

# Learning Sampling Distributions: A Practical Guide with R

Authored by  
**Mohammed loot**

November 6, 2025

## RECOMMENDED CITATION

Mohammed loot (2025). *Learning Sampling Distributions: A Practical Guide with R*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=11709>

Understanding the concept of a [sampling distribution](#) is absolutely fundamental to the field of **inferential statistics**. Formally, this distribution is defined as the probability distribution of a specific [statistic](#)--such as the sample mean, median, or proportion--which is derived by repeatedly drawing multiple random samples from a single, defined population.

When statisticians and data scientists work with extensive datasets or require complex Monte Carlo simulations, the open-source programming language [R](#) provides an unparalleled suite of tools for generating, visualizing, and rigorously analyzing these crucial distributions. This comprehensive tutorial will guide you through the essential steps required to construct and analyze a **sampling distribution of the mean** effectively using R code.

Generating the required sampling distribution through an iterative simulation process.

Visualizing the shape and properties of the resulting distribution using appropriate graphical methods.

Calculating the key measures of center and spread, specifically the mean and [standard deviation](#), which is commonly referred to as the **standard error**.

Calculating empirical probabilities related to specific outcomes within the distribution.

## Generating a Sampling Distribution in R

To accurately study and analyze the behavior of sample means, we must first simulate the entire process of drawing numerous random samples from a known population. This simulation approach forms the critical foundation for constructing the empirical [sampling distribution](#). The following R code block illustrates how to generate a distribution of sample means based on a simulated [normal distribution](#) using 10,000 iterations.

**#make this example reproducible**

**set.seed(0)**

#define number of samples

n = 10000

#create empty vector of length n

sample\_means = rep(NA, n)

#fill empty vector with means

for(i in 1:n){

sample\_means = mean(rnorm(20, mean=5.3, sd=9))

}

#view first six sample means

head(sample\_means)

5.283992 6.304845 4.259583 3.915274 7.756386 4.532656

In this simulation exercise, we deliberately set the total number of samples,  $N$ , to 10,000 to guarantee a robust and accurate calculation of the resulting [sampling distribution](#). We employed the R function `rnorm()` to efficiently draw samples of size 20 ( $n=20$ ) from a hypothetical population. This population is characterized by a true mean ( $\mu$ ) of 5.3 and a population [standard deviation](#) ( $\sigma$ ) of 9.

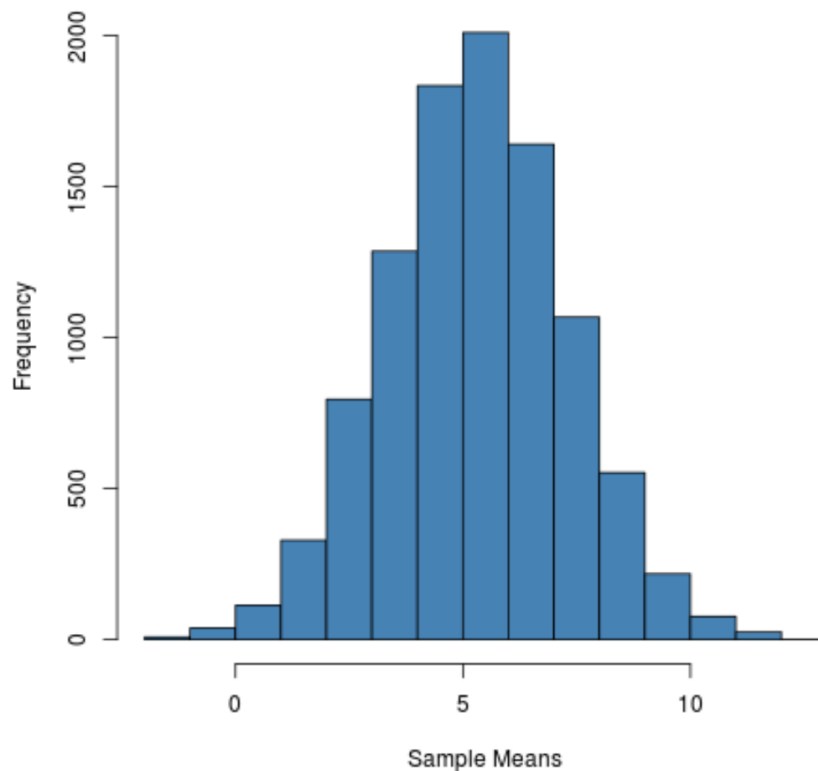
The resulting vector, `sample_means`, serves as the repository for the mean calculated from each of the 10,000 simulated samples. As evidenced by the `head()` function output, these individual sample means show natural variation; for instance, the first sample yielded a mean of 5.283992, while the second resulted in 6.304845. The collective assembly of these 10,000 values constitutes our complete simulated **empirical sampling distribution**.

## Visualizing the Sampling Distribution

After successfully generating the vector of sample means, the next essential phase involves visualizing the distribution. This visualization is critical for confirming its expected shape and intrinsic statistical properties. A cornerstone of statistical theory, the [Central Limit Theorem](#) dictates that the sampling distribution of the mean will inevitably tend toward a [normal distribution](#) (i.e., bell-shaped), irrespective of the underlying population's original shape, provided that the sample size ( $n$ ) is sufficiently large.

In R, we can employ a straightforward **histogram** to graphically represent this data. The following code uses the previously created `sample_means` vector to plot the frequency of means across various ranges, providing an immediate visual confirmation of the distribution's characteristics:

```
#create histogram to visualize the sampling distribution  
hist(sample_means, main = "", xlab = "Sample Means", col = "steelblue")
```



The resulting visualization clearly validates that the empirical [sampling distribution](#) is strikingly **bell-shaped**, with the highest concentration of sample means clustering near the expected population mean ( $\mu$ ) of 5.3. While the core of the distribution adheres to this central tendency, the tails demonstrate that some samples produced means significantly greater than 10 or less than 0. This spread is a natural reflection of the inherent **sampling variability** present in any statistical process.

## Analyzing Measures of Center and Spread

Analyzing the numerical properties of our simulated distribution is essential, as it allows us to compare the empirical results derived from the simulation against the established benchmarks of **statistical theory**. Theoretically, the mean of the sampling distribution should closely approximate the true population mean ( $\mu$ ), and its [standard deviation](#)--known specifically as the [standard error](#), denoted  $\sigma_{\bar{x}}$ --must conform to the formula  $\sigma / \sqrt{n}$ .

The following R code calculates these two critical statistical parameters directly from our generated `sample_means` data vector, providing us with the empirical mean and the empirical standard error for comparison:

**#mean of sampling distribution (Empirical Mean)**

```
mean(sample_means)
```

```
5.287195
```

```
#standard deviation of sampling distribution (Standard Error)
```

```
sd(sample_means)
```

```
2.00224
```

Our calculated empirical mean, 5.287195, shows remarkable proximity to the true population mean of 5.3. Similarly, we can assess the standard error. The theoretical standard error is calculated as  $9 / \sqrt{20}$ , which yields a value of approximately 2.012.

The actual calculated [standard deviation](#) of the sampling distribution, which is our **empirical standard error**, is  $\mathbf{2.00224}$ . This empirical finding aligns exceptionally well with the theoretical prediction of 2.012, thereby validating the accuracy of our simulation and serving as a practical demonstration of the power and reliability of the [Central Limit Theorem](#).

## Calculating Empirical Probabilities

One of the most significant practical applications of understanding the [sampling distribution](#) is its ability to calculate the probability of observing a specific [statistic](#), such as a given sample mean, based on known population parameters. Using our simulated distribution, we can easily calculate **empirical probabilities** by determining the proportion of sample means that successfully satisfy a defined condition.

For this specific example, we will calculate the probability that a sample mean drawn under our conditions is less than or equal to 6. Our known parameters are the population mean (5.3), population standard deviation (9), and a fixed sample size of 20. The following R code executes this calculation by summing the number of samples that meet the criteria and dividing that count by the total number of simulations (10,000):

```
#calculate probability that sample mean is less than or equal to 6
```

```
sum(sample_means <= 6) / length(sample_means)
```

Based on the results of our 10,000 simulations, the estimated probability that the sample mean is less than or equal to 6 is calculated as  $\mathbf{0.6417}$ . This empirical result shows strong consistency with the theoretical probability that would be derived using standard statistical methods or the R function `pnorm()`, further confirming the high accuracy of the simulation approach for generating reliable sampling distribution data.

This empirical probability value closely matches the probability calculated by a dedicated Sampling Distribution Calculator:

$\mu$  (population mean)

$\sigma$  (population standard deviation)

$n$  (sample size)

$X$  (random variable)

$P(\bar{X} \leq 6): 0.63602$

$P(\bar{X} \geq 6): 0.36398$

## The Complete R Simulation Script

For ease of use, verification, and complete reproducibility, the comprehensive [R](#) code that integrates every step outlined in this tutorial is presented below. This script covers the entire process: starting from setting the seed for simulation and generating the 10,000 sample means, moving through visualization (histogram), calculating the key parameters (mean and standard error), and concluding with the probability estimation. Running this consolidated script will allow you to recreate the entire **sampling distribution analysis** instantly.

```
#make this example reproducible  
set.seed(0)
```

```
#define number of samples
n = 10000

#create empty vector of length n
sample_means = rep(NA, n)

#fill empty vector with means
for(i in 1:n){
  sample_means = mean(rnorm(20, mean=5.3, sd=9))
}

#view first six sample means
head(sample_means)

#create histogram to visualize the sampling distribution
hist(sample_means, main = "", xlab = "Sample Means", col = "steelblue")

#mean of sampling distribution
mean(sample_means)

#standard deviation of sampling distribution
sd(sample_means)

#calculate probability that sample mean is less than or equal to 6
sum(sample_means <= 6) / length(sample_means)
```

## Additional Learning Resources

To further deepen your understanding of these core statistical concepts, we recommend exploring the following related resources:

[An Introduction to Sampling Distributions](#)

[Sampling Distribution Calculator](#)

[An Introduction to the Central Limit Theorem](#)