
Assignment III: Gauss and Beta Distribution

Exercises in Machine Learning (190.013), SS2022
Stefan Nehl¹

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

March 20, 2022

In the third assignment, I had to create the abstract class *ContinuousDistribution*. This class contains functions for importing and exporting the data to CSV files, compute the mean and standard deviation, generating samples and visualize the given data. Furthermore, I had to implement two additional classes, *GaussDistribution* and *BetaDistribution* which inheritances the class *ContinuousDistribution* and implement gauss and beta distribution. Finally, I plotted the result of this distributions.

1 Introduction

The gauss and beta distributions are often used in machine learning to generate random variables. The target of this assignment was to create a module *inference* which implements this two different distributions.

2 Implementation

For the implementation I created a module with the name *inference*. In this module i implemented the abstract class *ContinuousDistribution* with it's abstract methods.

3 Abstract class *ContinuousDistribution*

First the class was *ContinuousDistribution* created. This class is not completely abstract because some methods I reused in the other classes. For this I implemented first a constructor with the default properties. These properties are *dataSet*, *normalizedDataSet*, *mean*, *median*, *variance* and *standardDeviation*. This properties are reused in the *GaussDistribution* and *BetaDistribution*. The functions, *importCSV*, *exportCSV*, *calculateMean*, *calculateVariance*, *calculateStandardDeviation* and *normalizeDataSet*. Only the method *generateSamples* and

plotData is abstract, because every class needs there own implementation of this functions.

4 Gauss Distribution

The class *GaussDistribution* has a constructor with the parameters *dimension*, sets the dimension of this gauss distribution and the optional parameters *fileName*, for importing a CSV file, *numberOfSamplesToGenerate*, number of samples generated by the class, *mean* and *variance*. Important to mention here is, that the importing of a file and the generating of samples is excluding each other. Only one parameter can be set otherwise the class throws an exception. The construction also sets the data and calculate the needed values like mean and standard deviation and generates the gauss distribution.

4.1 Generating Samples

As already mentioned, the generating of samples needs to be implemented for each class separately. For the generation I used the *random()* function from the *numpy* library with the values of the mean and the variance for the generation and the dimension with the number of samples for the amount of data.

```
1
2 def generateSamples(self):
3     if len(self.dataSet) != 0:
4         raise Exception("Data
5             already added")
6
7     self.dataSet = np.random.
8         default_rng().normal(self.
9             mean, self.variance, size=(
10                 self.numberOfSamples, self.
11                     dimension))
```

4.2 Calculation

For the calculation I implemented two different methods. One for the one dimensional calculation, *generateGaussian1D*, and for the two dimensional calculation, *calculateGaussian2D*. For the one dimensional implementation I used the following formula.

$$N(x|\mu, \sigma) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp^{-\frac{1}{2\sigma^2}(x - \mu)^2}$$

Where μ stands for the mean and σ for the standard deviation. For the two dimensional implementation I used the following formula.

$$N(x|\mu, \sigma) = \frac{1}{(2\pi)^{D/2}|\Sigma|^{1/2}} \exp^{-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)}$$

Where D is the dimension, Σ the covariance and $|\Sigma|$ the determinant of the covariance. For the covariance I used created a vector with the mean and zeros. $\begin{pmatrix} \sigma & 0 \\ 0 & \sigma \end{pmatrix}$

Both formulas are from the book *An Introduction to Probabilistic Machine Learning*. (Rueckert, 2022)

4.3 Plotting

The implementation for the one dimensional plot was straightforward. I plotted a histogram of the generated data and the gauss distribution as a line above the histogram. Additionally, I add the raw the of the generated samples. With the two dimensional plots I had some issues. I was able to create a 3d model of the raw data and the 3d bar chart of the distribution. I wanted to plot the surface of the two dimensional gauss distribution following the paper (Roelants, 2018), but unfortunately I failed to create the surface.

5 Beta Distribution

5.1 Generating Samples

5.2 Calculation

5.3 Plotting

6 Results

7 Conclusion

APPENDIX

```

1 import csv
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import math
5
6 from abc import ABC, abstractmethod
7
8 class ContiniousDustribution():
9
10     def __init__(self):
11         self.dataSet = []
12         self.normalizedDataSet = []
13         self.mean = None
14         self.median = None
15         self.variance = None
16         self.standardDeviation = None
17
18     def importCsv(self, filename):
19         if len(self.dataSet) != 0:
20             raise Exception("Data already added")
21
22         with open(filename, mode="r") as file:
23             csvFile = csv.reader(file)
24
25             for row in csvFile:
26                 self.dataSet.append(row)
27
28     def exportCsv(self, filename):
29         if len(self.dataSet) == 0:
30             raise Exception("No Data added")
31
32         with open(filename, mode="w") as file:
33             csvWriter = csv.writer(file, delimiter = ";")
34             csvWriter.writerows(self.dataSet)
35
36     def calculateMean(self):
37         self.mean = np.mean(self.dataSet)
38
39     def calculateVariance(self):
40         length = len(self.dataSet)
41         mean = self.mean
42
43         squareDeviations = [(x - mean) ** 2 for x in self.dataSet]
44
45         # Bessel's correction (n-1) instead of n for better results
46         self.variance = sum(squareDeviations) / (length - 1)
47         return self.variance
48
49     def calculateStandardDeviation(self):
50         self.standardDeviation = math.sqrt(self.variance)
51         return self.standardDeviation
52
53     def normalizeDataSet(self):
54         self.dataSet = [(x - self.mean)/self.standardDeviation for x in self.dataSet]
55
56     @abstractmethod
57     def generateSampels(self):
58         pass

```

```
59
60     @abstractmethod
61     def plotData(self):
62         pass
```

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import math
4 from inference import ContinuousDistribution
5
6 class GaussDistribution(ContinuousDistribution):
7
8     def __init__(self, dimension, fileName = None, numberOfSamplesToGenerate =
9         None, mean = None, variance = None):
10         if((fileName is not None) & (numberOfSamplesToGenerate is not None)):
11             raise Exception("Can't load data and generate samples")
12
13         ContinuousDistribution.__init__(self)
14         self.dimension = dimension
15
16         if(fileName is not None):
17             self.importCsv(fileName)
18             self.numberOfSamples = len(self.dataSet)
19             self.calculateMean()
20             self.calculateVariance()
21
22         if(numberOfSamplesToGenerate is not None) & \
23             (mean is not None) & \
24             (variance is not None):
25             self.numberOfSamples = numberOfSamplesToGenerate
26             self.mean = np.array(mean)
27             self.variance = np.array(variance)
28             self.generateSamples()
29
30         if(len(self.dataSet) == 0):
31             raise Exception("Could not generate data, verify parameters")
32
33         self.calculateStandardDeviation()
34         self.gaussianDistribution = []
35         self.generateGaussian()
36
37     def generateSamples(self):
38         if len(self.dataSet) != 0:
39             raise Exception("Data already added")
40
41         self.dataSet = np.random.default_rng().normal(self.mean, self.variance,
42             size=(self.numberOfSamples, self.dimension))
43
44     def generateGaussian(self):
45         if len(self.dataSet) == 0:
46             raise Exception("No Data added")
47
48         if self.dimension == 1:
49             return self.generateGaussian1D()
50
51         return self.generateGaussian2D()
52
53     def calculateMean(self):
54         self.mean = np.mean(self.dataSet)
55
56     def calculateStandardDeviation(self):
57         if self.dimension == 1:
58             self.standardDeviation = np.std(self.dataSet)
59         else:

```

```

58         self.standardDeviation = np.array((math.sqrt(self.mean[0]), math.
59             sqrt(self.mean[1])))
60
61     def calculateGausse1D(self, x):
62         exponentialTerm = (-(1 / (2 * self.variance ** 2)) * (x - self.mean) **
63             2)
64         denominator = (2 * math.pi * self.variance ** 2) ** (0.5)
65         return (1 / denominator) * math.e ** (exponentialTerm)
66
67     def generateGausse1D(self):
68         vectorArray = np.array(self.dataSet)
69         self.gausse = []
70
71         for x in vectorArray:
72             self.gausse.append(self.calculateGausse1D(x))
73
74     def calculateGausse2D(self, vector):
75         xs = [self.mean[0], 0]
76         ys = [0, self.mean[1]]
77         covariance = [xs, ys]
78         inverseCovariance = np.linalg.inv(covariance)
79         determinantCovariance = np.linalg.det(covariance)
80
81         # exponentialTerm = (-0.5 * np.transpose(vector - self.mean)) *
82             inverseCovariance * (vector - self.mean)
83         exponentialTerm = -(np.linalg.solve(covariance, (vector - self.mean)).T.
84             dot((vector - self.mean))) / 2
85         denominator = ((2 * math.pi) ** (self.dimension / 2)) *
86             determinantCovariance ** (0.5)
87         result = (1 / denominator) * math.exp(exponentialTerm)
88         return result
89
90     def generateGausse2D(self):
91         vectorArray = np.array(self.dataSet)
92
93         self.gausseDistribution = []
94
95         for vector in vectorArray:
96             self.gausseDistribution.append(self.calculateGausse2D(vector))
97
98     def plotData(self):
99         if len(self.dataSet) == 0:
100             raise Exception("No Data added")
101
102         if self.dimension == 1:
103             self.plotData1D()
104         else:
105             self.plotData2D()
106
107     def plotData1D(self):
108         plotRange = range(len(self.dataSet))
109         x = np.linspace(min(self.dataSet), max(self.dataSet), self.
110             numberOfSamples)
111         y = self.calculateGausse1D(x)
112
113         plt.figure(figsize=(6, 6))
114
115         plt.subplot(2, 1, 1)
116         plt.hist(self.dataSet, bins=60, density=True, label="Histogram")

```

```

111     plt.plot(x, y, "r-", linewidth=1, label="Distribution")
112
113     plt.title("Distribution")
114     plt.xlabel("Value")
115     plt.ylabel("Frequency")
116     plt.legend(loc="upper right")
117
118     plt.subplot(2, 1, 2)
119     plt.scatter(plotRange, self.dataSet, label="Data", s=2)
120
121     plt.title(f"Raw data with n = {self.numberOfSamples} sample points")
122     plt.xlabel("Sample")
123     plt.ylabel("Value")
124     plt.legend(loc="best")
125
126     plt.suptitle(f"Gaus Distribution with  $\mu$  = {self.mean} and  $\sigma$  = {self.variance}")
127
128     plt.tight_layout()
129     plt.show()
130
131 def plotData2D(self):
132     plt.figure(figsize=(8, 12))
133
134     hist, xedges, yedges = np.histogram2d(self.dataSet[:,0], self.dataSet[:,1], bins=60)
135
136     # Construct arrays for the anchor positions of the 16 bars.
137     xpos, ypos = np.meshgrid(xedges[:-1] + 0.25, yedges[:-1] + 0.25,
138                               indexing="ij")
139     xpos = xpos.ravel()
140     ypos = ypos.ravel()
141     zpos = 0
142
143     # Construct arrays with the dimensions for the 16 bars.
144     dx = dy = 0.5 * np.ones_like(zpos)
145     dz = hist.ravel()
146
147     ax = plt.subplot(2, 1, 1, projection="3d")
148     # ax.bar3d(xpos, ypos, zpos, dx, dy, dz)
149     ax.set_zlabel("Frequency")
150
151     xy = np.linspace([min(self.dataSet[0]), min(self.dataSet[1])], [max(self.dataSet[0]), max(self.dataSet[1])], self.numberOfSamples)
152     # y = np.linspace(min(self.dataSet[1]), max(self.dataSet[1]), self.numberOfSamples)
153     z = np.array([self.calculateGaussien2D(v) for v in xy])
154     ax.plot_surface(xy[0], xy[1], z, )
155
156     plt.title("Distribution")
157     plt.xlabel("X")
158     plt.ylabel("Y")
159     plt.legend(loc="best")
160
161     ax = plt.subplot(2, 1, 2, projection="3d")
162     plotRange = range(self.numberOfSamples)
163     ax.scatter3D(plotRange, self.dataSet[:, 0], self.dataSet[:, 1], label="Data", s=2)

```

```
164 plt.title(f"Raw data with n = {self.numberOfSamples} sample points")
165 plt.xlabel("Sample")
166 plt.ylabel("Value")
167 plt.legend(loc="best")
168
169 plt.suptitle(f"Gaus Distribution with  $\mu$  = {self.mean} and  $\sigma$  =
    {self.variance}")
170
171 plt.tight_layout()
172 plt.show()
```



```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import math
4 from inference import ContinuousDistribution
5
6 class BetaDistribution(ContinuousDistribution):
7
8     def __init__(self, a, b, fileName = None, numberOfSamplesToGenerate = None):
9         if((fileName is not None) & (numberOfSamplesToGenerate is not None)):
10             raise Exception("Can't load data and generate samples")
11
12         if((fileName is None) & (numberOfSamplesToGenerate is None)):
13             raise Exception("No parameters for data")
14
15         ContinuousDistribution.__init__(self)
16         self.a = a
17         self.b = b
18         self.calculateAndSetBFromAAndB()
19         self.generatedBetaDistribution = []
20
21         if(numberOfSamplesToGenerate is not None):
22             self.numberOfSamples = numberOfSamplesToGenerate
23             self.generateSamples()
24
25         if (fileName is not None):
26             self.importCsv(fileName)
27             self.numberOfSamples = len(self.dataSet)
28
29         self.calculateMean()
30         self.calculateVariance()
31         self.calculateStandardDeviation()
32         self.generateBetaDistribution()
33
34     def generateSamples(self):
35         if len(self.dataSet) != 0:
36             raise Exception("Data already added")
37
38         self.dataSet = np.random.default_rng().beta(self.a, self.b, size=self.
39             numberOfSamples)
40
41     def calculateAndSetBFromAAndB(self):
42         self.bFromAAndB = (math.gamma(self.a + self.b) / (math.gamma(self.a) +
43             math.gamma(self.b)))
44
45     def calculateBeta(self, x):
46         return self.bFromAAndB * pow(x, (self.a - 1)) * pow((1 - x), (self.b -
47             1))
48
49     def generateBetaDistribution(self):
50         self.generatedBetaDistribution = []
51
52         for x in self.dataSet:
53             result = self.calculateBeta(x)
54             self.generatedBetaDistribution.append(result)
55
56     def plotData(self):
57         plotRange = range(len(self.dataSet))
58         x = np.linspace(0.01, 0.99, self.numberOfSamples)
59         y = self.calculateBeta(x)

```

```

57
58     plt.figure(figsize=(6, 6))
59
60     plt.subplot(2, 1, 1)
61     plt.hist(self.dataSet, bins=60, density=True, label="Histogram")
62     plt.plot(x, y, "r-", linewidth=1, label=f"alpha = {self.a}, beta = {self.b}")
63
64     plt.title("Distribution")
65     plt.xlabel("Value")
66     plt.ylabel("Frequency")
67     plt.legend(loc="upper right")
68
69     plt.subplot(2, 1, 2)
70     plt.scatter(plotRange, self.dataSet, label="Data", s=2)
71
72     plt.title("Raw Data")
73     plt.xlabel("Sample")
74     plt.ylabel("Value")
75     plt.legend(loc="upper right")
76
77     plt.suptitle(f"Beta Distribution with alpha = {self.a} and beta = {self.b}")
78
79     plt.tight_layout()
80     plt.show()
81
82     def plotDataWithDifferentAlphasAndBetas(self, alphaAndBetas):
83         if len(alphaAndBetas) == 0:
84             return
85
86         plt.figure(figsize=(4, 4))
87         plt.title("Distribution")
88         plt.xlabel("Value")
89         plt.ylabel("Frequency")
90
91         for alphaAndBeta in alphaAndBetas:
92             self.a = alphaAndBeta[0]
93             self.b = alphaAndBeta[1]
94             self.calculateAndSetBFromAAndB()
95
96             x = np.linspace(0.01, 0.99, self.numberOfSamples)
97             y = self.calculateBeta(x)
98             plt.plot(x, y, linewidth=1, label=f"alpha = {self.a}, beta = {self.b}")
99
100         plt.legend(loc = "best")
101
102         plt.tight_layout()
103         plt.show()

```

Bibliography

Roelants, Peter (2018). *Multivariate normal distribution*. URL: <https://peterroelants.github.io/posts/multivariate-normal-primer/>.

Rueckert, Elmar (2022). *An Introduction to Probabilistic Machine Learning*. Elmar Rueckert.