
Assignment II: Python Basics (Datafile Import, Plotting, Functions, Classes)

Exercises in Machine Learning (190.013), SS2022
Stefan Nehl¹

¹stefan-christopher.nehl@stud.unileoben.ac.at, MNR: 00935188, Montanuniversität Leoben, Austria

March 6, 2022

In the second assignment, I had to create basic statistic functions for calculating the mean, median, variance and standard deviation. Furthermore, I had to read in a CSV file, use the created functions and plot the values.

1 Introduction

The statistic functions I implemented were calculating the mean, median, variance, and standard deviation. Furthermore, I added a function for normalizing and standardizing a data set.

2 Implementation

For the implementation I created the class *BasicStatistics* with all the statistic functions. First, I implemented the mean.

2.1 Mean

The mean is the average of a collection of numbers. The formula for the calculation is the following (Institute, n.d.[a]):

$$mean = \frac{x_1 + x_2 + \dots + x_n}{n}$$

Where n is the number of elements in the collection and x_n the element on the position n . The calculation was implemented in the *getMean* function.

```
1 def getMean(self):
2     length = len(self.dataSet)
3     sumValue = sum(self.dataSet)
4     mean = sumValue / length
```

2.2 Median

The Median is the middle value of a collection of numbers. The steps of the implemented algorithm are (Institute, n.d.[b]):

- Sort the collection of numbers
- Calculate the mid index
- if the length of the collection is uneven: take the element with the mid index
- if the length is even calculate the median:

$$median = \frac{x_{mid} + x_{mid-1}}{2}$$

- return the result

```
1 def getMedian(self):
2     length = len(self.dataSet)
3     mid = length // 2
4
5     if length % 2 == 0:
6         median1 = self.sortedDataSet[mid]
7         median2 = self.sortedDataSet[mid - 1]
8         median = (median1 + median2) / 2
9
10    else:
11        median = self.sortedDataSet[mid]
12
13    return median
```

2.3 Variance

The variance is the expected variation between values in a collection of numbers. The formula for the

calculation of the variance it the following (Ramos, 2021):

$$\text{Variance } \sigma^2 = \frac{\sum_{i=0}^{n-1} (x_i - \bar{x})^2}{n - 1}$$

Where n is the number of elements, x_i the element on the index i and \bar{x} the mean. This formula uses the Bessel's correction for smaller numbers. Therefore, instead of dividing the aggregated values with n , I divided them with $n-1$ (Ramos, 2021).

```

1  def getVariance(self):
2      length = len(self.dataSet)
3      mean = self.getMean()
4
5      squareDeviations = [(x - mean)
6                          ** 2 for x in self.dataSet
7                          ]
8
9      #Bessel's correction (n-1)
      instead of n for better
      results
      variance = sum(
          squareDeviations) / (length
                                -1)
10     return variance

```

2.4 Standard Deviation

The standard deviation is the amount of the variation of a collection of numbers. The formula for the calculation of the standard deviation is (Ramos, 2021):

$$\text{Standard Deviation } \sigma = \sqrt{\sigma^2}$$

```

1  def getStandardDeviation(self):
2      variance = self.getVariance()
3      standardDeviation = math.sqrt(
4          variance)
5      return standardDeviation

```

2.5 Normalize Data

Standardization is a preprocessing step, to standardize the range of values of a collection of numbers. I used the following formula to standardize the data (Jaadi, 2019).

$$z = \frac{x_i - \bar{x}}{\sigma}$$

Where z is the standardize value, x_i the value in the collection on index i , \bar{x} the mean and σ the standard deviation.

```

1  def getNormalizeDataSet(self):
2      mean = self.getMean()
3      standardDeviation = self.
4          getStandardDeviation()
5      standardizeDataSet = [(x -
6                             mean)/standardDeviation)
7      for x in self.dataSet]

```

```

return standardizeDataSet

```

2.6 Testing the functions

I implemented a script file with the name *BasicStatistics_Test* which tests all the function with two data sets.

- 1, 2, 3, 4, 5
- 1, 2, 3, 4, 5, 6

2.7 Loading the data

The data which I should analyse is stored in the file *gauss.csv*. The data is one dimensional and I read the lines in the file with the function *reader()* of the imported csv package. The data was then stored in the array *dataSet*.

```

1  with open('gauss.csv', mode='r') as
      file:
2      csvFile = csv.reader(file)
3
4      for line in csvFile:
5          dataSet.append(float(line[0]))

```

2.8 Plotting the results

For plotting the results, I created a plot with the *figure()* function of *matplotlib* with the width of 8 inch and the height of 6 inch. I set the subtitle of the plot to *Data Distribution* and created a subplot. We need three different subplots in the plot. One histogram and two scatter plots. First, I created the histogram with the function *subplot(2,1,1)*. This subplot consumes two columns and one row and starts at the first position. I added the needed labels and plotted the values with the *hist()* function. I used 20 bins and the density property for the histogram. Also four vertical lines with the mean, median and the standard deviation with plus and minus was added to the subplot of the histogram. The other two subplots contain the raw and the standardized data, where each of the plots had three horizontal lines with the mean and the standard deviation.

3 Results

The charts in Figure 1 display the result of the basic statistic functions. The histogram at the top displays the data distribution of the normalized data with the mean, median and standard deviation, the scatter chart at the bottom left the raw data with the mean and standard deviation and the scatter chart at the bottom right the normalized data with mean and standard deviation.

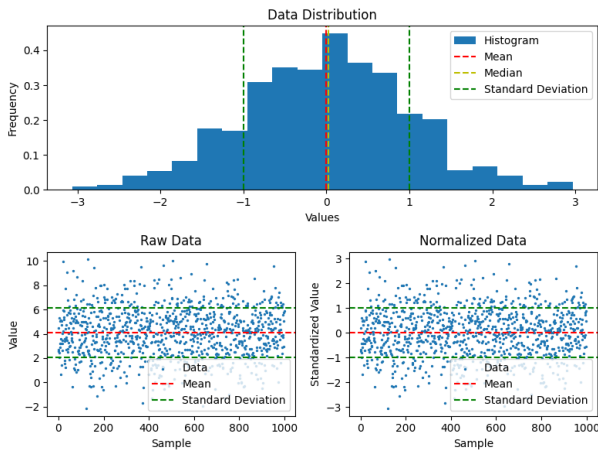


Figure 1: Illustrates the raw data in a Histogram and the raw and standardized data in a scatter chart

4 Conclusion

After some clarification of the tasks, implementing the statistic function was straightforward. The code needed for the implementation could keep clean and minimalistic. Only creating the plot functions needed some additional lines of code and the function `.tight_layout()` to improve the readability of the charts and labels.

APPENDIX

BasicStatistics class:

```

1  """
2  @author: Nehl Stefan
3  """
4
5  import math
6
7  class BasicStatistics:
8
9      def __init__(self, dataSet):
10         if self.checkDataSet(dataSet)
            == False:
11             return None
12
13         #clone dataSet
14         self.dataSet = dataSet[:]
15         self.sortedDataSet = dataSet
            [:]
16         self.sortedDataSet.sort()
17
18     def checkDataSet(self, dataSet):
19         if dataSet == None:
20             return False
21
22         if len(dataSet) == 0:
23             return False
24

```

```

25     return True
26
27     def getMean(self):
28         length = len(self.dataSet)
29         sumValue = sum(self.dataSet)
30         mean = sumValue / length
31         return mean
32
33     def getMedian(self):
34         length = len(self.dataSet)
35         mid = length // 2
36
37         if length % 2 == 0:
38             median1 = self.
                sortedDataSet[mid]
39             median2 = self.
                sortedDataSet[mid - 1]
40             median = (median1 +
                median2) / 2
41         else:
42             median = self.
                sortedDataSet[mid]
43
44         return median
45
46     def getVariance(self):
47         length = len(self.dataSet)
48         mean = self.getMean()
49
50         squareDeviations = [(x - mean)
51                               ** 2 for x in self.dataSet
52                               ]
53
54         #Bessel's correction (n-1)
55         #instead of n for better
56         #results
57         variance = sum(
58             squareDeviations) / (length
59             -1)
60         return variance
61
62     def getStandardDeviation(self):
63         variance = self.getVariance()
64         standardDeviation = math.sqrt(
65             variance)
66         return standardDeviation
67
68     def getMinValue(self):
69         return self.sortedDataSet[0]
70
71     def getMaxValue(self):
72         return self.sortedDataSet[len(
73             self.sortedDataSet) - 1]
74
75     def getNormalizeDataSetOld(self):
76         min = self.getMinValue()
77         max = self.getMaxValue()
78

```

```

71     dataNorm = [(i - min) / (max -
72                 min) for i in self.dataSet
73                 ]
74     return dataNorm
75
76 def getNormalizeDataSet(self):
77     mean = self.getMean()
78     standardDeviation = self.
79         getStandardDeviation()
80     standardizeDataSet = [((x -
81                             mean)/standardDeviation)
82                             for x in self.dataSet]
83
84     return standardizeDataSet

```

BasicStatistics_Test class:

```

1  """
2  @author: Nehl Stefan
3  """
4
5  from BasicStatistics import
6      BasicStatistics
7
8  dataSet1 = [1, 2, 3, 4, 5]
9  dataSet2 = [1, 2, 3, 4, 5, 6]
10
11  basicStatistics1 = BasicStatistics(
12      dataSet1)
13  basicStatistics2 = BasicStatistics(
14      dataSet2)
15  result = -1
16
17  #Mean
18  print('Mean')
19  expectedResult = 3
20  result = basicStatistics1.getMedian()
21  print('Expected: ' + str(
22      expectedResult) + ' Result: ' + str
23      (result))
24
25  expectedResult = 3.5
26  result = basicStatistics2.getMedian()
27  print('Expected: ' + str(
28      expectedResult) + ' Result: ' + str
29      (result))
30  print()
31
32  #Median
33  print('Median')
34  expectedResult = 3
35  result = basicStatistics1.getMedian()
36  print('Expected: ' + str(
37      expectedResult) + ' Result: ' + str
38      (result))
39
40  expectedResult = 3.5
41  result = basicStatistics2.getMedian()
42  print('Expected: ' + str(
43      expectedResult) + ' Result: ' + str
44      (result))
45  print()
46
47  #Variance
48  print('Variance')
49  expectedResult = 2.5
50  result = basicStatistics1.getVariance
51      ()
52  print('Expected: ' + str(
53      expectedResult) + ' Result: ' + str
54      (result))
55
56  expectedResult = 3.5
57  result = basicStatistics2.getVariance

```

```

    ()
44 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
45 print()
46
47 #StandardDefiation
48 print('Standard Defiation')
49 expectedResult = 1.5811388
50 result = basicStatistics1.
    getStandardDeviation()
51 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
52
53 expectedResult = 1.8708287
54 result = basicStatistics2.
    getStandardDeviation()
55 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
56 print()
57
58 #FindMinValue
59 print('Find MinValue')
60 expectedResult = 1
61 result = basicStatistics1.getMinValue
    ()
62 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
63
64 print('Find MinValue')
65 expectedResult = 1
66 result = basicStatistics2.getMinValue
    ()
67 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
68
69 #FindMaxValue
70 print('Find MaxValue')
71 expectedResult = 5
72 result = basicStatistics1.getMaxValue
    ()
73 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
74
75 print('Find MaxValue')
76 expectedResult = 6
77 result = basicStatistics2.getMaxValue
    ()
78 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
79
80 #Normalize Data
81 print('Standardized Data')

```

```

82 expectedResult = "??"
83 result = basicStatistics1.
    getNormalizeDataSetOld()
84 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
85
86 expectedResult = "??"
87 result = basicStatistics2.
    getNormalizeDataSetOld()
88 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
89 print()
90
91 #Standardized Data
92 print('Standardized Data')
93 expectedResult = "??"
94 result = basicStatistics1.
    getNormalizeDataSet()
95 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
96
97 expectedResult = "??"
98 result = basicStatistics2.
    getNormalizeDataSet()
99 print('Expected: ' + str(
    expectedResult) + ' Result: ' + str
    (result))
100 print()

```

Plotting of the values:

```

1  """
2  @author: Nehl Stefan
3  """
4
5  import csv
6
7  import matplotlib.pyplot as plt
8  import numpy
9  from BasicStatistics import
    BasicStatistics
10
11 dataSet = []
12
13
14 with open('gauss.csv', mode='r') as
    file:
15     csvFile = csv.reader(file)
16
17     for line in csvFile:
18         dataSet.append(float(line[0]))
19
20 basicStatistics = BasicStatistics(
    dataSet)
21
22 #dataSet = basicStatistics.
    normalizeDataSet(dataSet)
23 mean = basicStatistics.getMean()
24 median = basicStatistics.getMedian()
25 variance = basicStatistics.getVariance
    ()
26 standardDeviation = basicStatistics.
    getStandardDeviation()
27 standardizedDataSet = basicStatistics.
    getNormalizedDataSet()
28
29 standardizedStatistics =
    BasicStatistics(standardizedDataSet
    )
30 standardizedMean =
    standardizedStatistics.getMean()
31 standardizedMedian =
    standardizedStatistics.getMedian()
32 standardizedStandardDeviation =
    standardizedStatistics.
    getStandardDeviation()
33
34
35 length = len(standardizedDataSet)
36 plotRange = range(length)
37
38 plt.suptitle('Data Distribution')
39
40 plt.figure(figsize=(8,6))
41 # plot histogram (left, right, top)
42 plt.subplot(2, 1, 1)
43
44     density=True, label='Histogram')
45 plt.axvline(standardizedMean, label='
    Mean', color='r', ls='--')
46 plt.axvline(standardizedMedian, label=
    'Median', color='y', ls='--')
47 plt.axvline(standardizedMean -
    standardizedStandardDeviation,
    label='Standard Deviation', color='
    g', ls='--')
48 plt.axvline(standardizedMean +
    standardizedStandardDeviation,
    color='g', ls='--')
49
50 plt.title("Data Distribution")
51 plt.xlabel('Values')
52 plt.ylabel('Frequency')
53 plt.legend(loc="upper right")
54
55 # plot raw data (left, bottom)
56 plt.subplot(2, 2, 3)
57
58 plt.scatter(plotRange, dataSet, label=
    'Data', s=2)
59 plt.axhline(mean, label='Mean', color=
    'r', ls='--')
60 #plt.axhline(median, label='Median',
    color='y', ls='--')
61 plt.axhline(mean + standardDeviation,
    label='Standard Deviation', color='
    g', ls='--')
62 plt.axhline(mean - standardDeviation,
    color='g', ls='--')
63
64 plt.title('Raw Data')
65 plt.xlabel('Sample')
66 plt.ylabel('Value')
67 plt.legend(loc="best")
68
69 # plot standardizedData(left, bottom)
70 plt.subplot(2, 2, 4)
71 plt.title('Normalized Data')
72 plt.xlabel('Sample')
73 plt.ylabel('Standardized Value')
74
75 plt.scatter(plotRange,
    standardizedDataSet, label='Data',
    s=2)
76 plt.axhline(standardizedMean, label='
    Mean', color='r', ls='--')
77 #plt.axhline(standardizedMedian, label
    ='Median', color='y', ls='--')
78 plt.axhline(standardizedMean +
    standardizedStandardDeviation,
    label='Standard Deviation', color='
    g', ls='--')
79 plt.axhline(standardizedMean -
    standardizedStandardDeviation,
    color='g', ls='--')
80

```

```
81 plt.legend(loc="best")
82
83 plt.tight_layout()
84 plt.show()
```

Bibliography

Institute, Corporate Finance (n.d.[a]). *Mean*. URL: <https://corporatefinanceinstitute.com/resources/knowledge/other/mean/> (visited on 04/03/2022).

– (n.d.[b]). *Median*. URL: <https://corporatefinanceinstitute.com/resources/knowledge/other/median/> (visited on 04/03/2022).

Jaadi, Zakaria (2019). *Standardization*. URL: <https://builtin.com/data-science/when-and-why-standardize-your-data> (visited on 04/03/2022).

Ramos, Leodanis Pozo (2021). *Calculating Variance and Standard Deviation in Python*. URL: <https://stackabuse.com/calculating-variance-and-standard-deviation-in-python/> (visited on 04/03/2022).

Zach (2021). *How to Normalize Data in Python*. URL: <https://www.statology.org/normalize-data-in-python/> (visited on 04/03/2022).