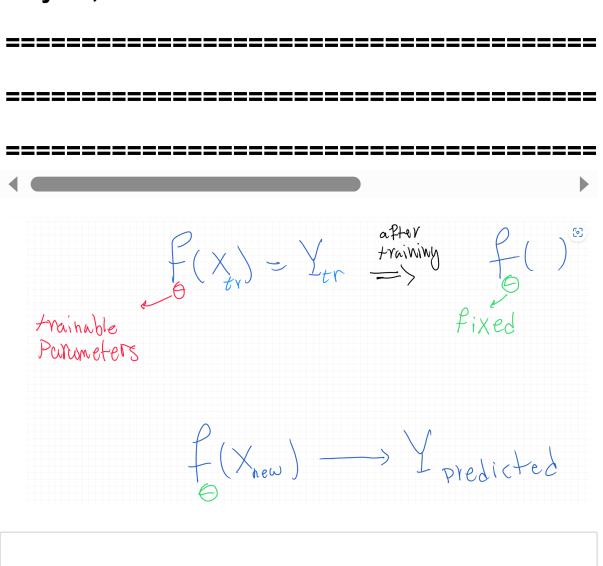
Course	Machine Learning Fundamentals (ML)
Instructor	Hossen Teimoorinia

### May 13, 2023

In [ ]:

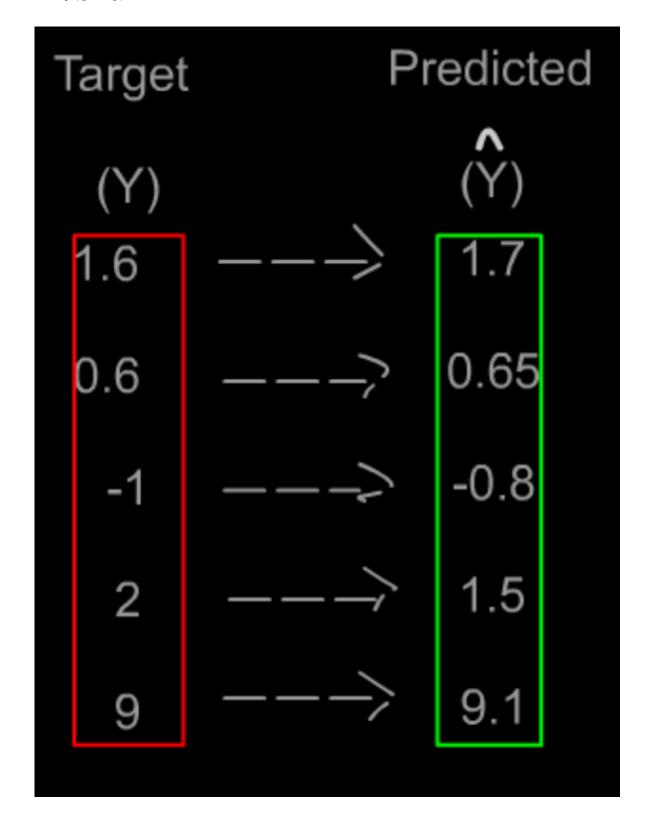


## A practical example in SKIrean

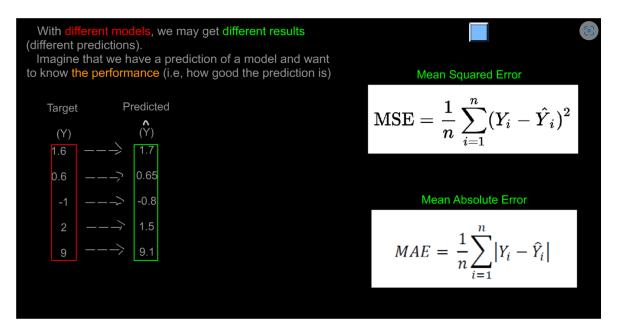
### In [1]:

from sklearn import linear\_model # choose your model
reg = linear\_model.Lasso(alpha=0.1) # set up the model hyper-paramaters
reg.fit([[0, 0], [1, 1]], [0, 1]) # train
reg.predict([[1, 1]]) # predict

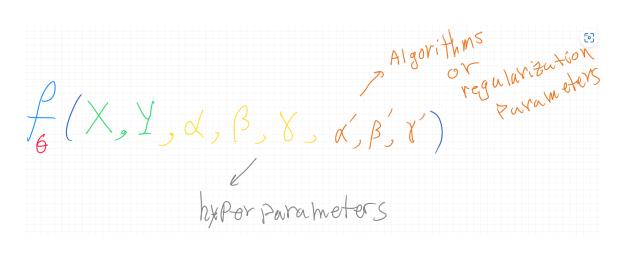
#### Out[1]: array([0.8])



### **Error / loss functions**



# A typical model (regression, classification or clustring)



# Classification, regression and clustring models: What are practical steps in predictions (what to do)?

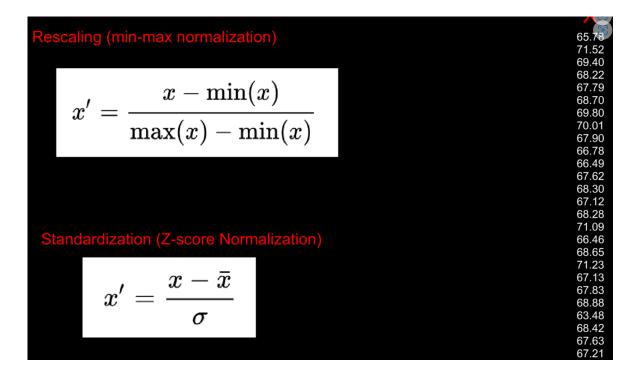
- 1- Define the problem: Clearly state the problem you want to solve, the inputs, and the desired outputs. Identify the type of machine learning task: classification, regression, clustering, etc.
- 2- Collect and preprocess data: Gather a dataset that is representative of the problem domain. Preprocess the data by cleaning, handling missing values, converting categorical data to numerical data, and normalizing or scaling the features if necessary.

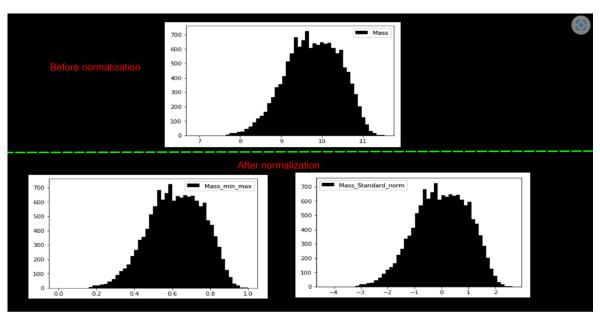
- 3- Split the data: Divide the dataset into training, validation (optional), and testing subsets. A common split ratio is 70% for training, 30% for validation.
- 4- Select a model: Choose a suitable machine learning algorithm based on the problem and data. Examples include linear regression, logistic regression, support vector machines, decision trees, random forests, and neural networks.
- 5- Train the model: Feed the training data to the algorithm, allowing it to learn the underlying patterns in the data. This usually involves minimizing a loss function through an optimization process, such as gradient descent or another optimization algorithm.
- 6- Tune hyperparameters: Optimize the hyperparameters of the model using the validation set. This can involve techniques like grid search, random search, or Bayesian optimization.
- 7- Evaluate the model: Assess the performance of the trained model on the test set using appropriate evaluation metrics (e.g., accuracy, precision, recall, F1 score, mean squared error, etc.). This step helps you understand how well the model generalizes to new, unseen data.
- 8- Iterate and improve: If the model's performance is not satisfactory, you may need to iterate by revisiting steps 2-7. This can involve collecting more data, trying different algorithms, or further tuning hyperparameters.
- 9- Deploy the model: Once you're satisfied with the model's performance, deploy it to a production environment, where it can be used to make predictions on new data

### In order to proceed, it is crucial to emphasize the significance of data preparation and normalization as the second step.

# 1- Cleaning the data can be done by Pandas or Numpy

### 2 - Normalizatin can by done by SKlearn





## A real and complex situation

the file loaded

In [3]: data\_exp # See the table

#### Out[3]:

	tar_1	tar_2	tar_3	tar_4	tar_5	tar_6	tar_7	tar_8	tar_9	tar_1
0	NaN	8.92791	8.54151	8.66548	8.66583	NaN	NaN	8.49927	8.44851	Nal
1	NaN	8.79355	NaN	Nal						
2	NaN	9.06602	NaN	Nal						
3	8.86975	8.88465	8.51746	8.64180	8.67329	8.96472	8.97729	8.48092	8.45596	8.8175
4	NaN	8.85254	NaN	8.73630	8.64105	NaN	NaN	8.54667	8.42411	Nal
158883	8.94060	8.96944	8.60891	8.66408	8.72856	9.06701	9.05309	8.49821	8.51291	8.9068
158884	NaN	9.08881	NaN	Nal						
158885	NaN	8.91223	8.60198	8.72192	8.69056	NaN	NaN	8.53811	8.47344	Nal
158886	NaN	9.07084	NaN	Nal						
158887	8.92745	8.95868	8.68229	8.68525	8.66979	8.98096	8.98990	8.51373	8.45245	8.8362

158888 rows × 33 columns

```
In [4]: data exp.info() # get information and check NaNs
        <class 'pandas.core.frame.DataFrame'>
        RangeIndex: 158888 entries, 0 to 158887
        Data columns (total 33 columns):
             Column
                        Non-Null Count
                                        Dtype
             ----
                        -----
                                         ----
         0
             tar 1
                        75852 non-null
                                         float64
         1
             tar 2
                        158852 non-null float64
         2
             tar_3
                        92991 non-null
                                        float64
         3
             tar 4
                        94951 non-null
                                        float64
         4
             tar_5
                        97650 non-null
                                        float64
         5
             tar_6
                        59447 non-null
                                        float64
         6
                                        float64
             tar_7
                        60146 non-null
         7
                                        float64
             tar 8
                        97634 non-null
         8
             tar 9
                        96436 non-null
                                        float64
         9
             tar 10
                        60306 non-null
                                        float64
         10 tar_11
                        95643 non-null
                                        float64
         11
            tar 12
                        97047 non-null
                                        float64
         12
                        158888 non-null float64
             Z
                        158852 non-null float64
         13
             mass
         14
                        158888 non-null float64
            agn_sn
         15 cf g
                        158888 non-null float64
                        158888 non-null float64
         16
            cf_r
         17 Ha
                        158888 non-null float64
         18
            Hb
                        158888 non-null float64
                        158888 non-null float64
         19 OII 1
         20 OII 2
                        158888 non-null float64
         21 OIII
                        158888 non-null float64
         22 NII
                        158888 non-null float64
         23
            SII 1
                        158888 non-null float64
         24 SII 2
                        158888 non-null float64
         25 err Ha
                        158888 non-null float64
                        158888 non-null float64
         26 err Hb
         27
            err_OII_1 158888 non-null float64
         28 err OII 2 158888 non-null float64
         29 err OIII
                        158888 non-null float64
         30 err NII
                        158888 non-null float64
            err SII 1 158888 non-null float64
         31
         32 err SII 2 158888 non-null float64
        dtypes: float64(33)
        memory usage: 40.0 MB
In [5]: data_exp.shape # get the shape
Out[5]: (158888, 33)
In [6]: print(type(data exp))
```

<class 'pandas.core.frame.DataFrame'>

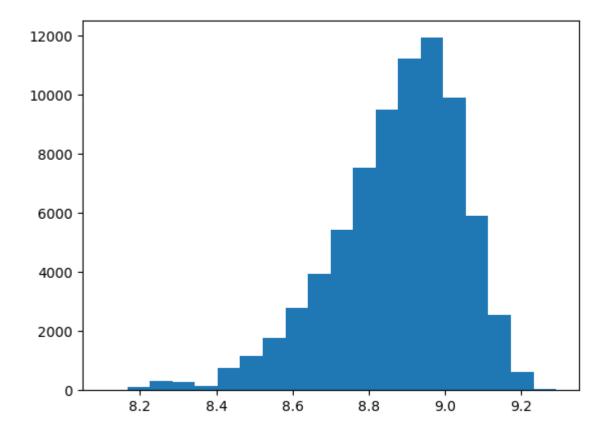
```
In [7]: data_clean_all=data_exp.dropna() # Choose a new name for the clean data set
In [8]: data clean all.shape # get shape
Out[8]: (56619, 33)
In [9]: # select a subset of data and then drop NaNs
        data_clean_selected= data_exp[['Ha', 'Hb', 'OII_1','NII','SII_2', 'tar_1']].d
        print (data clean selected)
                     На
                               Hb
                                      OII 1
                                                   NII
                                                          SII 2
                                                                   tar_1
        3
                514.525
                          180.6580
                                    268.702
                                              139.3550
                                                        73.0610 8.86975
        5
               1011.950 355.5870
                                    494.637
                                              470.7900 124.1400 9.10778
        6
                284.017 99.6942
                                    120.660
                                              119.6000
                                                       47.5061 9.11417
        8
                542.953
                          190.6650
                                    112.974
                                              213.9970
                                                        69.0888 9.08641
               1176.250 413.0610 520.896
                                              299.1270 158.9760 8.86924
        11
                    . . .
                                        . . .
                                                   . . .
                                                            . . .
        . . .
        158879 3185.680 1118.5100 2230.150
                                              522.5390 359.2150 8.65184
                                   205.794 42.4103 38.3080 8.80593
        158880
               311.140 109.2260
        158882 3832.080 1346.6800 1474.790 1308.8700 408.8190 9.00067
        158883 687.306
                          241.4110
                                   213.749
                                              202.9140
                                                       85.3614 8.94060
        158887
                321.273
                          112.7690
                                    150.629
                                              102,4770
                                                        41.4276 8.92745
        [75852 rows x 6 columns]
```

# To utilize supervised methods like classification or regression, you can split the data into input and target sets.

```
In [13]: print(type(input data))
         <class 'numpy.ndarray'>
In [14]: import numpy as np # to work with numpy files
In [15]: input_data # just vlaues
Out[15]: array([[ 514.525 , 180.658 , 268.702 , 139.355 , 73.061 ],
                [1011.95 , 355.587 , 494.637 , 470.79 , 124.14 ],
                [ 284.017 , 99.6942 , 120.66 , 119.6 ,
                                                           47.5061],
                . . . ,
                [3832.08 , 1346.68 , 1474.79 , 1308.87 , 408.819 ],
               [ 687.306 , 241.411 , 213.749 , 202.914 ,
                                                            85.3614],
                [ 321.273 , 112.769 , 150.629 , 102.477 , 41.4276]])
In [16]: input_data[:,1] # the content of the second column
Out[16]: array([ 180.658 , 355.587 , 99.6942, ..., 1346.68 , 241.411 ,
                112.769 ])
In [17]: | target_data[:,0] # the content of the first column of the target data
Out[17]: array([8.86975, 9.10778, 9.11417, ..., 9.00067, 8.9406 , 8.92745])
```

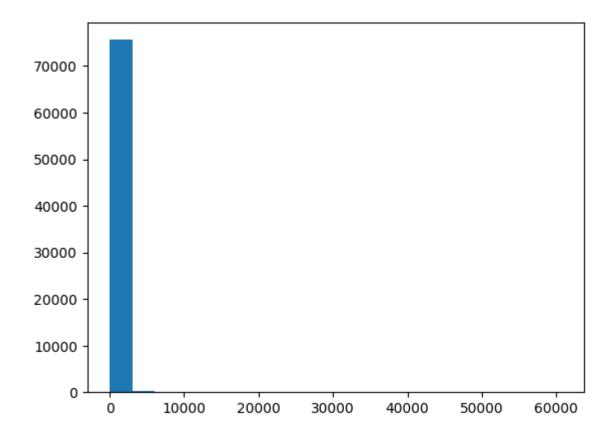
In [18]: # histogram can give us for information
 from matplotlib import pyplot as plt
 plt.hist(target\_data[:,0], bins=20)

```
Out[18]: (array([
                                           267.,
                                                                         1781.,
                    23.,
                           104.,
                                   329.,
                                                   160.,
                                                          761., 1146.,
                  2774.,
                          3942.,
                                 5419.,
                                         7531., 9507., 11204., 11919.,
                                                                         9880.,
                                            43.]),
                  5914.,
                          2539.,
                                   609.,
          array([8.10717
                         , 8.1664485, 8.225727 , 8.2850055, 8.344284 , 8.4035625,
                 8.462841 , 8.5221195, 8.581398 , 8.6406765, 8.699955 , 8.7592335,
                 8.818512 , 8.8777905, 8.937069 , 8.9963475, 9.055626 , 9.1149045,
                 9.174183 , 9.2334615, 9.29274 ]),
          <BarContainer object of 20 artists>)
```



```
In [19]: plt.hist(input_data[:,1],bins=20) # the distribution of the second column
print (len(input_data[:,1])) # the number of samples
```

75852



```
In [20]: print ('max= ', np.max(input_data[:,1])) # maximum of the column
print ('min= ',np.min(input_data[:,1])) # maximum of the column
print ('mean= ',np.mean(input_data[:,1])) # maximum of the column
print ('median= ', np.median(input_data[:,1])) # maximum of the column

max= 60645.8
```

```
max= 60645.8
min= 14.4689
mean= 348.8803033565364
median= 226.672
```

Measures of Central Tendency: These are ways to describe the center of a data set. Commonly used measures include:

Mean: The average of all data points. Median: The middle value in a sorted dataset. Mode: The most frequently occurring value in the dataset. Measures of Dispersion or Variability: These are ways to describe the spread of a data set. Commonly used measures include:

Range: The difference between the maximum and minimum values in the dataset. Variance: The average of the squared differences from the mean. Standard Deviation: The square root of the variance. It is a measure of the amount of variation or dispersion in a set of values.

Interquartile Range (IQR): The range between the first quartile (25th percentile) and the third quartile (75th percentile). Measures of Position: These tell you where a specific data point is within the data set. Common measures include percentiles and quartiles.

Measures of Shape: These describe the shape of the data distribution.

Skewness: Measures the asymmetry of the data. Kurtosis: Measures the "tailedness" of the data distribution. Descriptive statistics are often presented in tables, histograms, bar charts, pie charts, frequency distributions, and other graphical or pieterial methods to visualize the

In [21]: # descriptive analysis
data\_exp.describe()

#### Out[21]:

	tar_1	tar_2	tar_3	tar_4	tar_5	tar_6	
count	75852.000000	158852.000000	92991.000000	94951.000000	97650.000000	59447.000000	(
mean	8.874875	7.707753	8.546300	8.646865	8.643520	8.897636	
std	0.169721	11.217666	0.241618	0.100803	0.106223	0.168795	
min	8.107170	-99.900000	7.762050	7.962620	7.735390	8.352340	
25%	8.778530	8.744627	8.379140	8.593565	8.587332	8.779070	
50%	8.900750	8.941810	8.572460	8.665100	8.666935	8.916170	
75%	8.996460	9.049900	8.734620	8.715075	8.722340	9.030810	
max	9.292740	9.466640	9.049880	8.849830	8.843860	9.278110	

8 rows × 33 columns



1- Measures of Central Tendency: These are ways to describe the center of a data set. Commonly used measures include:

Mean: The average of all data points.

Median: The middle value in a sorted dataset.

Mode: The most frequently occurring value in the dataset.

2- Measures of Dispersion or Variability: These are ways to describe the spread of a data set. Commonly used measures include:

Range: The difference between the maximum and minimum values in the dataset.

Variance: The average of the squared differences from the mean.

Standard Deviation: The square root of the variance. It is a measure of the amount of variation or dispersion in a set of values.

Interquartile Range (IQR): The range between the first quartile (25t h percentile) and the third quartile (75th percentile).

- 3- Measures of Position: These tell you where a specific data point is within the data set. Common measures include percentiles and quartiles.
- 4- Measures of Shape: These describe the shape of the data distribution.

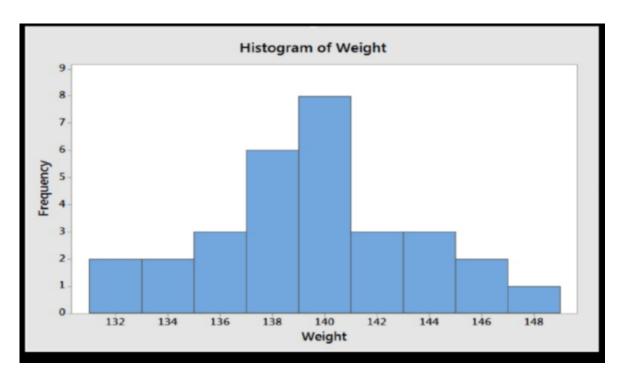
Skewness: Measures the asymmetry of the data.

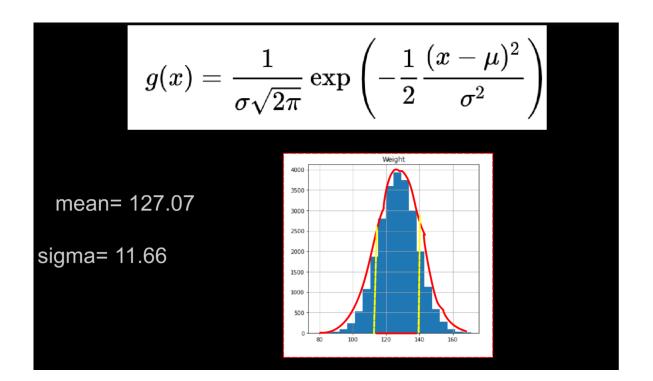
Kurtosis: Measures the "tailedness" of the data distribution.

Descriptive statistics are often presented in tables, histograms, bar charts, pie charts, frequency distributions, and other graphical or pictorial methods to visualize the data. They provide a foundational understanding for a dataset before proceeding to more complex statistical analyses or predictive modeling.

### How to remove outliers?

# A little about distributions and the information theory



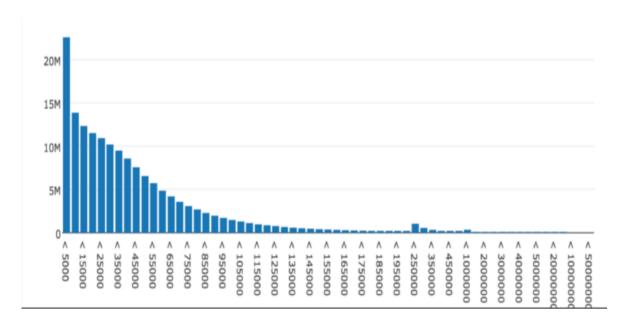


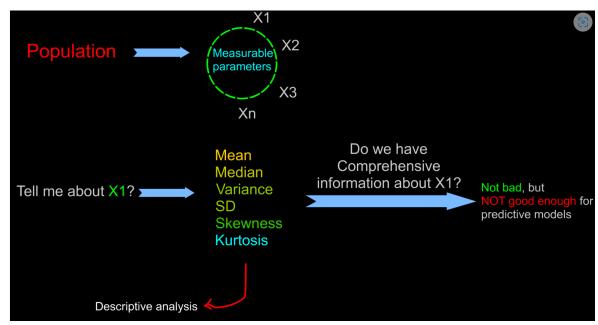
------For Height and Weight file:

chaose a feature (weight)

Find Mo (mean) and Go (standard diviation)

$$g(x) = \frac{1}{G_0 \text{ Fix}} \text{ exp} \left[-\frac{1}{2} \frac{(x-u_0)^2}{G_0^2}\right]$$





## HT: Draw an example of outliers; how to make a decision to remove them

## what is the difference between ML, statistics and data science

```
In [22]: lst =[0, 2,3,4 ,8, 1000, 11]
lst = np.array(lst)
```

### **Percentile by Numpy**

# Draw to show the case regarding input\_data and target\_data

```
In [29]: input_good = input_data[idx_good]
    target_good = target_data[idx_good]
    print('size of input_good= ', np.shape(input_good))
    print('size of target_good=', np.shape(target_good))

size of input_good= (64486, 5)
    size of target_good= (64486, 1)
```

## >>>>> The class activity (number 2)

Remove ouliers for percentile= 95 and create new input\_data and target\_data without ouliers and 0 for all cloumns

# speliting input\_data and target data to the training set (70%) and the validation set (30%)

```
In [32]: # we want to normalize the input using SKlearn
    # import two models "StandardScaler" and "MinMaxScaler"
    from sklearn.preprocessing import StandardScaler, MinMaxScaler #line #1

# choose one model with proper name
    scaler_MinMax= MinMaxScaler()

# Find the trainable paramters of the model

scaler_MinMax.fit(X_train)

X_train_MinMax= scaler_MinMax.transform(X_train) # Line #4

X_test_MinMax= scaler_MinMax.transform(X_test) # Line #4

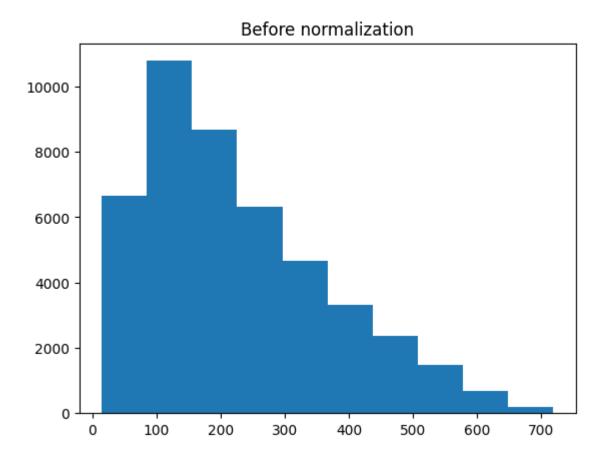
print('Done!')
```

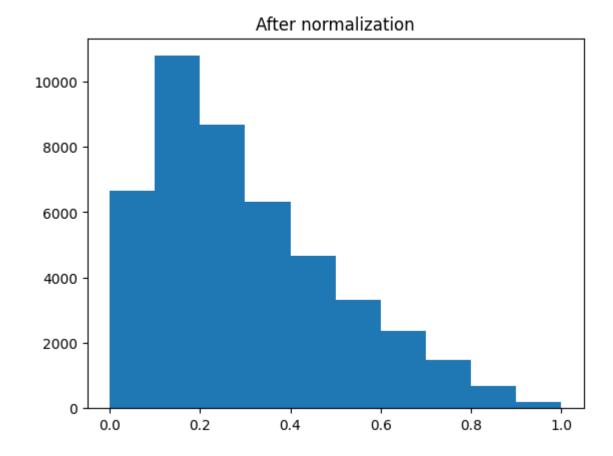
Done!

```
In [33]: plt.figure(1)
    plt.hist(X_train[:,1])
    plt.title('Before normalization')

plt.figure(2)
    plt.hist(X_train_MinMax[:,1])
    plt.title('After normalization')
```

Out[33]: Text(0.5, 1.0, 'After normalization')





>>>> Activity #3

Choose weight-height file and split the data to training (75%) and test set(25%) and normalize the split data

In [ ]:	