

In [2]:

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.preprocessing import StandardScaler
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 sns.set(style="white")
8 sns.set(style="whitegrid", color_codes=True)
9 import warnings
10 warnings.simplefilter(action='ignore')
```

In [3]:

```
1 df=pd.read_csv(r"C:\Users\teppa\Downloads\Heart Disease.csv")
2 df
```

Out[3]:

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp
0	1	39	4.0	0	0.0	0.0	0	0
1	0	46	2.0	0	0.0	0.0	0	0
2	1	48	1.0	1	20.0	0.0	0	0
3	0	61	3.0	1	30.0	0.0	0	1
4	0	46	3.0	1	23.0	0.0	0	0
...	...	...	...	...	...	...	...	...
4233	1	50	1.0	1	1.0	0.0	0	1
4234	1	51	3.0	1	43.0	0.0	0	0
4235	0	48	2.0	1	20.0	NaN	0	0
4236	0	44	1.0	1	15.0	0.0	0	0
4237	0	52	2.0	0	0.0	0.0	0	0

4238 rows × 9 columns



In [4]:

```
1 df.head()
```

Out[4]:

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp	d
0	1	39	4.0	0	0.0	0.0	0	0	
1	0	46	2.0	0	0.0	0.0	0	0	
2	1	48	1.0	1	20.0	0.0	0	0	
3	0	61	3.0	1	30.0	0.0	0	1	
4	0	46	3.0	1	23.0	0.0	0	0	



In [5]: 1 df.tail()

Out[5]:

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp
4233	1	50	1.0	1	1.0	0.0	0	1
4234	1	51	3.0	1	43.0	0.0	0	0
4235	0	48	2.0	1	20.0	NaN	0	0
4236	0	44	1.0	1	15.0	0.0	0	0
4237	0	52	2.0	0	0.0	0.0	0	0


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]: 1 df.describe()

Out[7]:

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevale
count	4238.000000	4238.000000	4133.000000	4238.000000	4209.000000	4185.000000	423
mean	0.429212	49.584946	1.978950	0.494101	9.003089	0.029630	
std	0.495022	8.572160	1.019791	0.500024	11.920094	0.169584	
min	0.000000	32.000000	1.000000	0.000000	0.000000	0.000000	
25%	0.000000	42.000000	1.000000	0.000000	0.000000	0.000000	
50%	0.000000	49.000000	2.000000	0.000000	0.000000	0.000000	
75%	1.000000	56.000000	3.000000	1.000000	20.000000	0.000000	
max	1.000000	70.000000	4.000000	1.000000	70.000000	1.000000	



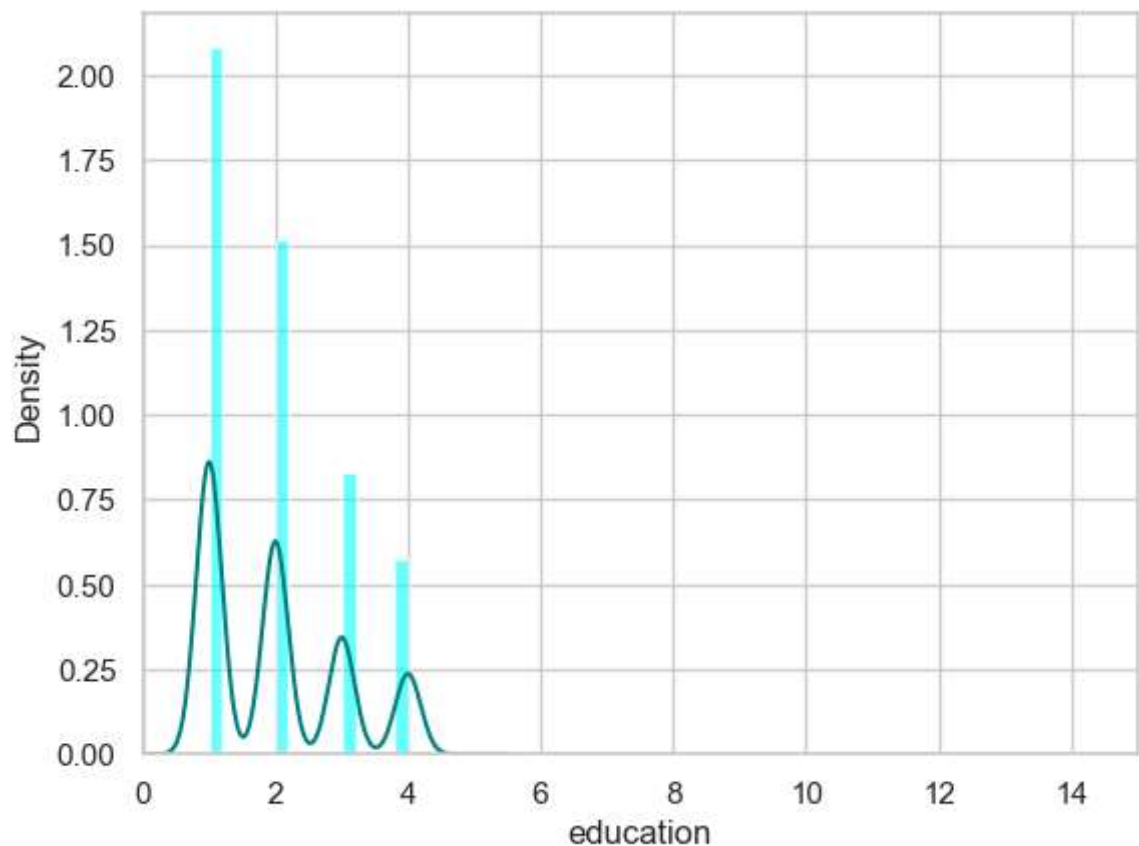
In [8]: 1 df.isnull().sum()

Out[8]: 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 [9]: 1 df.describe().any()
```

```
Out[9]: male                True
age                True
education          True
currentSmoker      True
cigsPerDay         True
BPMeds             True
prevalentStroke    True
prevalentHyp       True
diabetes           True
totChol            True
sysBP             True
diaBP             True
BMI               True
heartRate          True
glucose            True
TenYearCHD         True
dtype: bool
```

```
In [10]: 1 ax=df["education"].hist(bins=15,density=True,stacked=True,color='cyan',al
2 df["education"].plot(kind='density',color='teal')
3 ax.set(xlabel='education')
4 plt.xlim(-0,15)
5 plt.show()
```

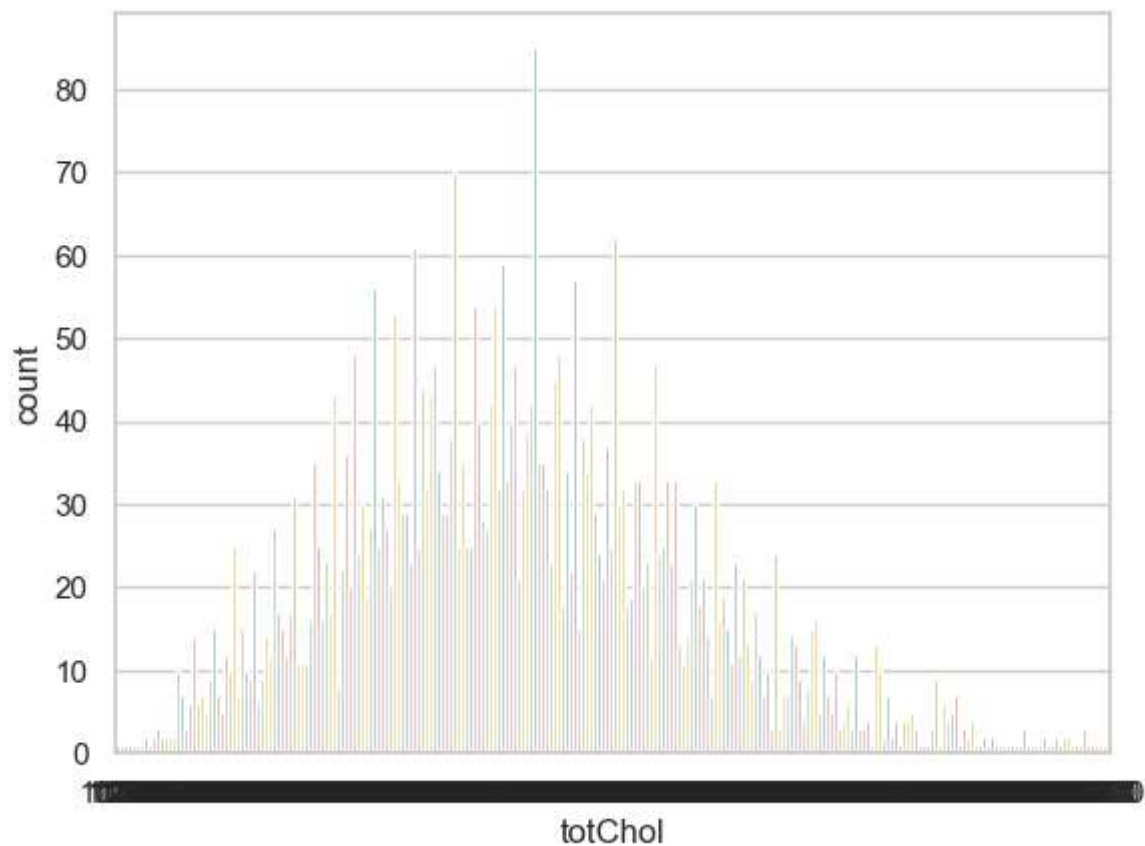


```
In [11]: 1 print(df["education"].mean(skipna=True))
          2 print(df["education"].median(skipna=True))
```

```
1.9789499153157513
2.0
```

```
In [12]: 1 print(df['totChol'].value_counts())
          2 sns.countplot(x='totChol',data=df,palette='Set2')
          3 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]: 1 print(df['totChol'].value_counts().idxmax())
```

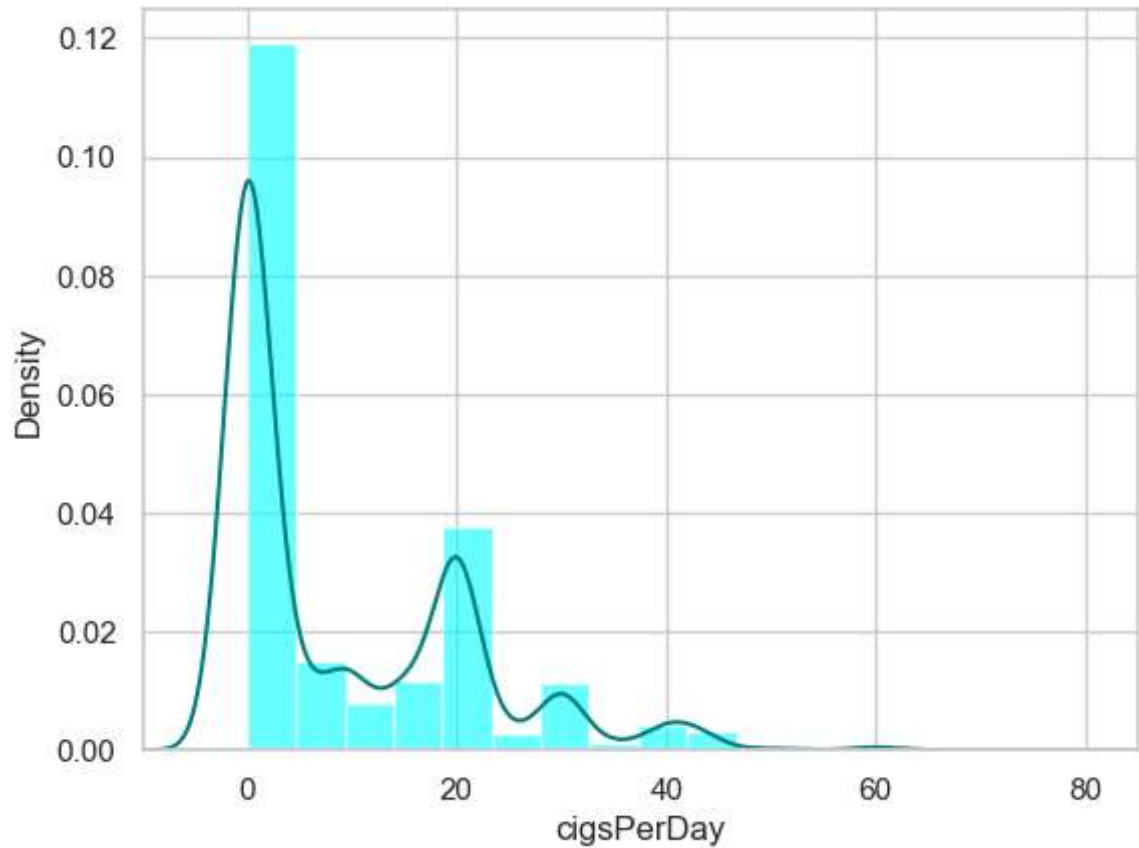
```
240.0
```

```
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()
```

```
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',a
2 df["cigsPerDay"].plot(kind='density',color='teal')
3 ax.set(xlabel='cigsPerDay')
4 plt.xlim(-10,85)
5 plt.show()
6
7
```



```
In [17]: 1 print(df["cigsPerDay"].mean(skipna=True))
2 print(df["cigsPerDay"].median(skipna=True))
```

```
9.003088619624615
0.0
```

```
In [18]: 1 print((df['BPMeds'].isnull().sum()/df.shape[0]*100))
```

```
1.2505899008966492
```

```
In [19]: 1 print((df['BMI'].isnull().sum()/df.shape[0]*100))
```

```
0.4483246814535158
```

```
In [20]: 1 print((df['heartRate'].isnull().sum()/df.shape[0]*100))
```

```
0.023596035865974516
```

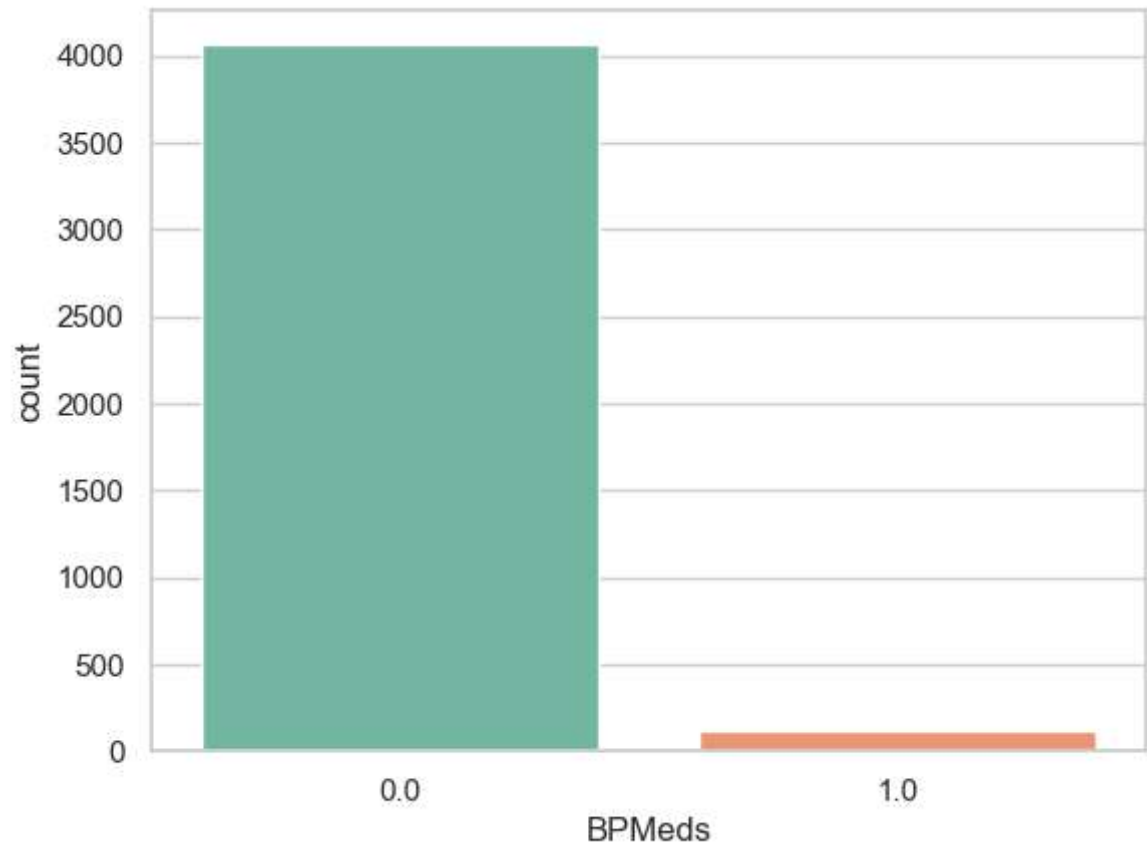
```
In [21]: 1 print(df['BPMeds'].value_counts())
2 sns.countplot(x='BPMeds',data=df,palette='Set2')
3 plt.show()
```

BPMeds

0.0 4061

1.0 124

Name: count, dtype: int64

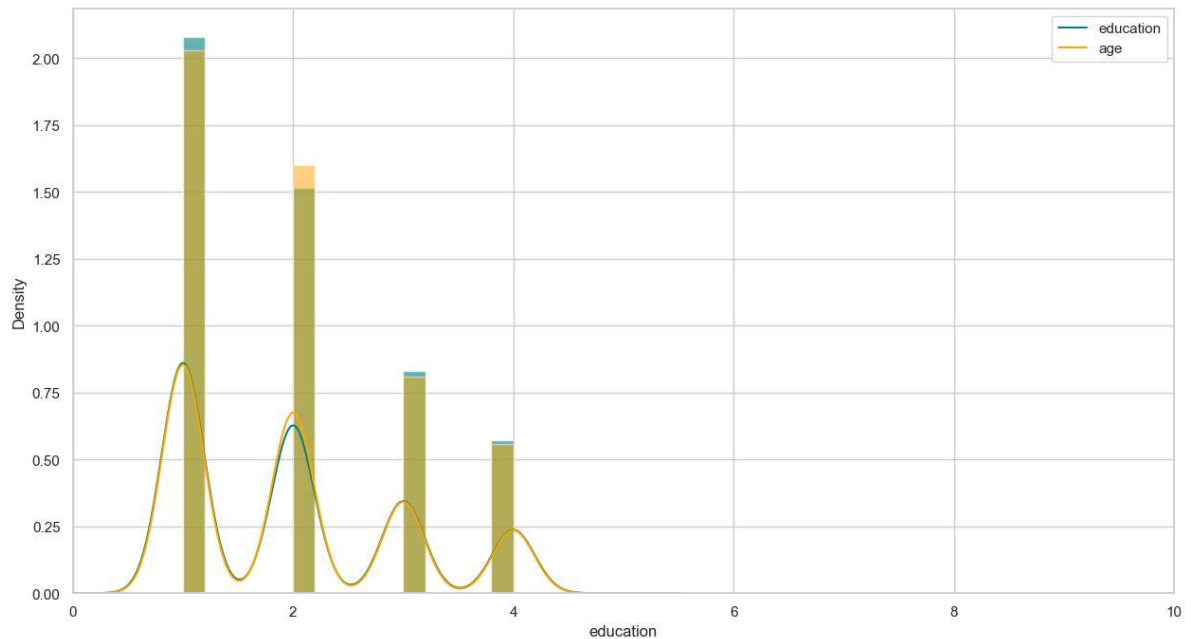


```
In [22]: 1 print(df['heartRate'].value_counts().idxmax())
```

75.0



```
In [23]: 1 plt.figure(figsize=(15,8))
2 ax=df["education"].hist(bins=15,density=True,stacked=True,color='teal',al
3 df["education"].plot(kind='density',color='teal')
4 ax=data["education"].hist(bins=15,density=True,stacked=True,color='orange
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()
```



```
In [27]: 1 df['Disease']=np.where((df["prevalentHyp"]+df["prevalentStroke"])>0,0,1)
2 df.drop('prevalentHyp',axis=1,inplace=True)
3 df.drop('prevalentStroke',axis=1,inplace=True)
```

```
In [28]: 1 training=pd.get_dummies(df,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()
```

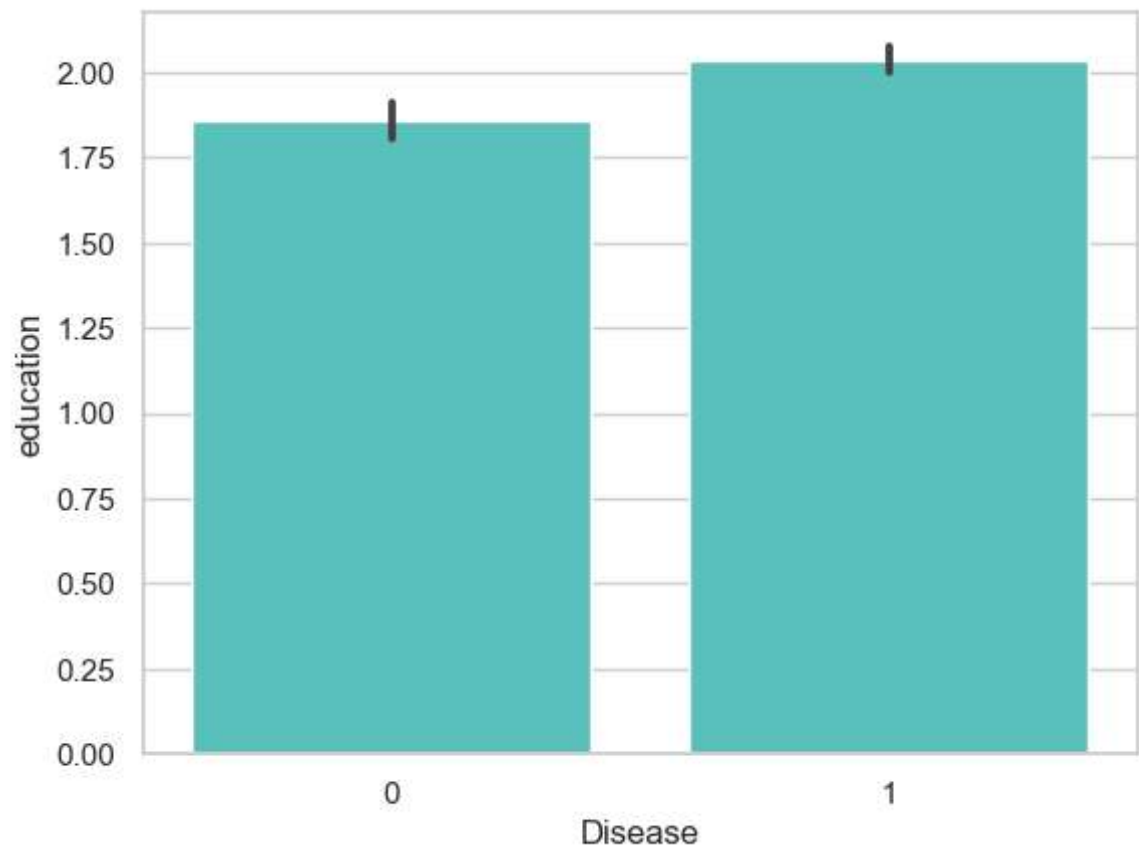
```
Out[28]:
```

	age	education	cigsPerDay	BPMeds	diabetes	BMI	heartRate	glucose	Disease	currentS
0	39	4.0	0.0	0.0	0	26.97	80.0	77.0	1	
1	46	2.0	0.0	0.0	0	28.73	95.0	76.0	1	
2	48	1.0	20.0	0.0	0	25.34	75.0	70.0	1	
3	61	3.0	30.0	0.0	0	28.58	65.0	103.0	0	
4	46	3.0	23.0	0.0	0	23.10	85.0	85.0	1	

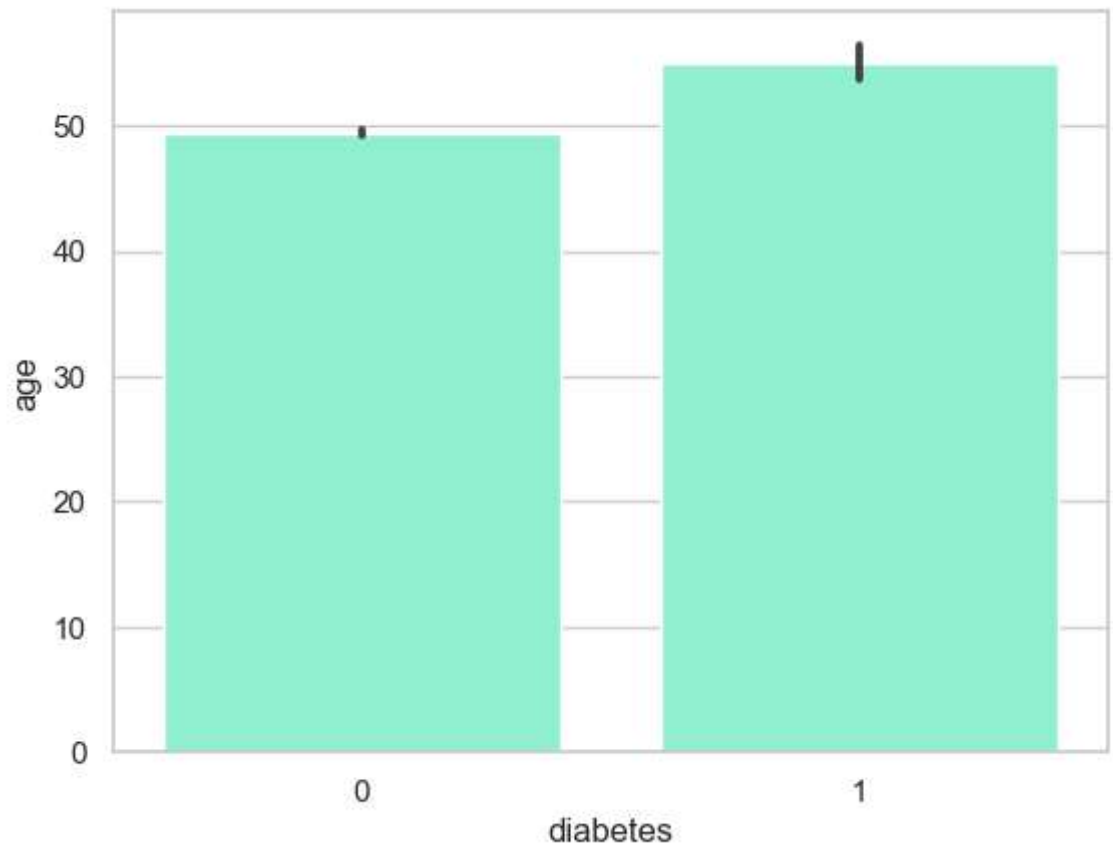
5 rows × 493 columns



```
In [31]: 1 sns.barplot(x='Disease', y='education', data=final_train, color="mediumturquoise")
2         plt.show()
```



```
In [32]: 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
```