

---

# Assignment III: Gauss and Beta Distribution

Exercises in Machine Learning (190.013), SS2022  
Stefan Nehl<sup>1</sup>

<sup>1</sup>stefan-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 self.  
2     calculateStandardDeviation  
3     ()  
4 self.gaussianDistribution = []  
5 self.generateGaussian()  
6  
7 def generateSamples(self):  
    if len(self.dataSet) != 0:
```

### 4.2 Calculation

For the calculation I implemented two different methods. One for the one dimensional calculation, *generate-*

*Gaussen1D*, and for the two dimensional calculation, *calculateGaussen2D*. 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  $\mu$  stands for the mean and  $\sigma$  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,  $\Sigma$  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. Additional, 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

The class *BetaDistribution* has analogue to the class *GaussDistribution* also a constructor for handling the initialization for the parameters. Also the limitation for file name for CSV or generation of samples is the same. The difference is, that the beta distribution needs the parameter  $a$  and  $b$  and not the dimension for the distribution.

### 5.1 Generating Samples

The generation for the samples was created again with the *random()* function from the *numpy* library. However, I changed the distribution to the beta distribution.

```

1         self.calculateMean()
2         self.calculateVariance()
3         self.
           calculateStandardDeviation
           ()
4         self.generateBetaDistribution
           ()
5
6         def generateSampels(self):
7             if len(self.dataSet) != 0:
```

### 5.2 Calculation

For the calculation I used the following formula.

$$Beta(x|a, b) = B(a, b)x^{a-1}(1-x)^{b-1}$$

where  $a, b$  are in the constructor given parameters as a scalar and  $B(a, b)$  the Beta function.

$$B(a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$$

where  $\Gamma$  is the gamma function. The formulas are again from book *An Introduction to Probabilistic Machine Learning*. (Rueckert, 2022)

### 5.3 Plotting

For plotting the results I used a histogram with the distribution and the raw data. In addition I created a plot with the beta distribution with different values for the parameters *alpha* and *beta*. The values are displayed in Table 1. For the plotting of the parameters with

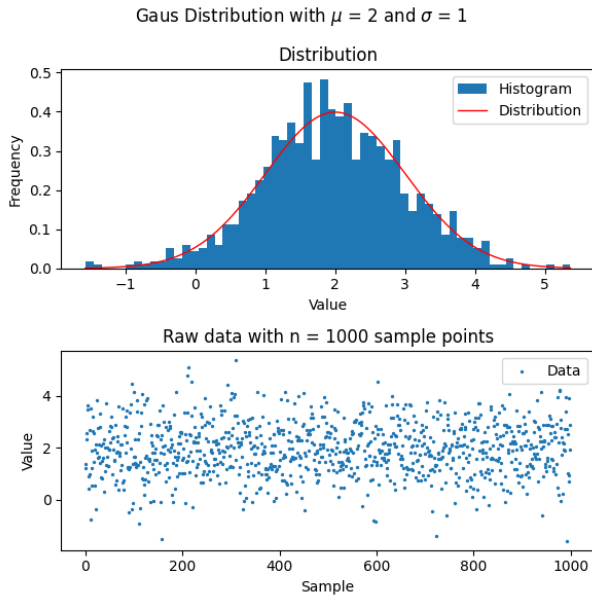
**Table 1:** Values for alpha and beta

alpha	beta
0.5	0.5
5	2
2	5

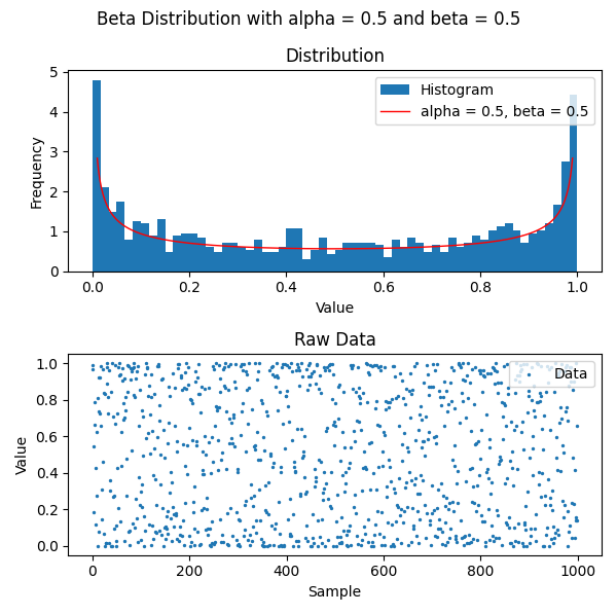
different values I added a method named *plotDataWithDifferentAlphasAndBetas* with a list of the different settings as a parameter. These function sets the values and generates the plots.

## 6 Results

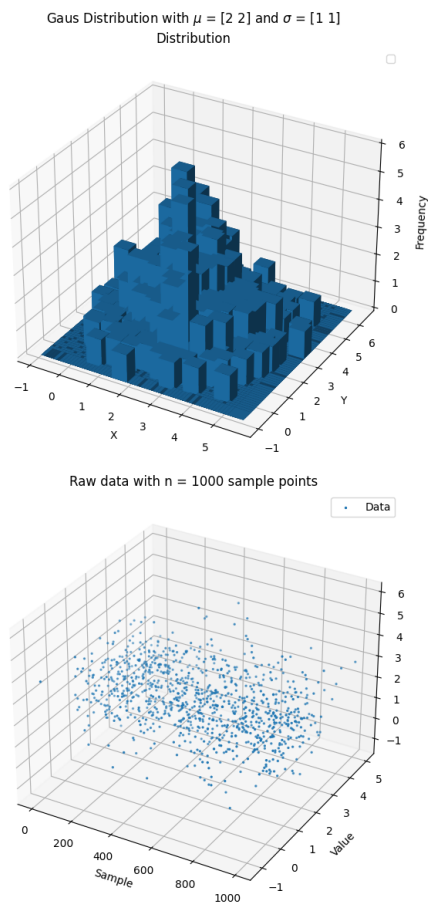
Figure 1 shows the gauss distribution for one dimension. It displays the values distribution and the frequencies of those values. The orange line displays the gauss distribution itself. Figure 2 shows the gauss distribution for two dimensions in a three dimensional plot. The beta distribution is displayed in figure 3. This plot uses an alpha with 0.5 and a beta with 0.5. The result displays a higher frequency of the values around 0 and 1. This is also displayed in the scatter chart with the raw values of this distribution. Figure 4 displays the beta distribution with alpha 2 and beta 5. The frequency is now higher at the start of the plot with a maximum at around 0.2. The beta distribution in figure 5 displays the frequency with the value for alpha 5 and beta 2. Therefore, the frequency is higher at the end of the plot near 0.8. The last figure, Figure 5, displays the three settings for the beta distribution.



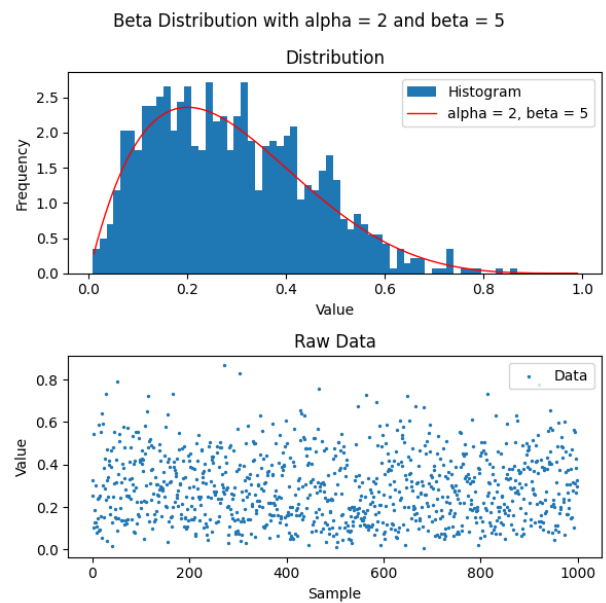
**Figure 1:** Gauss Distribution for one dimension with the raw data



**Figure 3:** Beta Distribution with  $\alpha = 0.5$  and  $\beta = 0.5$



**Figure 2:** Gauss Distribution for two dimensions with the raw data

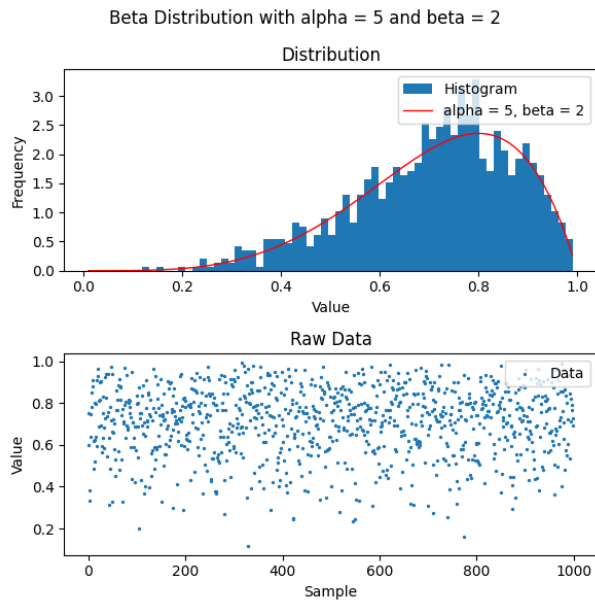


**Figure 4:** Beta Distribution with  $\alpha = 2$  and  $\beta = 5$

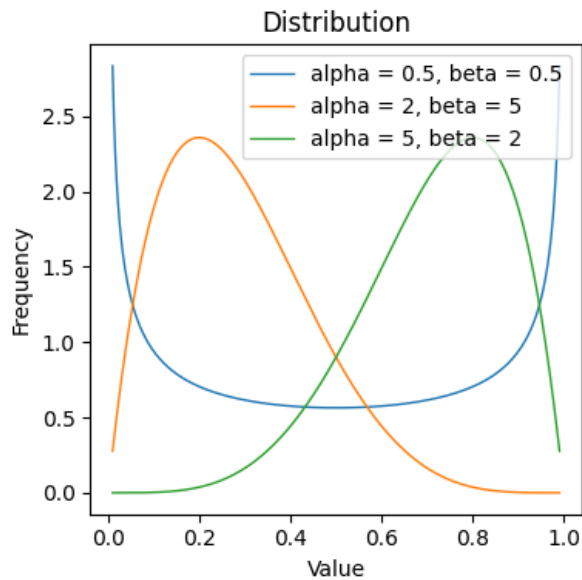
## 7 Conclusion

The implementation of the abstract class was straightforward. For the other classes I made some changes. I moved the classes to separate files to handle them better. Unfortunately, I was not able to create a satisfying plot for the gauss distribution for two dimensions. The whole code is in the appendix of this paper.

## APPENDIX



**Figure 5:** Beta Distribution with  $\alpha = 5$  and  $\beta = 2$



**Figure 6:** Beta Distribution with the different values for  $\alpha$  and  $\beta$

```
1  """
2  @author: Nehl Stefan
3  """
4
5  import sys
6  sys.path.insert(0, '../modules')
7  import inference
8  from BetaDistribution import BetaDistribution
9  from GaussDistribution import GaussDistribution
10
11  gaussDistr = GaussDistribution(1, numberOfSamplesToGenerate = 1000, mean = 2,
12                                variance = 1)
13  gaussDistr.plotData()
14
15  gaussDistr = GaussDistribution(2, numberOfSamplesToGenerate = 1000, mean = (2,
16                                2), variance = (1, 1))
17  gaussDistr.plotData()
18
19  betaDistr = BetaDistribution(0.5,0.5, numberOfSamplesToGenerate = 1000)
20  betaDistr.plotData()
21
22  betaDistr = BetaDistribution(2,5, numberOfSamplesToGenerate = 1000)
23  betaDistr.plotData()
24
25  betaDistr = BetaDistribution(5,2, numberOfSamplesToGenerate = 1000)
26  betaDistr.plotData()
27
28  settings = [(0.5, 0.5), (2, 5), (5, 2)]
29  betaDistr.plotDataWithDifferentAlphasAndBetas(settings)
```

```

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

```



```

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

```

```

111 plt.figure(figsize=(6, 6))
112
113 plt.subplot(2, 1, 1)
114 plt.hist(self.dataSet, bins=60, density=True, label="Histogram")
115 plt.plot(x, y, "r-", linewidth=1, label="Distribution")
116
117 plt.title("Distribution")
118 plt.xlabel("Value")
119 plt.ylabel("Frequency")
120 plt.legend(loc="upper right")
121
122 plt.subplot(2, 1, 2)
123 plt.scatter(plotRange, self.dataSet, label="Data", s=2)
124
125 plt.title(f"Raw data with n = {self.numberOfSamples} sample points")
126 plt.xlabel("Sample")
127 plt.ylabel("Value")
128 plt.legend(loc="best")
129
130 plt.suptitle(f"Gaus Distribution with  $\mu$  = {self.mean} and  $\sigma$  =
    {self.variance}")
131
132 plt.tight_layout()
133 plt.show()
134
135 def plotData2D(self):
136     plt.figure(figsize=(8, 12))
137
138     hist, xedges, yedges = np.histogram2d(self.dataSet[:,0], self.dataSet
139        [:,1], bins=60)
140
141     # Construct arrays for the anchor positions of the 16 bars.
142     xpos, ypos = np.meshgrid(xedges[:-1] + 0.25, yedges[:-1] + 0.25,
143         indexing="ij")
144     xpos = xpos.ravel()
145     ypos = ypos.ravel()
146     zpos = 0
147
148     # Construct arrays with the dimensions for the 16 bars.
149     dx = dy = 0.5 * np.ones_like(zpos)
150     dz = hist.ravel()
151
152     ax = plt.subplot(2, 1, 1, projection="3d")
153     ax.bar3d(xpos, ypos, zpos, dx, dy, dz)
154     ax.set_zlabel("Frequency")
155
156     xy = np.linspace([min(self.dataSet[0]), min(self.dataSet[1])], [max(self.
157         dataSet[0]), max(self.dataSet[1])], self.numberOfSamples)
158     # y = np.linspace(min(self.dataSet[1]), max(self.dataSet[1]), self.
159         numberOfSamples)
160     z = np.array([self.calculateGaussen2D(v) for v in xy])
161     # does not work
162     # ax.plot_surface(xy[0], xy[1], z, )
163
164     plt.title("Distribution")
165     plt.xlabel("X")
166     plt.ylabel("Y")
167     plt.legend(loc="best")

```

```
165 ax = plt.subplot(2, 1, 2, projection="3d")
166 plotRange = range(self.numberOfSamples)
167 ax.scatter3D(plotRange, self.dataSet[:, 0], self.dataSet[:, 1], label="
    Data", s=2)
168
169 plt.title(f"Raw data with n = {self.numberOfSamples} sample points")
170 plt.xlabel("Sample")
171 plt.ylabel("Value")
172 plt.legend(loc="best")
173
174 plt.suptitle(f"Gaus Distribution with  $\mu$  = {self.mean} and  $\sigma$  =
    {self.variance}")
175
176 plt.tight_layout()
177 plt.show()
```

```

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

```

```

57 def plotData(self):
58     plotRange = range(len(self.dataSet))
59     x = np.linspace(0.01, 0.99, self.numberOfSamples)
60     y = self.calculateBeta(x)
61
62     plt.figure(figsize=(6, 6))
63
64     plt.subplot(2, 1, 1)
65     plt.hist(self.dataSet, bins=60, density=True, label="Histogram")
66     plt.plot(x, y, "r-", linewidth=1, label=f"alpha = {self.a}, beta = {self.b}")
67
68     plt.title("Distribution")
69     plt.xlabel("Value")
70     plt.ylabel("Frequency")
71     plt.legend(loc="upper right")
72
73     plt.subplot(2, 1, 2)
74     plt.scatter(plotRange, self.dataSet, label="Data", s=2)
75
76     plt.title("Raw Data")
77     plt.xlabel("Sample")
78     plt.ylabel("Value")
79     plt.legend(loc="upper right")
80
81     plt.suptitle(f"Beta Distribution with alpha = {self.a} and beta = {self.b}")
82
83     plt.tight_layout()
84     plt.show()
85
86 def plotDataWithDifferentAlphasAndBetas(self, alphaAndBetas):
87     if len(alphaAndBetas) == 0:
88         return
89
90     plt.figure(figsize=(4, 4))
91     plt.title("Distribution")
92     plt.xlabel("Value")
93     plt.ylabel("Frequency")
94
95     for alphaAndBeta in alphaAndBetas:
96         self.a = alphaAndBeta[0]
97         self.b = alphaAndBeta[1]
98         self.calculateAndSetBFromAAndB()
99
100     x = np.linspace(0.01, 0.99, self.numberOfSamples)
101     y = self.calculateBeta(x)
102     plt.plot(x, y, linewidth=1, label=f"alpha = {self.a}, beta = {self.b}")
103
104     plt.legend(loc = "best")
105
106     plt.tight_layout()
107     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.