# Assignment III: Probability Theory

## Introduction to Machine Learning Lab (190.013), SS2023
## Björn Ellensohn[1]

[1]*m01435615, bjoern.ellensohn@stud.unileoben.ac.at*, *Montanuniversität Leoben, Austria*
April 27, 2023

This document guides through the process of solving Assignment 3.

## 1  Introduction

The Task of 3rd assignment was coming up with an abstract class called `ContinuousDistribution` that provides the outline for two subclasses `GaussDistribution` and `BetaDistribution`. Again, provided was a .csv file with a dataset. The assignment is divided into three main parts and the bonus part.

## 2  Part I - Abstract Class

### 2.1  Preparation

When searching online for information on Python's abstract classes, you may come across the module known as `ABC`. This module includes a specific decorator, `@abstractmethod`, that can be used to define abstract methods. Essentially, this allows you to create a class that serves as a template or blueprint for other classes, without actually implementing the methods. To create a functional child class, you must define and implement the necessary methods there.

### 2.2  Defining the Abstract Class

The class `ContinuousDistribution` is created that provides the basic outline for the main classes of this exercise. The following functions should be inlcuded:

- Data Import and Export using csv files.
- Computation of the mean based on the samples from the csv.
- Computation of the standard deviation based on the samples from the csv.
- Visualization of the distribution, the raw data or the generated samples.
- Generating/Drawing Samples from the distribution.

So in the end, the abstract class `ContinuousDistribution` should look like this:

```python
import abc

class
↪ ContinuousDistribution(metaclass=abc.ABCMeta):
@abc.abstractmethod
def import_data(self, file_path):
    pass

@abc.abstractmethod
def export_data(self, data, file_path):
    pass

@abc.abstractmethod
def compute_mean(self, data):
    pass

@abc.abstractmethod
def compute_standard_deviation(self, data):
    pass

@abc.abstractmethod
def visualize(self, data=None):
    pass

@abc.abstractmethod
def generate_samples(self, n_samples):
    pass
```

# 3 Part II - Plot and Sample Gaussian Distributions

In the second part of this exercise, it is required to initiate a child class which is responsible for dealing with Gaussian distributions.

Explicitly, the following should be implemented:

- Implement the functions defined in "ContinousDistribution".
- Implement a constructor that optionally takes the dimension of the multivariate distribution.
- Implement a visualization for Multivariate Gaussians up to 3 dimensions.
- Find the empirical parameters of the distribution that created the samples in the 'MGD.csv' file.
- Plot the samples of the 'MGD.csv' file and the sample from the learned distribution in two subfigures.

So the main goal of this class is to be able to compute the statistical parameters from a provided dataset and being able to sample a new dataset from this distribution. In the end, the datapoints should be plotted.

At first, let us declare how the Gaussian dsitribution is defined. For all normal distributions, 68.2% of the observations will appear within plus or minus one standard deviation of the mean; 95.4% of the observations will fall within +/- two standard deviations; and 99.7% within +/- three standard deviations. This fact is sometimes referred to as the "empirical rule," a heuristic that describes where most of the data in a normal distribution will appear. This means that data falling outside of three standard deviations ("3-sigma") would signify rare occurrences.

In mathematical terms this is explained in Equation 1.

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)} \qquad (1)$$

For sampling from the Gaussian Distribution, a Box-Muller Transform could be used (Keng, 2015). But for simplicity, I will use a module called `multivariate_normal` from the `scipy` package.

The constructor of `GaussianDistribution` then looks like this:

```
class GaussDistribution(ContinuousDistribution):
    def __init__(self, dim=1):
```

```
        self.dim = dim
        self.mean = np.zeros(dim)
        self.covariance = np.eye(dim)
        self.data = pd.DataFrame()
        self.samples = None
```

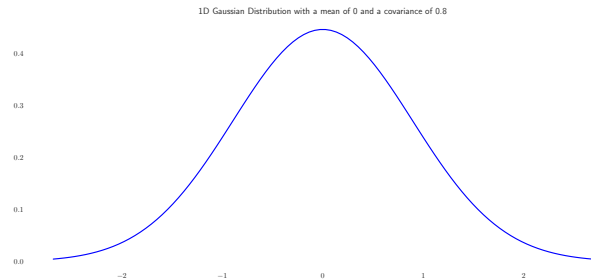Figure 1 to Figure 3 are showing the end results. In the appendix you will find the full code.



**Figure 1:** *Plot of one dimensional Gaussian distribution.*
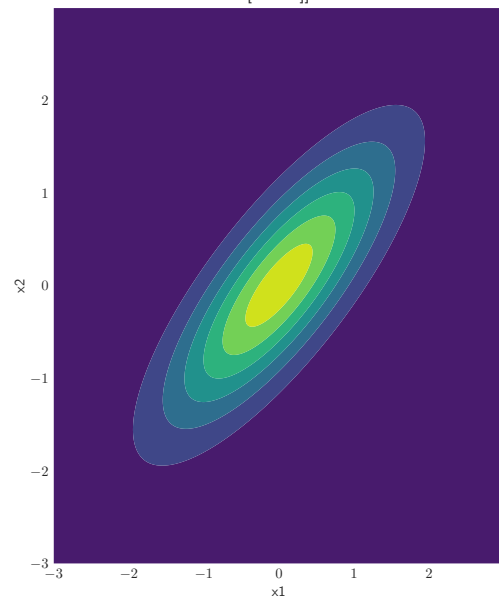


**Figure 2:** *Contour plot of two dimensional Gaussian distribution.*

# 4 Part III - Plot and Sample Beta Distributions

Next, a class `BetaDistribution` shouldbe inherited from the meta class. It should implement:

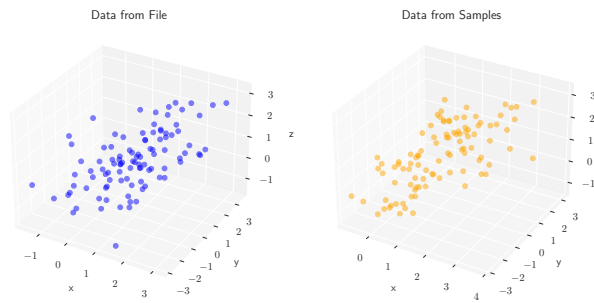- Generate beta distributed samples and plot the distribution giving the parameters a and b.

**Figure 3:** *Plots of three dimensional Gaussian distribution.*

- The constructor should take the parameters a and b as arguments.
- A visualization for Beta distributions, including the mean and the standard deviation lines.

Again, there is a module in the `scipy` package that does the hard work here `beta`. The constructor of `BetaDistribution` then looks like this:

```
class
↪   BetaDistribution(ContinuousDistribution):
def __init__(self, a, b):
    self.a = a
    self.b = b
    self.data = None
```

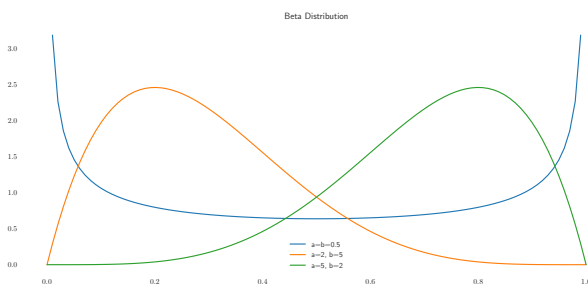Provided is a plot showing the end results in Figure 4.



**Figure 4:** *Plot of three different Beta Distributions.*

# 5  Part IV - Bonus Exercise

We were provided with a fictional card game with the aim to better understand the Bayes' Rule. A deck of four known cards is shuffled, a card is chosen until the deck is empty. To solve this, missing functions need to be filled. Those functions are:

- `update_probabilities`
- `get_posterior`
- `get_likelihood`
- `get_marginal_likelihood`

I tried to solve the exercise but in the end I ran out of time. Still, I think I was close. The code will be attached in the appendix. Here I will try to explain my functions. The function `update_probabilities` is just there to call the `get_posterior` function. The function `get_posterior` returns the posterior, based on the likelyhood, the prior and the marginal likelyhood. Likewise, `get_likelihood` returns the likelyhood of drawing a certain card (red or black) given the color of the drawn card. Last but not least `get_marginal_likelihood` returns the marginal likelyhood of drawing a certain card (red or black) given the color of the drawn card.

# APPENDIX

This section houses the code.

## Code for Part I

```
1   import abc
2
3   class ContinuousDistribution(metaclass=abc.ABCMeta):
4       @abc.abstractmethod
5       def import_data(self, file_path):
6           pass
7
8       @abc.abstractmethod
9       def export_data(self, data, file_path):
10          pass
11
12      @abc.abstractmethod
13      def compute_mean(self, data):
14          pass
15
16      @abc.abstractmethod
17      def compute_standard_deviation(self, data):
18          pass
19
20      @abc.abstractmethod
21      def visualize(self, data=None):
22          pass
23
24      @abc.abstractmethod
25      def generate_samples(self, n_samples):
26          pass
```

## Code for Part II

```
1   import numpy as np
2   import matplotlib
3   import matplotlib.pyplot as plt
4   from mpl_toolkits.mplot3d import Axes3D
5   from scipy.stats import multivariate_normal
6   import pandas as pd
7
8
9   class GaussDistribution(ContinuousDistribution):
10      def __init__(self, dim=1):
11          self.dim = dim
12          self.mean = np.zeros(dim)
13          self.covariance = np.eye(dim)
14          self.data = pd.DataFrame()
15          self.samples = None
16
17      def import_data(self, file_path):
18          # implementation to import data from file
19          self.data = pd.read_csv(file_path)
20
21      def export_data(self, data, file_path):
22          # implementation to export data to file
23          df = pd.DataFrame(data)
24          df.to_csv(file_path)
25
26      def compute_mean(self, data):
```

```
27              self.mean = np.mean(data, axis=0)
28
29          def compute_standard_deviation(self, data):
30              self.covariance = np.cov(data, rowvar=False)
31
32          def visualize(self, data=None):
33              if data is None:
34                  data = multivariate_normal.rvs(mean=self.mean, cov=self.covariance, size=1000)
35
36              if self.dim == 1:
37                  mean = 0
38                  covariance = 0.8
39                  x = np.linspace(mean - 3*np.sqrt(covariance), mean + 3*np.sqrt(covariance), 100)
40                  plt.plot(x, multivariate_normal.pdf(x, mean=mean, cov=covariance), color = 'blue')
41                  plt.title(f'1D Gaussian Distribution with a mean of {mean} and a covariance of
                  ↪  {covariance}')
42
43                  plt.savefig('gaussian1D.pdf', bbox_inches='tight', transparent=True)
44                  plt.show()
45
46              elif self.dim == 2:
47                  covariance = np.array([[1, 0.8],
48                                          [0.8, 1]])
49                  mean = np.array([0, 0])
50                  x, y =
                  ↪  np.mgrid[mean[0]-3*np.sqrt(covariance[0,0]):mean[0]+3*np.sqrt(covariance[0,0]):.01,
51                              mean[1]-3*np.sqrt(covariance[1,1]):mean[1]+3*np.sqrt(covariance[1,1]):⌋
                              ↪  .01]
52                  pos = np.empty(x.shape + (2,))
53                  pos[:, :, 0] = x
54                  pos[:, :, 1] = y
55                  rv = multivariate_normal(mean, covariance)
56
57                  # Generating the density function
58                  # for each point in the meshgrid
59                  pdf = np.zeros(x.shape)
60                  for i in range(x.shape[0]):
61                      for j in range(x.shape[1]):
62                          pdf[i,j] = rv.pdf([x[i,j], y[i,j]])
63
64                  pdf_list = []
65                  fig = plt.figure()
66
67                  # Plotting the density function values
68                  bx = fig.add_subplot(131, projection = '3d')
69                  bx.plot_surface(x, y, pdf, cmap = 'viridis')
70                  plt.xlabel("x1")
71                  plt.ylabel("x2")
72                  plt.title(f'2D Gaussian Distribution with a mean of {mean} and a covariance of
                  ↪  {covariance}')
73                  pdf_list.append(pdf)
74                  bx.axes.zaxis.set_ticks([])
75
76                  plt.tight_layout()
77
78                  plt.savefig('gaussian2D_surface.pdf', bbox_inches='tight', transparent=True)
79                  plt.show()
80
81                  # Plotting contour plots
82                  for idx, val in enumerate(pdf_list):
83                      plt.subplot(1,3,idx+1)
84                      plt.contourf(x, y, val, cmap='viridis')
85                      plt.xlabel("x1")
86                      plt.ylabel("x2")
```

```
87              plt.tight_layout()
88              plt.title(f'2D Gaussian Distribution Contour with a mean of {mean} and a covariance
   ↪   of {covariance}')

90              plt.savefig('gaussian2D_contour.pdf', bbox_inches='tight', transparent=True)
91              plt.show()

93          elif self.dim == 3:
94              fig = plt.figure(figsize=(10, 5))
95              ax1 = fig.add_subplot(1, 2, 1, projection='3d')
96              x, y, z = np.mgrid[self.mean[0]-3*np.sqrt(self.covariance[0,0]):self.mean[0]+3*np.sqrt
   ↪   (self.covariance[0,0]):.1,
97                          self.mean[1]-3*np.sqrt(self.covariance[1,1]):self.mean[1]+3*np.sqrt
   ↪   (self.covariance[1,1]):.1,
98                          self.mean[2]-3*np.sqrt(self.covariance[2,2]):self.mean[2]+3*np.sqrt
   ↪   (self.covariance[2,2]):.1]
99              #Plot the samples from the file

101             ax1.scatter(data.iloc[:, 0], data.iloc[:, 1], data.iloc[:, 2], c='blue', alpha=0.5)
102             ax1.set_xlabel('x')
103             ax1.set_ylabel('y')
104             ax1.set_zlabel('z')
105             ax1.set_title('Data from File')

107             ax2 = fig.add_subplot(1, 2, 2, projection='3d')
108             ax2.scatter(self.samples[:, 0], self.samples[:, 1], self.samples[:, 2], c='orange',
   ↪   alpha=0.5)
109             ax2.set_xlabel('x')
110             ax2.set_ylabel('y')
111             ax2.set_zlabel('z')
112             ax2.set_title('Data from Samples')

114             plt.savefig('gaussian3D.pdf', bbox_inches='tight', transparent=True)
115             plt.show()

117     def generate_samples(self, n_samples):
118         self.samples = multivariate_normal.rvs(mean=self.mean, cov=self.covariance, size=n_samples)
```

## Code for Part III

```
1  from scipy.stats import beta

3  class BetaDistribution(ContinuousDistribution):
4      def __init__(self, a, b):
5          self.a = a
6          self.b = b
7          self.data = None

9      def import_data(self, file_path):
10         self.data = pd.read_csv(file_path)

12     def export_data(self, data, file_path):
13         pass

15     def compute_mean(self, data):
16         pass

18     def compute_standard_deviation(self, data):
19         pass

21     def visualize(self, data=None):
22         # create a range of x values
```

```
23          x = np.linspace(0, 1, 100)
24
25          # calculate the beta PDF for the given parameters a and b
26          y = beta.pdf(x, self.a, self.b)
27
28          # plot the beta PDF
29          plt.plot(x, y, label='Beta PDF')
30
31          # plot the mean and standard deviation lines
32          mean = beta.mean(self.a, self.b)
33          std = beta.std(self.a, self.b)
34          plt.axvline(mean, color='red', label=f'Mean={mean:.2f}')
35          plt.axvline(mean - std, linestyle='--', color='green', label=f'Std Dev={std:.2f}')
36          plt.axvline(mean + std, linestyle='--', color='green')
37
38          # set the plot title and legend
39          plt.title(f'Beta Distribution (a={self.a}, b={self.b})')
40          plt.legend()
41
42          # show the plot
43          plt.savefig('beta.pdf', bbox_inches='tight', transparent=True)
44
45          plt.show()
46
47      def visualize_book(self, data=None):
48          # create a range of x values
49          x = np.linspace(0, 1, 100)
50
51          # calculate the beta PDF for the given parameters a and b
52          y1 = beta.pdf(x, 0.5, 0.5)
53          y2 = beta.pdf(x, 2, 5)
54          y3 = beta.pdf(x, 5, 2)
55
56          # plot the beta PDF
57          plt.plot(x, y1, label='a=b=0.5')
58          plt.plot(x, y2, label='a=2, b=5')
59          plt.plot(x, y3, label='a=5, b=2')
60
61          # set the plot title and legend
62          plt.title(f'Beta Distribution')
63          plt.legend()
64
65          # show the plot
66          plt.savefig('beta.pdf' , bbox_inches='tight', transparent=True)
67          plt.show()
68
69      def generate_samples(self, n_samples):
70          # generate beta distributed samples using the given parameters a and b
71          return beta.rvs(self.a, self.b, size=n_samples)
```

## Code for Part IV

```
1   import random
2
3   class CardDeck:
4       def __init__(self):
5           # Initialize the deck of cards
6           self.cards = ['R', 'R', 'B', 'B']
7           # Shuffle the deck of cards
8           random.shuffle(self.cards)
9           # Initialize the prior probability
10          self.prior_red = 0.5
```

```python
11              self.likelihood = 0.5


14      def draw_card(self, index):
15          """
16          Draw a card from the deck
17          :param index: the index of the card to be drawn
18          :return: the card and its color
19          """
20          card = self.cards.pop(index)
21          color = 'Red' if card == 'R' else 'Black'
22          return (card, color)

24      def get_card_color(self, index):
25          """
26          Get the color of a card
27          :return: the color of the card
28          """
29          card = self.cards[index]
30          color = 'Red' if card == 'R' else 'Black'
31          return color

33      def get_remainder(self):
34          """
35          Get the remainder deck
36          :return: the remainder deck
37          """
38          return self.cards

40      """
41      Update the prior probability based on Bayes' rule
42      function name: update_probabilities
43      """
44      def update_probabilities(self, card_color):
45          """
46          Update the prior probability based on Bayes' rule
47          :param card_color: the color of the card drawn
48          """
49          if card_color == 'Red':
50              self.prior_red = self.get_posterior('Red')
51          elif card_color == 'Black':
52              self.prior_red = self.get_posterior('Black')

54      """
55      Calculate the posterior probability based on the prior probability, likelihood, and marginal
    ↪   likelihood
56      function name: get_posterior
57      """
58      def get_posterior(self, card_color):
59          """
60          Calculate the posterior probability based on the prior probability, likelihood, and
    ↪   marginal likelihood
61          :param card_color: the color of the card drawn
62          :return: the posterior probability of drawing a red card or a black card
63          """
64          likelihood = self.get_likelihood(card_color)
65          marginal_likelihood = self.get_marginal_likelihood(card_color)
66          posterior = (self.prior_red * likelihood) / marginal_likelihood
67          return posterior

69      def get_likelihood(self, card_color):
70          """
71          Calculate the likelihood of drawing a certain card (red or black) given the color of the
    ↪   card already drawn
```

```python
72              :param card_color: the color of the card drawn
73              :return: the likelihood of drawing a certain card (red or black) given the color of the
   ↪   card already drawn
74              """
75          if card_color == 'Red':
76              num_red_cards = self.cards.count('R')
77              num_cards = len(self.cards)
78              return num_red_cards / num_cards
79          elif card_color == 'Black':
80              num_black_cards = self.cards.count('B')
81              num_cards = len(self.cards)
82              return num_black_cards / num_cards
83
84      def get_marginal_likelihood(self, card_color):
85          """
86          Calculate the marginal likelihood of drawing a certain card (red or black) regardless of
   ↪   the color of the card already drawn
87              :param card_color: the color of the card drawn
88              :return: the marginal likelihood of drawing a certain card (red or black) regardless of
   ↪   the color of the card already drawn
89          """
90          if card_color == 'Red':
91              num_red_cards = self.cards.count('R')
92              num_black_cards = self.cards.count('B')
93              return (self.prior_red * num_red_cards / len(self.cards)) + ((1 - self.prior_red) *
   ↪   num_black_cards / len(self.cards))
94          elif card_color == 'Black':
95              num_red_cards = self.cards.count('R')
96              num_black_cards = self.cards.count('B')
97              return ((1 - self.prior_red) * num_red_cards / len(self.cards)) + (self.prior_red *
   ↪   num_black_cards / len(self.cards))
98
99
100     def play_game(self):
101         while self.cards:
102             for i, card in enumerate(self.cards):
103                 print(f'{i}: Unknown')
104              # only one card left
105             if len(self.cards) == 1:
106                 # code missing here
107                 self.cards.pop()
108             else:
109                 index = int(input(f"Remainder deck: {self.get_remainder()}\nEnter the index of
   ↪   the card you want to draw:"))
110                 card, color = self.draw_card(index-1)
111                 print(f'You drew a {color} card with value {card}')
112
113                 # Calculate the likelihood and marginal likelihood based on the color of the drawn
   ↪   card
114                 self.get_likelihood(color)
115                 self.get_marginal_likelihood(color)
116
117             # Update the prior probability based on Bayes' rule
118             self.update_probabilities(color)
119             posterior_red = self.get_posterior('Red')   # you must calculate this value
120             posterior_black = self.get_posterior('Black') # you must calculate this value
121             # Print the posterior probabilities
122             print(f"After drawing a {card} card:")
123             print(f"The probability that the other card is red: {posterior_red:.2f}")
124             print(f"The probability that the other card is black: {posterior_black:.2f}")
125         print("The deck is empty.")
```

**Bibliography**

Keng, Brian. "Sampling from a Normal Distribution | Bounded Rationality." Sampling from a Normal Distribution, November 28, 2015. https://bjlkeng.github.io/posts/sampling-from-a-normal-distribution/.