

```
import pandas as pd
```

```
a=pd.read_csv("/content/drive/MyDrive/DSBDA/NHANES Weight and Height.csv")
```

```
a.head
```

```
<bound method NDFrame.head of
0      0      0      97.1      160.2      37.8      BMI(kg/m**2)
1      1      1      98.8      182.3      29.7
2      2      2      74.3      184.2      21.9
3      3      3     103.7      185.3      30.2
4      4      4      83.3      177.1      26.6
...     ...     ...     ...     ...     ...
8383    8383    8383     94.3     178.8     29.5
8384    8384    8384     82.8     147.8     37.9
8385    8385    8385    108.8     168.7     38.2
8386    8386    8386     79.5     176.4     25.5
8387    8387    8387     59.7     167.5     21.3
```

```
[8388 rows x 4 columns]>
```

```
a.columns
```

```
Index(['Unnamed: 0', 'Weight (kg)', 'Height (cm)', 'BMI(kg/m**2)'], dtype='object')
```

```
a.shape
```

```
(8388, 4)
```

```
a=a.drop("Weight (kg)",axis=1)
```

```
a.shape
```

```
(8388, 3)
```

```
a.isnull().sum()
```

```
Unnamed: 0      0
Height (cm)      0
BMI(kg/m**2)     0
dtype: int64
```

```
a['Height (cm)']=a['Height (cm)'].fillna(a['Height (cm)'].mean())
```

```
a.isnull().sum()
```

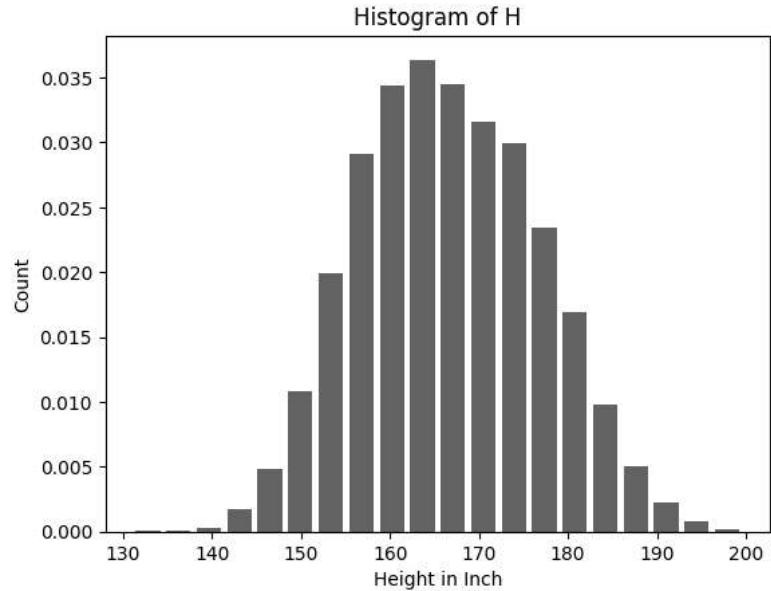
```
Unnamed: 0      0
Height (cm)      0
BMI(kg/m**2)     0
dtype: int64
```

```
import matplotlib.pyplot as plt
```

```
from scipy.stats import norm
```

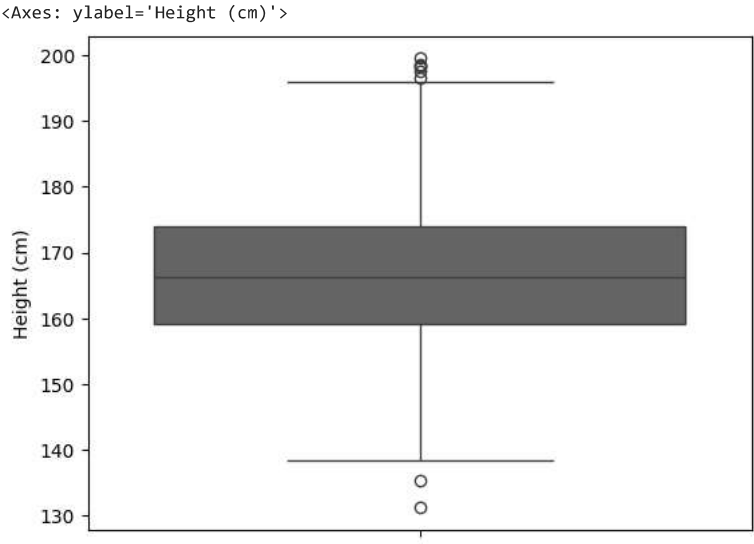
```
plt.hist(a['Height (cm)'], bins=20, rwidth=0.8, density=True)
plt.title("Histogram of H")
plt.xlabel("Height in Inch")
plt.ylabel("Count")
```

```
Text(0, 0.5, 'Count')
```



```
import seaborn as sb

sb.boxplot(a['Height (cm)'])
```



```
a.describe()
```

	Unnamed: 0	Height (cm)	BMI(kg/m**2)
count	8388.000000	8388.000000	8388.000000
mean	4193.500000	166.641190	30.034859
std	2421.551362	10.079013	7.565376
min	0.000000	131.100000	14.200000
25%	2096.750000	159.100000	24.900000
50%	4193.500000	166.200000	28.800000
75%	6290.250000	173.900000	33.800000
max	8387.000000	199.600000	92.300000

```
#####3333Z-Score#####
#upper limit
ul=a['Height (cm)'].mean()+3*a['Height (cm)'].std()

ll=a['Height (cm)'].mean()-3*a['Height (cm)'].std()

print(ul)
```

```
196.87823017566316

print(l1)

136.4041494142272

a.loc[(a['Height (cm)']>=u1) | (a['Height (cm)']<=l1)]
```

Unnamed: 0	Height (cm)	BMI(kg/m**2)	
60	60	198.7	27.1
1906	1906	135.3	29.4
2165	2165	131.1	35.1
3379	3379	197.7	24.9
4026	4026	198.4	23.8
5815	5815	198.3	27.7
7576	7576	199.6	29.5

```
#trimming
a1=a.loc[(a['Height (cm)']<=u1) & (a['Height (cm)']>=l1)]

a1
```

Unnamed: 0	Height (cm)	BMI(kg/m**2)	
0	0	160.2	37.8
1	1	182.3	29.7
2	2	184.2	21.9
3	3	185.3	30.2
4	4	177.1	26.6
...
8383	8383	178.8	29.5
8384	8384	147.8	37.9
8385	8385	168.7	38.2
8386	8386	176.4	25.5
8387	8387	167.5	21.3

```
8381 rows × 3 columns

print("Before Trim :",len(a))

Before Trim : 8388

print("After Trim:",len(a1))

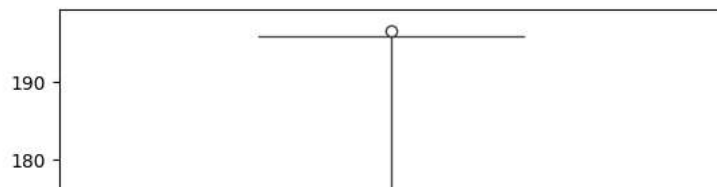
After Trim: 8381

print("No of outliers :",len(a)-len(a1))

No of outliers : 7

sb.boxplot(a1['Height (cm)'])
```

<Axes: ylabel='Height (cm) '>



#capping

a2=a.copy()

a2.loc[(a2['Height (cm)']>=u1),'Height (cm)']=u1

|

|

|

a2.loc[(a2['Height (cm)']<=l1),'Height (cm)']=l1

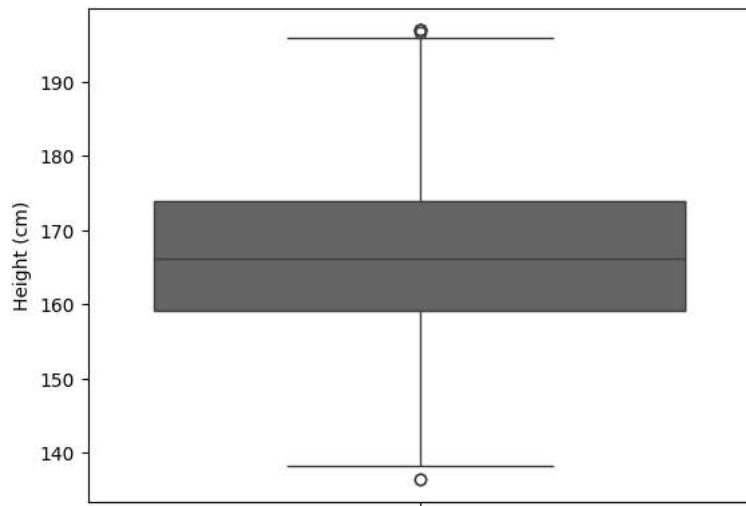
|

|

|

sb.boxplot(a2['Height (cm)'])

<Axes: ylabel='Height (cm) '>



print("Before Trim :",len(a))

Before Trim : 8388

print("After Trim :",len(a2))

After Trim : 8388