

# How to Calculate Sample & Population Variance in R?

Authored by  
**stats writer**

December 17, 2025

## RECOMMENDED CITATION

stats writer (2025). *How to Calculate Sample & Population Variance in R?*.

PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=107704>

The initial approach to calculating variance in the R Programming Language relies primarily on the built-in `var()` function. It is crucial for users to understand that, by default, the `var()` function calculates the **sample variance**, which is the most common requirement in statistical analysis. Historically, discussions often mention a theoretical `var.test()` function for population variance; however, in standard R practice, the `var()` function is used, and the population variance must be derived through a simple mathematical adjustment, as R does not include a direct, single-function implementation for population variance calculation comparable to the sample function.

To successfully execute the `var()` function for **sample variance**, the user simply needs to provide the dataset--typically a numeric vector or a column from a data frame--as the primary input. The result returned is a single numerical value representing the dispersion of the data points around the mean. For instance, if attempting to calculate the **population variance**, one must understand the subtle but vital difference in the underlying statistical formulas, which necessitates knowing the size of the dataset (N) and applying the appropriate correction factor to ensure accuracy in the calculated measure of spread.

Both functions are integral to descriptive statistics, providing essential information about the variability inherent in the collected data. While `var()` handles the majority of inferential tasks by providing the sample estimate, the methodology for calculating the true population parameter requires careful consideration of the statistical context and the application of manual algebraic steps within the R environment, especially when working with complete census data rather than a limited sample.

## Introduction to Statistical Variance in R

Understanding statistical variance is fundamental to mastering data analysis, as it quantifies the spread or dispersion of data points within a distribution. Specifically, **variance** measures how far a set of numbers is spread out from their average value, offering a key metric of variability. A high variance indicates that data points are widely dispersed, while a low variance suggests that the data points tend to be clustered closely around the mean. This measure is essential for researchers and statisticians using R, as it provides crucial context for interpreting standard deviation and performing subsequent hypothesis testing, making it a cornerstone of both descriptive and inferential statistics.

In the context of the R Programming Language, calculating this measure of dispersion is straightforward, although the distinction between calculating variance for a population versus a sample remains paramount. The standard R function `var()` is optimized to calculate sample variance, recognizing the reality that most statistical work involves inferential statistics based on subsets of the total population. This optimization reflects a core principle of statistical inference, where results from a small group are generalized to a larger, unknown group, requiring the use of

unbiased estimators to maintain statistical rigor.

When working with real-world datasets in R, one must always prioritize clarity regarding whether the dataset represents a complete census (a population) or a subset chosen for observation (a sample). This distinction dictates which formula--and thus, which calculation method--is appropriate. Failing to account for this difference can lead to skewed analyses and inaccurate inferential conclusions regarding the true variability of the underlying system being studied. The subsequent sections will detail how R handles both these critical calculations effectively, ensuring that analysts can accurately quantify variability regardless of the dataset type.

## Theoretical Foundations: Sample vs. Population Variance

Statistically, a clear differentiation exists between a population and a sample, and this difference directly impacts how variance is calculated. A **population** encompasses all possible observations or subjects relevant to a study, while a **sample** is merely a manageable subset drawn from that larger population. Because a sample is only an estimate of the overall distribution, its variance calculation requires a specific adjustment known as Bessel's Correction to ensure the resulting statistic is an unbiased estimator of the true population variance.

The core issue revolves around the calculation of the Degrees of Freedom. When calculating sample variance, the sum of squared differences is divided by  $(n-1)$ , where 'n' is the sample size. This reduction by one degree of freedom accounts for the fact that the sample mean, used in the calculation, is itself derived from the sample data. If we were to use the sample mean and divide by 'n' instead of  $(n-1)$ , the resulting variance estimate would systematically underestimate the true population variance, particularly in scenarios involving smaller sample sizes where this bias is most pronounced.

Conversely, when calculating **population variance**, we assume that the dataset contains all relevant observations, and therefore, the true population mean ( $\mu$ ) is known, or at least treated as known for the purpose of the calculation. Consequently, no correction is required, and the sum of squared differences is simply divided by the total population size (N). Understanding the necessity of the Bessel's Correction is crucial for anyone performing statistical analysis in R, as the standard `var()` function inherently uses the sample formula, requiring division by the denominator  $(n-1)$  to produce an unbiased estimate.

## Mathematical Formulas for Dispersion Measurement

The fundamental definition of variance remains consistent--it is the average of the squared deviations from the mean. However, the exact mathematical implementation diverges based on the dataset type being analyzed. The formula to find the variance of a **population** (often denoted by the Greek letter sigma squared,  $\sigma^2$ ) explicitly uses the population mean ( $\mu$ ) and the total size (N) in

its denominator, representing the true average squared distance from the center point of the entire data distribution.

The formula to find the variance of a population is:

$$\sigma^2 = \sum (x_i - \mu)^2 / N$$

In this formula,  $\mu$  represents the **population mean**,  $x_i$  is the  $i$ th element from the population,  $N$  denotes the population size, and  $\sum$  is the standard mathematical symbol signifying summation across all elements. This calculation assumes perfect knowledge of the entire dataset, which is rare outside of academic exercises or specific, bounded datasets where a complete census has been performed. The result provides the true dispersion of the entire set of observations, not just an estimate.

In contrast, the formula required to find the variance of a **sample** (denoted by  $s^2$ ) relies on the sample mean ( $\bar{x}$ ) and incorporates the necessary adjustment for unbiased estimation. This critical distinction is the primary reason the standard R function `var()` is designed to calculate this value, as statistical inference based on samples is far more common in empirical research than full population analysis. The inclusion of the degrees of freedom correction ensures that the sample variance does not systematically underreport the true variability.

The formula to find the variance of a sample is:

$$s^2 = \sum (x_i - \bar{x})^2 / (n-1)$$

Here,  $\bar{x}$  is the **sample mean**,  $x_i$  is the  $i$ th element in the sample, and  $n$  is the sample size. The denominator  $(n-1)$  represents the Degrees of Freedom, applying Bessel's Correction to provide a statistically sound estimate of the population variance based on the limited sample data. This mathematical refinement ensures that our estimates are as accurate and unbiased as possible given the constraints of working with a subset of the total observations.

### Example: Calculating Sample Variance using the R var() Function

When working within the R Programming Language environment, calculating the **sample variance** is remarkably straightforward due to the dedicated `var()` function. This function is designed explicitly for analyzing sampled data, adhering to the industry standard of using the  $(n-1)$  denominator correction, making it the preferred method for assessing data variability when the dataset is not exhaustive of the entire population. This simplicity allows analysts to quickly obtain a crucial measure of data spread with minimal coding effort.

Suppose we introduce a small, representative dataset into R to demonstrate this functionality. This dataset, consisting of ten numeric observations, allows us to clearly illustrate the output of the

`var()` command and set the baseline for comparison with the subsequent population variance calculation. We first define the numeric vector, which is the standard input format for this function:

```
# Define a representative dataset in R
```

```
data <- c(2, 4, 4, 7, 8, 12, 14, 15, 19, 22)
```

Once the data is defined and assigned to the variable `data`, calculating the **sample variance** requires only a single, simple command. The R engine automatically determines the sample size ( $n=10$ ), calculates the sample mean, and applies the sample variance formula (dividing by  $n-1 = 9$ ). The output clearly shows the calculated measure of dispersion, which is the sample variance ( $s^2$ ), providing immediate insight into the dataset's volatility.

```
# Calculate the sample variance using the var() function
```

```
var(data)
```

```
46.01111
```

This result, 46.01111, represents the estimated variance of the underlying population based on the input sample data. It is critical to always remember that the `var()` function returns the sample variance by default; no additional arguments are needed to achieve this standard calculation, making it highly efficient for exploratory data analysis and initial statistical reporting.

## Deriving Population Variance through Manual Adjustment

Since the base R Programming Language does not offer a dedicated, single-step function like `var()` for calculating **population variance** ( $\sigma^2$ ), statisticians must perform a simple manual adjustment to the result obtained from the sample variance calculation. This necessity arises because the population formula requires division by  $N$  (the total count) rather than  $(n-1)$ , effectively removing Bessel's Correction from the formula. This procedure ensures the calculation reflects the true population parameter when the dataset is complete.

The mathematical relationship between the two variances is straightforward: Population Variance ( $\sigma^2$ ) equals Sample Variance ( $s^2$ ) multiplied by the scaling factor  $(n-1)/n$ . Therefore, we can leverage the already calculated sample variance from the `var()` function and apply this scaling factor directly. To achieve this in R, we first need to determine the sample size ( $n$ ) using the highly useful `length()` function before performing the necessary multiplication.

```
# Determine the length of the data vector (n)
```

```
n <- length(data)
```

```
# Calculate population variance by adjusting sample variance
```

```
var(data) * (n-1)/n
```

```
41.41
```

The resulting value, 41.41, is the true **population variance** for this specific dataset, assuming it constitutes the complete population of interest. This manual calculation consistently yields a lower value than the sample variance (46.01111 vs. 41.41). This difference is expected, as the divisor  $N$  is always larger than  $(n-1)$ , thus providing a slightly smaller measure of dispersion when the entire population is considered, reflecting the reduced uncertainty compared to estimating from a sample.

It is an important statistical note that the population variance will invariably be smaller than the sample variance for any given dataset (where  $N > 1$ ). In practical applications, however, most research involves calculating **sample variances** for datasets, as collecting data for an entire theoretical population is extremely rare. Therefore, the standard `var()` output is often the relevant statistic for reporting and subsequent inference, making the ability to switch between the two calculations a fundamental skill for R users.

### Advanced Application: Calculating Variance Across Data Frame Columns

Beyond single vectors, real-world data analysis in R frequently involves working with complex data frames containing multiple variables (columns). A common task is to efficiently calculate the sample variance for every single variable within the data frame. R offers vectorized solutions that streamline this process, preventing the need for repetitive coding loops which can be slow and error-prone, especially with large datasets.

To illustrate this advanced technique, let us define a sample data frame containing three different variables (columns a, b, and c). This setup mimics a common scenario in data analysis where multiple features or metrics are collected for the same set of observations:

```
# Create a sample data frame with multiple variables
```

```
data <- data.frame(a=c(1, 3, 4, 4, 6, 7, 8, 12),
```

```
b=c(2, 4, 4, 5, 5, 6, 7, 16),
```

```
c=c(6, 6, 7, 8, 8, 9, 9, 12))
```

```
# View the structure of the data frame
```

```
data
```

```
a b c
```

```
1 1 2 6
```

```
2 3 4 6
```

```
4 4 5 8
```

```
5 6 5 8
6 7 6 9
7 8 7 9
8 12 16 12
```

To calculate the **sample variance** for each column simultaneously, we employ the powerful `sapply()` function. The `sapply()` function applies a specified function (in this case, `var`) across all elements (columns) of an input list or data frame, simplifying the calculation immensely. This method is highly preferred over iterative loops for its speed, conciseness, and superior readability in R programming.

```
# Find the sample variance of each column using sapply()
sapply(data, var)
```

```
a b c
11.696429 18.125000 3.839286
```

The resulting output provides three variance measures, corresponding to columns a, b, and c, respectively. This demonstrates the efficiency of R's functional programming capabilities when dealing with multivariate analysis. For instance, column 'b' exhibits the highest variance (18.125), indicating its values are the most spread out, while column 'c' has the lowest variance (3.839), showing its values are tightly clustered around its mean, suggesting lower overall variability for that specific feature.

## Connecting Variance to Standard Deviation in R

While variance provides an essential mathematical measure of dispersion, its utility for direct intuitive interpretation is often limited because it is expressed in squared units, which do not align with the original data measurements. This limitation leads analysts to frequently utilize the **standard deviation**, which is simply the positive square root of the variance. Standard deviation restores the measure of spread back into the original units of measurement, making it highly interpretable and easier to report to non-statistical audiences.

In R, the `sd()` function is used to calculate the **sample standard deviation**, and similar to `var()`, it incorporates Bessel's Correction by default, ensuring an unbiased estimate. If we wish to calculate the standard deviation for the same columns used in the previous example, we can substitute `sd` for `var` within the `sapply()` function. This demonstrates the seamless integration and flexibility of the R environment for computing related descriptive statistics quickly and efficiently across multiple variables.

Using the same data frame containing columns a, b, and c, the command structure remains simple, showcasing the consistency of R's base functions when transitioning between variance and standard deviation calculations:

**# Find the sample standard deviation of each column**

**sapply(data, sd)**

```
a b c
```

```
3.420004 4.257347 1.959410
```

These results confirm the earlier findings regarding dispersion: the standard deviation for column 'b' (4.257) is the highest, and for column 'c' (1.959) is the lowest. Note that 4.257 is the square root of 18.125, confirming the inherent statistical link between the two measures. In summary, while **variance** is essential for theoretical computations and inferential statistics (especially in ANOVA or linear modeling), **standard deviation** is often preferred for summarizing data due to its intuitive scale and direct comparability to the mean.

*For analysts looking to deepen their expertise in statistical methods and data manipulation techniques within the R environment, extensive additional tutorials and guides are readily available.*