

How do I perform systematic sampling in R?

Authored by
stats writer

December 23, 2025

RECOMMENDED CITATION

stats writer (2025). *How do I perform systematic sampling in R?*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=108470>

Welcome to this comprehensive guide on executing **systematic sampling** efficiently using the R programming environment. Systematic sampling is a straightforward and widely utilized statistical method for selecting a representative subset from a larger population. Unlike simple random sampling, where every element has an equal chance of selection at every draw, systematic sampling introduces a fixed interval after a random start, ensuring coverage across the entire ordered list. This methodology is particularly powerful when the population elements are already ordered in a meaningful way, such as chronologically or alphabetically.

While R provides powerful built-in functions like `sample()` for general selection, performing true systematic sampling requires a careful combination of functions to establish the fixed interval selection rule. This tutorial will walk through the conceptual steps, illustrate how to define a custom function to handle the selection interval (k), and apply it to a practical data scenario, ensuring clarity and ease of implementation. We will emphasize the importance of setting a seed for reproducibility, which is a cornerstone of reliable data analysis.

Understanding Systematic Sampling

In statistical research, it is often impractical or impossible to study every member of a large population. Therefore, researchers frequently draw a smaller, manageable sample from the population and use the data gathered from this subset to make statistically valid inferences about the population as a whole. The method by which this subset is selected is critical to ensuring the validity of these conclusions.

Systematic sampling is categorized as a probability sampling technique, meaning that every element in the population has a known, non-zero chance of being selected. Its primary advantage lies in its simplicity and efficiency, especially when dealing with very large populations where physical ordering is feasible. The technique minimizes the risk of researcher bias while ensuring that the sample is spread evenly throughout the population list, assuming the underlying order is not cyclical or biased itself.

Researchers often utilize this method because it offers the convenience of simple random sampling but provides better coverage across the population structure, potentially yielding a more precise estimate than simple random sampling when elements are ordered randomly or sequentially.

The Mechanics of Systematic Sampling

The process of systematic sampling is remarkably simple and relies on two core actions. Before beginning the selection process in R, it is vital to understand these foundational steps, as they directly translate into the code we will write. The simplicity of this approach makes it a favored

method in contexts like quality control, public surveys, and administrative data analysis.

The implementation of systematic sampling involves the following step-by-step procedure:

Ordering the Population: Every member of a population must first be placed into a specific, predetermined order. This order might be based on alphabetical criteria (like student names), numerical criteria (like transaction IDs), or chronological criteria (like time stamps).

Determining the Sampling Interval (k): The sampling interval, denoted as 'k', is calculated by dividing the total population size (N) by the required sample size (n). If the result is not a whole number, it is standard practice to use the `ceiling()` function to ensure the interval is maintained appropriately for index selection.

Selecting a Random Start (r): A random starting point (r) is chosen within the range of 1 to k. This initial random selection prevents any bias that might arise from always starting at the beginning of the list.

Iterative Selection: Beginning with the random start (r), select every kth member thereafter until the desired sample size (n) is achieved.

This structured approach ensures that the sample is evenly dispersed, providing geographic or temporal representation across the dataset, depending on the initial ordering criteria chosen.

Implementing Systematic Sampling in R: Core Concepts

While R's base functions are powerful, there isn't a single function dedicated solely to systematic sampling. Instead, we utilize a combination of functions, including `sample()`, `ceiling()`, and `seq()`, within a customized function structure. This approach demonstrates the flexibility of R for handling complex statistical procedures.

The initial step involves defining the population size (N) and the desired sample size (n). In R, if our data is stored in a data frame, N will correspond to the number of rows, which can be easily retrieved using `nrow(df)`. We then calculate the sampling interval *k* using the `ceiling()` function in conjunction with the division N/n to ensure *k* is an integer that can be used for indexing. The truly random part of the selection process is handled by choosing the starting index (r) using `sample(1:k, 1)`.

Finally, the sequence of selected indices is generated using the `seq()` function, which takes the starting point (r), the endpoint ($r + k*(n-1)$), and the interval (k). These generated indices are then used to subset the original data frame, extracting the systematic sample. This method ensures that the selection is purely systematic and based strictly on positional indices, which is crucial for statistical integrity.

Setting Up the Example Data Frame in R

To illustrate this concept practically, let us consider a scenario frequently encountered in educational administration. Suppose a school superintendent wishes to obtain a sample of 100 students from a larger student population of 500 total students ($N=500$, $n=100$). The superintendent decides to employ systematic sampling. The students are first placed in order alphabetically according to their last names. Since the ratio N/n is $500/100 = 5$, the sampling interval (k) will be 5. She will randomly choose a starting point between 1 and 5, and then pick every 5th student to be included in the final sample.

We begin by creating a simulated data frame in R to mimic this administrative list. This data structure will contain 500 observations, each representing a student with a randomly generated last name and a corresponding Grade Point Average (GPA). We use the `rnorm()` function to generate realistic, normally distributed GPA scores.

The crucial element here is the initial ordering. By default, when we create the data frame in R, the observations are indexed from 1 to 500. For systematic sampling to be valid in this example, we assume that this default row order corresponds to the alphabetical ordering by last name, fulfilling the first requirement of the methodology.

Ensuring Reproducibility: The Role of `set.seed()`

Before executing any code that involves random generation--such as the creation of random names or the selection of the random starting point--it is absolutely essential to set a seed. Using the `set.seed()` function ensures that the results obtained are fully reproducible. If you or another researcher run the code at a later date, setting the same seed guarantees that the random number generation sequence remains identical, yielding the exact same output. This is fundamental to good scientific practice and data auditing, especially when relying on random processes like the `sample()` function.

The following code block demonstrates how to set the seed, create a simple helper function to generate random last names (since the ordering of the population is important in systematic sampling), and construct the main data frame, `df`, containing 500 student records with simulated GPA scores normally distributed around a mean of 82.

#make this example reproducible

```
set.seed(1)
```

```
#create simple function to generate random last names
```

```
randomNames <- function(n = 5000) {
```

```
do.call(paste0, replicate(5, sample(LETTERS, n, TRUE), FALSE))
```

```
}  
  
#create data frame  
df <- data.frame(last_name = randomNames(500),  
gpa = rnorm(500, mean=82, sd=3))  
  
#view first six rows of data frame  
head(df)  
  
last_name gpa  
1 GONBW 82.19580  
2 JRRWZ 85.10598  
3 ORJFW 88.78065  
4 XRYNL 85.94409  
5 FMDCE 79.38993  
6 XZBJC 80.49061
```

Developing the Systematic Sampling Function

The next crucial step is defining a reusable function in R that encapsulates the logic of systematic sampling. This function, which we name `obtain_sys`, takes two arguments: `N` (the total population size, usually `nrow(df)`) and `n` (the desired sample size).

Inside the function, we first calculate the sampling interval `k` using `k = ceiling(N/n)`. The `ceiling()` function is vital here as it ensures `k` is an integer large enough to cover the entire population when generating the indices. Next, we determine the random starting point `r` by randomly selecting one number from the range 1 to `k` using `r = sample(1:k, 1)`. This single random selection satisfies the probability requirement of systematic sampling.

The core of the selection process is handled by the `seq()` function. It generates a sequence of indices starting at `r`, stepping by `k`, and stopping at `r + k*(n-1)`. This sequence guarantees that we obtain exactly `n` members, each separated by the fixed interval `k`. This sequence is then returned by the function, ready to be used for data frame subsetting.

Executing the Sample Selection and Analyzing Results

With our custom function defined, applying the systematic sampling procedure is straightforward. We call `obtain_sys` using the total number of rows in our data frame (`N = 500`) and the desired sample size (`n = 100`). The returned indices are then used to subset the original data frame `df`, creating `sys_sample_df`, which holds our final systematic sample.

The subsequent code demonstrates this execution and shows the initial results. Since $N=500$ and $n=100$, the interval k is 5. By reviewing the output, we can observe that the randomly selected starting index (r), determined by the seed and the `sample()` function within our custom definition, was 3, meaning the selection sequence begins at row 3 and increases by 5 for each subsequent member (3, 8, 13, 18, and so on). This consistent interval spacing is the hallmark of systematic sampling.

#define function to obtain systematic sample

```
obtain_sys = function(N,n){
```

```
  k = ceiling(N/n)
```

```
  r = sample(1:k, 1)
```

```
  seq(r, r + k*(n-1), k)
```

```
}
```

```
#obtain systematic sample
```

```
sys_sample_df = df
```

```
#view first six rows of data frame
```

```
head(sys_sample_df)
```

```
last_name gpa
```

```
3 ORJFW 88.78065
```

```
8 RWPSB 81.96988
```

```
13 RACZU 79.21433
```

```
18 ZOHKA 80.47246
```

```
23 QJETK 87.09991
```

```
28 JTHWB 83.87300
```

```
#view dimensions of data frame
```

```
dim(sys_sample_df)
```

```
100 2
```

After reviewing the initial six rows of the resulting data frame, we confirm that the systematic selection was successful. The first member included in the sample was indeed located in row 3 of the original data frame. Crucially, each subsequent member in the sample is located exactly 5 rows after the previous member (8, 13, 18, etc.). This adherence to the fixed interval $k=5$ confirms that the selection process was purely systematic.

To finalize the verification, we utilized the `dim()` function, which returns the dimensions of the generated data structure. The output `100 2` confirms two important facts: first, the systematic

sample we obtained is precisely the desired size of 100 observations; and second, it retains the two variables (columns) from the original data frame (`last_name` and `gpa`). This confirms the successful use of R for generating a highly structured and representative systematic sample.

Conclusion: Advantages and Applications

The implementation of systematic sampling in R, as demonstrated through a custom function, provides a reliable and transparent method for researchers. Its primary advantage over simple random sampling is its ability to spread the sample uniformly across the entire population, minimizing the chance of clustered selection and ensuring that all strata of the ordered population are represented.

However, analysts must be cautious of potential pitfalls. If the ordering of the population contains a hidden periodicity or cycle that aligns with the sampling interval (k), the resulting sample may be highly biased. For instance, if $k=7$ and the list is ordered by day of the week, one might consistently sample only observations from the same day, leading to skewed results. In cases where hidden periodic patterns are suspected, alternative methods like stratified sampling or cluster sampling might be more appropriate.

Nonetheless, for ordered administrative data, manufacturing quality checks, or simple temporal sampling, systematic sampling remains an effective and powerful tool when properly implemented using robust coding practices, ensuring reproducibility and statistical rigor.

For those interested in exploring other methodologies for statistical data selection, resources are available detailing alternative sampling strategies:

[Types of Sampling Methods](#)

[Stratified Sampling in R](#)

[Cluster Sampling in R](#)