

# Assignment 1 - Q4: Simulating Gaussian Distributions

Welcome to Pluto Notebooks! In this course you will be spending quite a bit of time working within these notebooks. These notebooks are an HTML/CSS/Javascript built interface for interacting with and working with Julia. They were inspired by Jupyter, and you can learn more about the notebooks through [their github](#). The notebooks are lightweight (with files that can be used by the normal julia interpreter), and track how notebook cells depend on each other to re-run cells when their dependencies change.

The only code in this notebook which is not found in julia's base is the plotting code provided by [StatsPlots](#) and [Plots](#). Another package we use is [PlutoUI](#) which contains utilities for building interfaces in pluto notebooks. To complete this assignment you will not need to import any other packages.

For all the cells, some are hidden by default. To see hidden cells click on the eye to the top-left of the cell of interest.

In this assignment you will:

1. Learn about using Pluto Notebooks and Stats Plots for visualizing and analysing data.
2. Get experience with calculating the mean, variance, and other metrics of sample data.
3. Build intuition about how different Gaussian parameters impact our estimators.

```
1 # Import the packages and export their exported functions to the main namespace.  
2 using StatsPlots, PlutoUI, Random
```

## !!!IMPORTANT!!!

Insert your details below. You should see a green checkmark.

```
student =
  (name = "Manav Patel", email = "mdpatel1@ualberta.ca", ccid = "mdpatel1", idnumber :
1 student = (name="Manav Patel", email="mdpatel1@ualberta.ca", ccid="mdpatel1",
  idnumber=1707001)
```

Welcome Manav Patel! 

## Gaussian Distribution

A Gaussian distribution has mean  $\mu$  and standard deviation  $\sigma$ . We will want to sample from Gaussian distribution. We provide an implementation below. We also discuss that implementation, to help you better understand Julia syntax that will be useful for your own implementation.

### GaussianDistribution

GaussianDistribution( $\mu::\text{Float64}$ ,  $\sigma::\text{Float64}$ )

A Gaussian distribution with mean  $\mu$ , standard deviation  $\sigma$ . You can sample data from this distribution using `sample(gd, n)` to get  $n$  samples. You can get the mean using `mean(gd)`, the standard deviation using `stddev(gd)`, and the variance using `var(gd)`.

```
1 # The block of text below add documentation to julia struct or function.
2 # Check out the live docs to the right when your cursor
3 # is in GaussianDistribution.
4 """
5     GaussianDistribution( $\mu::\text{Float64}$ ,  $\sigma::\text{Float64}$ )
6
7     A Gaussian distribution with mean  $\mu$ , standard deviation  $\sigma$ . You can sample
8     data from this distribution using `sample(gd, n)` to get `n` samples. You
9     can get the mean using `mean(gd)`, the standard deviation using
10    `stddev(gd)`, and the variance using `var(gd)`.
11    """
12 struct GaussianDistribution
13      $\mu::\text{Float64}$  # mean
14      $\sigma::\text{Float64}$  # standard deviation
15 end
```

mean (generic function with 1 method)

```
1 mean(gd::GaussianDistribution) = gd. $\mu$ 
```

stddev (generic function with 1 method)

```
1 stddev(gd::GaussianDistribution) = gd.σ
```

var (generic function with 1 method)

```
1 var(gd::GaussianDistribution) = gd.σ^2
```

sample (generic function with 2 methods)

```
1 function sample(gd::GaussianDistribution, n = 1)
2     gd.σ*randn(n) .+ gd.μ
3 end
```

## Understanding the Julia code for the Gaussian distribution

Note that in Julia we use `+` for scalar addition and `.+` to add two vectors. If we have two vectors  $u$  and  $v$ , both  $d > 1$  dimensional, then we would write `u .+ v` to add these elementwise. If we have a scalar  $s$ , then `u .+ s` adds  $s$  to every element of  $u$ .

Let us look more carefully at the `sample` function. First note that to generate a Gaussian sample with mean  $\mu$  and variance  $\sigma^2$ , we 1) call `randn(1)` to generate a sample from a zero-mean, unit variance Gaussian (a normal distribution) and 2) scale it by  $\sigma$  and shift it by the mean  $\mu$ . We can either call this function `n` times, to get  $n$  samples. Or, we can leverage the fact that `randn(n)` returns  $n$  samples from a normal distribution. `randn(n)` is a vector of size  $n$  of samples from a normal distribution. Multiplying this vector by the scalar `gd.σ` rescales every element in the vector and then we `.+` the scalar `gd.μ` to shift every element in the vector.

Equivalently, we could have used a for loop and written

```
dataset = zeros(n)
for i in 1:n
    dataset[i] = gd.σ*randn(1)[1] + gd.μ
end
return dataset
```

where `randn(1)` returns a vector of dimension 1, so we have to further index this vector to return this single scalar. You may wonder why in our vector-based implementation, we did not explicitly have a `return`. In Julia the last value computed in the function is returned when there is no explicit `return`.

Note that Julia is 1-indexed, rather than 0-indexed. This means indexing an array or vector starts at 1, rather than 0. This contrasts Python and C, where indexing starts at 0, and matches Matlab, where indexing starts at 1. It's possible 1-indexing was chosen for Julia to help make it a suitable replacement for Matlab, which is (or was) a popular numerical computing language.

One other note. We have `n=1` as an argument to `sample`. This means that  $n$  defaults to 1 if it is not provided.

# Implementing basic statistics from data

Below you will be implementing the sample mean, variance, and standard deviation of a dataset  $D$ . This dataset is guaranteed to be a vector of floating point numbers, where the  $i$ th entry corresponds to  $X_i$ . You need to fill in the code between `#### BEGIN SOLUTION` and `#### END SOLUTION`.

The sample mean is

$$\text{sample-mean}(D) = \frac{1}{n} \sum_{i=1}^n X_i$$

To implement the sample variance, we want you to use the unbiased sample variance formula

$$\text{sample-variance}(D) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \text{sample-mean}(D))^2$$

Finally, the sample standard deviation is the square root of the sample variance.

A few useful functions for this section include the following, where  $v$  is a vector and  $s$  is a scalar. `sum(v)` takes the sum of the elements in  $v$  and `length(v)` returns the length of the vector  $v$ . `sqrt(s)` returns the square root of the scalar  $s$  and `s^2` squares the scalar  $s$ . You can also call `sqrt` and squaring on a vector, by calling `sqrt(v)` and `v.^2`. As mentioned above, most basic operations on scalars like `+`, `-`, `^2` can be turned into elementwise operations on a vector by adding a `.`, namely using `.+`, `.-`, `.^2`.

Finally, you might want to use a `for` loop, which was explained above when discussing the `sample` function. Note that you can get away with simply using the above vector operations, but it can be mentally simpler and just as correct to use a `for` loop, so it is up to you.

Great job!  

mean (generic function with 2 methods)

```
1 function mean(D::Vector{Float64})
2     #### BEGIN SOLUTION
3     return sum(D) / length(D)
4
5
6     #### END SOLUTION
7 end
```

var (generic function with 2 methods)

```
1 function var(D::Vector{Float64})
2     ##### BEGIN SOLUTION
3     n = length(D)
4     mean_value = mean(D)
5     sample_variance = sum((D .- mean_value).^2) / (n - 1)
6
7
8     ##### END SOLUTION
9 end
```

stddev (generic function with 2 methods)

```
1 function stddev(D::Vector{Float64})
2     ##### BEGIN SOLUTION
3     return sqrt(var(D))
4
5     ##### END SOLUTION
6 end
```

## Simulating sample variance

Use the below let blocks to complete question 4(bcde). You will be graded on your written work, not on the code in these cells. Here we have given you one example of how you might call the above functions, to avoid getting hung up on Julia syntax.

```
1 let # example
2     n = 3
3     gd = GaussianDistribution(6.2, 0.1)
4     dataset = sample(gd, n)
5     for i in 1:n
6         println(dataset[i])
7     end
8 end
```

```
6.234456591737808
6.4044063465144045
6.0844980562994655
```



```
1 let
2
3 # Define parameters
4  $\mu$  = 0.0
5  $\sigma^2$  = 1.0
6 num_samples = 10
7 num_estimates = 5
8
9 # Generate 5 sets of sample means
10 sample_means = [mean( $\mu$  .+ sqrt( $\sigma^2$ ) .* randn(num_samples)) for _ in
11 1:num_estimates]
12
13 # Calculate the sample variance of the sample means
14 sample_variance_of_means = var(sample_means)
15
16 # Print the results
17 println("Sample Means: ", sample_means)
18 println("Sample Variance of Means: ", sample_variance_of_means)
19
20
21
end
```

```
Sample Means: [0.16993710729604267, 0.28534399387815645, -0.1601554055
3575816, 0.14440337640375292, 0.22322431622615135]
Sample Variance of Means: 0.02969564723500241
```



```

1 let # 4c
2
3 # Function to generate samples
4 function generate_samples( $\mu$ ,  $\sigma^2$ , num_samples)
5     return  $\mu$  .+ sqrt( $\sigma^2$ ) .* randn(num_samples)
6 end
7
8 # Define parameters
9  $\mu$  = 0.0
10  $\sigma^2$  = 1.0
11 num_samples_10 = 10
12 num_samples_100 = 100
13 num_estimates = 5
14
15 # Initialize arrays to store sample variances
16 sample_variances_10 = Float64[]
17 sample_variances_100 = Float64[]
18
19 # Generate 5 sets of sample variances with 10 samples
20 for _ in 1:num_estimates
21     samples = generate_samples( $\mu$ ,  $\sigma^2$ , num_samples_10)
22     mean_value = mean(samples)
23     push!(sample_variances_10, var(samples) / (num_samples_10 - 1))
24 end
25
26 # Generate 5 sets of sample variances with 100 samples
27 for _ in 1:num_estimates
28     samples = generate_samples( $\mu$ ,  $\sigma^2$ , num_samples_100)
29     mean_value = mean(samples)
30     push!(sample_variances_100, var(samples) / (num_samples_100 - 1))
31 end
32
33 # Print the results
34 println("Sample Variances with 10 samples: ", sample_variances_10)
35 println("Average Sample Variance with 10 samples: ", mean(sample_variances_10))
36
37 println("Sample Variances with 100 samples: ", sample_variances_100)
38 println("Average Sample Variance with 100 samples: ",
39 mean(sample_variances_100))
40
end

```

```

Sample Variances with 10 samples: [0.052714405762473354, 0.07113856909
565892, 0.03760666574949808, 0.13208194398891482, 0.11750081645160589]
Average Sample Variance with 10 samples: 0.0822084802096302
Sample Variances with 100 samples: [0.011922614059955945, 0.01213790100341
0759, 0.008416827022228807, 0.009091817519650795, 0.008708323737103169]
Average Sample Variance with 100 samples: 0.010055496668469895

```



```

1 let # 4d, 4e
2 n=30
3 gd = GaussianDistribution(0,sqrt(10))
4 dataset = sample(gd,n)
5 sample_mean = mean(dataset)
6 println("Sample mean ",sample_mean)
7 lower_95 = sample_mean - 1.96(sqrt(10/30))
8 upper_95 = sample_mean + 1.96(sqrt(10/30))
9 println("95% confidence interval ::")
10 println("Lower Bound : ",lower_95)
11 println("Upper Bound : ",upper_95)
12 s=0.05
13 e = sqrt(10/(s*n))
14 lower_cheb = sample_mean - e
15 upper_cheb = sample_mean + e
16 println("Chebychev's Inequality ::")
17 println("Lower Bound : ",lower_cheb)
18 println("Upper Bound : ",upper_cheb)
19 end

```

```

Sample mean 0.06497363294997256
95% confidence interval ::
Lower Bound : -1.0666328946616939
Upper Bound : 1.196580160561639
Chebychev's Inequality ::
Lower Bound : -2.5170152645216386
Upper Bound : 2.646962530421584

```

## Plotting

In this section we help you plot samples from the Gaussian, to get a better intuition for the impact of the underlying mean, variance and the number of samples. There are no explicit questions related to this section, but it might help you better understand your answers to question 4.

Below you will see some example plotting code `plot_density` and `plot_box_and_violin`.

`plot_density` plots a histogram of the data `D` and a density estimated through a kernel density algorithm (see implementation on [github](#) for more details).

`plot_box_and_violin` plots a box plot over a violin plot. A violin plot shows the density of the sampled data (same as the density function), while the overlaid box plot shows the first quartile, median, and third quartile. More information can be found on [github](#) about these plotting utilites.

plot\_density (generic function with 1 method)

```
1 function plot_density(D)
2     histogram(
3         # data/transform parameters
4         D, norm=true,
5         # make plot pretty parameters
6         grid=false, # removes background grid
7         tickdir=:out, # changes tick direction to be out
8         lw=1, # makes line width thicker
9         color=RGB(87/255, 123/255, 181/255), # Changes fill color of histogram
10        legend=nothing, # removes legend
11        fillalpha=0.6) # makes the histogram transparent
12
13    density!(D, color=:black, lw=2)
14 end
```

plot\_box\_and\_violin (generic function with 1 method)

```
1 function plot_box_and_violin(D)
2     plt = violin(
3         ["data"], #The label for the data on the x-axis
4         D, # the data
5         grid=false, # remove the background grid
6         tickdir=:out, # set the ticks to be out
7         lw=0, # set the line width to be zero
8         color=RGB(87/255, 123/255, 181/255), # set color
9         legend=nothing)
10    boxplot!(
11        plt, # explicitly pass in plt object
12        ["data"], #The label for the data on the x-axis
13        D, # the data
14        fillalpha=0.5, # make transparent
15        lw=3) # emphasize the lines
16
17 end
```

# Visualizing the distribution from samples

Below are some sliders you can use to visualize different normal distributions interactively. The data is then plotted using the above plotting functions.

Mean:  -0.1

Std Dev:  1.0

number of samples:  100

Resample

```

1 # This is a markdown block.
2 md"""
3 ## Visualizing the distribution from samples
4
5 Below are some sliders you can use to visualize different normal distributions
  interactively. The data is then plotted using the above plotting functions.
6
7 Mean: $(@bind mu Slider(-100.0 : 0.1 : 100.0; default=0.0, show_value=true))
8
9 Std Dev: $(@bind sigma Slider(0.1 : 0.1 : 100.0; default=1.0, show_value=true))
10
11 number of samples: $(@bind n Slider(10 : 10 : 10000; default=100,
12 show_value=true))
13
14 $(@bind resample_btn Button("Resample"))
  """

```

```
gd = GaussianDistribution(-0.1, 1.0)
```

```
1 gd = GaussianDistribution(mu, sigma)
```

-0.1

```
1 gd. :μ
```

```

1 begin
2     resample_btn
3     D = sample(gd, n)
4 end;

```

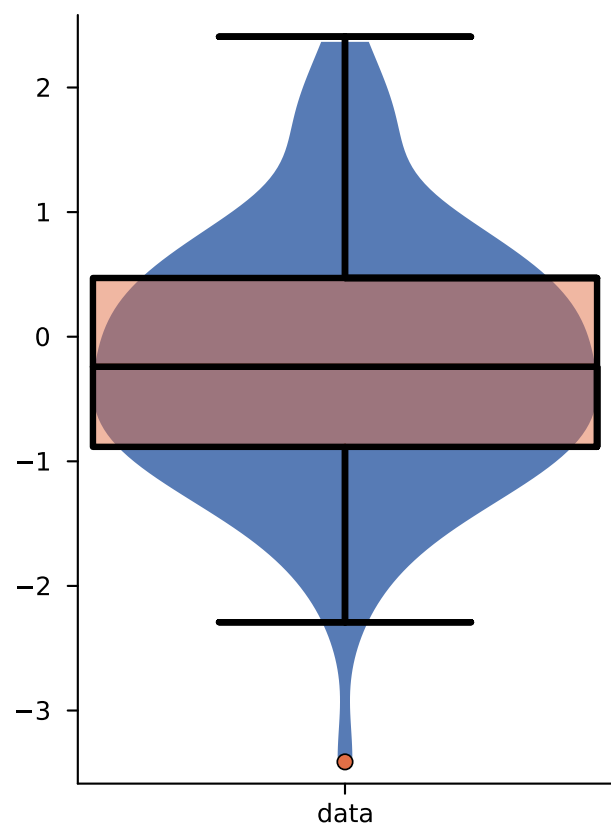
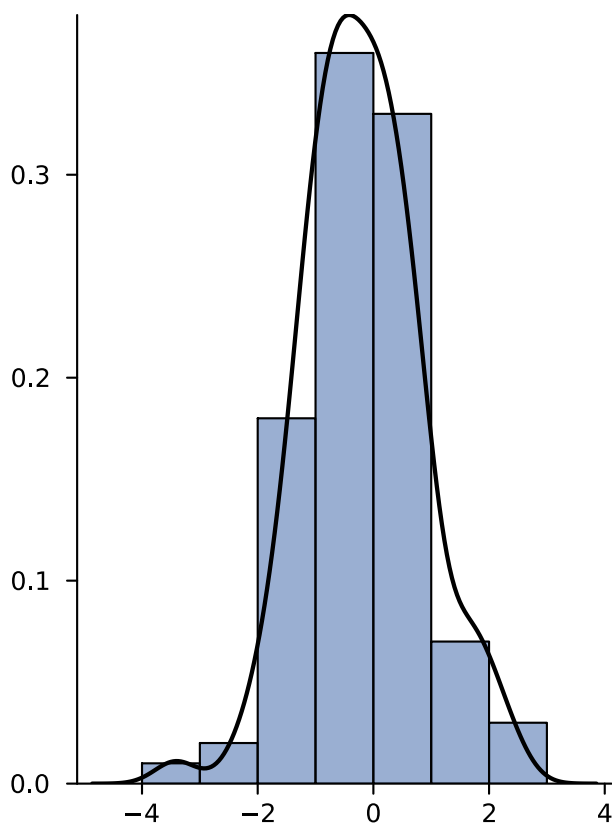
**True mean: -0.1, True Standard Deviation: 1.0**

**Sample mean: -0.185048448530157, Sample Standard Deviation: 1.0104641693264647**

```

1 let
2    $\bar{x}$  = mean(D)
3    $\bar{\sigma}$  = stddev(D)
4
5   md"""
6     **True mean:** $(gd.μ), **True Standard Deviation:** $(gd.σ)
7
8     **Sample mean:** $( $\bar{x}$ ), **Sample Standard Deviation:** $( $\bar{\sigma}$ )"""
9 end

```



```

1 let
2   plt_1 = plot_density(D)
3   plt_2 = plot_box_and_violin(D)
4   plot(plt_1, plt_2)
5 end

```

