

# Problem Statement

A real estate agent want help to predict the house price for regions in USA.He gave us the dataset to work on to use linear regression model.Create a model that helps him to estimate of what the house would sell for

## Import libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: # To import dataset
df=pd.read_csv('16 Sleep csv')
df
```

Out[2]:

|     | Person ID | Gender | Age | Occupation           | Sleep Duration | Quality of Sleep | Physical Activity Level | Stress Level | BMI Category | Blo Pressure |
|-----|-----------|--------|-----|----------------------|----------------|------------------|-------------------------|--------------|--------------|--------------|
| 0   | 1         | Male   | 27  | Software Engineer    | 6.1            | 6                | 42                      | 6            | Overweight   | 126          |
| 1   | 2         | Male   | 28  | Doctor               | 6.2            | 6                | 60                      | 8            | Normal       | 125          |
| 2   | 3         | Male   | 28  | Doctor               | 6.2            | 6                | 60                      | 8            | Normal       | 125          |
| 3   | 4         | Male   | 28  | Sales Representative | 5.9            | 4                | 30                      | 8            | Obese        | 140          |
| 4   | 5         | Male   | 28  | Sales Representative | 5.9            | 4                | 30                      | 8            | Obese        | 140          |
| ... | ...       | ...    | ... | ...                  | ...            | ...              | ...                     | ...          | ...          | ...          |
| 369 | 370       | Female | 59  | Nurse                | 8.1            | 9                | 75                      | 3            | Overweight   | 140          |
| 370 | 371       | Female | 59  | Nurse                | 8.0            | 9                | 75                      | 3            | Overweight   | 140          |
| 371 | 372       | Female | 59  | Nurse                | 8.1            | 9                | 75                      | 3            | Overweight   | 140          |
| 372 | 373       | Female | 59  | Nurse                | 8.1            | 9                | 75                      | 3            | Overweight   | 140          |
| 373 | 374       | Female | 59  | Nurse                | 8.1            | 9                | 75                      | 3            | Overweight   | 140          |

374 rows × 13 columns



```
In [3]: # To display top 10 rows
df.head(10)
```

Out[3]:

|   | Person ID | Gender | Age | Occupation           | Sleep Duration | Quality of Sleep | Physical Activity Level | Stress Level | BMI Category | Blood Pressure |
|---|-----------|--------|-----|----------------------|----------------|------------------|-------------------------|--------------|--------------|----------------|
| 0 | 1         | Male   | 27  | Software Engineer    | 6.1            | 6                | 42                      | 6            | Overweight   | 126/80         |
| 1 | 2         | Male   | 28  | Doctor               | 6.2            | 6                | 60                      | 8            | Normal       | 125/80         |
| 2 | 3         | Male   | 28  | Doctor               | 6.2            | 6                | 60                      | 8            | Normal       | 125/80         |
| 3 | 4         | Male   | 28  | Sales Representative | 5.9            | 4                | 30                      | 8            | Obese        | 140/90         |
| 4 | 5         | Male   | 28  | Sales Representative | 5.9            | 4                | 30                      | 8            | Obese        | 140/90         |
| 5 | 6         | Male   | 28  | Software Engineer    | 5.9            | 4                | 30                      | 8            | Obese        | 140/90         |
| 6 | 7         | Male   | 29  | Teacher              | 6.3            | 6                | 40                      | 7            | Obese        | 140/90         |
| 7 | 8         | Male   | 29  | Doctor               | 7.8            | 7                | 75                      | 6            | Normal       | 120/80         |
| 8 | 9         | Male   | 29  | Doctor               | 7.8            | 7                | 75                      | 6            | Normal       | 120/80         |
| 9 | 10        | Male   | 29  | Doctor               | 7.8            | 7                | 75                      | 6            | Normal       | 120/80         |

## Data Cleaning and Pre-Processing

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 374 entries, 0 to 373
Data columns (total 13 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Person ID                            374 non-null    int64
1   Gender                               374 non-null    object
2   Age                                   374 non-null    int64
3   Occupation                           374 non-null    object
4   Sleep Duration                       374 non-null    float64
5   Quality of Sleep                     374 non-null    int64
6   Physical Activity Level               374 non-null    int64
7   Stress Level                         374 non-null    int64
8   BMI Category                         374 non-null    object
9   Blood Pressure                       374 non-null    object
10  Heart Rate                           374 non-null    int64
11  Daily Steps                          374 non-null    int64
12  Sleep Disorder                       374 non-null    object
dtypes: float64(1), int64(7), object(5)
memory usage: 38.1+ KB
```

In [5]: `df.describe()`

Out[5]:

|              | Person ID  | Age        | Sleep Duration | Quality of Sleep | Physical Activity Level | Stress Level | Heart Rate | Daily Steps |
|--------------|------------|------------|----------------|------------------|-------------------------|--------------|------------|-------------|
| <b>count</b> | 374.000000 | 374.000000 | 374.000000     | 374.000000       | 374.000000              | 374.000000   | 374.000000 | 374.000000  |
| <b>mean</b>  | 187.500000 | 42.184492  | 7.132086       | 7.312834         | 59.171123               | 5.385027     | 70.165775  | 681.125000  |
| <b>std</b>   | 108.108742 | 8.673133   | 0.795657       | 1.196956         | 20.830804               | 1.774526     | 4.135676   | 161.125000  |
| <b>min</b>   | 1.000000   | 27.000000  | 5.800000       | 4.000000         | 30.000000               | 3.000000     | 65.000000  | 300.000000  |
| <b>25%</b>   | 94.250000  | 35.250000  | 6.400000       | 6.000000         | 45.000000               | 4.000000     | 68.000000  | 560.000000  |
| <b>50%</b>   | 187.500000 | 43.000000  | 7.200000       | 7.000000         | 60.000000               | 5.000000     | 70.000000  | 700.000000  |
| <b>75%</b>   | 280.750000 | 50.000000  | 7.800000       | 8.000000         | 75.000000               | 7.000000     | 72.000000  | 800.000000  |
| <b>max</b>   | 374.000000 | 59.000000  | 8.500000       | 9.000000         | 90.000000               | 8.000000     | 86.000000  | 1000.000000 |

In [6]: `df.columns`

Out[6]: Index(['Person ID', 'Gender', 'Age', 'Occupation', 'Sleep Duration', 'Quality of Sleep', 'Physical Activity Level', 'Stress Level', 'BMI Category', 'Blood Pressure', 'Heart Rate', 'Daily Steps', 'Sleep Disorder'], dtype='object')

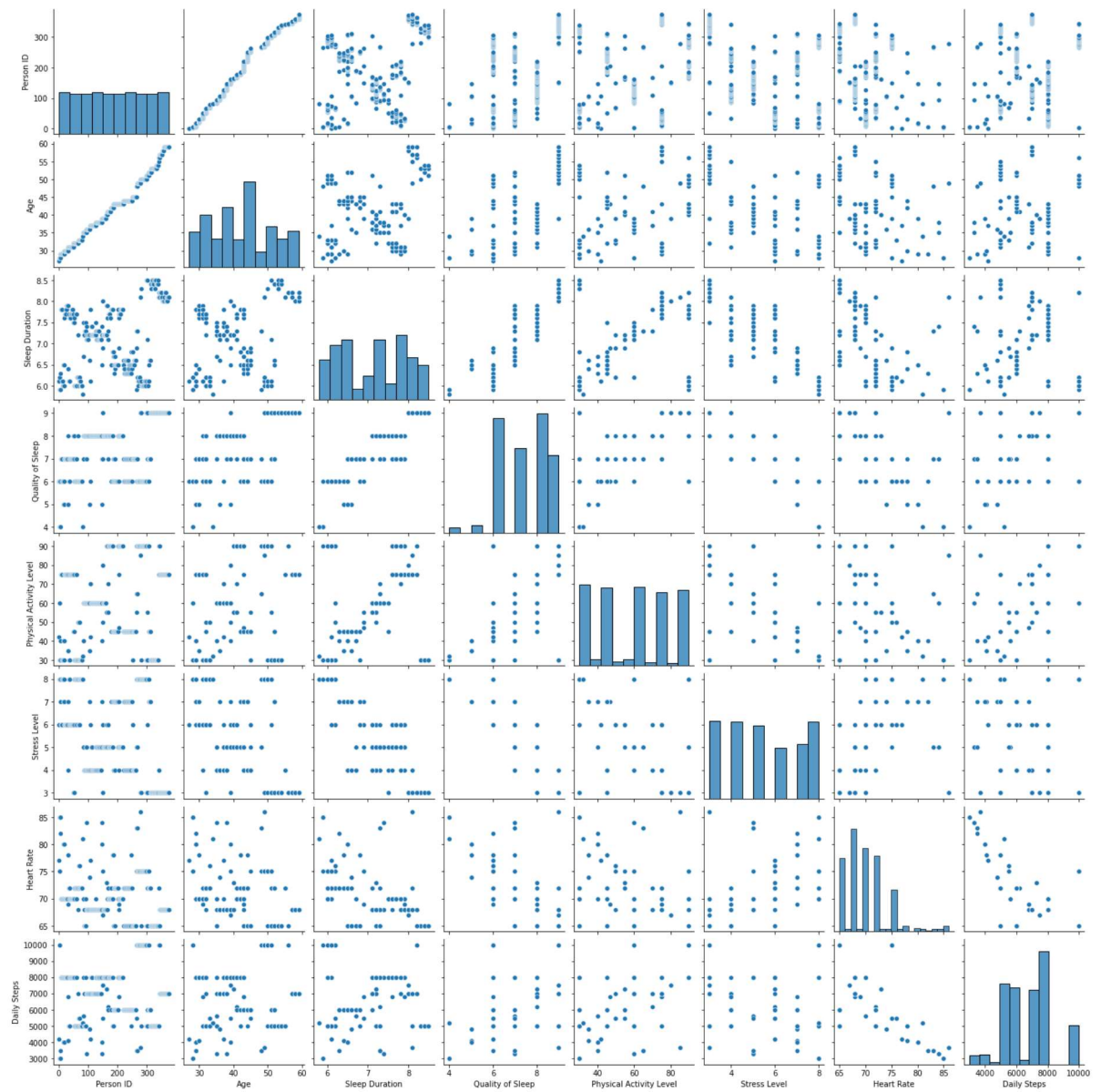
In [7]: `a = df.dropna(axis='columns')`  
`a.columns`

Out[7]: Index(['Person ID', 'Gender', 'Age', 'Occupation', 'Sleep Duration', 'Quality of Sleep', 'Physical Activity Level', 'Stress Level', 'BMI Category', 'Blood Pressure', 'Heart Rate', 'Daily Steps', 'Sleep Disorder'], dtype='object')

## EDA and Visualization

```
In [8]: sns.pairplot(a)
```

```
Out[8]: <seaborn.axisgrid.PairGrid at 0x22b36c741c0>
```

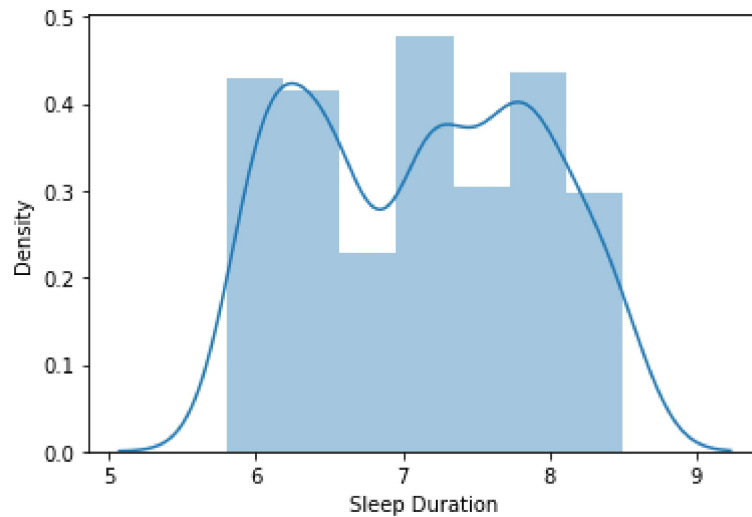


In [9]: `sns.distplot(a['Sleep Duration'])`

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

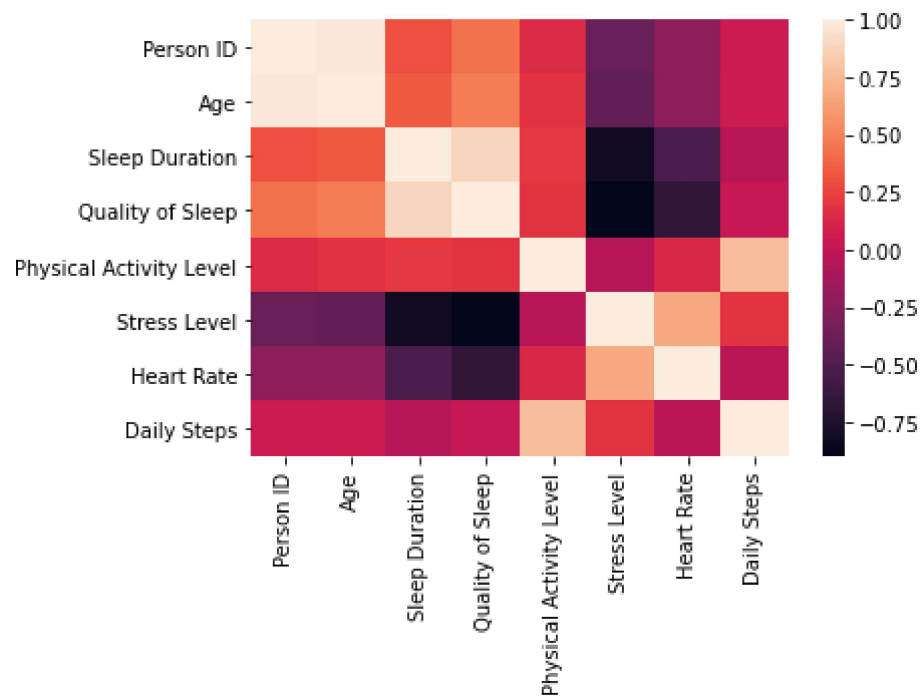
Out[9]: `<AxesSubplot:xlabel='Sleep Duration', ylabel='Density'>`



In [10]: `a1=a[['Person ID', 'Age', 'Sleep Duration', 'Quality of Sleep', 'Physical Activity Level', 'Stress Level', 'Heart Rate', 'Daily Steps']]`

In [11]: `sns.heatmap(a1.corr())`

Out[11]: `<AxesSubplot:>`



## To Train the Model - Model Building

We are going to train Linear Regression model; We need to split out data into two variables x and y where x is independent variable (input) and y is dependent on x (output). We could ignore address column as it is not required for our model.

```
In [12]: x=a1[['Person ID', 'Age', 'Sleep Duration','Quality of Sleep', 'Physical Activity Level']]
        y=a1['Sleep Duration']
```

## To split my dataset into training and test data

```
In [13]: from sklearn.model_selection import train_test_split
        x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3)
```

```
In [14]: from sklearn.linear_model import LinearRegression
        lr=LinearRegression()
        lr.fit(x_train,y_train)
```

Out[14]: LinearRegression()

```
In [15]: print(lr.intercept_)
5.595524044110789e-14
```

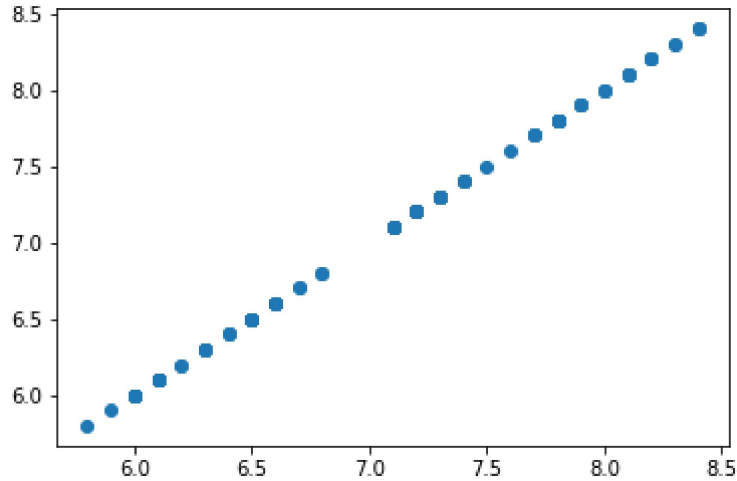
```
In [16]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
        coeff
```

Out[16]:

|                         | Co-efficient  |
|-------------------------|---------------|
| Person ID               | 1.695337e-16  |
| Age                     | -1.275859e-15 |
| Sleep Duration          | 1.000000e+00  |
| Quality of Sleep        | -6.431170e-16 |
| Physical Activity Level | -5.409609e-18 |
| Stress Level            | -6.733586e-16 |
| Heart Rate              | 1.955594e-16  |
| Daily Steps             | -5.447393e-18 |

```
In [17]: prediction=lr.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[17]: <matplotlib.collections.PathCollection at 0x22b3c265550>



```
In [18]: print(lr.score(x_test,y_test))

1.0
```

## ACCURACY

```
In [19]: from sklearn.linear_model import Ridge,Lasso
```

```
In [20]: rr=Ridge(alpha=10)
rr.fit(x_train,y_train)
rr.score(x_test,y_test)
rr.score(x_train,y_train)
```

Out[20]: 0.9903183765111357

```
In [21]: rr.score(x_test,y_test)
```

Out[21]: 0.9877786833377165

```
In [22]: la=Lasso(alpha=10)
la.fit(x_train,y_train)
```

Out[22]: Lasso(alpha=10)

```
In [23]: la.score(x_test,y_test)
```

Out[23]: -0.008487295791862692

```
In [ ]:
```

