

```
print('The target matrix has %d rows and %d columns'%(np.array(target_vector).reshape(-1,1).shape))
```

The target matrix has 351 rows and 1 columns

In [8]:

```
features_matrix_standardized=StandardScaler().fit_transform(features_matrix)
```

In [14]:

```
algorithm=LogisticRegression(penalty='l2',dual=False,tol=1e-4,C=1.0,fit_intercept=True,intercept_scaling=1,class_weight=None,random_s
```

In [15]:

```
Logistic_Regression_Model=algorithm.fit(features_matrix_standardized,target_vector)
```

In [16]:

```
observation=[[1,0,0.99539,-0.05889,0.8524299999999999,0.02306,0.8339799999999999,-0.37708,1.0,0.0376,0.8524299999999999,-0.17755,0.59
```

In [17]:

```
predictions=Logistic_Regression_Model.predict(observation)
```

In [19]:

```
print('The model predicted the observation to belong to class %s'%(predictions))
```

The model predicted the observation to belong to class ['g']

In [20]:

```
print('The algorithm was trained to predict one of the two classes :%s'%(algorithm.classes_))
```

The algorithm was trained to predict one of the two classes :['b' 'g']

In [21]:

```
print("""The model says the probability of the observation we passed belonging to class['b'] is %s""%(algorithm.predict_proba(observ
```

The model says the probability of the observation we passed belonging to class['b'] is 0.00777393160013784

In [22]:

```
print()
```

In [23]:

```
print("""The model says the probability of the observation we passed belonging to class['g'] is %s""%(algorithm.predict_proba(observ
```

The model says the probability of the observation we passed belonging to class['g'] is 0.9922260683998622

In []:

In [86]:

```
#2 DIGITS
import re
from sklearn.datasets import load_digits
from sklearn.model_selection import train_test_split
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import metrics
%matplotlib inline
digits=load_digits()
```

In [87]:

```
print("Image Data Shape",digits.data.shape)
print("Label Data Shape",digits.target.shape)
```

Image Data Shape (1797, 64)
Label Data Shape (1797,)

In [88]:

```
plt.figure(figsize=(20,4))
```

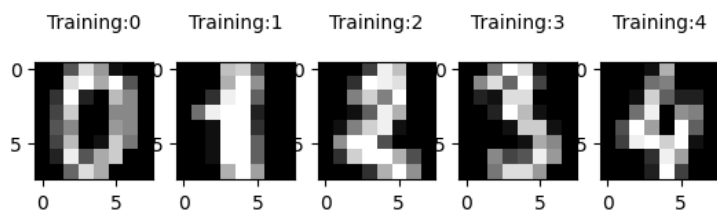
Out[88]:

<Figure size 2000x400 with 0 Axes>

<Figure size 2000x400 with 0 Axes>

In [89]:

```
for index,(image,label)in enumerate(zip(digits.data[0:5],digits.target[0:5])):
    plt.subplot(1,5,index+1)
    plt.imshow(np.reshape(image,(8,8)),cmap=plt.cm.gray)
    plt.title('Training:%i\n'%label,fontsize=10)
```



In [90]:

```
x_train,x_test,y_train,y_test=train_test_split(digits.data,digits.target,test_size=0.30,random_state=2)
```

In [91]:

```
print(x_train.shape)
```

(1257, 64)

In [92]:

```
print(y_train.shape)
```

(1257,)

In [93]:

```
print(x_test.shape)
```

(540, 64)

In [94]:

```
print(y_test.shape)
```

(540,)

In [95]:

```
from sklearn.linear_model import LogisticRegression
logisticRegr=LogisticRegression(max_iter=10000)
```

In [96]:

```
logisticRegr.fit(x_train,y_train)
```

Out[96]:

LogisticRegression(max_iter=10000)

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [97]:

```
print(logisticRegr.predict(x_test))
```

```
[4 0 9 1 8 7 1 5 1 6 6 7 6 1 5 5 8 6 2 7 4 6 4 1 5 2 9 5 4 6 5 6 3 4 0 9 9
 8 4 6 8 8 5 7 9 8 9 6 1 7 0 1 9 7 3 3 1 8 8 8 9 8 5 8 4 9 3 5 8 4 3 1 3 8
 7 3 3 0 8 7 2 8 5 3 8 7 6 4 6 2 2 0 1 1 5 3 5 7 1 8 2 2 6 4 6 7 3 7 3 9 4
 7 0 3 5 1 5 0 3 9 2 7 3 2 0 8 1 9 2 1 5 1 0 3 4 3 0 8 3 2 2 7 3 1 6 7 2 8
 3 1 1 6 4 8 2 1 8 4 1 3 1 1 9 5 4 8 7 4 8 9 5 7 6 9 4 0 4 0 0 9 0 6 5 8 8
 3 7 9 2 0 8 2 7 3 0 2 1 9 2 7 0 6 9 3 1 1 3 5 2 5 5 2 1 2 9 4 6 5 5 5 9 7
 1 5 9 6 3 7 1 7 5 1 7 2 7 5 5 4 8 6 6 2 8 7 3 7 8 0 9 5 7 4 3 4 1 0 3 3 5
 4 1 3 1 2 5 1 4 0 3 1 5 5 7 4 0 1 0 9 5 5 5 4 0 1 8 6 2 1 1 1 7 9 6 7 9 7
 0 4 9 6 9 2 7 2 1 0 8 2 8 6 5 7 8 4 5 7 8 6 4 2 6 9 3 0 0 8 0 6 6 7 1 4 5
 6 9 7 2 8 5 1 2 4 1 8 8 7 6 0 8 0 6 1 5 7 8 0 4 1 4 5 9 2 2 3 9 1 3 9 3 2
 8 0 6 5 6 2 5 2 3 2 6 1 0 7 6 0 6 2 7 0 3 2 4 2 3 6 9 7 7 0 3 5 4 1 2 2 1
 2 7 7 0 4 9 8 5 6 1 6 5 2 0 8 2 4 3 3 2 9 3 8 9 9 5 9 0 3 4 7 9 8 5 7 5 0
 5 3 5 0 2 7 3 0 4 3 6 6 1 9 6 3 4 6 4 6 7 2 7 6 3 0 3 0 1 3 6 1 0 4 3 8 4
 3 3 4 8 6 9 6 3 3 0 5 7 8 9 1 5 3 2 5 1 7 6 0 6 9 5 2 4 4 7 2 0 5 6 2 0 8
 4 4 4 7 1 0 4 1 9 2 1 3 0 5 3 9 8 2 6 0 0 4]
```

In [98]:

```
score=logisticRegr.score(x_test,y_test)
print(score)
```

```
0.9537037037037037
```

In []:

In [2]:

```
#3 GENDER SUBMISSION
import re
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn import metrics
%matplotlib inline
```

In [3]:

```
df=pd.read_csv(r"C:\Users\G S R KARTHIK\Downloads\gender_submission.csv")
print(df)
```

	PassengerId	Survived
0	892	0
1	893	1
2	894	0
3	895	0
4	896	1
..
413	1305	0
414	1306	1
415	1307	0
416	1308	0
417	1309	0

```
[418 rows x 2 columns]
```

In [4]:

```
plt.figure(figsize=(20,2))
```

Out[4]:

```
<Figure size 2000x200 with 0 Axes>
```

```
<Figure size 2000x200 with 0 Axes>
```

In [5]:

```
df.describe()
```

Out[5]:

	PassengerId	Survived
count	418.000000	418.000000
mean	1100.500000	0.363636
std	120.810458	0.481622
min	892.000000	0.000000
25%	996.250000	0.000000
50%	1100.500000	0.000000
75%	1204.750000	1.000000
max	1309.000000	1.000000

In [6]:

```
df.isnull().any()
```

Out[6]:

PassengerId False
Survived False
dtype: bool

In [7]:

```
pd.set_option('display.max_rows',1000000000)  
pd.set_option('display.max_columns',1000000000)  
pd.set_option('display.width',95)  
print('This dataframe has %d rows and %d columns'%(df.shape))
```

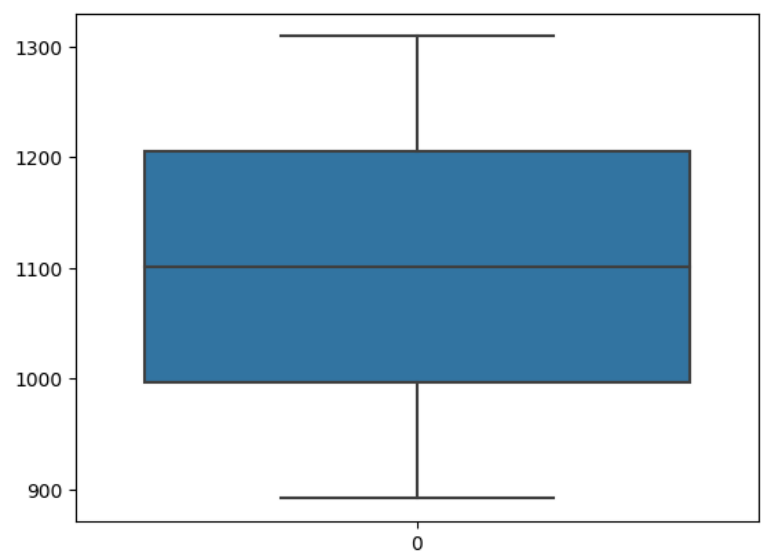
This dataframe has 418 rows and 2 columns

In [8]:

```
sns.boxplot(df[ 'PassengerId' ])
```

Out[8]:

<Axes: >



In [9]:

```
df.head(10)
```

Out[9]:

	PassengerId	Survived
0	892	0
1	893	1
2	894	0
3	895	0
4	896	1
5	897	0
6	898	1
7	899	0
8	900	1
9	901	0

In [13]:

```
features_matrix=df.iloc[:,0:1]
```

In [14]:

```
target_vector=df.iloc[:,-1]
```

In [15]:

```
print('The features matrix has %d rows and %d columns'%(features_matrix.shape))
```

The features matrix has 418 rows and 1 columns

In [16]:

```
print('The target matrix has %d rows and %d columns'%(np.array(target_vector).reshape(-1,1).shape))
```

The target matrix has 418 rows and 1 columns

In [18]:

```
from sklearn.preprocessing import StandardScaler  
features_matrix_standardized=StandardScaler().fit_transform(features_matrix)
```

In [19]:

```
ss_weight=None,random_state=None,solver='lbfgs',max_iter=100,multi_class='auto',verbose=0,warm_start=False,n_jobs=None,l1_ratio=None)
```

In [20]:

```
Logistic_Regression_Model=algorithm.fit(features_matrix_standardized,target_vector)
```

In [23]:

```
observation=[[1]]
```

In [24]:

```
predictions=Logistic_Regression_Model.predict(observation)
```

In [25]:

```
print('The model predicted the observation to belong to class %d'%(predictions))
```

The model predicted the observation to belong to class 0

In [26]:

```
print('The algorithm was trained to predict one of the two classes :%d'%(algorithm.classes_))
```

The algorithm was trained to predict one of the two classes :[0 1]

In [31]:

```
print("""The model says the probability of the observation we passed belonging to class['0'] is %s""%(algorithm.predict_proba(observ
```

The model says the probability of the observation we passed belonging to class['0'] is 0.6474324251144166

In [32]:

```
print("""The model says the probability of the observation we passed belonging to class['1'] is %s""%(algorithm.predict_proba(observ
```

The model says the probability of the observation we passed belonging to class['1'] is 0.35256757488558343

In []:

PROBLEM STATEMENT: TO PREDICT AND ANALYZE WHICH GENDER HAS A HIGH CHANCE OF SURVIVAL AT THE TIME OF DISASTER

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

In [2]:

```
train_df = pd.read_csv(r"C:\Users\G S R KARTHIK\Downloads\train.gender_submission (1).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 x 12 columns

In [3]:

```
test_df = pd.read_csv(r"C:\Users\G S R KARTHIK\Downloads\test.gender_submission (1).csv")
test_df
```

Out[3]:

	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 [4]:

```
train_df.shape
```

Out[4]:

(891, 12)

In [5]:

```
test_df.shape
```

Out[5]:

(418, 11)

In [6]:

train_df.info

Out[6]:

```
<bound method DataFrame.info 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 [7]:

test_df.info

Out[7]:

```
<bound method DataFrame.info of      PassengerId  Pclass      Name
0             892        3      Kelly, Mr. James  \
1             893        3  Wilkes, Mrs. James (Ellen Needs)
2             894        2      Myles, Mr. Thomas Francis
3             895        3      Wirz, Mr. Albert
4             896        3  Hirvonen, Mrs. Alexander (Helga E Lindqvist)
..          ...     ...
413          1305        3      Spector, Mr. Woolf
414          1306        1  Oliva y Ocana, Dona. Fermina
415          1307        3  Saether, Mr. Simon Sivertsen
416          1308        3      Ware, Mr. Frederick
417          1309        3  Peter, Master. Michael J

      Sex  Age  SibSp  Parch      Ticket    Fare Cabin Embarked
0    male  34.5     0     0      330911   7.8292   NaN        Q
1  female  47.0     1     0      363272   7.0000   NaN        S
2    male  62.0     0     0      240276   9.6875   NaN        Q
3    male  27.0     0     0      315154   8.6625   NaN        S
4  female  22.0     1     1      3101298  12.2875   NaN        S
..      ...     ...     ...     ...     ...     ...     ...
413    male  NaN     0     0      A.5. 3236   8.0500   NaN        S
414  female  39.0     0     0      PC 17758  108.9000  C105        C
415    male  38.5     0     0  SOTON/O.Q. 3101262   7.2500   NaN        S
416    male  NaN     0     0      359309   8.0500   NaN        S
417    male  NaN     1     1         2668  22.3583   NaN        C

[418 rows x 11 columns]>
```

In [8]:

```
#to find missing values  
train_df.isnull().sum()
```

Out[8]:

```
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 [9]:

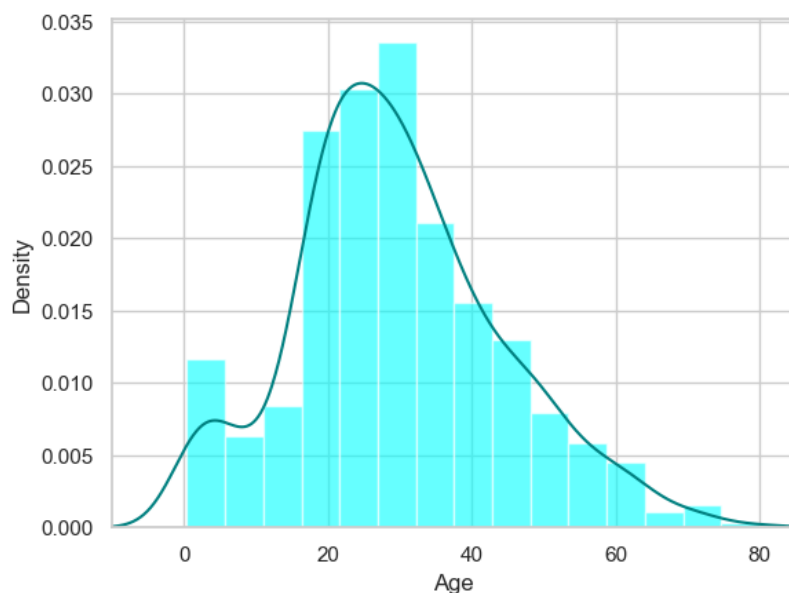
```
test_df.isnull().sum()
```

Out[9]:

```
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 [10]:

```
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 [11]:

```
print(train_df["Age"].mean(skipna=True))  
print(train_df["Age"].median(skipna=True))
```

```
29.69911764705882  
28.0
```

In [12]:

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

77.10437710437711

In [13]:

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

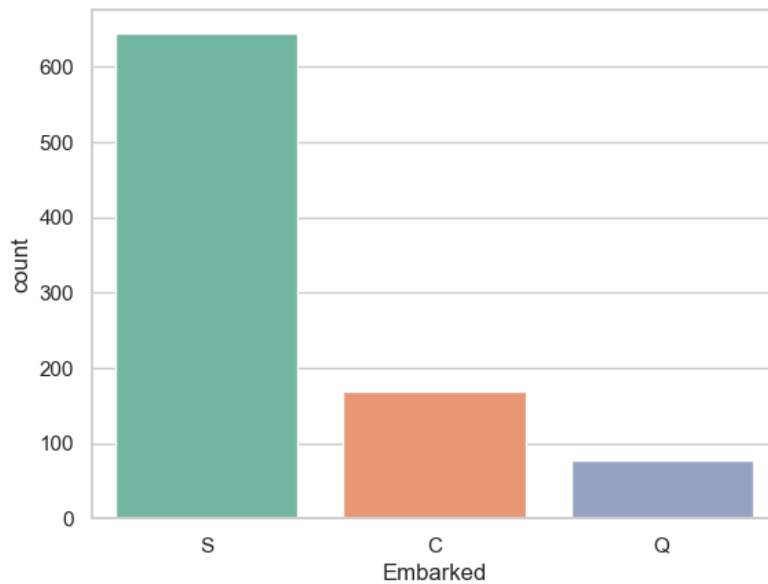
0.22446689113355783

In [14]:

```
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 [15]:

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

S

In [16]:

```
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 [17]:

```
train_data.isnull().sum()
```

Out[17]:

```
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 [18]:

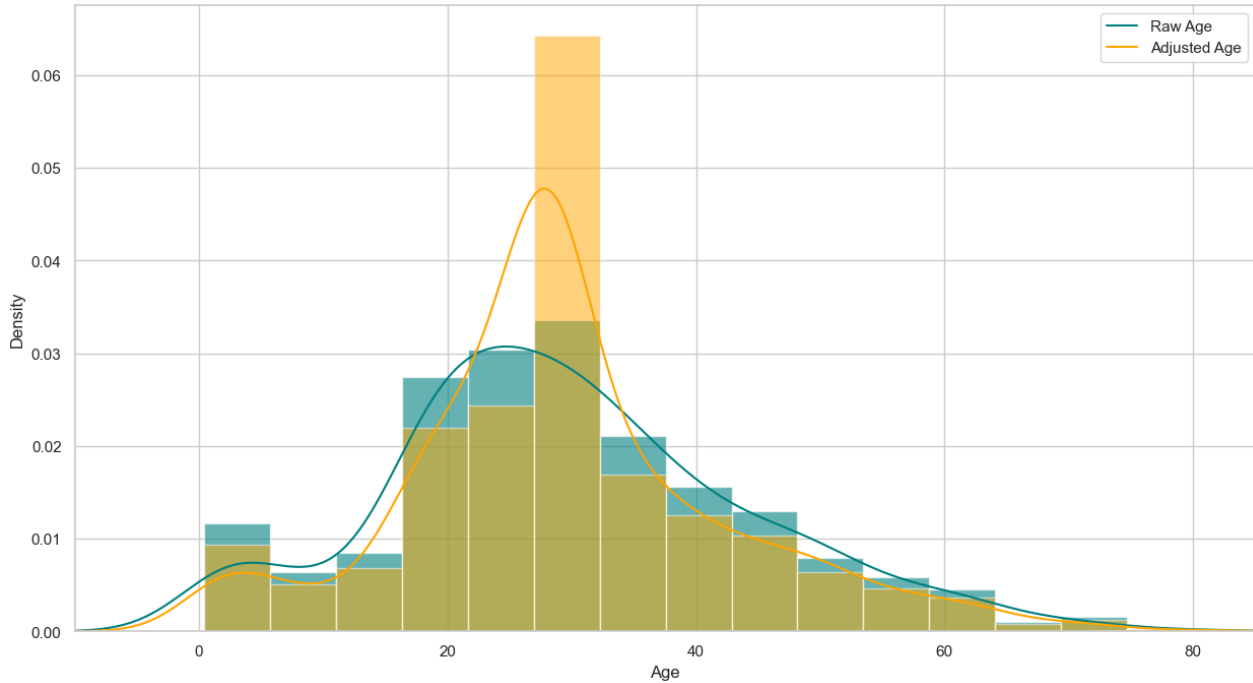
```
train_data.head()
```

Out[18]:

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 [19]:

```
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=0.5)
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 [20]:

```
## 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 [21]:

```
#create categorical 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[21]:

	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

In [22]:

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

Out[22]:

```
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 [23]:

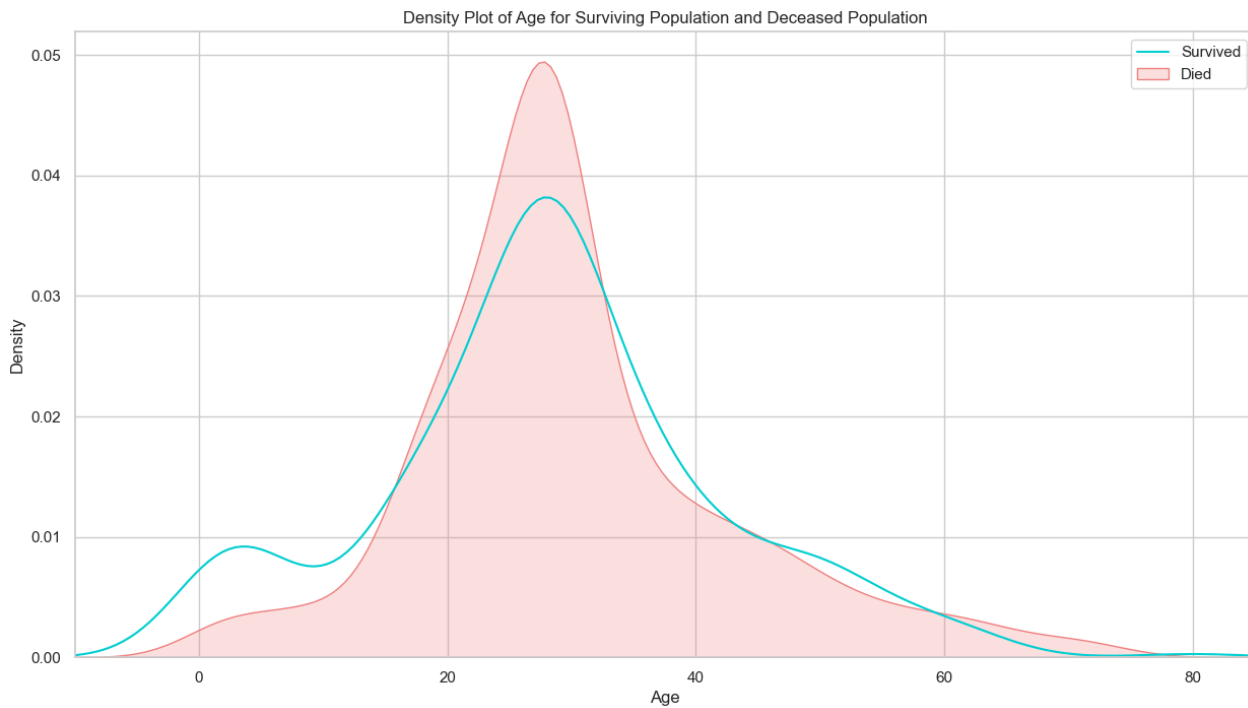
```
test_data = test_df.copy()
test_data["Age"].fillna(train_df["Age"].median(skipna=True), inplace=True)
test_data["Fare"].fillna(train_df["Fare"].median(skipna=True), inplace=True)
test_data.drop('Cabin', axis=1, inplace=True)
test_data['TravelAlone']=np.where((test_data["SibSp"]+test_data["Parch"])>0, 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()
```

Out[23]:

	Age	Fare	TravelAlone	Pclass_1	Pclass_2	Pclass_3	Embarked_C	Embarked_Q	Embarked_S	Sex_male
0	34.5	7.8292	1	False	False	True	False	True	False	True
1	47.0	7.0000	0	False	False	True	False	False	True	False
2	62.0	9.6875	1	False	True	False	False	True	False	True
3	27.0	8.6625	1	False	False	True	False	False	True	True
4	22.0	12.2875	0	False	False	True	False	False	True	False

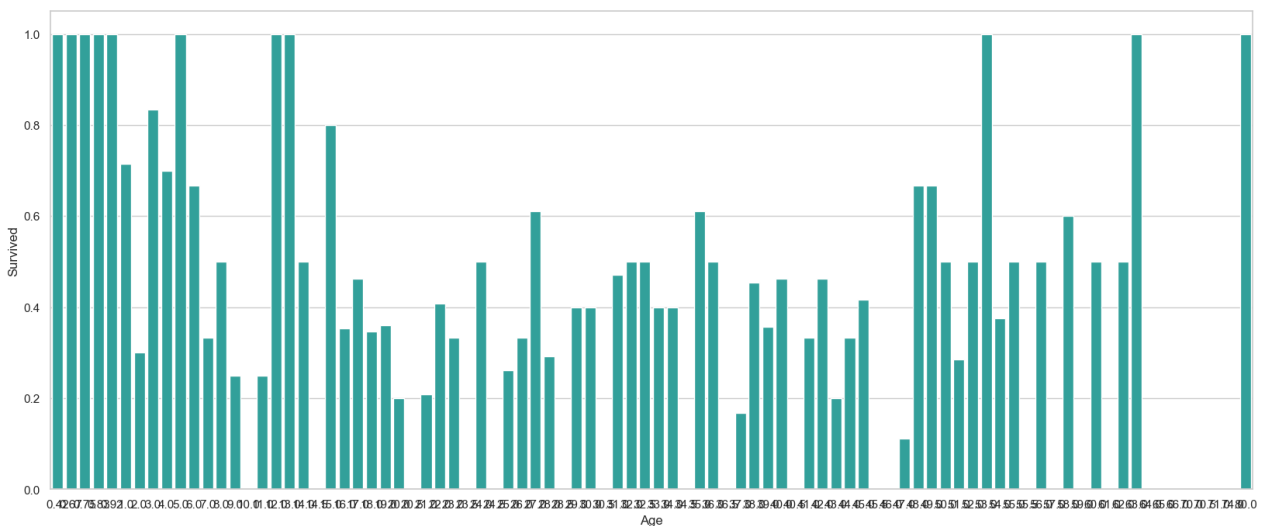
In [24]:

```
#EDA
plt.figure(figsize=(15,8))
ax = sns.kdeplot(final_train["Age"][final_train.Survived == 1], color="darkturquoise")
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()
```



In [27]:

```
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 [28]:

```
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
```

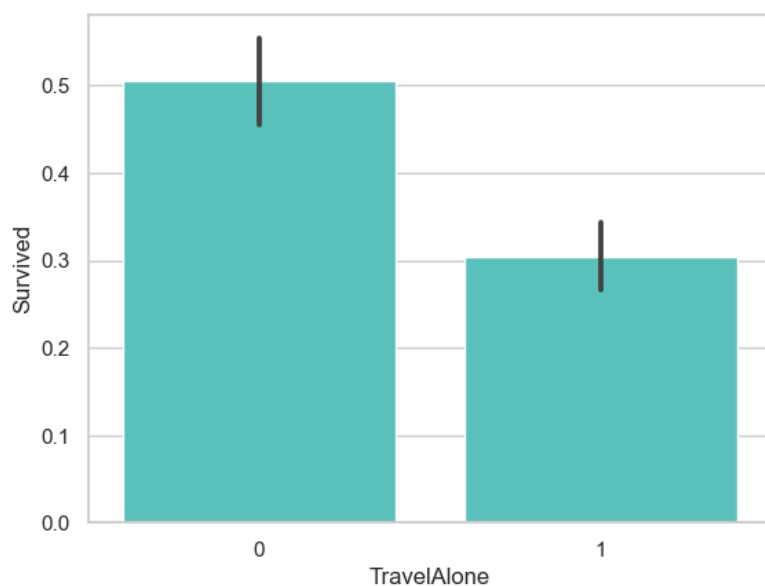
In [29]:

```
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
```

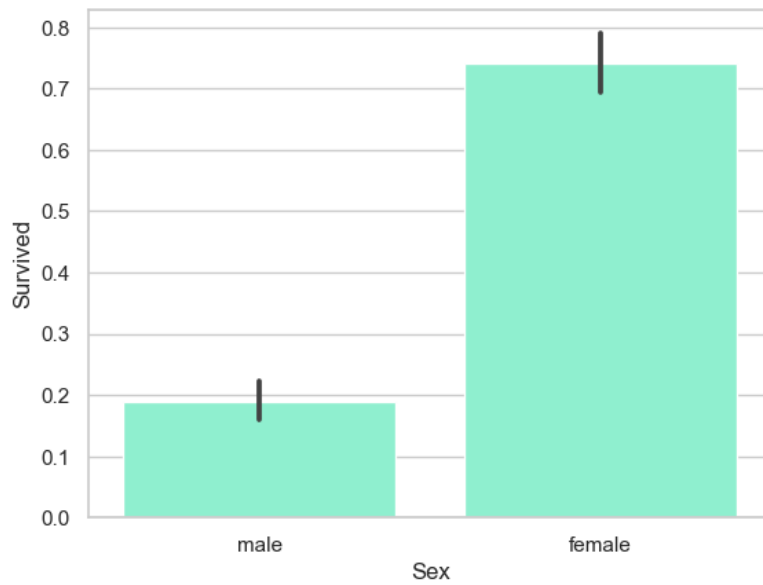
In [30]:

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



In [31]:

```
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()
```



In []: