

Start coding or [generate](#) with AI.

```
**AGABA LUCKY**
**2024/HD05/21913U **
```

RIGOROUS EXPLORATORY DATA ANALYSIS OF DATA CONCERNING HEALTHCARE, AND HEALTH POLICY ISSUES AFFECTING AMERICANS AGED 50 AND OLDER UTILIZING PYTHON LIBRARIES INCLUDING PANDAS, MATPLOTLIB, AND SEABORN

Start coding or [generate](#) with AI.

Double-click (or enter) to edit

```
#import the necessary python libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
# Read the csv file
df = pd.read_csv('/content/NPHA-doctor-visits.csv')
```

```
#check for the first five rows of the dataset
df.head()
```

|   | Number<br>of<br>Doctors<br>Visited | Age | Physical<br>Health | Mental<br>Health | Dental<br>Health | Employment | Stress<br>Keeps<br>Patient<br>from<br>Sleeping | Medication<br>Keeps<br>Patient<br>from<br>Sleeping | Pain<br>Keeps<br>Patient<br>from<br>Sleeping | Bathroom<br>Needs<br>Keeps<br>Patient<br>from<br>Sleeping | Unknown<br>Keeps<br>Patient<br>from<br>Sleeping | Trouble<br>Sleeping | Prescription<br>Sleep<br>Medication | Ra |
|---|------------------------------------|-----|--------------------|------------------|------------------|------------|--|--|--|---|---|---------------------|-------------------------------------|----|
| 0 | 3                                  | 2   | 4                  | 3                | 3                | 3          | 0  | 0  | 0  | 0   | 1   | 2                   | 3                                   |    |
| 1 | 2                                  | 2   | 4                  | 2                | 3                | 3          | 1  | 0  | 0  | 1   | 0   | 3                   | 3                                   |    |
| 2 | 3                                  | 2   | 3                  | 2                | 3                | 3          | 0  | 0  | 0  | 0   | 1   | 3                   | 3                                   |    |

Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)

```
#check for the last five rows of the dataset
df.tail()
```

|     | Number<br>of<br>Doctors<br>Visited | Age | Physical<br>Health | Mental<br>Health | Dental<br>Health | Employment | Stress<br>Keeps<br>Patient<br>from<br>Sleeping | Medication<br>Keeps<br>Patient<br>from<br>Sleeping | Pain<br>Keeps<br>Patient<br>from<br>Sleeping | Bathroom<br>Needs<br>Keeps<br>Patient<br>from<br>Sleeping | Unknown<br>Keeps<br>Patient<br>from<br>Sleeping | Trouble<br>Sleeping | Prescription<br>Sleep<br>Medication |  |
|-----|------------------------------------|-----|--------------------|------------------|------------------|------------|--|--|--|---|---|---------------------|-------------------------------------|--|
| 709 | 2                                  | 2   | 2                  | 2                | 2                | 3          | 0  | 0  | 0  | 1   | 0   | 3                   | 3                                   |  |
| 710 | 3                                  | 2   | 2                  | 2                | 2                | 2          | 1  | 0  | 0  | 0   | 1   | 2                   | 3                                   |  |
| 711 | 3                                  | 2   | 4                  | 2                | 3                | 3          | 0  | 0  | 0  | 0   | 0   | 3                   | 3                                   |  |
| 712 | 3                                  | 2   | 3                  | 1                | 3                | 3          | 1  | 0  | 1  | 1   | 1   | 3                   | 3                                   |  |

```
#check for the number of rows and columns in the dataset
df.shape
```

(714, 15)

```
#check for the columns together with their datatypes
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 714 entries, 0 to 713
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Number of Doctors Visited             714 non-null    int64
1   Age                                    714 non-null    int64
2   Physical Health                       714 non-null    int64
3   Mental Health                         714 non-null    int64
```

```

4  Dental Health          714 non-null  int64
5  Employment            714 non-null  int64
6  Stress Keeps Patient from Sleeping  714 non-null  int64
7  Medication Keeps Patient from Sleeping  714 non-null  int64
8  Pain Keeps Patient from Sleeping      714 non-null  int64
9  Bathroom Needs Keeps Patient from Sleeping  714 non-null  int64
10 Unknown Keeps Patient from Sleeping  714 non-null  int64
11 Trouble Sleeping        714 non-null  int64
12 Prescription Sleep Medication  714 non-null  int64
13 Race                    714 non-null  int64
14 Gender                  714 non-null  int64
dtypes: int64(15)
memory usage: 83.8 KB

```

```
#check for the number of unique elements in the dataset
df.nunique()
```

|   |   |
|---|---|
|   | 0 |
| <b>Number of Doctors Visited</b>                  | 3 |
| <b>Age</b>  | 1 |
| <b>Physiscal Health</b>                           | 6 |
| <b>Mental Health</b>                              | 6 |
| <b>Dental Health</b>                              | 7 |
| <b>Employment</b>                                 | 4 |
| <b>Stress Keeps Patient from Sleeping</b>         | 2 |
| <b>Medication Keeps Patient from Sleeping</b>     | 2 |
| <b>Pain Keeps Patient from Sleeping</b>           | 2 |
| <b>Bathroom Needs Keeps Patient from Sleeping</b> | 2 |
| <b>Unknown Keeps Patient from Sleeping</b>        | 2 |
| <b>Trouble Sleeping</b>                           | 4 |
| <b>Prescription Sleep Medication</b>              | 4 |
| <b>Race</b>                                       | 5 |
| <b>Gender</b>                                     | 2 |

```
#obtain a summary of the dataset/ descriptive statistics
df.describe().T
```

|   | count | mean     | std      | min  | 25% | 50% | 75% | max |  |
|---|-------|----------|----------|------|-----|-----|-----|-----|--|
| <b>Number of Doctors Visited</b>                  | 714.0 | 2.112045 | 0.683441 | 1.0  | 2.0 | 2.0 | 3.0 | 3.0 |  |
| <b>Age</b>  | 714.0 | 2.000000 | 0.000000 | 2.0  | 2.0 | 2.0 | 2.0 | 2.0 |  |
| <b>Physiscal Health</b>                           | 714.0 | 2.794118 | 0.900939 | -1.0 | 2.0 | 3.0 | 3.0 | 5.0 |  |
| <b>Mental Health</b>                              | 714.0 | 1.988796 | 0.939928 | -1.0 | 1.0 | 2.0 | 3.0 | 5.0 |  |
| <b>Dental Health</b>                              | 714.0 | 3.009804 | 1.361117 | -1.0 | 2.0 | 3.0 | 4.0 | 6.0 |  |
| <b>Employment</b>                                 | 714.0 | 2.806723 | 0.586582 | 1.0  | 3.0 | 3.0 | 3.0 | 4.0 |  |
| <b>Stress Keeps Patient from Sleeping</b>         | 714.0 | 0.247899 | 0.432096 | 0.0  | 0.0 | 0.0 | 0.0 | 1.0 |  |
| <b>Medication Keeps Patient from Sleeping</b>     | 714.0 | 0.056022 | 0.230126 | 0.0  | 0.0 | 0.0 | 0.0 | 1.0 |  |
| <b>Pain Keeps Patient from Sleeping</b>           | 714.0 | 0.218487 | 0.413510 | 0.0  | 0.0 | 0.0 | 0.0 | 1.0 |  |
| <b>Bathroom Needs Keeps Patient from Sleeping</b> | 714.0 | 0.504202 | 0.500333 | 0.0  | 0.0 | 1.0 | 1.0 | 1.0 |  |
| <b>Unknown Keeps Patient from Sleeping</b>        | 714.0 | 0.417367 | 0.493470 | 0.0  | 0.0 | 0.0 | 1.0 | 1.0 |  |
| <b>Trouble Sleeping</b>                           | 714.0 | 2.407563 | 0.670349 | -1.0 | 2.0 | 3.0 | 3.0 | 3.0 |  |
| <b>Prescription Sleep Medication</b>              | 714.0 | 2.829132 | 0.546767 | -1.0 | 3.0 | 3.0 | 3.0 | 3.0 |  |
| <b>Race</b>                                       | 714.0 | 1.425770 | 1.003896 | 1.0  | 1.0 | 1.0 | 1.0 | 5.0 |  |
| <b>Gender</b>                                     | 714.0 | 1.550420 | 0.497800 | 1.0  | 1.0 | 2.0 | 2.0 | 2.0 |  |

```
df.describe()
```




|       | Number of<br>Doctors<br>Visited | Age   | Phyiscal<br>Health | Mental<br>Health | Dental<br>Health | Employment | Stress<br>Keeps<br>Patient<br>from<br>Sleeping | Medication<br>Keeps<br>Patient<br>from<br>Sleeping | Pain<br>Keeps<br>Patient<br>from<br>Sleeping | Bathroom<br>Needs<br>Keeps<br>Patient<br>from<br>Sleeping | Uknown<br>Keeps<br>Patient<br>from<br>Sleeping |            |
|-------|---------------------------------|-------|--------------------|------------------|------------------|------------|--|--|--|---|--|------------|
| count | 714.000000                      | 714.0 | 714.000000         | 714.000000       | 714.000000       | 714.000000 | 714.000000                                     | 714.000000   | 714.000000                                   | 714.000000  | 714.000000                                     | 714.000000 |
| mean  | 2.112045                        | 2.0   | 2.794118           | 1.988796         | 3.009804         | 2.806723   | 0.247899                                       | 0.056022   | 0.218487                                     | 0.504202  | 0.417367                                       | 0.417367   |
| std   | 0.683441                        | 0.0   | 0.900939           | 0.939928         | 1.361117         | 0.586582   | 0.432096                                       | 0.230126   | 0.413510                                     | 0.500333  | 0.493470                                       | 0.493470   |
| min   | 1.000000                        | 2.0   | -1.000000          | -1.000000        | -1.000000        | 1.000000   | 0.000000                                       | 0.000000   | 0.000000                                     | 0.000000  | 0.000000                                       | 0.000000   |
| 25%   | 2.000000                        | 2.0   | 2.000000           | 1.000000         | 2.000000         | 3.000000   | 0.000000                                       | 0.000000   | 0.000000                                     | 0.000000  | 0.000000                                       | 0.000000   |
| 50%   | 2.000000                        | 2.0   | 3.000000           | 2.000000         | 3.000000         | 3.000000   | 0.000000                                       | 0.000000   | 0.000000                                     | 1.000000  | 0.000000                                       | 0.000000   |
| 75%   | 3.000000                        | 2.0   | 3.000000           | 3.000000         | 4.000000         | 3.000000   | 0.000000                                       | 0.000000   | 0.000000                                     | 1.000000  | 1.000000                                       | 1.000000   |
| max   | 3.000000                        | 2.0   | 5.000000           | 5.000000         | 6.000000         | 4.000000   | 1.000000                                       | 1.000000   | 1.000000                                     | 1.000000  | 1.000000                                       | 1.000000   |

```
#checking for missing values  
df.isnull().sum()
```



|  | 0 |
|--|---|
| Number of Doctors Visited                  | 0 |
| Age  | 0 |
| Phyiscal Health                            | 0 |
| Mental Health                              | 0 |
| Dental Health                              | 0 |
| Employment                                 | 0 |
| Stress Keeps Patient from Sleeping         | 0 |
| Medication Keeps Patient from Sleeping     | 0 |
| Pain Keeps Patient from Sleeping           | 0 |
| Bathroom Needs Keeps Patient from Sleeping | 0 |
| Uknown Keeps Patient from Sleeping         | 0 |
| Trouble Sleeping                           | 0 |
| Prescription Sleep Medication              | 0 |
| Race                                       | 0 |
| Gender                                     | 0 |

```
#checking for missing values  
df.isnull().sum()/df.shape[0]*100
```




|  | 0   |
|--|-----|
| Number of Doctors Visited                  | 0.0 |
| Age  | 0.0 |
| Phyiscal Health                            | 0.0 |
| Mental Health                              | 0.0 |
| Dental Health                              | 0.0 |
| Employment                                 | 0.0 |
| Stress Keeps Patient from Sleeping         | 0.0 |
| Medication Keeps Patient from Sleeping     | 0.0 |
| Pain Keeps Patient from Sleeping           | 0.0 |
| Bathroom Needs Keeps Patient from Sleeping | 0.0 |
| Unknown Keeps Patient from Sleeping        | 0.0 |
| Trouble Sleeping                           | 0.0 |
| Prescription Sleep Medication              | 0.0 |
| Race                                       | 0.0 |
| Gender                                     | 0.0 |

```
#checking for duplicates in the dataset
df.duplicated().sum()
```

 42

```
#dropping or eliminating the duplicated values
df.drop_duplicates()
```

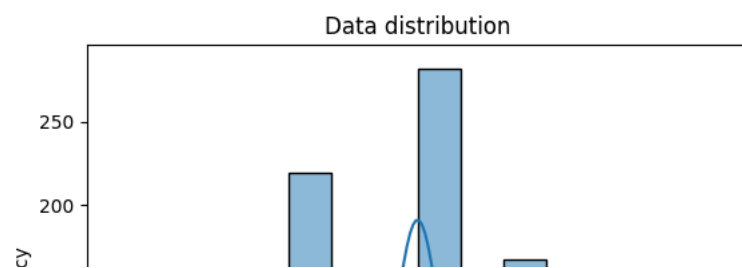
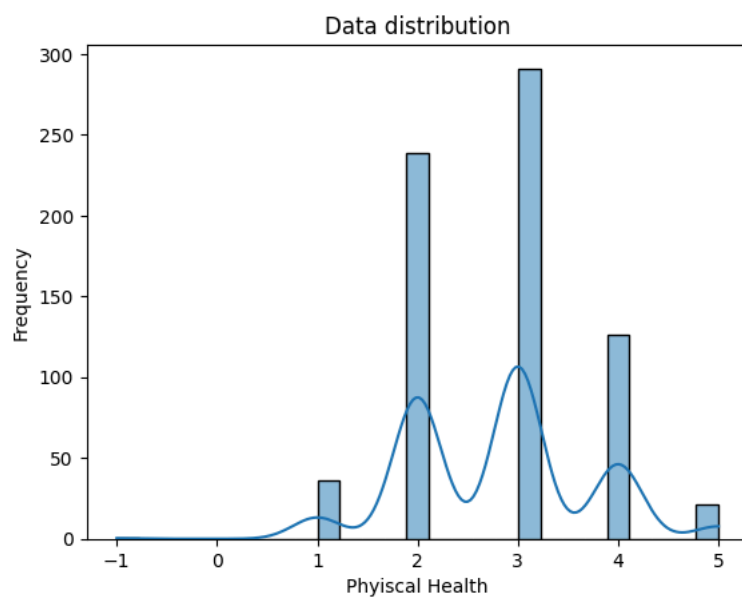
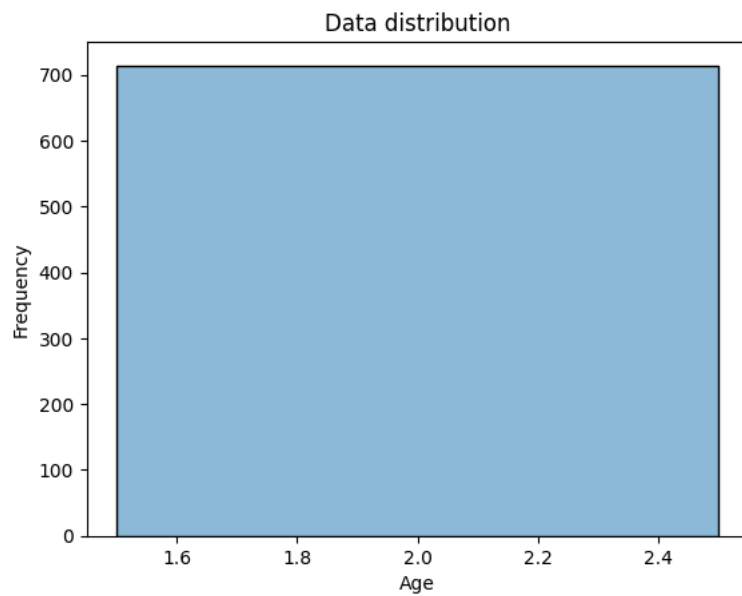
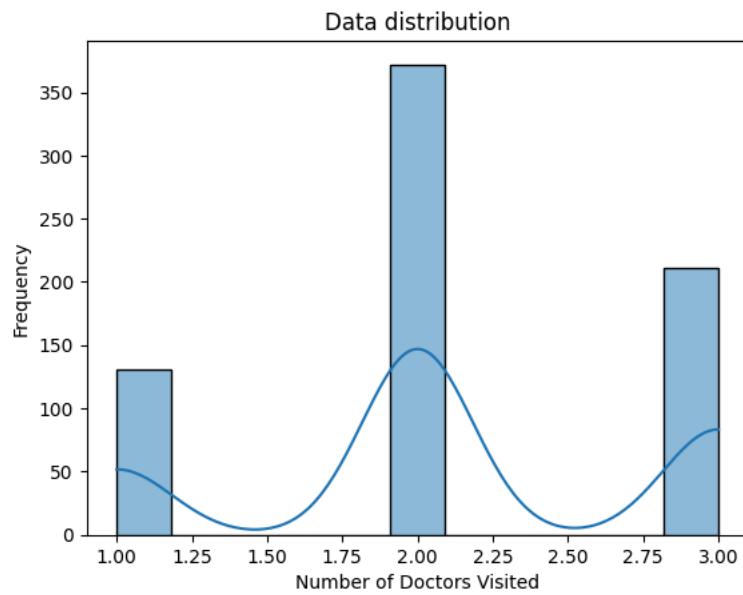


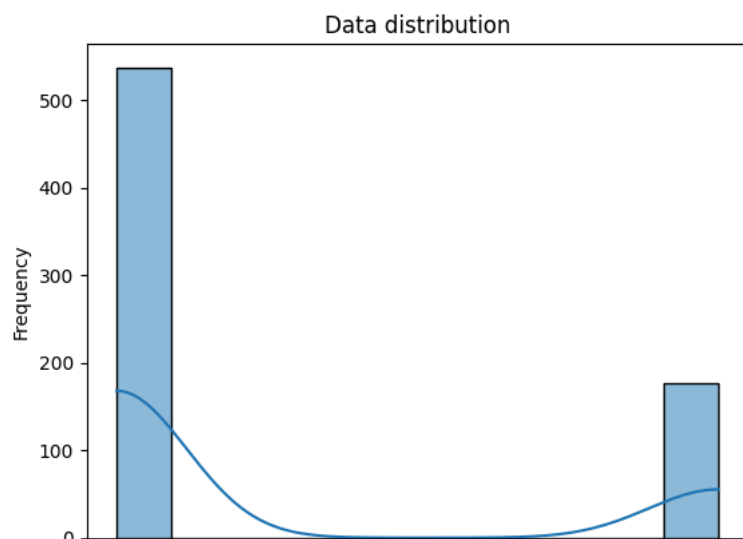
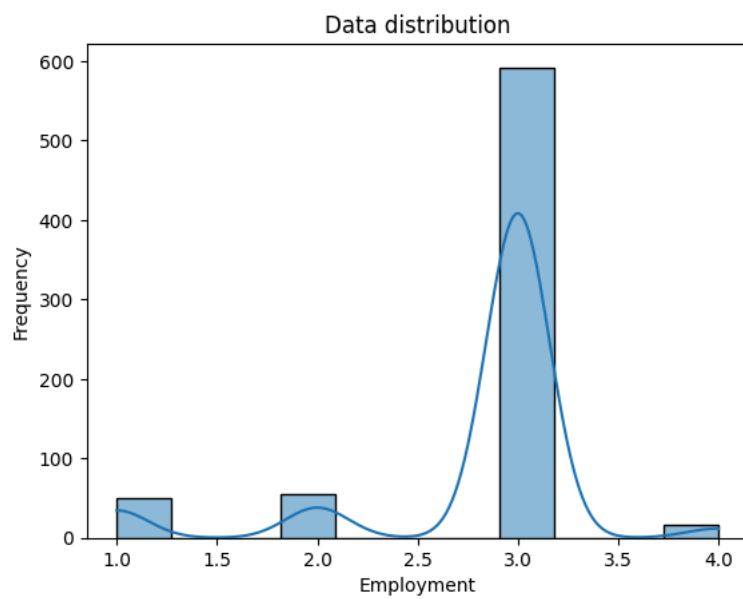
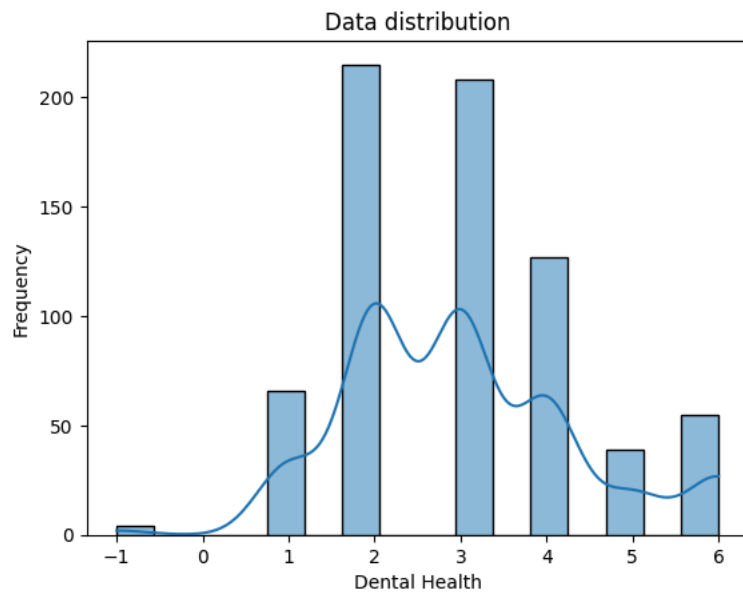
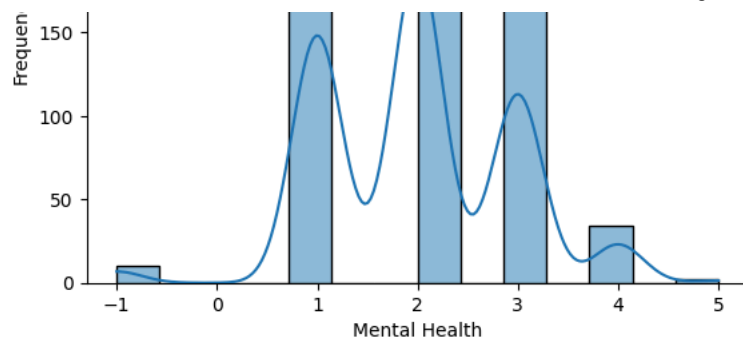
|     | Number<br>of<br>Doctors<br>Visited | Age | Phyiscal<br>Health | Mental<br>Health | Dental<br>Health | Employment | Stress<br>Keeps<br>Patient<br>from<br>Sleeping | Medication<br>Keeps<br>Patient<br>from<br>Sleeping | Pain<br>Keeps<br>Patient<br>from<br>Sleeping | Bathroom<br>Needs<br>Keeps<br>Patient<br>from<br>Sleeping | Unknown<br>Keeps<br>Patient<br>from<br>Sleeping | Trouble<br>Sleeping | Prescription<br>Sleep<br>Medication |
|-----|------------------------------------|-----|--------------------|------------------|------------------|------------|--|--|--|---|---|---------------------|-------------------------------------|
| 0   | 3                                  | 2   | 4                  | 3                | 3                | 3          | 0  | 0  | 0  | 0   | 1   | 2                   | 3                                   |
| 1   | 2                                  | 2   | 4                  | 2                | 3                | 3          | 1  | 0  | 0  | 1   | 0   | 3                   | 3                                   |
| 2   | 3                                  | 2   | 3                  | 2                | 3                | 3          | 0  | 0  | 0  | 0   | 1   | 3                   | 3                                   |
| 3   | 1                                  | 2   | 3                  | 2                | 3                | 3          | 0  | 0  | 0  | 1   | 0   | 3                   | 3                                   |
| 4   | 3                                  | 2   | 3                  | 3                | 3                | 3          | 1  | 0  | 0  | 0   | 0   | 2                   | 3                                   |
| ... | ...                                | ... | ...                | ...              | ...              | ...        | ...  | ...  | ...  | ...   | ...   | ...                 | ...                                 |
| 706 | 3                                  | 2   | 4                  | 2                | 2                | 3          | 0  | 0  | 1  | 1   | 0   | -1                  | 3                                   |
| 710 | 3                                  | 2   | 2                  | 2                | 2                | 2          | 1  | 0  | 0  | 0   | 1   | 2                   | 3                                   |
| 711 | 3                                  | 2   | 4                  | 2                | 3                | 3          | 0  | 0  | 0  | 0   | 0   | 3                   | 3                                   |
| 712 | 3                                  | 2   | 3                  | 1                | 3                | 3          | 1  | 0  | 1  | 1   | 1   | 3                   | 3                                   |
| 713 | 3                                  | 2   | 3                  | 2                | 2                | 3          | 1  | 0  | 1  | 1   | 0   | 3                   | 3                                   |

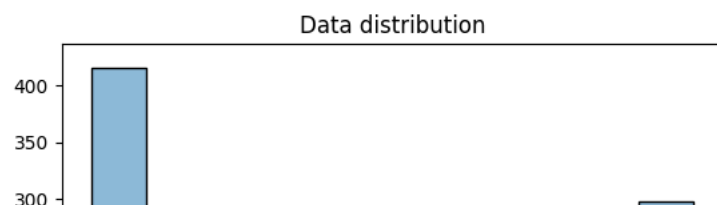
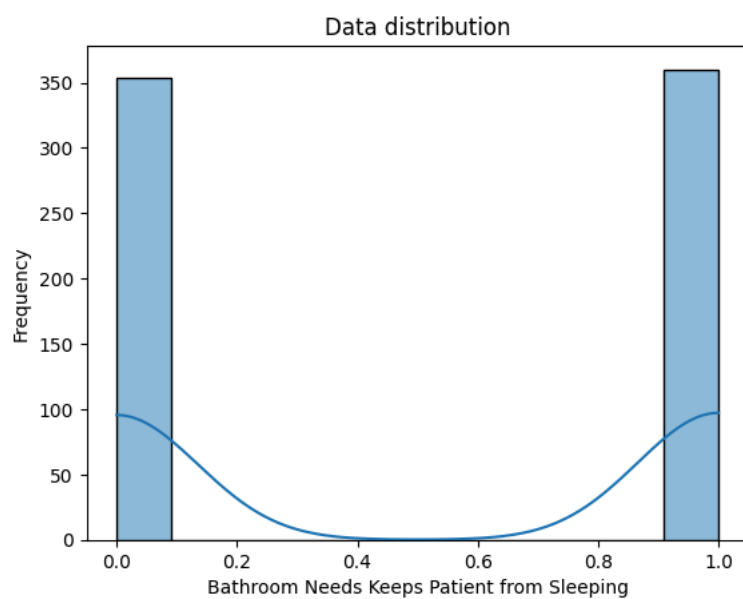
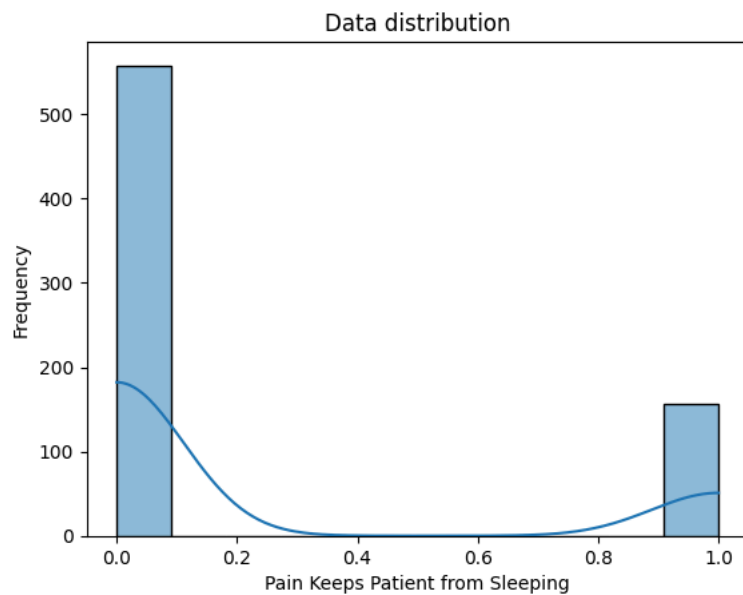
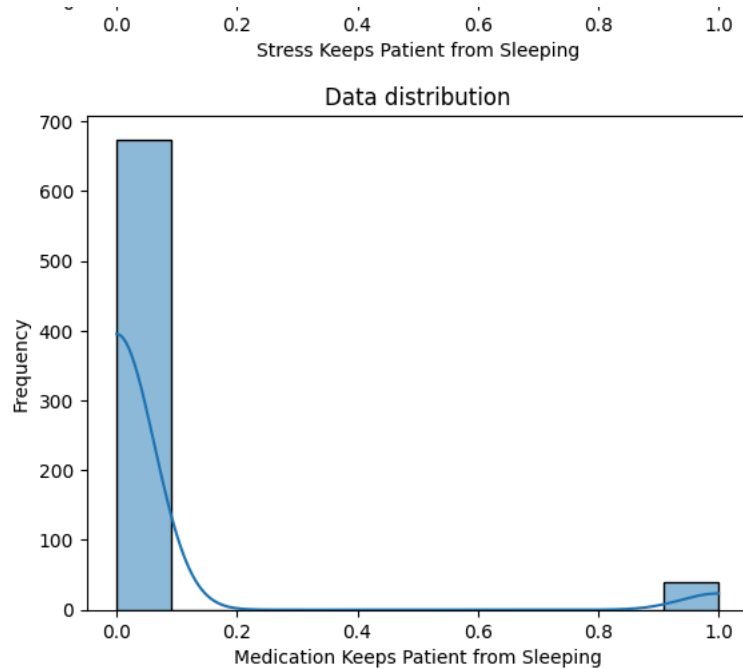
```
#checking for garbage values
for i in df.select_dtypes(include='object').columns:
    print(df[i].value_counts())
print("****10)
```

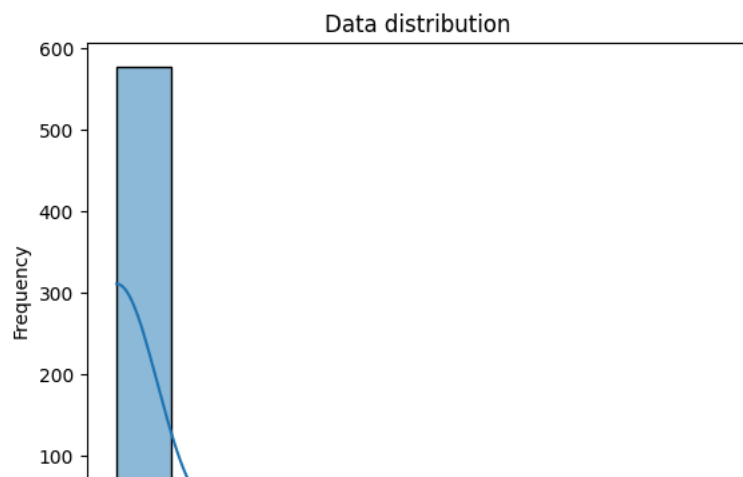
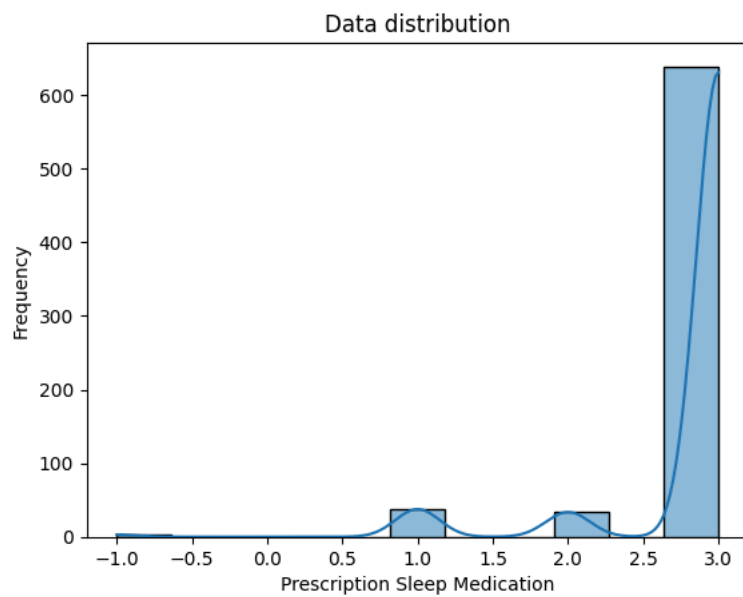
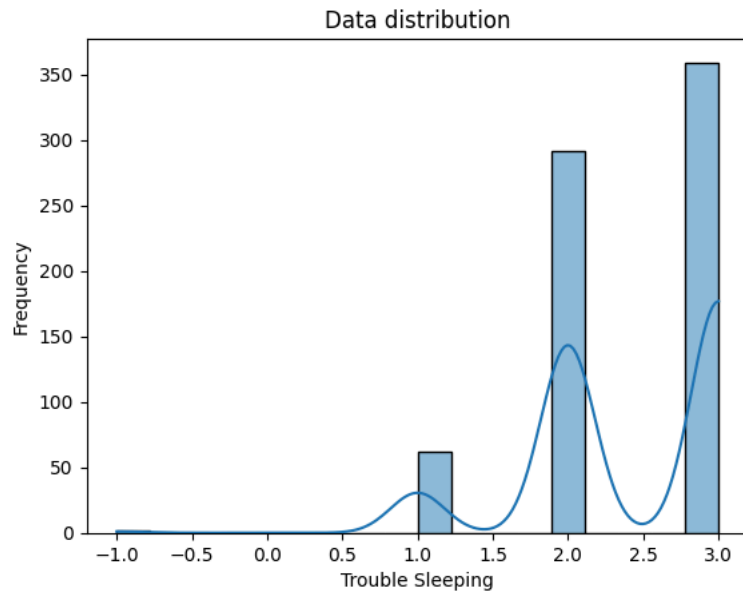
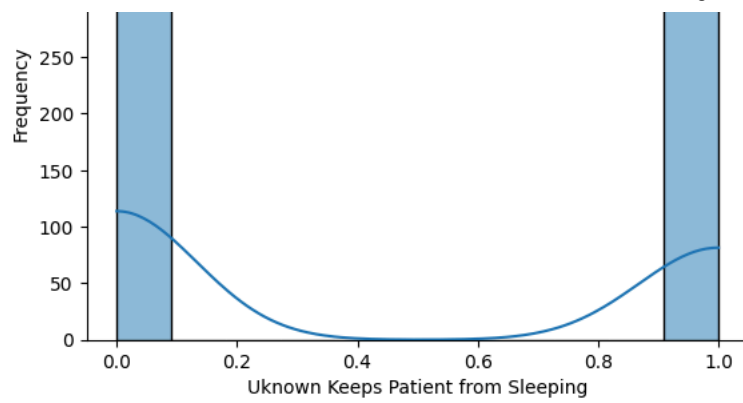
### Exploratory Data Analysis (EDA)

```
#understanding the distribution of the data for each numerical column using a histogram
for i in df.select_dtypes(include='number').columns:
    sns.histplot(data=df, x=i, kde=True)
    plt.ylabel('Frequency')
    plt.title('Data distribution')
    plt.show()
```

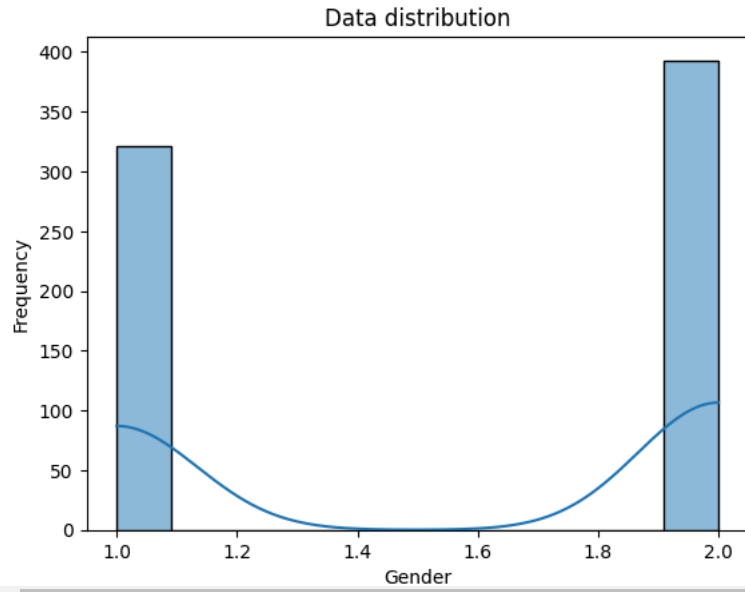
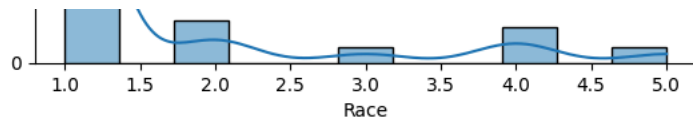




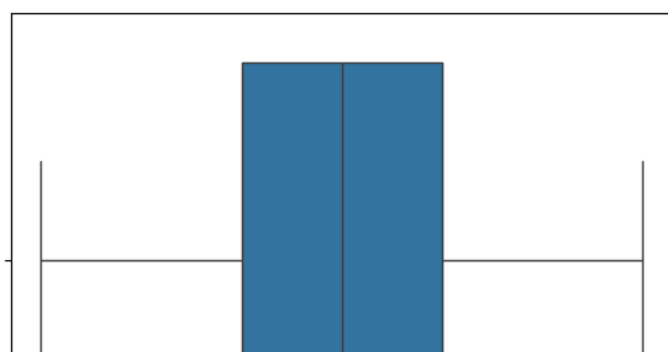
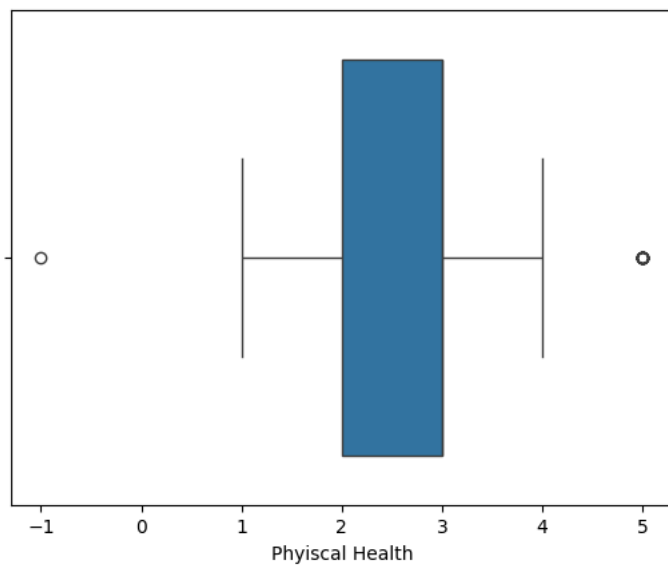
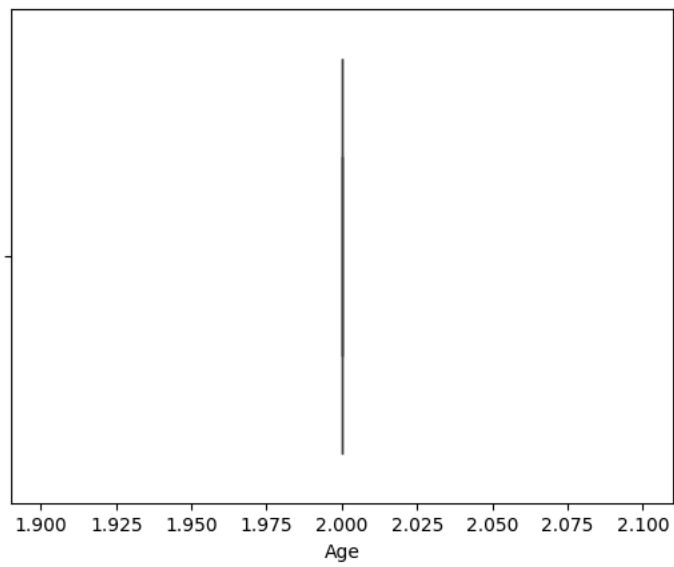
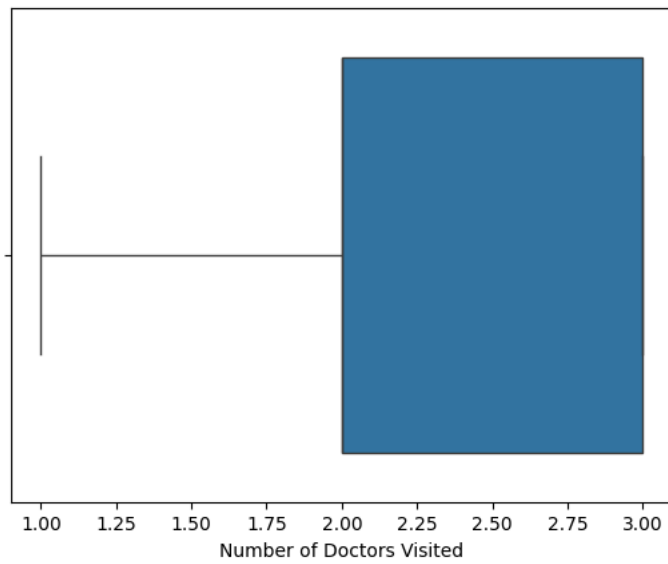


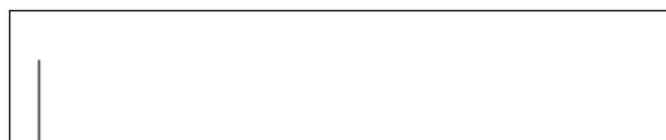
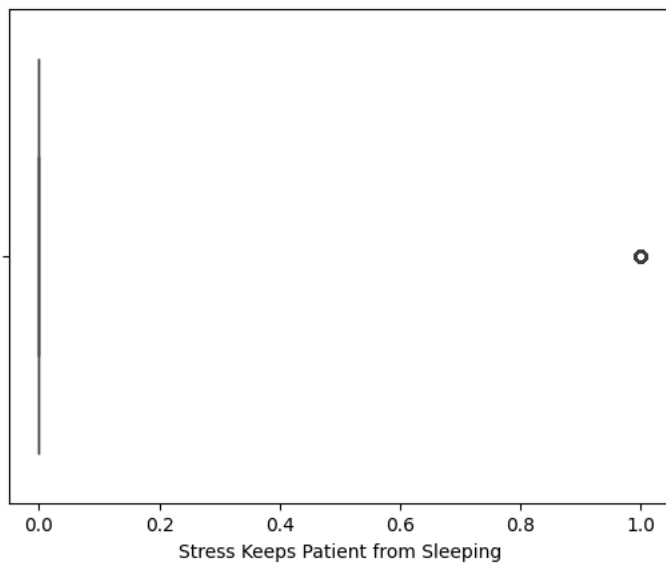
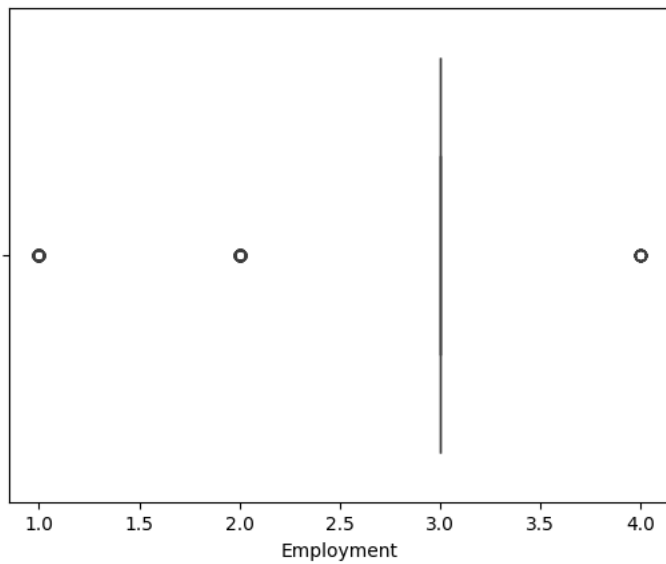
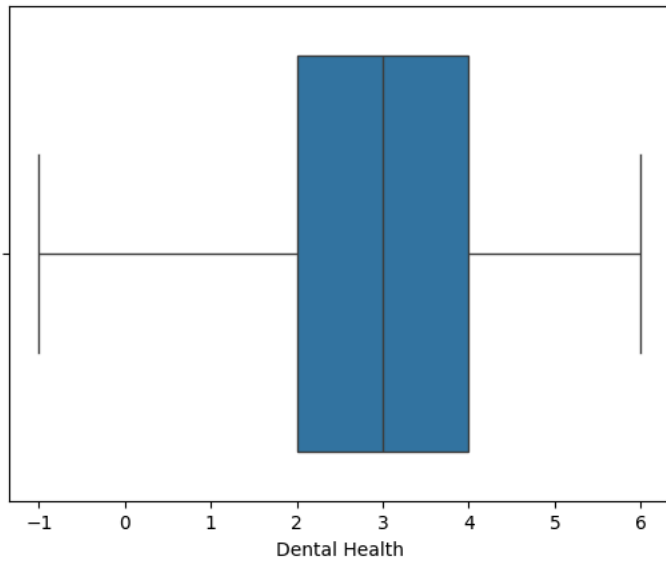
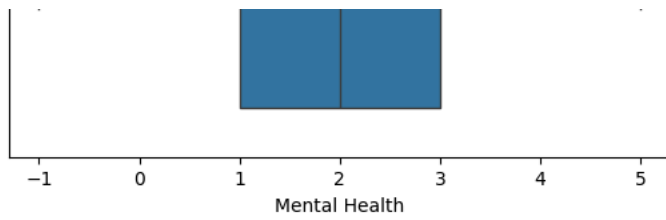


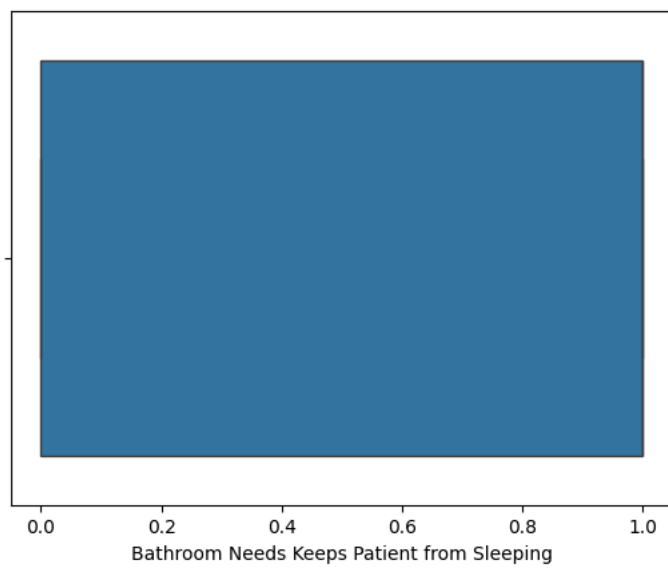
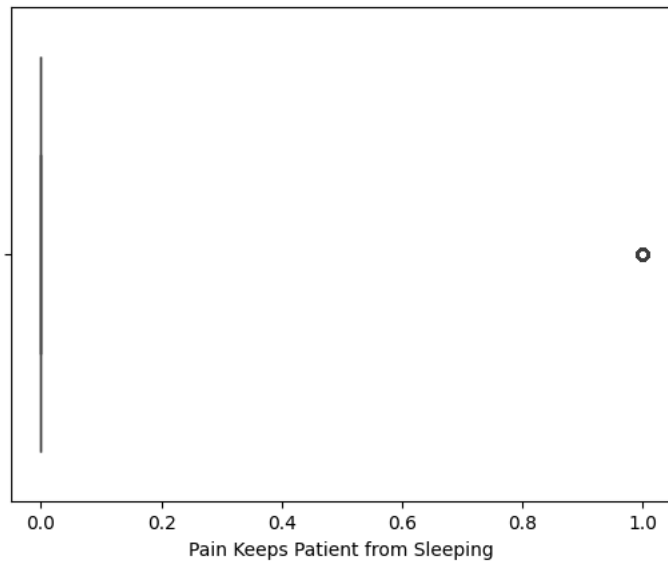
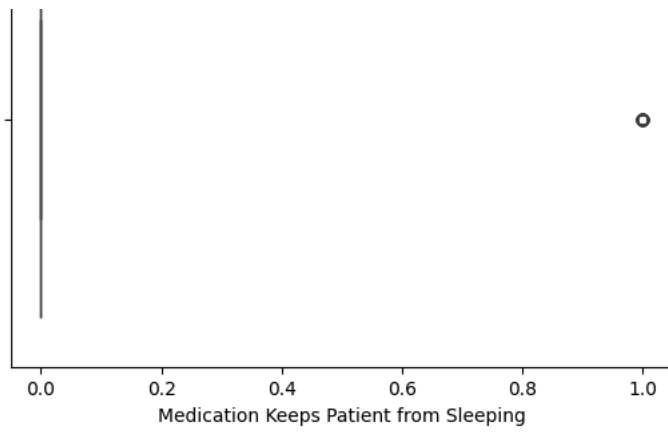


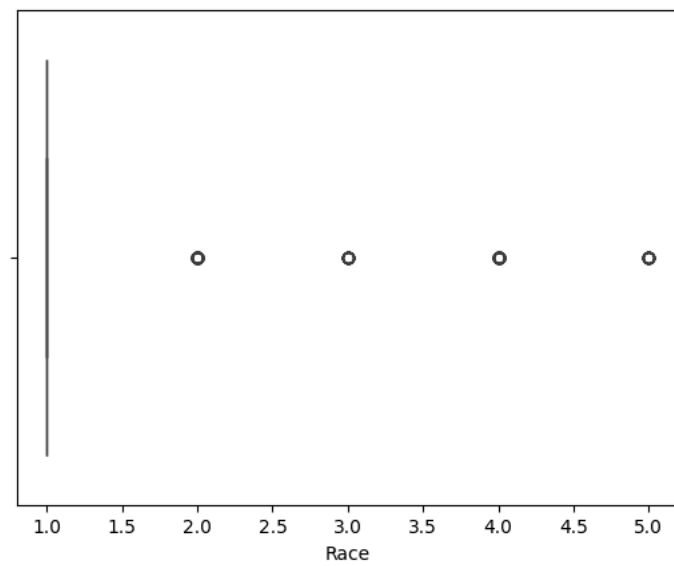
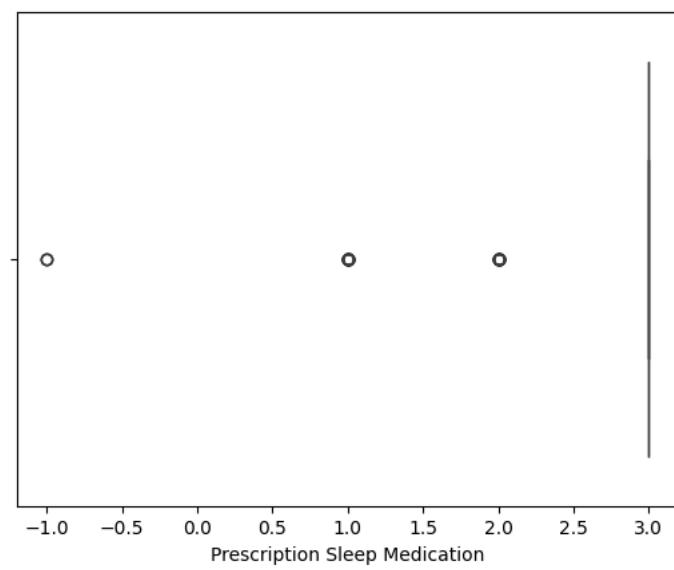
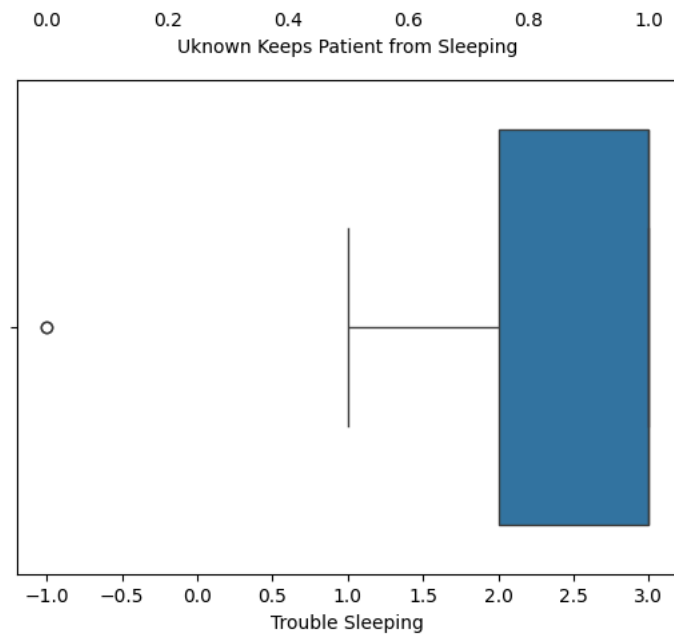


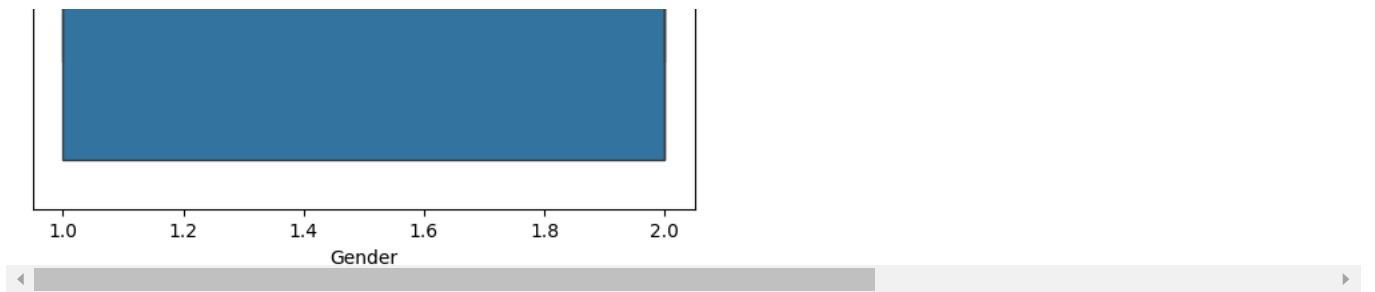
```
#identifying outliers in the dataset
for i in df.select_dtypes(include='number').columns:
    sns.boxplot(data=df, x=i)
    plt.show()
```











```
#dealing with the outliers
#def out_liers(col):
#    q1,q3=np.percentile(col,[25,75])
#    iqr=q3-q1
#    upper_bound=q3+(1.5*iqr)
#    lower_bound=q1-(1.5*iqr)
#    return upper_bound,lower_bound
```

```
#for i in ['Phyiscal Health', 'Employment', 'Stress Keeps Patient from Sleeping', 'Medication Keeps Patient from Sleeping',
#         'Pain Keeps Patient from Sleeping', 'Trouble Sleeping', 'Prescription Sleep Medication', 'Race']:
#    upper_bound,lower_bound=out_liers(df[i])
#    df[i]=np.where(df[i]>upper_bound,upper_bound,df[i])
#    df[i]=np.where(df[i]<lower_bound,lower_bound,df[i])
```

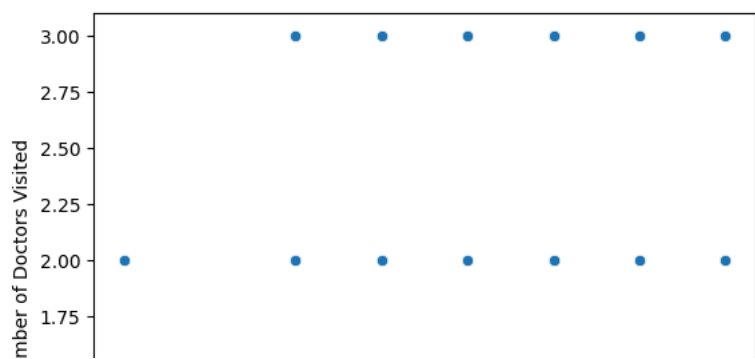
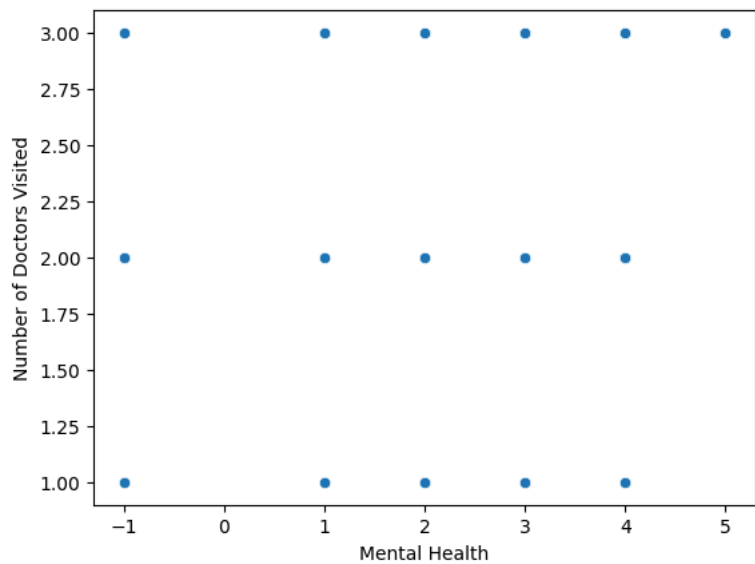
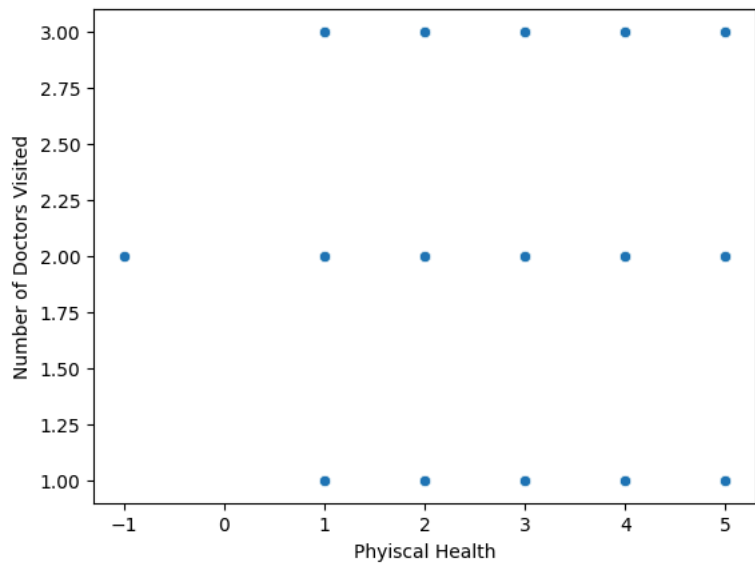
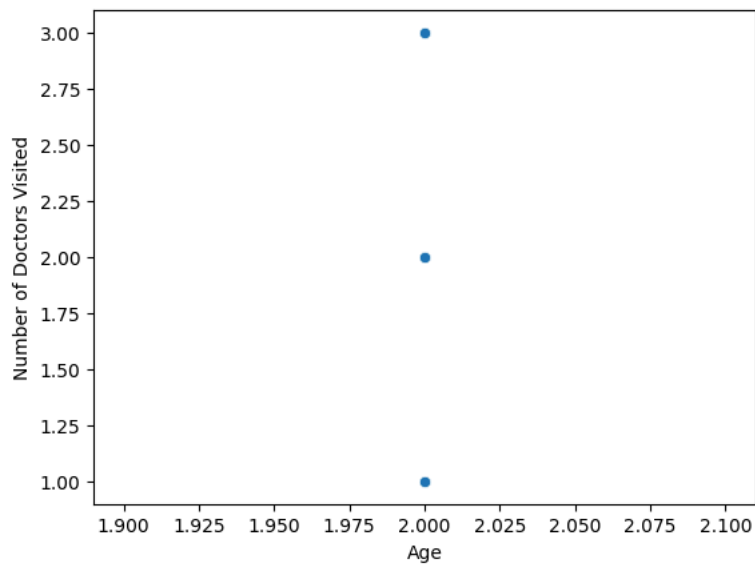
```
#for i in ['Phyiscal Health', 'Employment', 'Stress Keeps Patient from Sleeping', 'Medication Keeps Patient from Sleeping',
#         'Pain Keeps Patient from Sleeping', 'Trouble Sleeping', 'Prescription Sleep Medication', 'Race']:
#    sns.boxplot(df[i])
#    plt.show()
```

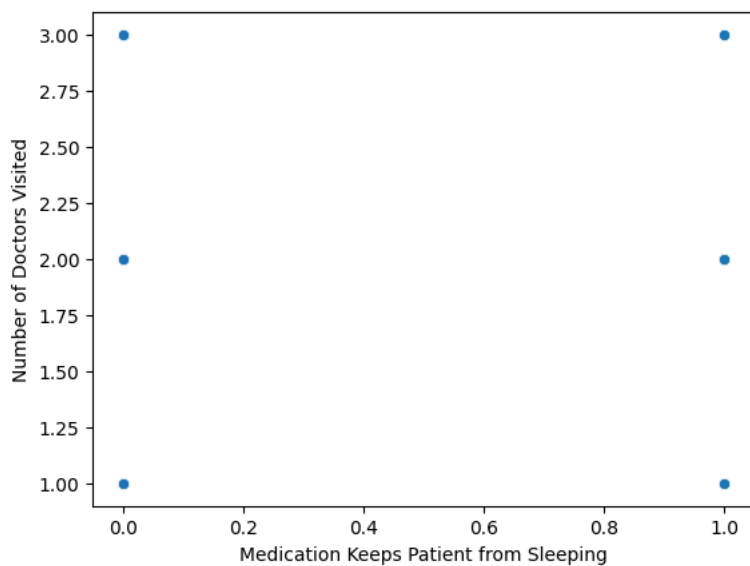
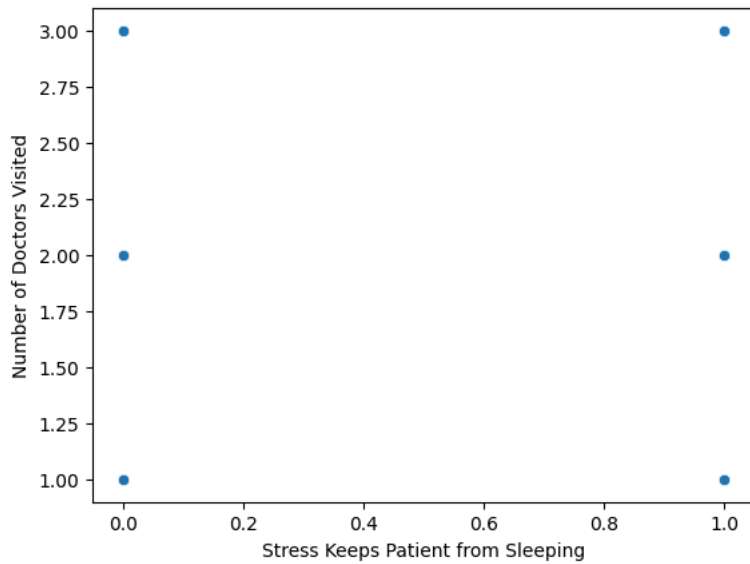
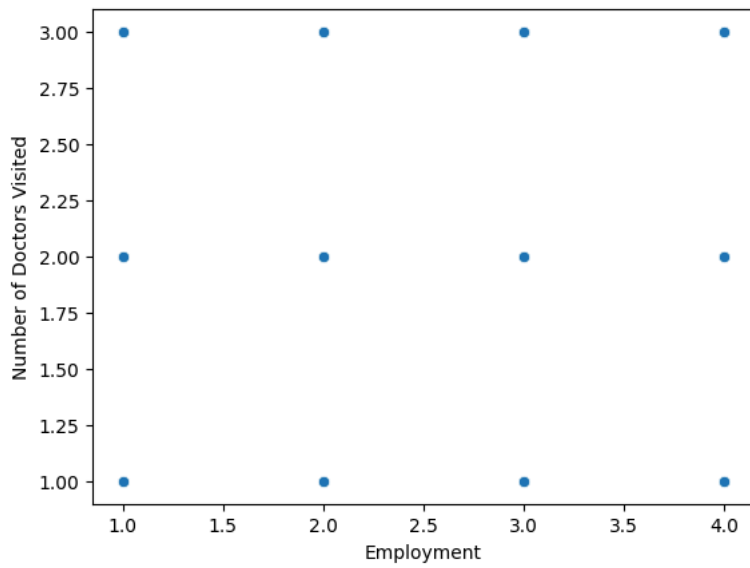
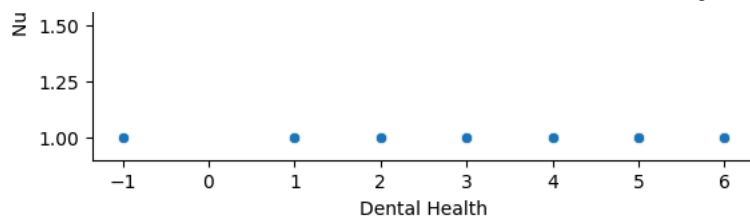
```
df.select_dtypes(include='number').columns
```

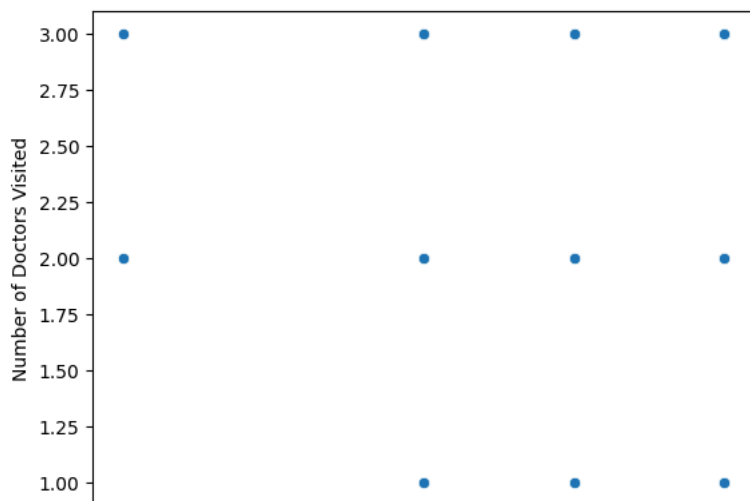
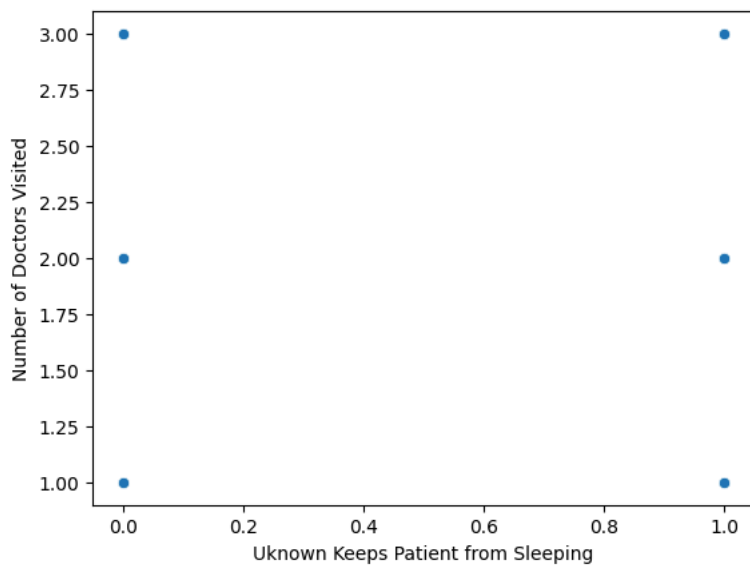
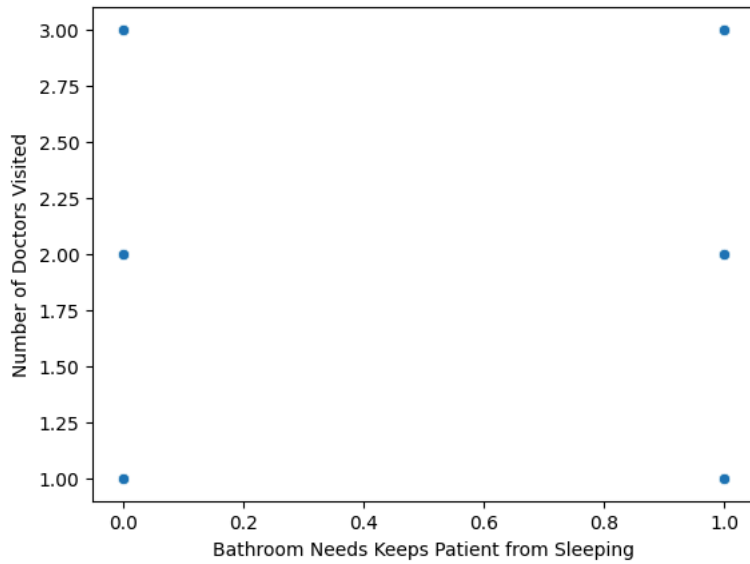
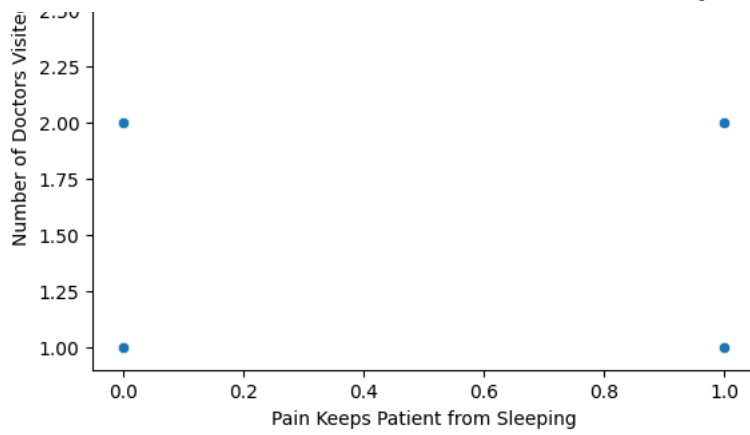
```
Index(['Number of Doctors Visited', 'Age', 'Phyiscal Health', 'Mental Health',
      'Dental Health', 'Employment', 'Stress Keeps Patient from Sleeping',
      'Medication Keeps Patient from Sleeping',
      'Pain Keeps Patient from Sleeping',
      'Bathroom Needs Keeps Patient from Sleeping',
      'Uknown Keeps Patient from Sleeping', 'Trouble Sleeping',
      'Prescription Sleep Medication', 'Race', 'Gender'],
      dtype='object')
```

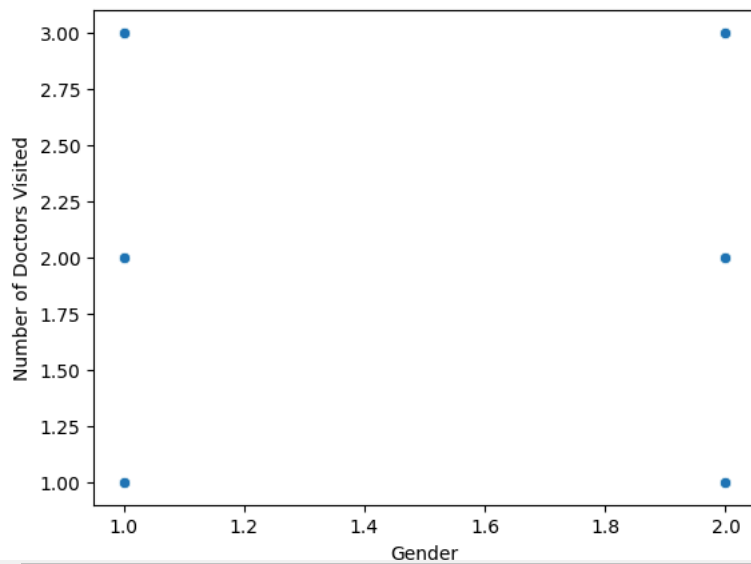
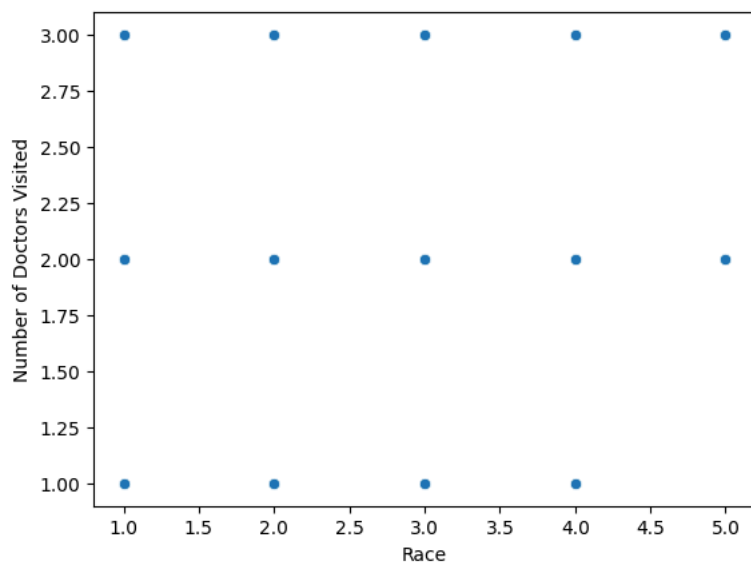
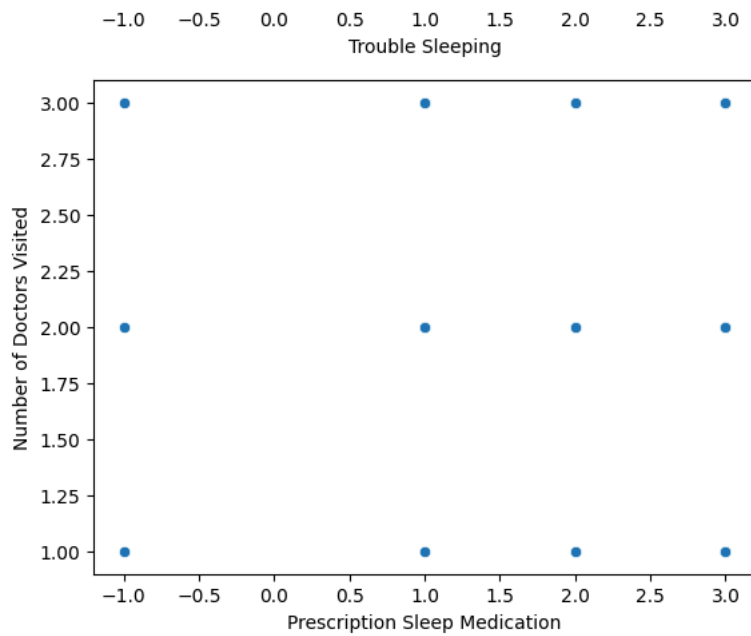
```
#scatter plot to understand the relationship between my target variable and other variables
for i in ['Age', 'Phyiscal Health', 'Mental Health',
        'Dental Health', 'Employment', 'Stress Keeps Patient from Sleeping',
        'Medication Keeps Patient from Sleeping',
        'Pain Keeps Patient from Sleeping',
        'Bathroom Needs Keeps Patient from Sleeping',
        'Uknown Keeps Patient from Sleeping', 'Trouble Sleeping',
        'Prescription Sleep Medication', 'Race', 'Gender']:
    sns.scatterplot(data=df, x=i, y='Number of Doctors Visited')
    plt.show()
```











```
#checking for the correlation with heatmap to interpret the relation and multicollinearity
df.select_dtypes(include='number').corr()
```



|  | Number<br>of<br>Doctors<br>Visited | Age | Physcal<br>Health | Mental<br>Health | Dental<br>Health | Employment | Stress<br>Keeps<br>Patient<br>from<br>Sleeping | Medication<br>Keeps<br>Patient<br>from<br>Sleeping | Pain<br>Keeps<br>Patient<br>from<br>Sleeping | Bathroom<br>Needs<br>Keeps<br>Patient<br>from<br>Sleeping | Uknown<br>Keeps<br>Patient<br>from<br>Sleeping | Trouble<br>Sleeping |
|--|------------------------------------|-----|-------------------|------------------|------------------|------------|--|--|--|---|--|---------------------|
| Number of<br>Doctors<br>Visited                        | 1.000000                           | NaN | 0.169629          | 0.049990         | 0.009371         | 0.092578   | 0.053040                                       | 0.120549   | 0.081990                                     | 0.056043  | -0.014095                                      | -0.063079           |
| Age  | NaN                                | NaN | NaN               | NaN              | NaN              | NaN        | NaN  | NaN  | NaN  | NaN   | NaN  | NaN                 |
| Physcal<br>Health                                      | 0.169629                           | NaN | 1.000000          | 0.404705         | 0.404238         | 0.147526   | 0.034014                                       | 0.109827   | 0.275266                                     | 0.003477  | -0.011505                                      | -0.213855           |
| Mental<br>Health                                       | 0.049990                           | NaN | 0.404705          | 1.000000         | 0.269770         | 0.077469   | 0.138074                                       | 0.139072   | 0.121780                                     | 0.044835  | -0.038285                                      | -0.110718           |
| Dental<br>Health                                       | 0.009371                           | NaN | 0.404238          | 0.269770         | 1.000000         | 0.076156   | -0.018446                                      | 0.029588   | 0.080913                                     | -0.007269   | -0.014453                                      | -0.044351           |
| Employment   | 0.092578                           | NaN | 0.147526          | 0.077469         | 0.076156         | 1.000000   | -0.043106                                      | 0.059546   | -0.004908                                    | 0.012329  | 0.017427                                       | -0.116836           |
| Stress<br>Keeps<br>Patient from<br>Sleeping            | 0.053040                           | NaN | 0.034014          | 0.138074         | -0.018446        | -0.043106  | 1.000000                                       | 0.029395   | 0.136015                                     | -0.001581   | -0.314897                                      | -0.150775           |
| Medication<br>Keeps<br>Patient from<br>Sleeping        | 0.120549                           | NaN | 0.109827          | 0.139072         | 0.029588         | 0.059546   | 0.029395                                       | 1.000000   | 0.165965                                     | 0.071039  | -0.119734                                      | -0.148217           |
| Pain Keeps<br>Patient from<br>Sleeping                 | 0.081990                           | NaN | 0.275266          | 0.121780         | 0.080913         | -0.004908  | 0.136015                                       | 0.165965   | 1.000000                                     | 0.144695  | -0.213823                                      | -0.235680           |
| Bathroom<br>Needs<br>Keeps<br>Patient from<br>Sleeping | 0.056043                           | NaN | 0.003477          | 0.044835         | -0.007269        | 0.012329   | -0.001581                                      | 0.071039   | 0.144695                                     | 1.000000  | -0.382029                                      | 0.038795            |
| Uknown<br>Keeps<br>Patient from<br>Sleeping            | -0.014095                          | NaN | -0.011505         | -0.038285        | -0.014453        | 0.017427   | -0.314897                                      | -0.119734  | -0.213823                                    | -0.382029   | 1.000000                                       | -0.023123           |
| Trouble<br>Sleeping                                    | -0.063079                          | NaN | -0.213855         | -0.110718        | -0.044351        | -0.116836  | -0.150775                                      | -0.148217  | -0.235680                                    | 0.038795  | -0.023123                                      | 1.000000            |

```
plt.figure(figsize=(10,10))
sns.heatmap(df.select_dtypes(include='number').corr(), annot=True)
```

Number of Doctors Visited

## Phyiscal Health

Mental Health

Dental Health

## Employment

## Stress Keeps Patient from Sleeping

### Medication Keeps Patient from Sleeping

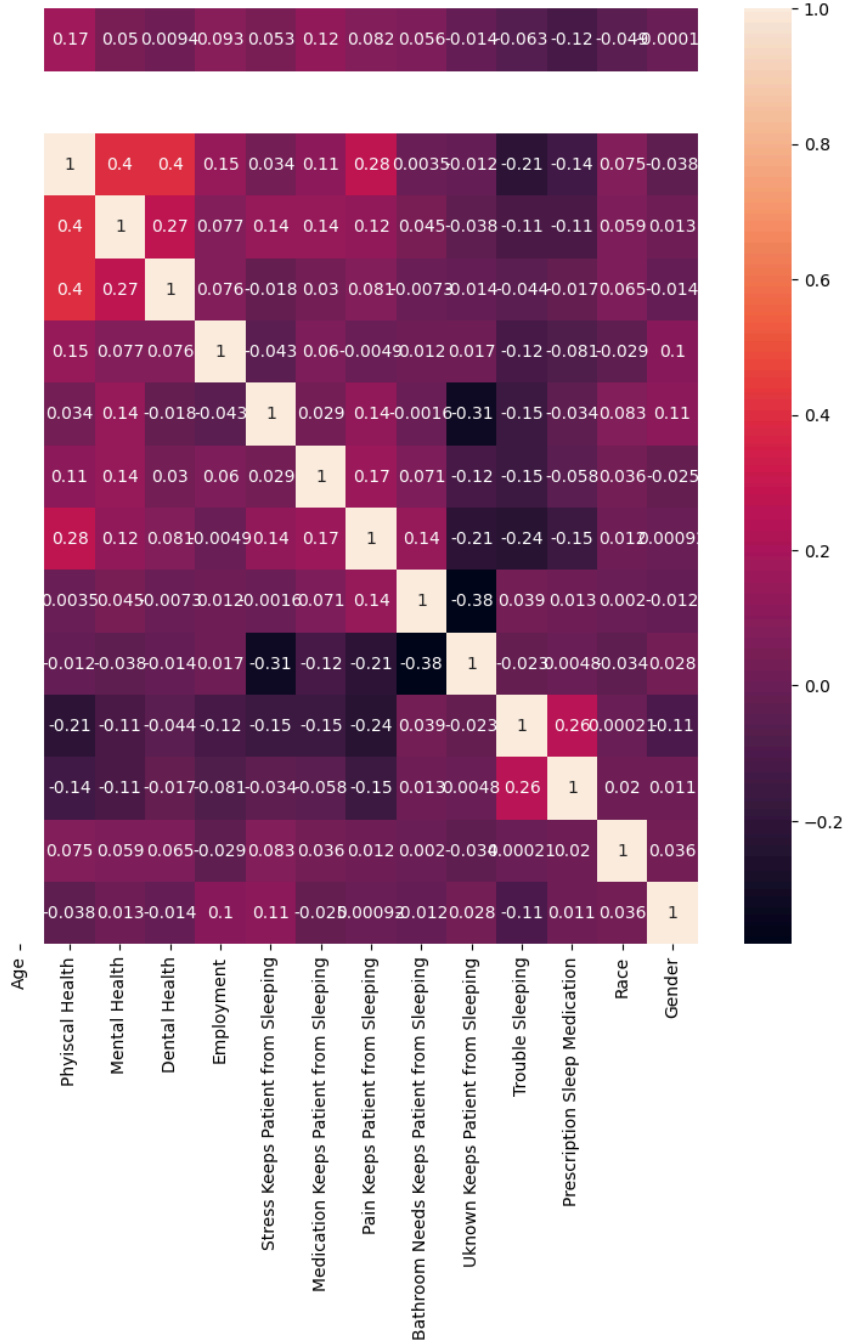
### Pain Keeps Patient from Sleeping

### Bathroom Needs Keeps Patient from Sleeping

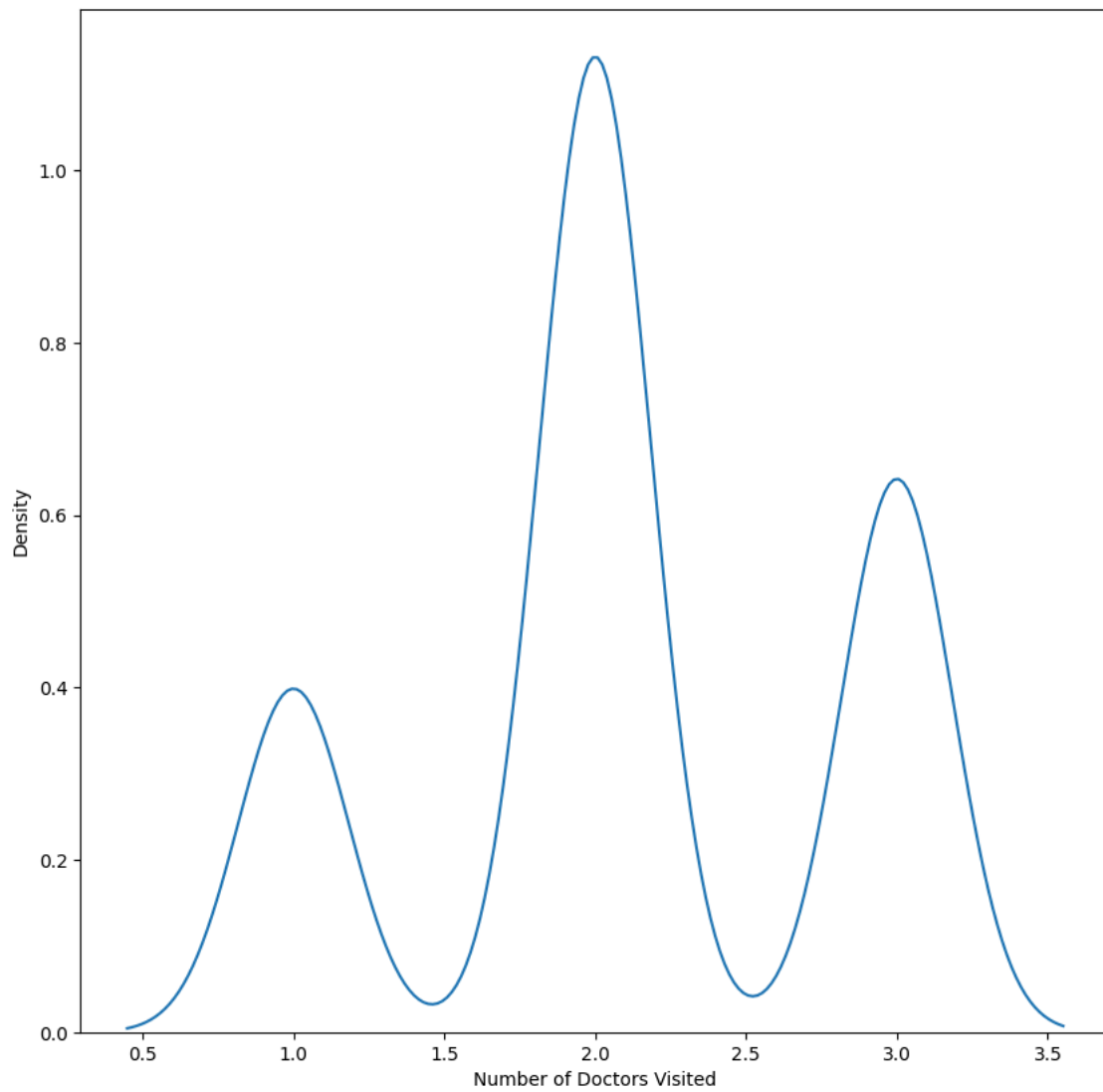
### Unknown Keeps Patient from Sleeping

multiple cleaning

### Results and discussion



22/25



```
<ipython-input-15-8369a3293a54>:3: UserWarning: Dataset has 0 variance; skipping density estimate. Pass `warn_singular=False` to  
sns.kdeplot(data=df, x=i)
```

