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

In [2]:
```python
df=pd.read_csv(r"C:\Users\P. VIJAY KUMAR\Downloads\heart disease (1).csv")
df
```

Out[2]:

| | male | age | education | currentSmoker | cigsPerDay | BPMeds | prevalentStroke | prevalentHyp | diabetes | totChol | s |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 39 | 4.0 | 0 | 0.0 | 0.0 | 0 | 0 | 0 | 195.0 | |
| 1 | 0 | 46 | 2.0 | 0 | 0.0 | 0.0 | 0 | 0 | 0 | 250.0 | |
| 2 | 1 | 48 | 1.0 | 1 | 20.0 | 0.0 | 0 | 0 | 0 | 245.0 | |
| 3 | 0 | 61 | 3.0 | 1 | 30.0 | 0.0 | 0 | 1 | 0 | 225.0 | |
| 4 | 0 | 46 | 3.0 | 1 | 23.0 | 0.0 | 0 | 0 | 0 | 285.0 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 4233 | 1 | 50 | 1.0 | 1 | 1.0 | 0.0 | 0 | 1 | 0 | 313.0 | |
| 4234 | 1 | 51 | 3.0 | 1 | 43.0 | 0.0 | 0 | 0 | 0 | 207.0 | |
| 4235 | 0 | 48 | 2.0 | 1 | 20.0 | NaN | 0 | 0 | 0 | 248.0 | |
| 4236 | 0 | 44 | 1.0 | 1 | 15.0 | 0.0 | 0 | 0 | 0 | 210.0 | |
| 4237 | 0 | 52 | 2.0 | 0 | 0.0 | 0.0 | 0 | 0 | 0 | 269.0 | |

4238 rows × 16 columns

In [3]:
```python
df.head()
```

Out[3]:

| | male | age | education | currentSmoker | cigsPerDay | BPMeds | prevalentStroke | prevalentHyp | diabetes | totChol | sysE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 39 | 4.0 | 0 | 0.0 | 0.0 | 0 | 0 | 0 | 195.0 | 106 |
| 1 | 0 | 46 | 2.0 | 0 | 0.0 | 0.0 | 0 | 0 | 0 | 250.0 | 121 |
| 2 | 1 | 48 | 1.0 | 1 | 20.0 | 0.0 | 0 | 0 | 0 | 245.0 | 127 |
| 3 | 0 | 61 | 3.0 | 1 | 30.0 | 0.0 | 0 | 1 | 0 | 225.0 | 150 |
| 4 | 0 | 46 | 3.0 | 1 | 23.0 | 0.0 | 0 | 0 | 0 | 285.0 | 130 |

In [4]:    1  df.describe()

Out[4]:

|       | male        | age         | education   | currentSmoker | cigsPerDay  | BPMeds      | prevalentStroke | prevalentHy |
|-------|-------------|-------------|-------------|---------------|-------------|-------------|-----------------|-------------|
| count | 4238.000000 | 4238.000000 | 4133.000000 | 4238.000000   | 4209.000000 | 4185.000000 | 4238.000000     | 4238.00000  |
| mean  | 0.429212    | 49.584946   | 1.978950    | 0.494101      | 9.003089    | 0.029630    | 0.005899        | 0.31052     |
| std   | 0.495022    | 8.572160    | 1.019791    | 0.500024      | 11.920094   | 0.169584    | 0.076587        | 0.46276     |
| min   | 0.000000    | 32.000000   | 1.000000    | 0.000000      | 0.000000    | 0.000000    | 0.000000        | 0.00000     |
| 25%   | 0.000000    | 42.000000   | 1.000000    | 0.000000      | 0.000000    | 0.000000    | 0.000000        | 0.00000     |
| 50%   | 0.000000    | 49.000000   | 2.000000    | 0.000000      | 0.000000    | 0.000000    | 0.000000        | 0.00000     |
| 75%   | 1.000000    | 56.000000   | 3.000000    | 1.000000      | 20.000000   | 0.000000    | 0.000000        | 1.00000     |
| max   | 1.000000    | 70.000000   | 4.000000    | 1.000000      | 70.000000   | 1.000000    | 1.000000        | 1.00000     |

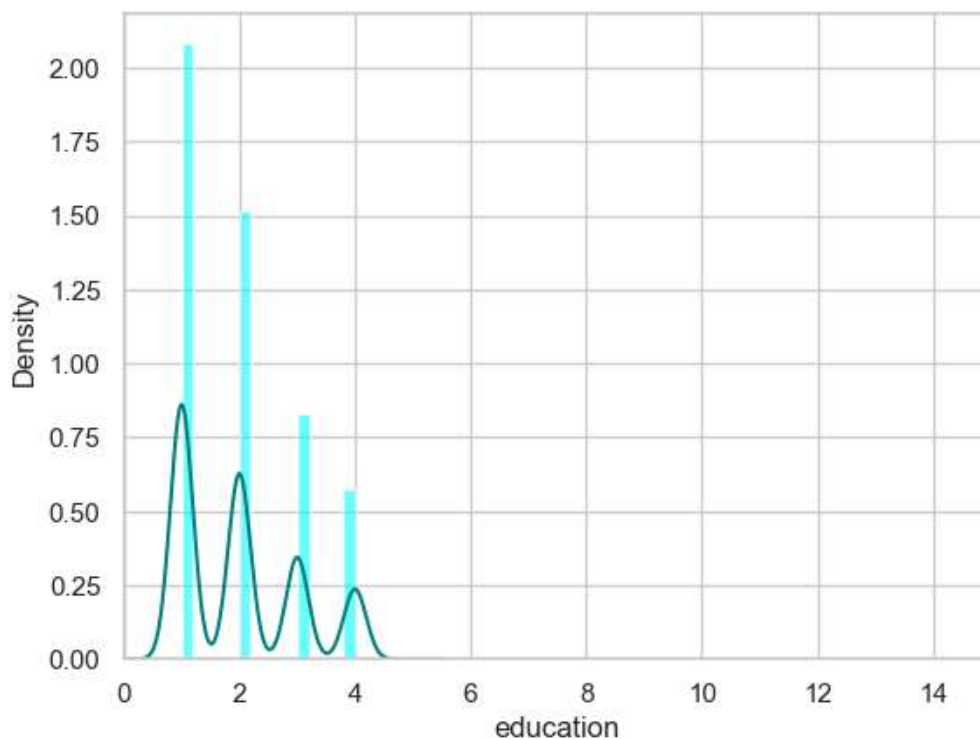In [6]:    1  df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4238 entries, 0 to 4237
Data columns (total 16 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   male             4238 non-null   int64
 1   age              4238 non-null   int64
 2   education        4133 non-null   float64
 3   currentSmoker    4238 non-null   int64
 4   cigsPerDay       4209 non-null   float64
 5   BPMeds           4185 non-null   float64
 6   prevalentStroke  4238 non-null   int64
 7   prevalentHyp     4238 non-null   int64
 8   diabetes         4238 non-null   int64
 9   totChol          4188 non-null   float64
 10  sysBP            4238 non-null   float64
 11  diaBP            4238 non-null   float64
 12  BMI              4219 non-null   float64
 13  heartRate        4237 non-null   float64
 14  glucose          3850 non-null   float64
 15  TenYearCHD       4238 non-null   int64
dtypes: float64(9), int64(7)
memory usage: 529.9 KB
```

In [7]:
```python
1  df.isnull().sum()
```

Out[7]:
```
male               0
age                0
education        105
currentSmoker      0
cigsPerDay        29
BPMeds            53
prevalentStroke    0
prevalentHyp       0
diabetes           0
totChol           50
sysBP              0
diaBP              0
BMI               19
heartRate          1
glucose          388
TenYearCHD         0
dtype: int64
```

In [8]:
```python
1  ax = df["education"].hist(bins=15, density=True, stacked=True, color='cyan', alpha=0.6)
2  df["education"].plot(kind='density', color='teal')
3  ax.set(xlabel='education')
4  plt.xlim(-0,15)
5  plt.show()
```



In [9]:
```python
1  print(df["education"].mean(skipna=True))
2  print(df["education"].median(skipna=True))
3
```

```
1.9789499153157513
2.0
```

In [10]:
```python
print((df['glucose'].isnull().sum()/df.shape[0])*100)
```

9.155261915998112

In [11]:
```python
print((df['totChol'].isnull().sum()/df.shape[0])*100)
```

1.1798017932987257

In [12]:
```python
print(df['totChol'].value_counts())
sns.countplot(x='totChol', data=df, palette='Set2')
plt.show()
```

```
totChol
240.0    85
220.0    70
260.0    62
210.0    61
232.0    59
         ..
392.0     1
405.0     1
359.0     1
398.0     1
119.0     1
Name: count, Length: 248, dtype: int64
```



In [13]:
```python
print(df['totChol'].value_counts().idxmax())
```
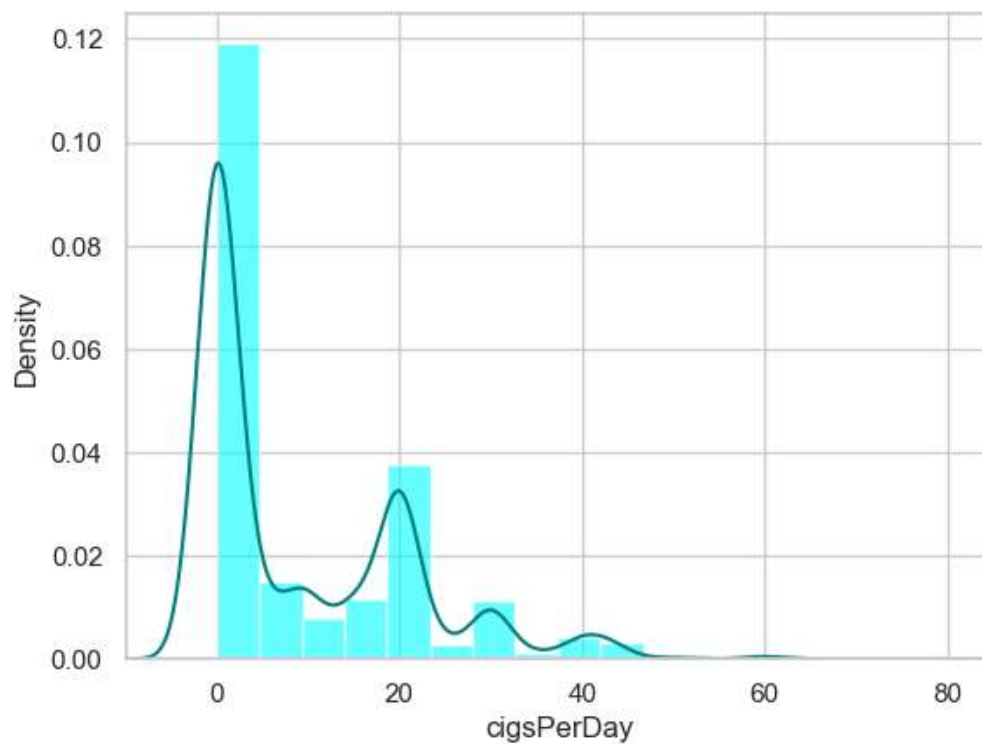
240.0

```
In [14]:   1  data = df.copy()
           2  data["education"].fillna(df["education"].median(skipna=True), inplace=True)
           3  data["totChol"].fillna(df['totChol'].value_counts().idxmax(), inplace=True)
           4  data.drop('glucose', axis=1, inplace=True)
```

```
In [15]:   1  data.isnull().sum()
           2
```

```
Out[15]:  male                0
          age                 0
          education           0
          currentSmoker       0
          cigsPerDay         29
          BPMeds             53
          prevalentStroke     0
          prevalentHyp        0
          diabetes            0
          totChol             0
          sysBP               0
          diaBP               0
          BMI                19
          heartRate           1
          TenYearCHD          0
          dtype: int64
```

```
In [16]:   1  ax = df["cigsPerDay"].hist(bins=15, density=True, stacked=True, color='cyan', alpha=0.6)
           2  df["cigsPerDay"].plot(kind='density', color='teal')
           3  ax.set(xlabel='cigsPerDay')
           4  plt.xlim(-10,85)
           5  plt.show()
```

In [17]:
```python
print(df["cigsPerDay"].mean(skipna=True))
print(df["cigsPerDay"].median(skipna=True))
```

```
9.003088619624615
0.0
```

In [18]:
```python
print((df['BPMeds'].isnull().sum()/df.shape[0])*100)
```

```
1.2505899008966492
```

In [19]:
```python
print((df['BMI'].isnull().sum()/df.shape[0])*100)
```
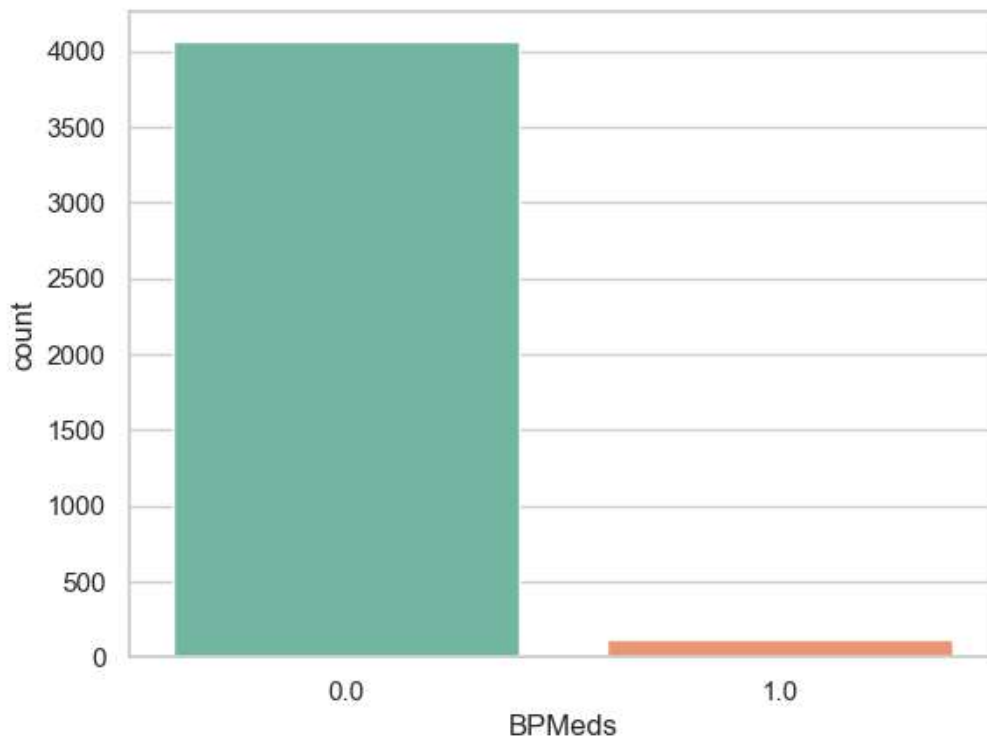
```
0.4483246814535158
```

In [20]:
```python
print((df['heartRate'].isnull().sum()/df.shape[0])*100)
```

```
0.023596035865974516
```

In [21]:
```python
print(df['BPMeds'].value_counts())
sns.countplot(x='BPMeds', data=df, palette='Set2')
plt.show()
```

```
BPMeds
0.0    4061
1.0     124
Name: count, dtype: int64
```



In [22]:
```python
print(df['heartRate'].value_counts().idxmax())
```

```
75.0
```

In [23]:
```python
1  data = df.copy()
2  data["cigsPerDay"].fillna(df["cigsPerDay"].median(skipna=True), inplace=True)
3  data["BPMeds"].fillna(df['BPMeds'].value_counts().idxmax(), inplace=True)
4  data["education"].fillna(df["education"].median(skipna=True), inplace=True)
5  data["totChol"].fillna(df['totChol'].value_counts().idxmax(), inplace=True)
6  data.drop('glucose', axis=1, inplace=True)
7  data.drop('BMI', axis=1, inplace=True)
8  data.drop('heartRate', axis=1, inplace=True)
```

In [24]:
```python
1  data.isnull().sum()
```

Out[24]:
```
male               0
age                0
education          0
currentSmoker      0
cigsPerDay         0
BPMeds             0
prevalentStroke    0
prevalentHyp       0
diabetes           0
totChol            0
sysBP              0
diaBP              0
TenYearCHD         0
dtype: int64
```
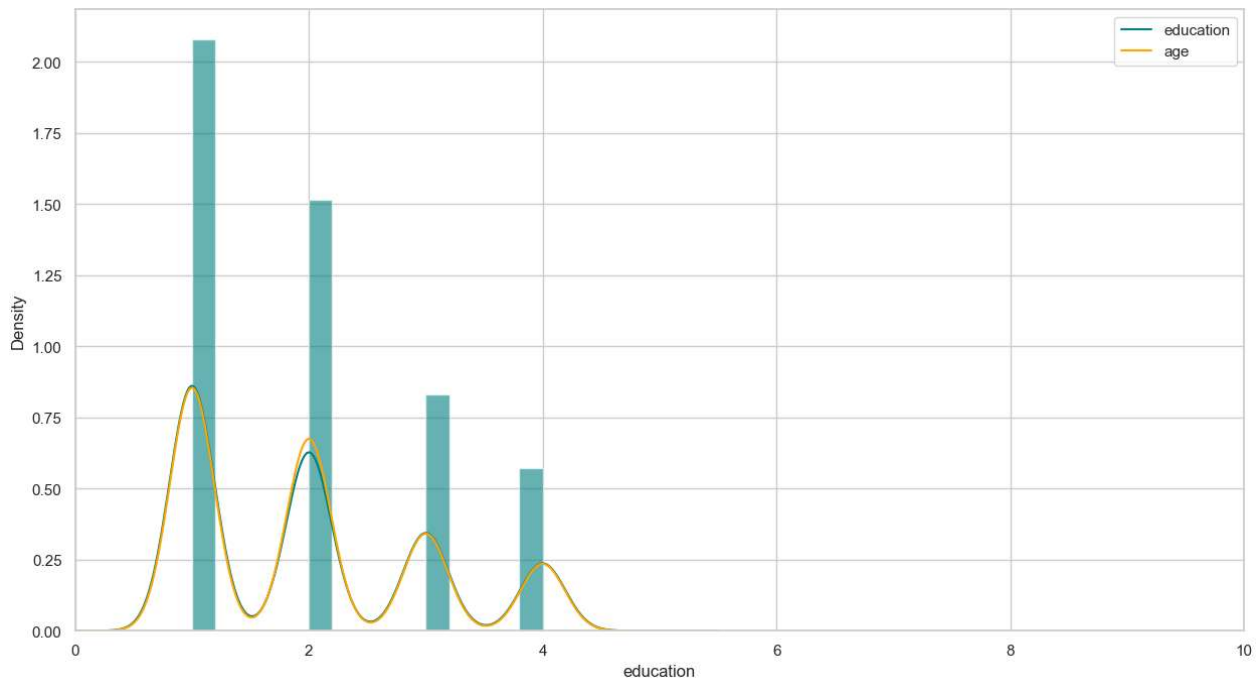
In [25]:
```python
1  data.head()
```

Out[25]:

|   | male | age | education | currentSmoker | cigsPerDay | BPMeds | prevalentStroke | prevalentHyp | diabetes | totChol | sysB |
|---|------|-----|-----------|---------------|------------|--------|-----------------|--------------|----------|---------|------|
| 0 | 1 | 39 | 4.0 | 0 | 0.0 | 0.0 | 0 | 0 | 0 | 195.0 | 106 |
| 1 | 0 | 46 | 2.0 | 0 | 0.0 | 0.0 | 0 | 0 | 0 | 250.0 | 121 |
| 2 | 1 | 48 | 1.0 | 1 | 20.0 | 0.0 | 0 | 0 | 0 | 245.0 | 127 |
| 3 | 0 | 61 | 3.0 | 1 | 30.0 | 0.0 | 0 | 1 | 0 | 225.0 | 150 |
| 4 | 0 | 46 | 3.0 | 1 | 23.0 | 0.0 | 0 | 0 | 0 | 285.0 | 130 |

In [27]:
```python
1  plt.figure(figsize=(15,8))
2  ax = df["education"].hist(bins=15, density=True, stacked=True, color='teal', alpha=0.6)
3  df["education"].plot(kind='density', color='teal')
4  ax = data["education"].hist(bins=15, density=True, stacked=True, color='orange', alpha=0)
5  data["education"].plot(kind='density', color='orange')
6  ax.legend(['education', 'age'])
7  ax.set(xlabel='education')
8  plt.xlim(-0,10)
9  plt.show()
10
```



In [28]:
```python
1  data['Disease']=np.where((data["prevalentHyp"]+data["prevalentStroke"])>0, 0, 1)
2  data.drop('prevalentHyp', axis=1, inplace=True)
3  data.drop('prevalentStroke', axis=1, inplace=True)
```

In [29]:
```python
1  training=pd.get_dummies(data, columns=["currentSmoker","totChol","sysBP"])
2  training.drop('TenYearCHD', axis=1, inplace=True)
3  training.drop('male', axis=1, inplace=True)
4  training.drop('diaBP', axis=1, inplace=True)
5  final_train = training
6  final_train.head()
7
```

Out[29]:

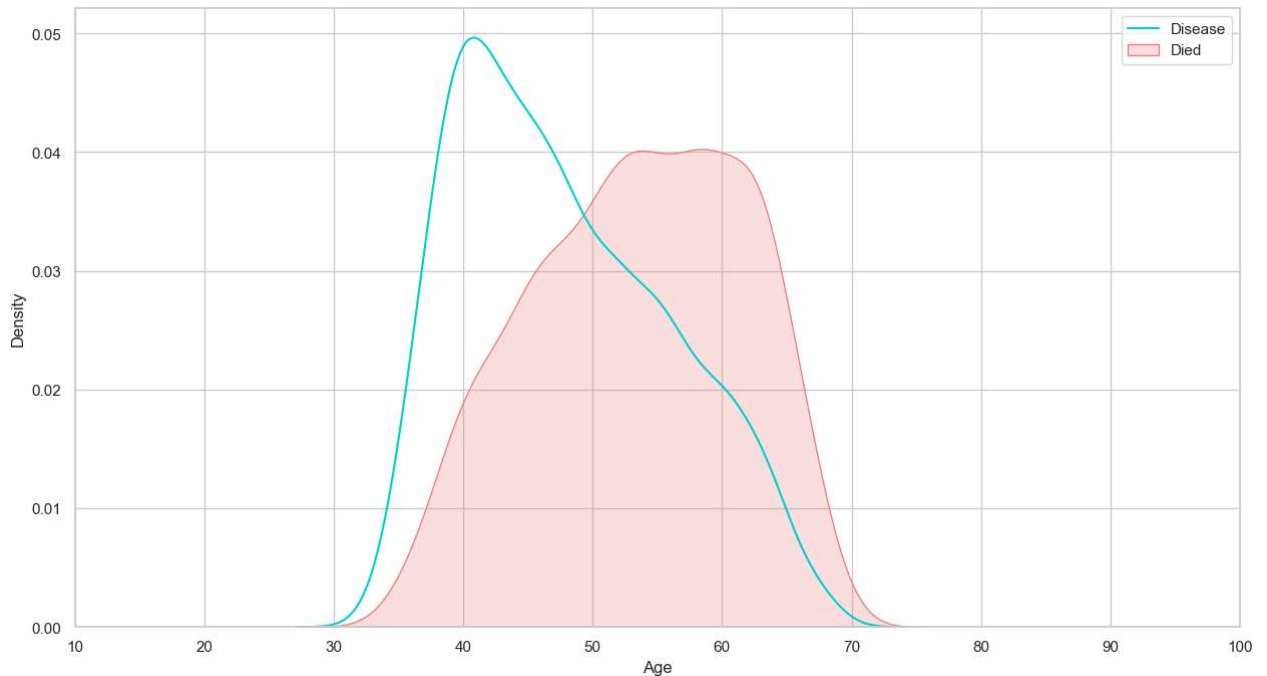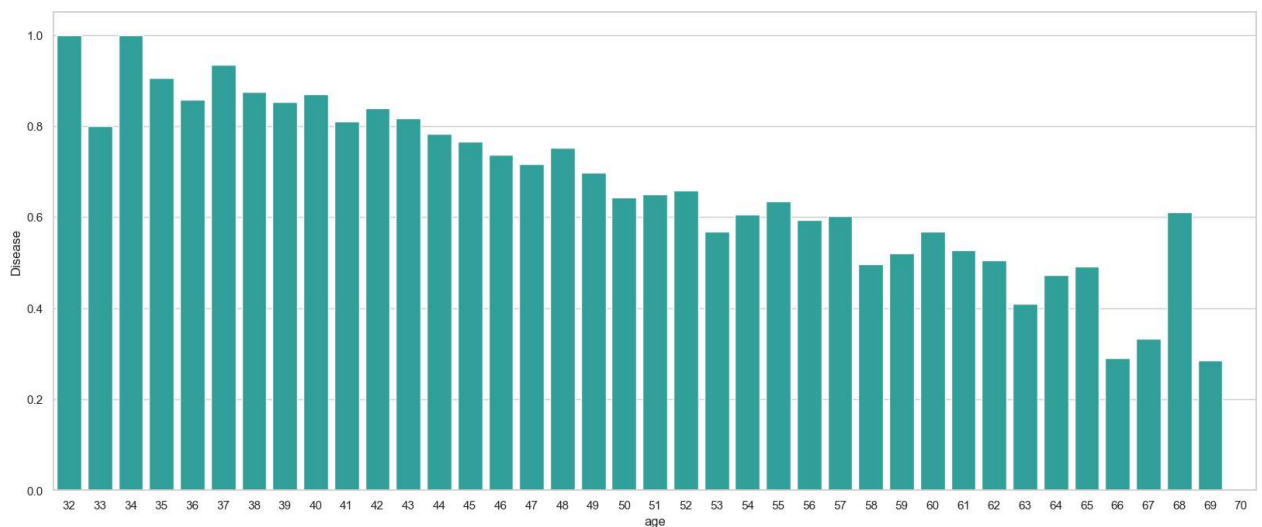| | age | education | cigsPerDay | BPMeds | diabetes | Disease | currentSmoker_0 | currentSmoker_1 | totChol_107.0 | totCho |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 39 | 4.0 | 0.0 | 0.0 | 0 | 1 | True | False | False | |
| 1 | 46 | 2.0 | 0.0 | 0.0 | 0 | 1 | True | False | False | |
| 2 | 48 | 1.0 | 20.0 | 0.0 | 0 | 1 | False | True | False | |
| 3 | 61 | 3.0 | 30.0 | 0.0 | 0 | 0 | False | True | False | |
| 4 | 46 | 3.0 | 23.0 | 0.0 | 0 | 1 | False | True | False | |

5 rows × 490 columns

In [31]:
```python
#EDA
plt.figure(figsize=(15,8))
ax = sns.kdeplot(final_train["age"][final_train.Disease == 1], color="darkturquoise")
sns.kdeplot(final_train["age"][final_train.Disease == 0], color="lightcoral", shade=True)
plt.legend(['Disease', 'Died'])
ax.set(xlabel='Age')
plt.xlim(10,100)
plt.show()
```
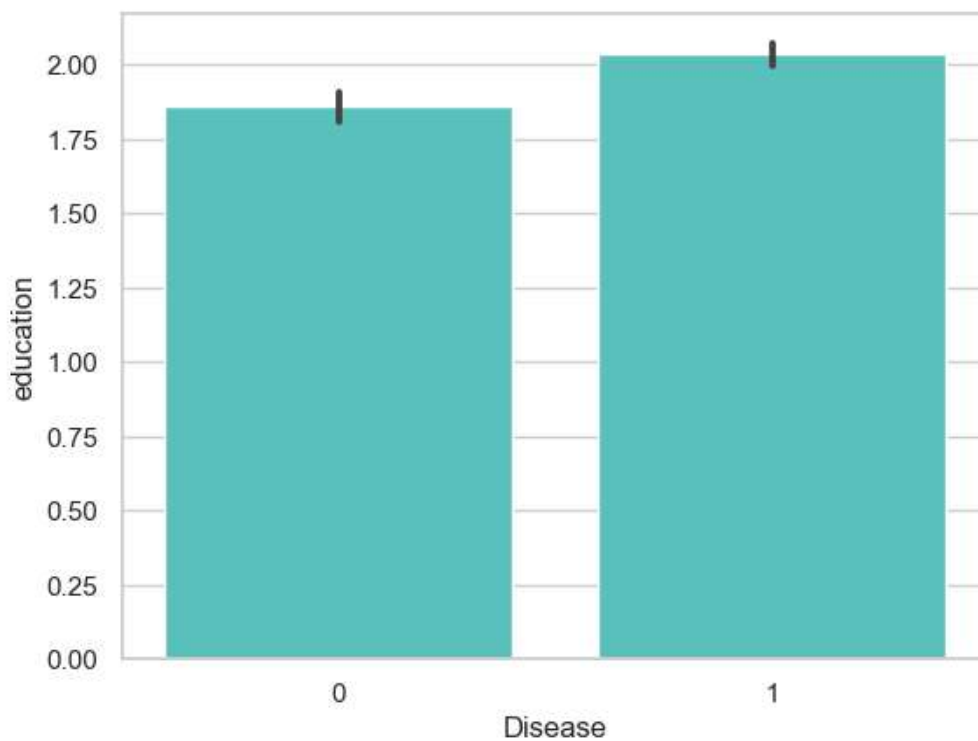


In [33]:
```python
plt.figure(figsize=(20,8))
avg_survival_byage = final_train[["age", "Disease"]].groupby(['age'], as_index=False).mean
g = sns.barplot(x='age', y='Disease', data=avg_survival_byage, color="LightSeaGreen")
plt.show()
```

In [34]:
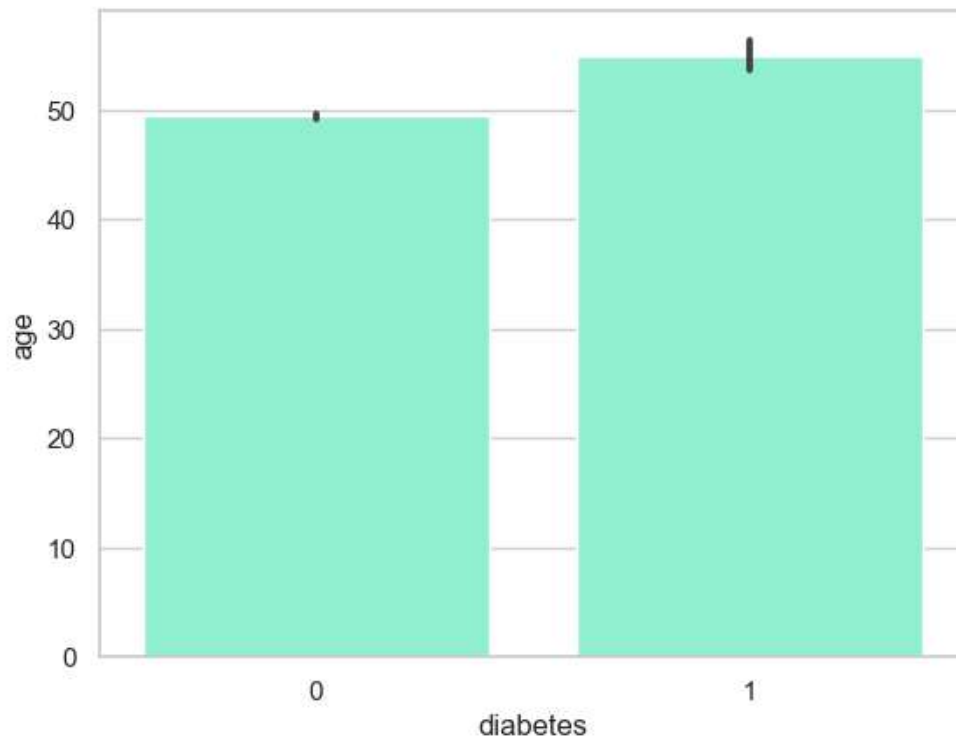```python
final_train['IsMinor']=np.where(final_train['age']<=16, 1, 0)
print(final_train['IsMinor'])
```

```
0       0
1       0
2       0
3       0
4       0
       ..
4233    0
4234    0
4235    0
4236    0
4237    0
Name: IsMinor, Length: 4238, dtype: int32
```

In [35]:
```python
sns.barplot(x='Disease', y='education', data=final_train, color="mediumturquoise")
plt.show()
```

```
In [36]:  1  import seaborn as sns
          2  import matplotlib.pyplot as plt
          3  # Assuming 'train_df' is your DataFrame containing the data
          4  sns.barplot(x='diabetes', y='age', data=df, color='aquamarine')
          5  plt.show()
```



```
In [ ]:  1
```