

Problem Statement

To predict the risk of heart diseases using Logistic Regression

In [1]:

```
1 import numpy as np
2 import pandas as pd
3
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 import warnings
8 warnings.filterwarnings("ignore")
9
```

In [2]:

```
1 df=pd.read_csv(r"C:\Users\91949\Downloads\framingham.csv")
2 df
```

Out[2]:

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

4238 rows × 16 columns



In [4]:

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 [5]:

1 df.describe()

Out[5]:

	male	age	education	currentSmoker	cigsPerDay	BPMeds	pre
count	4238.000000	4238.000000	4133.000000	4238.000000	4209.000000	4185.000000	
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 [6]:

1 df.shape

Out[6]:

(4238, 16)

In [7]:

```
1 df.isna().any()
```

Out[7]:

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

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]:

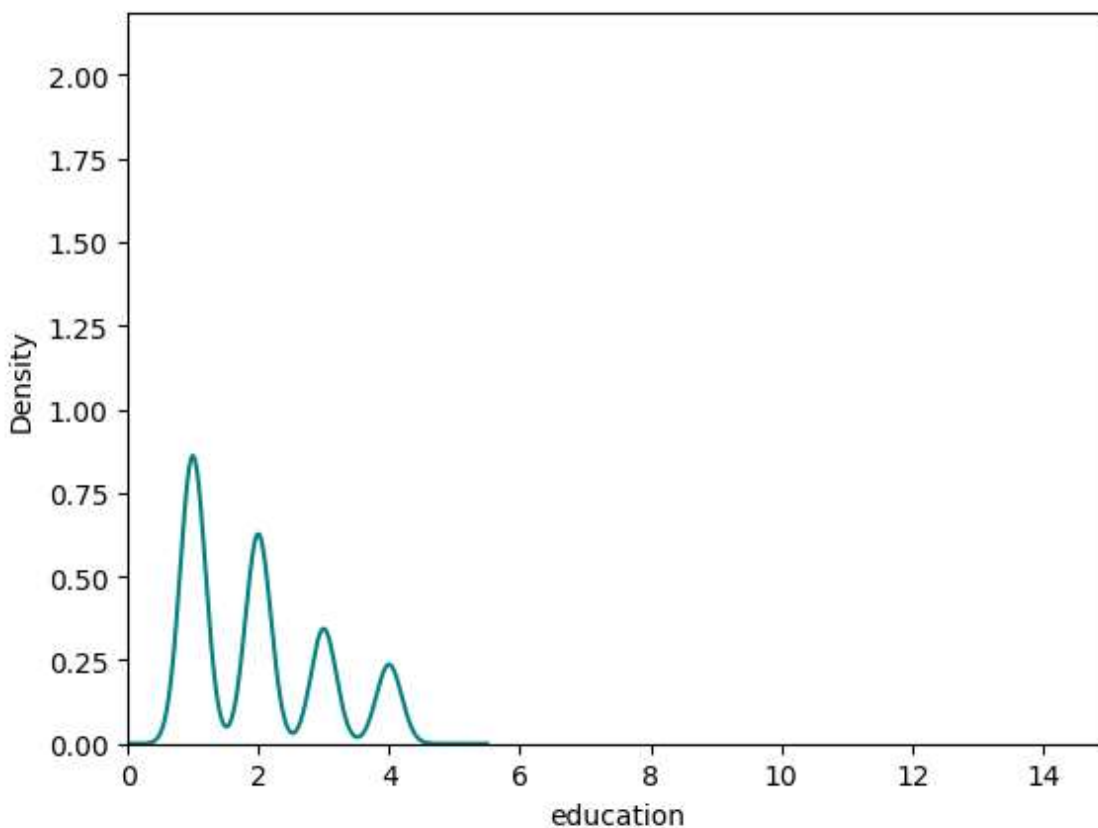
```
1 df.describe().any()
```

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', alpha=0
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['glucose'].isnull().sum()/df.shape[0]*100))
```

9.155261915998112

In [13]:

```
1 print((df['totChol'].isnull().sum()/df.shape[0]*100))
```

1.1798017932987257

In [14]:

```

1 print(df['totChol'].value_counts())
2 sns.countplot(x='totChol', data=df,palette='Set2')
3 plt.show()

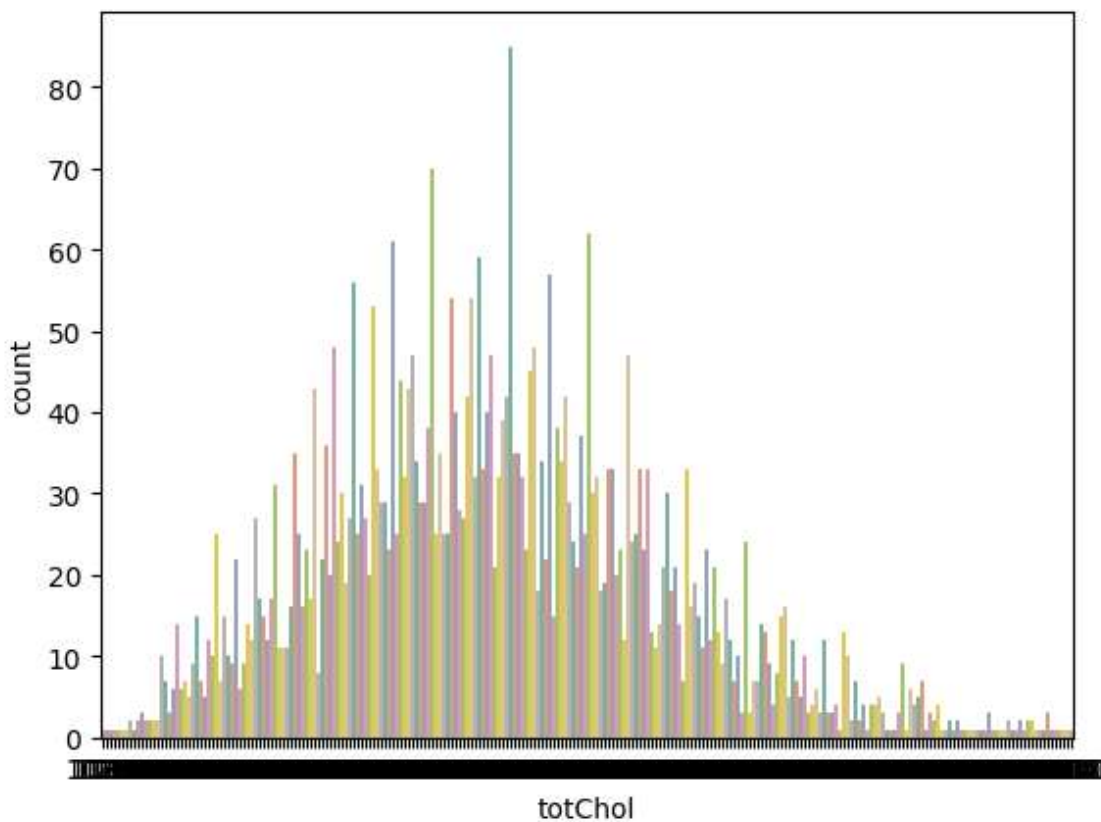
```

```

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: totChol, Length: 248, dtype: int64



In [15]:

```
1 print(df['totChol'].value_counts().idxmax())
```

240.0

In [16]:

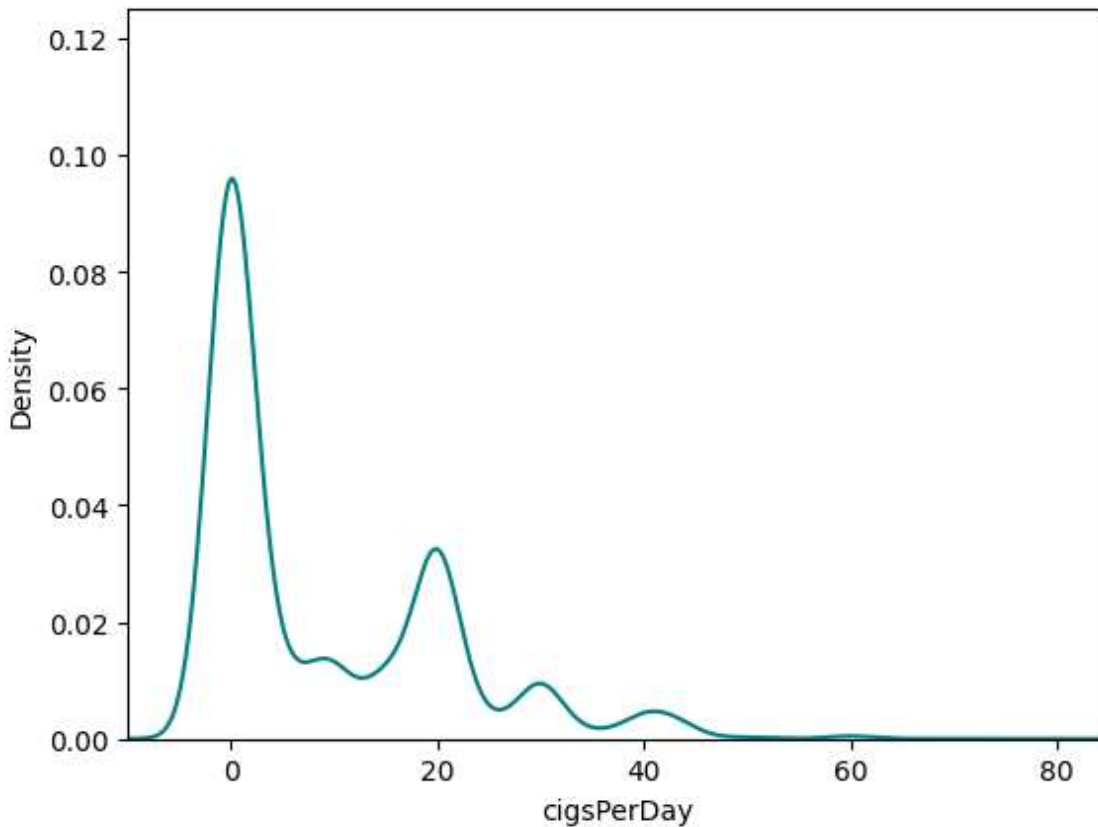
```

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 [17]:

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



In [18]:

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

```
9.003088619624615
0.0
```

In [19]:

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

```
1.2505899008966492
```

In [20]:

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

```
0.4483246814535158
```

In [21]:

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

0.023596035865974516

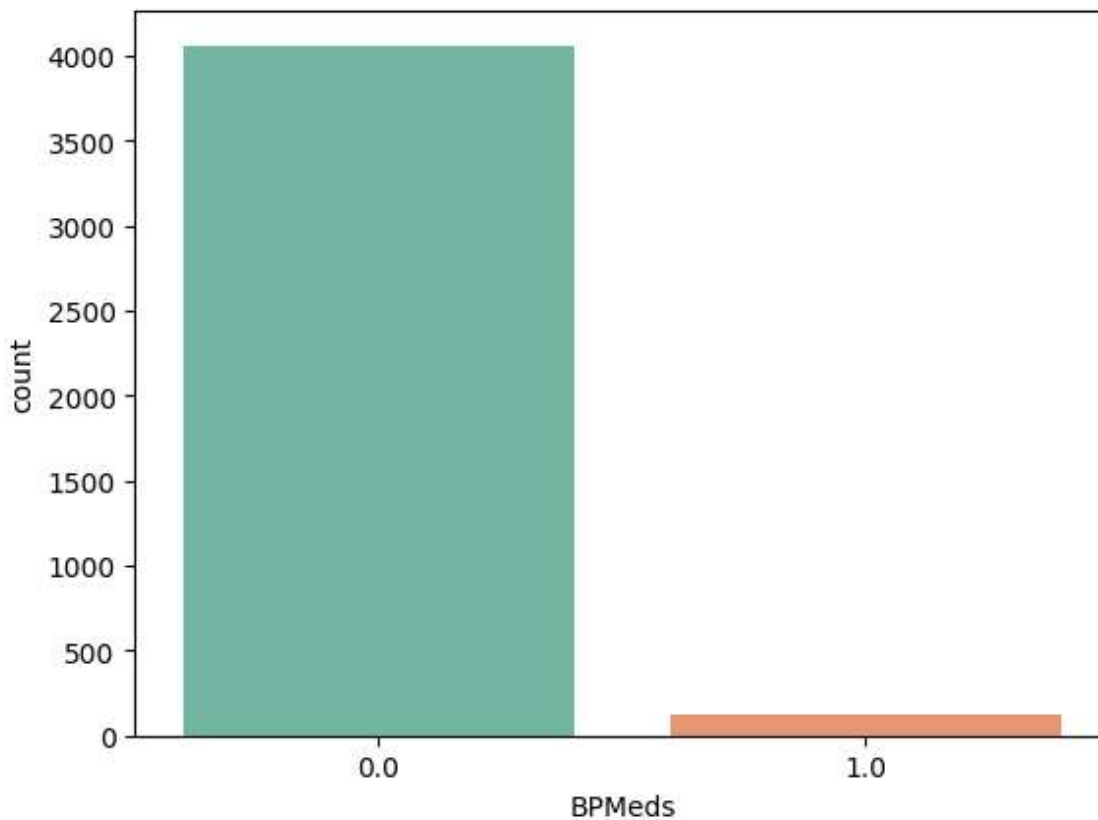
In [22]:

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

0.0 4061

1.0 124

Name: BPMeds, dtype: int64



In [23]:

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

75.0

In [24]:

```

1 data=df.copy()
2 data["cigsPerDay"].fillna(df["cigsPerDay"].median(skipna=True), inplace=True)
3 data["BPMeds"].fillna(df["BPMeds"].median(skipna=True), 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)

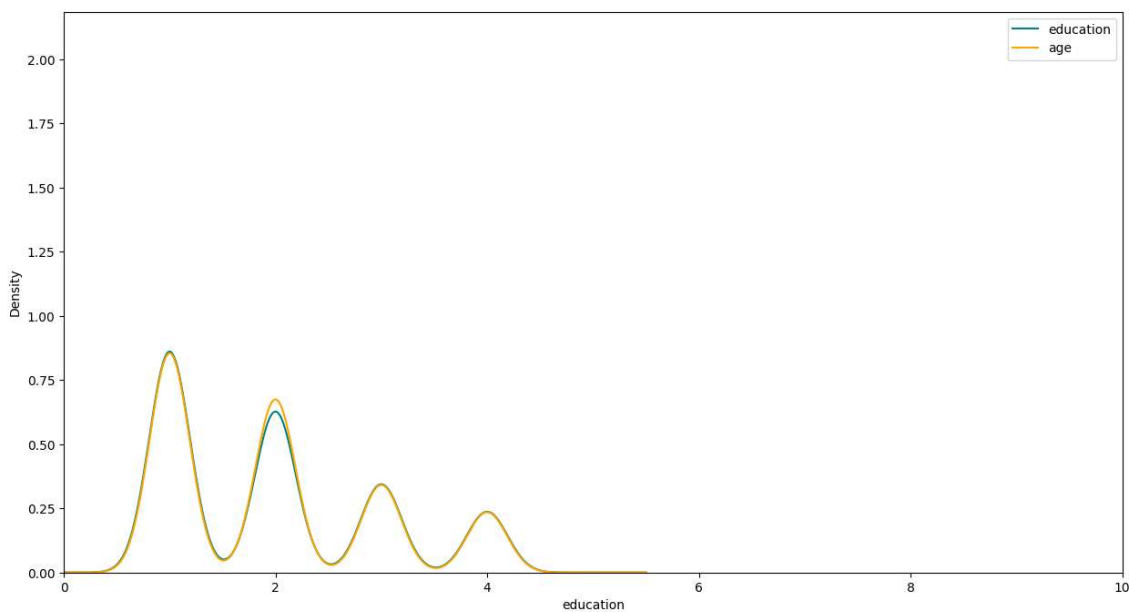
```

In [25]:

```

1 plt.figure(figsize=(15,8))
2 ax=df["education"].hist(bins=15, density=True, stacked=True, color='teal', alpha=0
3 df["education"].plot(kind='density', color='teal')
4 ax=data["education"].hist(bins=15, density=True, stacked=True, color='orange', alp
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 [26]:

```

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 [27]:

```

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)

```

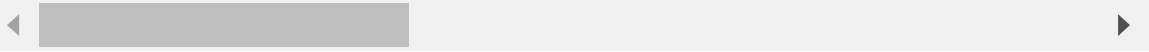
In [28]:

```
1 training.drop("diaBP",axis=1,inplace=True)
2 final_train=training
3 final_train.head()
```

Out[28]:

	age	education	cigsPerDay	BPMeds	diabetes	Disease	currentSmoker_0	currentSmoke
0	39	4.0	0.0	0.0	0	1	1	
1	46	2.0	0.0	0.0	0	1	1	
2	48	1.0	20.0	0.0	0	1	0	
3	61	3.0	30.0	0.0	0	0	0	
4	46	3.0	23.0	0.0	0	1	0	

5 rows × 490 columns



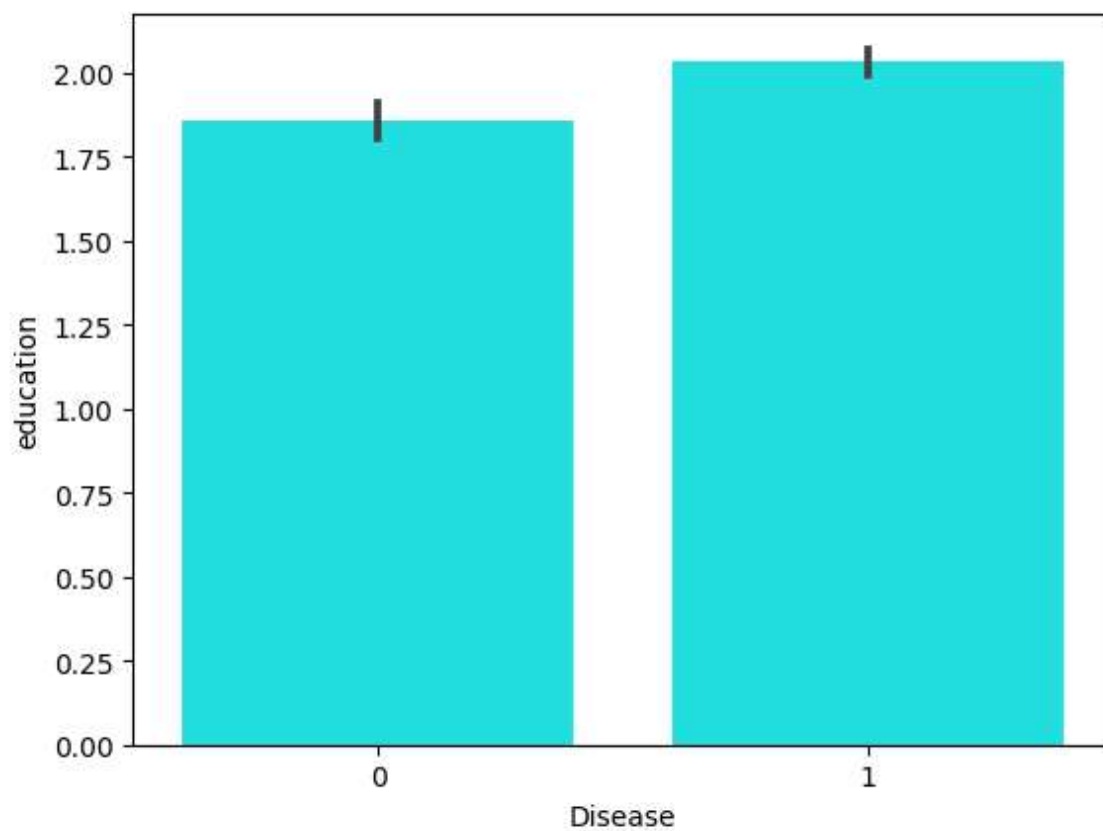
In [29]:

```
1 final_train['IsMinor']=np.where(final_train['age']<=16, 1, 0)
2 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 [30]:

```
1 sns.barplot (x= 'Disease', y='education', data=final_train, color="cyan")  
2 plt.show()
```



In [31]:

```
1 sns.barplot(x='diabetes', y='age', data=df, color='aquamarine')  
2 plt.show()
```

