

How to Generate a Random Sample in R Using the `sample()` Function

Authored by
stats writer

March 2, 2026

RECOMMENDED CITATION

stats writer (2026). *How to Generate a Random Sample in R Using the `sample()` Function*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=133514>

The Foundations of Random Sampling within the R Programming Environment

In the expansive field of **data science** and **statistical analysis**, the ability to extract a representative subset from a larger population is an essential skill. The **R programming language**, a premier tool for statistical computing, provides a robust solution for this requirement through its built-in **sample()** function. This function is designed to facilitate **random sampling**, a process that ensures every element within a given **dataset** has a measurable chance of being selected. By focusing on a subset rather than the entire population, researchers can perform complex calculations and **exploratory data analysis** more efficiently, significantly reducing the computational overhead and time required to derive insights.

The importance of utilizing the **sample()** function extends beyond mere convenience; it is a critical component in ensuring the validity of **statistical inference**. When a sample is truly random, it minimizes **selection bias**, allowing the characteristics of the sample to be generalized to the population from which it was drawn. Whether you are working with a simple numeric **vector** or a multi-dimensional **data frame**, understanding the mechanics of sampling is vital for any professional working with **big data** or **predictive modeling**. This guide will explore the intricacies of the **sample()** function, providing clear examples and best practices for its implementation.

As datasets grow in complexity and scale, the strategic use of sampling becomes even more pertinent. High-performance computing often relies on **sampling techniques** to prototype algorithms or validate **machine learning** models before deploying them across full-scale production environments. The **R programming language** simplifies this workflow, offering a syntax that is both powerful for experts and accessible for beginners. By mastering the **sample()** function, users can take full advantage of R's capabilities to manage data with precision and statistical rigor.

Deconstructing the Syntax and Parameters of the Sample Function

To effectively utilize the **sample()** function, one must first understand its formal **syntax** and the specific **parameters** that govern its behavior. The flexibility of the function is derived from its four primary arguments, which allow the user to customize the sampling process to meet specific analytical needs. In R, the function signature is defined as **sample(x, size, replace = FALSE, prob = NULL)**. Each of these components plays a distinct role in determining how the elements are selected and whether the resulting output is suitable for the intended **statistical model**.

The first argument, **x**, represents the source data, which is typically a **vector** or a **dataset** containing the elements from which the sample will be drawn. The **size** parameter is an integer specifying the number of items to be returned in the sample. It is important to ensure that the

requested size is compatible with the length of the input vector, especially when performing sampling without **replacement**. The **replace** parameter is a **Boolean** value--defaulting to **FALSE**--that determines if an element can be selected more than once. Finally, the **prob** argument allows for **weighted sampling** by providing a vector of **probability** weights, giving the user granular control over the likelihood of each element's selection.

Understanding the default settings of these parameters is crucial for avoiding common errors in **data manipulation**. For instance, because the **replace** argument is set to **FALSE** by default, attempting to draw a sample larger than the original vector will result in an error. Conversely, setting **replace = TRUE** enables **bootstrap sampling**, a popular technique in **non-parametric statistics**. By carefully adjusting these arguments, a data scientist can tailor the **sample()** function to perform everything from simple random draws to complex simulations and **Monte Carlo methods**.

The complete **documentation** for **sample()** can be found [here](#).

Practical Implementation: Generating Samples from Numeric Vectors

One of the most common applications of the **sample()** function is extracting a random subset from a numeric **vector**. This process is frequently used in **hypothesis testing** and **randomized controlled trials** where subjects or data points must be assigned to different groups without **human bias**. To demonstrate this, consider a simple vector containing ten sequential integers. By applying the function, we can quickly generate a subset that maintains the **randomness** required for high-quality **data analysis**.

Suppose we have vector *a* with 10 elements in it:

```
#define vector a with 10 elements in it  
a <- c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)
```

To generate a random sample of 5 elements from vector *a* without replacement, we can use the following syntax:

```
#generate random sample of 5 elements from vector a  
sample(a, 5)  
  
# 3 1 4 7 5
```

When executing the code above, the **R interpreter** selects five distinct values from the vector. It is a fundamental characteristic of **random number generation** that each execution of this command will likely yield a different result. This inherent **variability** is what makes random sampling a

powerful tool for simulating real-world uncertainty. However, this same variability can pose challenges when a researcher needs to document their steps or share their findings with others for **peer review**, as the exact results may be difficult to replicate without further intervention.

```
#generate another random sample of 5 elements from vector a
```

```
sample(a, 5)
```

```
# 1 8 7 4 2
```

As shown in the second execution, the output differs from the first. This demonstrates the function's ability to produce **non-deterministic** results, which is essential for **stochastic modeling**. For users who require the same sample for multiple iterations of an analysis--perhaps for debugging purposes or for consistency in a published **report**--the **R programming language** offers a specific mechanism to control the **pseudorandom number generator**, which we will discuss in the subsequent section.

The Role of Randomness and Reproducibility with `set.seed()`

In the realm of **scientific computing**, **reproducibility** is a cornerstone of credible research. While the `sample()` function relies on **randomness**, there are many scenarios where a data scientist must be able to recreate the exact same "random" sample. This is particularly important when developing **algorithms**, conducting **audits**, or sharing code with colleagues. To achieve this in R, we use the `set.seed()` function. This function initializes the **pseudorandom number generator** (PRNG) with a specific starting point, ensuring that the sequence of "random" numbers generated is identical every time the code is run with that same seed.

```
#set.seed(some random number) to ensure that we get the same sample each time
```

```
set.seed(122)
```

```
#define vector a with 10 elements in it
```

```
a <- c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)
```

```
#generate random sample of 5 elements from vector a
```

```
sample(a, 5)
```

```
# 10 9 2 1 4
```

```
#generate another random sample of 5 elements from vector a
```

```
sample(a, 5)
```

```
# 10 9 2 1 4
```

By defining **set.seed(122)**, we anchor the sampling process. As demonstrated in the code block above, subsequent calls to **sample()** now yield the exact same sequence of numbers. This technique is indispensable for **data science** workflows where consistency is required across different **computing environments**. It allows other researchers to run your script and obtain the same **statistical results**, which is a vital part of the **open science** movement and general **quality assurance** in data-driven industries.

It is important to note that the specific number passed to **set.seed()** is arbitrary; the key is to use the same number consistently throughout your project. While the output appears random, it is technically **deterministic** based on the seed. This balance between **stochasticity** and control is what makes the **R programming language** such a versatile environment for **quantitative analysis**. Whether you are performing a **t-test** or training a **neural network**, managing your seeds is a hallmark of professional-grade **programming**.

Comparative Analysis: Sampling With vs. Without Replacement

A critical decision when using the **sample()** function is whether to perform sampling with or without **replacement**. This choice significantly impacts the **probability distribution** of your results and the **statistical validity** of your model. By default, R performs sampling without replacement (**replace = FALSE**), meaning that once an element is selected, it is removed from the pool and cannot be chosen again. This is analogous to drawing names from a hat; once a name is drawn, the total number of names remaining in the hat decreases.

In contrast, sampling with **replacement** (**replace = TRUE**) allows elements to be selected multiple times. This is a fundamental requirement for **bootstrapping**, a **resampling method** used to estimate the **sampling distribution** of an estimator. When replacement is enabled, the selection of one element does not affect the probability of selecting that same element in the next draw. This is useful for simulating processes like **coin flips** or **rolling dice**, where each event is independent of the previous one.

```
#generate random sample of 5 elements from vector a using sampling with replacement  
sample(a, 5, replace = TRUE)
```

```
# 10 10 2 1 6
```

In the example provided, notice that the number 10 appears twice in the resulting sample. This would be impossible under the default settings. Choosing the correct approach depends entirely on the **research design**. For **randomized assignments** in experiments, you generally sample without replacement to ensure each subject is only assigned once. However, if you are conducting **statistical simulations** or working with **Bayesian inference**, sampling with replacement might be

the more appropriate **methodology**. Understanding these mathematical nuances ensures that your **data analysis** is theoretically sound.

Manipulating Complex Data Structures and Data Frames

While sampling from vectors is useful for basic tasks, real-world **data science** projects usually involve **data frames**. A data frame in R is a two-dimensional structure where each column can contain different types of data (e.g., numeric, character, factor). The **sample()** function is frequently used to select a random subset of rows from a **dataset**, which is an essential step in **data preprocessing** and **model validation**. This allows the user to create **training datasets** and **test datasets** for **machine learning** applications.

For the following example, we will generate a random sample of 10 rows from the built-in R dataset **iris**, which is a famous dataset in **statistics** and **pattern recognition**. The **iris dataset** contains 150 total rows representing different flower measurements. By sampling rows, we can analyze a manageable portion of the data while maintaining the relationship between different variables like sepal length and petal width.

```
#view first 6 rows of iris dataset
```

```
head(iris)# Sepal.Length Sepal.Width Petal.Length Petal.Width Species
```

```
#1 5.1 3.5 1.4 0.2 setosa
```

```
#2 4.9 3.0 1.4 0.2 setosa
```

```
#3 4.7 3.2 1.3 0.2 setosa
```

```
#4 4.6 3.1 1.5 0.2 setosa
```

```
#5 5.0 3.6 1.4 0.2 setosa
```

```
#6 5.4 3.9 1.7 0.4 setosa#set seed to ensure that this example is replicable  
set.seed(100)#choose a random vector of 10 elements from all 150 rows in iris dataset
```

```
sample_rows <- sample(1:nrow(iris), 10)
```

```
sample_rows
```

```
# 47 39 82 9 69 71 117 53 78 25
```

```
#choose the 10 rows of the iris dataset that match the row numbers above
```

```
sample <- iris
```

```
sample
```

```
# Sepal.Length Sepal.Width Petal.Length Petal.Width Species
```

```
#47 5.1 3.8 1.6 0.2 setosa
```

```
#39 4.4 3.0 1.3 0.2 setosa
```

```
#82 5.5 2.4 3.7 1.0 versicolor
```

```
#9 4.4 2.9 1.4 0.2 setosa
#69 6.2 2.2 4.5 1.5 versicolor
#71 5.9 3.2 4.8 1.8 versicolor
#117 6.5 3.0 5.5 1.8 virginica
#53 6.9 3.1 4.9 1.5 versicolor
#78 6.7 3.0 5.0 1.7 versicolor
#25 4.8 3.4 1.9 0.2 setosa
```

The logic here involves two steps: first, we use **sample()** to generate a vector of random row indices (from 1 to the total number of rows in the **iris** dataset). Second, we use those indices to subset the **data frame**. This method is highly efficient and is a standard practice in **exploratory data analysis**. By selecting a small, random portion of a large **database**, researchers can visualize trends and identify **outliers** without the need for excessive **computational power**.

Advanced Probability Weighting and Bias Control

One of the most powerful yet underutilized features of the **sample()** function is the **prob** argument. This allows for **non-uniform sampling**, where certain elements have a higher **probability** of being selected than others. In many real-world scenarios, data is not uniformly distributed, and **stratified sampling** or weighted draws are necessary to reflect the true nature of the **population**. For instance, in a **marketing survey**, you might want to weight responses from a specific demographic more heavily to correct for **underrepresentation**.

To implement this, you provide a **numeric vector** to the **prob** parameter that is the same length as the input vector **x**. These weights do not necessarily need to sum to one, as R will automatically normalize them. This capability is vital in **econometrics** and **social sciences**, where **weighted averages** and **sampling weights** are used to produce accurate **statistical estimates**. By mastering the **prob** argument, you can move beyond simple draws and conduct sophisticated **simulations** that account for complex real-world variables.

Ultimately, the **sample()** function in R is a multifaceted tool that supports a wide range of **analytical workflows**. From simple random draws to complex **reproducible research** and weighted **data frame** subsetting, it provides the flexibility required for modern **data science**. By integrating these techniques into your **R programming** toolkit, you ensure that your **data analysis** is not only efficient but also statistically rigorous and aligned with industry best practices. Whether you are a student or a professional **data analyst**, the ability to generate reliable samples is a foundational skill that will enhance the quality of your insights.