

# How to Perform Univariate Analysis in R (With Examples)

Authored by  
**stats writer**

December 8, 2025

## RECOMMENDED CITATION

stats writer (2025). *How to Perform Univariate Analysis in R (With Examples)*.  
PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=106677>

The initial step in any rigorous data analysis project is understanding the fundamental characteristics of the variables involved. This process is known as univariate analysis, a crucial technique used in statistics to examine the distribution of a single variable. The term itself is derived from the prefix "uni," meaning "one," highlighting its focus on isolating and studying one dataset column at a time. This methodology provides foundational insights into the data's shape, central location, and variability, helping analysts identify potential outliers, skewness, and the overall reliability of the measurements before moving to more complex multivariate techniques. Mastering univariate analysis is indispensable for data cleaning, preliminary modeling, and descriptive reporting, establishing the context for all subsequent statistical exploration.

When working within the R programming language, analysts have access to a robust suite of tools for conducting thorough univariate examinations. R is particularly adept at handling both categorical and numerical data distributions, offering functions that yield quick, yet powerful, descriptive results. For numerical data, such as measurements, counts, or scores, univariate analysis relies on three primary methods that collectively paint a complete picture of the variable's behavior. These methods are essential for transforming raw data into meaningful statistical information, ensuring that assumptions for future inferential tests are met and that the data is well understood by the researcher.

## The Three Pillars of Univariate Data Exploration

To perform a complete analysis of a single variable, statisticians rely on a triad of complementary techniques. Each technique addresses a different aspect of the data's distribution, ensuring that no important features are overlooked. The integration of these methods--numerical summaries, tabulation, and visualization--provides a holistic view necessary for informed decision-making. Failing to utilize all three approaches can lead to an incomplete understanding of the data's underlying structure, potentially resulting in flawed interpretations during later stages of analysis.

The first pillar involves calculating summary statistics. These numerical descriptors quantify the characteristics of the distribution, providing precise measures of where the data is centered and how widely it is spread. Key statistics include the mean, median, mode (for location), and standard deviation, variance, and range (for dispersion). These statistics offer an immediate, quantitative assessment of the variable, often serving as the primary output in descriptive reporting. Analysts rely on these numbers to quickly compare the variable against established benchmarks or industry standards.

The second method utilizes the frequency table, which is particularly useful for discrete or categorical data. A frequency table systematically counts how often each unique value appears within the dataset. While less common for truly continuous data (where grouping or binning is required), it is invaluable for understanding the distribution of nominal or ordinal variables, such as

counts of occurrences, ratings, or categories. This tabulation provides granularity that summary statistics often mask, revealing exact value counts and helping to identify the most frequent observations.

Finally, the third and perhaps most intuitive method involves creating specialized charts and visualizations. Graphical representations, such as histograms, boxplots, and density curves, are indispensable for visualizing the overall distribution of values. These charts allow analysts to quickly assess the shape of the data--including symmetry, skewness, and modality--and detect outliers that might skew numerical summaries. A visual inspection often confirms or challenges the conclusions drawn from the numerical statistics, providing robust evidence for the variable's distributional properties.

## Setting Up Data for Analysis in R

Before delving into the specific analytical functions, we must first establish the dataset we intend to examine. For the purpose of this tutorial, we will focus on a simple numerical vector created in R. This vector, labeled `x`, consists of fifteen values designed to represent a typical small sample of quantitative observations. This variable allows us to demonstrate how R functions handle basic numerical data and how descriptive statistics are derived from this input.

The creation of the variable `x` utilizes R's concatenation function, `c()`, which combines the individual values into a single vector object. This is a fundamental operation in R data handling. It is good practice to ensure the data type is appropriate for the intended analysis; since these are numerical values, R recognizes them as either integers or doubles, suitable for calculation of means and standard deviations. The following code block illustrates the setup and creation of our sample variable, which will be the basis for all subsequent univariate explorations:

```
#create variable with 15 numerical values for analysis  
x <- c(1, 1, 2, 3.5, 4, 4, 4, 5, 5, 6.5, 7, 7.4, 8, 13, 14.2)
```

Once the vector `x` is successfully loaded into the R environment, we can proceed directly to the calculation of its descriptive statistics. The subsequent sections will utilize this variable to demonstrate the necessary R commands and interpret the output, providing a clear, step-by-step guide to univariate analysis.

## Calculating and Interpreting Measures of Central Tendency and Dispersion

The first step in analyzing the vector `x` is to calculate key summary statistics, which provide numerical benchmarks for the center and spread of the data. Understanding the central location helps us determine the typical value, while measuring dispersion shows us how much the data

varies around that center. R provides specific, intuitive functions for calculating these statistics rapidly. These measures are foundational to statistical inference, as they quantify the primary characteristics of the variable's distribution.

We begin by examining the Measures of central tendency, specifically the arithmetic mean and the median. The **mean** (calculated using `mean(x)`) is the average value, often sensitive to outliers. For our dataset, the mean is approximately 5.71. The **median** (calculated using `median(x)`) is the middle value when the data is ordered, offering a measure of center that is resistant to extreme values. The median of our dataset is 5. The fact that the mean (5.71) is slightly higher than the median (5) suggests a minor positive skew in the data distribution, meaning there are a few higher values pulling the average upward.

Next, we quantify the data's dispersion or variability. The simplest measure of spread is the **range**, which is the difference between the maximum and minimum values (calculated by `max(x) - min(x)`). Our range is 13.2 (14.2 - 1). While useful, the range is highly sensitive to extreme observations. A more robust measure of spread is the **Interquartile Range (IQR)**, calculated using `IQR(x)`, which measures the spread of the middle 50% of the data. Our IQR of 3.45 indicates the variability concentrated around the median. Finally, the Standard deviation (calculated using `sd(x)`) provides the average distance between each data point and the mean. A standard deviation of approximately 3.86 indicates a substantial spread relative to the mean of 5.71, suggesting heterogeneity in the sample values.

The following R code demonstrates the calculation of these key statistics and the corresponding output generated by the R environment, allowing for a precise numerical description of the variable `x`:

```
#find mean (average value)
mean(x)
5.706667

#find median (middle value)
median(x)

5

#find range (max minus min)
max(x) - min(x)

13.2

#find interquartile range (spread of middle 50% of values)
IQR(x)
```

3.45

```
#find standard deviation (average deviation from the mean)
```

```
sd(x)
```

3.858287

## Analyzing Discrete Observations with Frequency Tables

While summary statistics offer a high-level overview of centrality and dispersion, they do not reveal how often specific values appear. For a more granular view, especially when dealing with discrete data or when examining repeated values in continuous data, a frequency table is essential. This table systematically organizes the data by counting the occurrences of each unique observation, providing the absolute frequency distribution.

In R, the `table()` function is used to easily generate a frequency distribution for a given vector. When applied to our variable `x`, which contains several repeated values (like 1, 4, and 5), the function counts these repetitions and presents the results in a concise table format. This step is crucial for identifying modes (the most frequent values) and understanding the overall concentration of observations across the measured spectrum. This tabulation helps confirm or refine hypotheses about the distribution shape suggested by the mean and median comparison.

Executing the `table(x)` command provides the exact count for every unique value present in the vector `x`. The output shows the values listed in the top row and their corresponding frequencies (counts) in the bottom row. Observing these counts allows the analyst to see immediately where the data clusters. For instance, the value 4 has the highest frequency (3 times), indicating it is the mode of this specific sample, whereas many higher values, such as 7.4, 8, 13, and 14.2, occur only once.

```
#produce frequency table
```

```
table(x)
```

```
1 2 3.5 4 5 6.5 7 7.4 8 13 14.2
```

```
2 1 1 3 2 1 1 1 1 1 1
```

This detailed breakdown tells us explicitly about the composition of the sample:

The value **1** occurs 2 times, showing a cluster at the low end of the range.

The value **2** occurs 1 time.

The value **3.5** occurs 1 time.

The value **4** occurs 3 times, making it the most frequent observation.

The value **5** occurs 2 times.

And so forth, providing counts for all eleven unique values within the fifteen-observation sample. This confirms that the variable is not uniformly distributed and exhibits definite points of concentration, which we will further explore visually.

## Using Boxplots for Visualizing Data Spread and Outliers

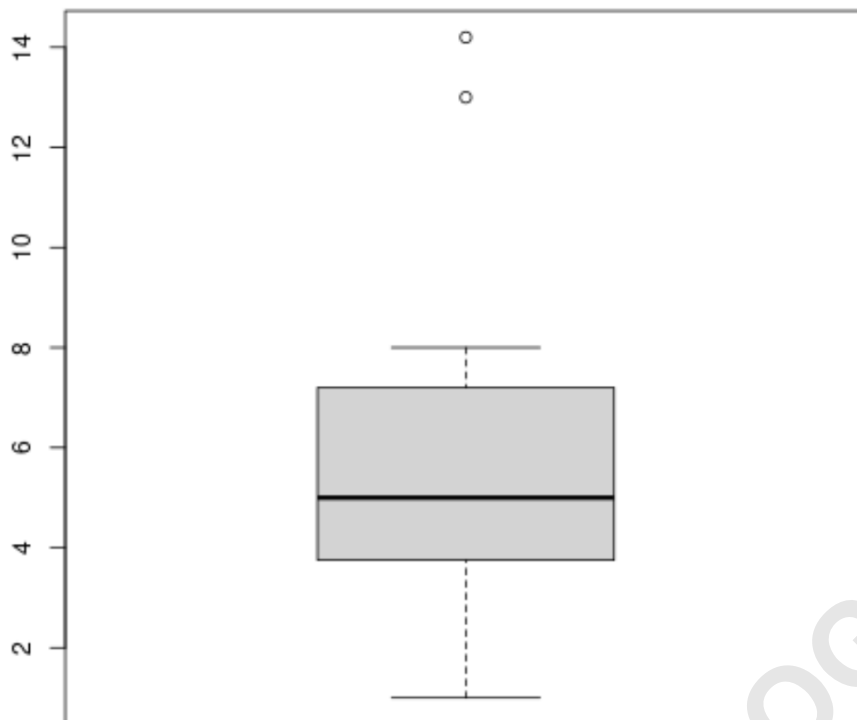
While numerical statistics are precise, they cannot replace the power of visualization in univariate analysis. Charts offer an immediate, intuitive understanding of the data's distribution characteristics. One of the most effective visualizations for summarizing key numerical properties is the **boxplot**, or box-and-whisker plot. This chart succinctly displays the five-number summary: minimum, first quartile (Q1), median (Q2), third quartile (Q3), and maximum.

The boxplot provides excellent insight into the variable's central location, skewness, and the presence of outliers. The central box represents the Interquartile Range (IQR), encompassing the middle 50% of the data. The line inside the box marks the median, giving a clear visual reference for the center. The whiskers extend from the box to the minimum and maximum values that are not considered outliers (typically within 1.5 times the IQR distance from the quartiles). If points fall outside these whiskers, they are plotted individually, flagging them as potential outliers for further investigation.

To generate the boxplot in R for our variable `x`, we use the simple `boxplot()` function. This function automatically calculates the necessary quartiles and determines the appropriate range for the whiskers based on standard statistical conventions:

```
#produce boxplot  
boxplot(x)
```

The resulting plot visually confirms the distribution characteristics identified numerically. For example, if the median line is closer to Q1 than to Q3, it suggests a positive (right) skew. If the whiskers are unequal in length, it further supports the finding of asymmetry. Reviewing the generated image allows for a quick assessment of the data's stability and dispersion across its range.



In the boxplot above, the structure confirms the slight positive skew noted earlier (mean > median). The plot visually emphasizes the spread of the data, particularly how the top 25% (Q3 to Max) occupies a wider range than the bottom 25% (Min to Q1), illustrating the influence of the higher values (13 and 14.2) on the overall distribution.

## Histograms and Density Curves for Shape Assessment

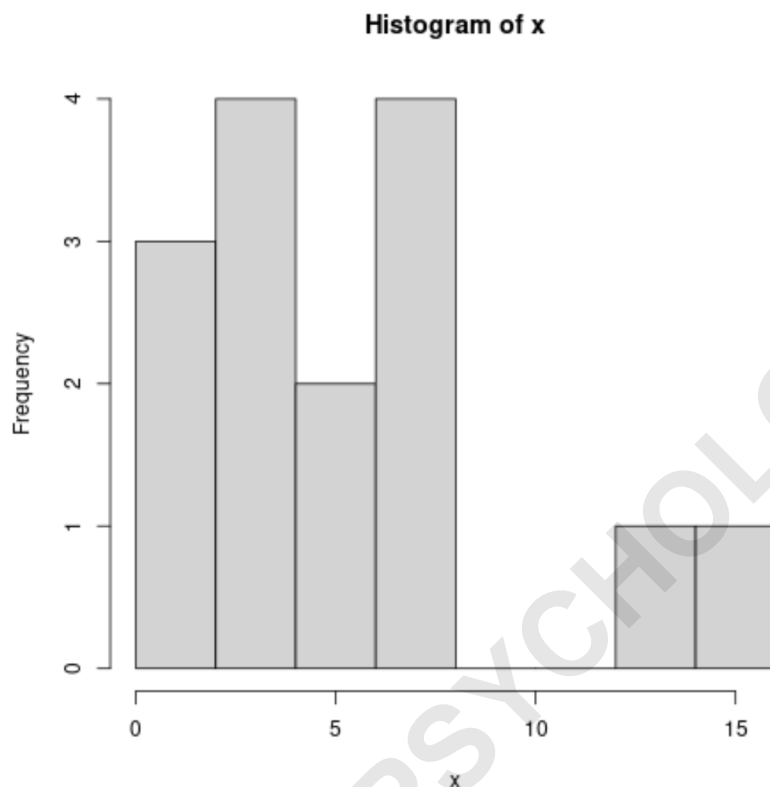
While boxplots excel at summarizing five key statistics and identifying outliers, they do not provide detailed information about the overall shape or modality of the distribution. For this purpose, the **histogram** is the visualization of choice. A histogram groups continuous data into bins (intervals) and plots the frequency (or count) of observations falling into each bin, thereby illustrating the probability distribution of the data.

The histogram is crucial for determining if the data follows a recognizable theoretical distribution, such as a normal distribution, or if it exhibits significant skewness or multiple modes. Generating a histogram in R is straightforward, using the `hist()` function on the variable `x`. R automatically selects a suitable number of bins based on the sample size and range, although this can be customized for specific analytical needs.

**#produce histogram**

**hist(x)**

The resulting histogram clearly displays the frequency of values across different intervals. If the bars are taller on the left and trail off to the right, it indicates a positive skew, consistent with our earlier numerical findings. If the distribution shows two distinct peaks, it suggests a bimodal distribution, potentially signaling that the sample combines data from two different underlying populations. For our specific variable, the histogram visually reinforces the concentration of data in the lower range (around 4-5) and the sparsity in the higher ranges.



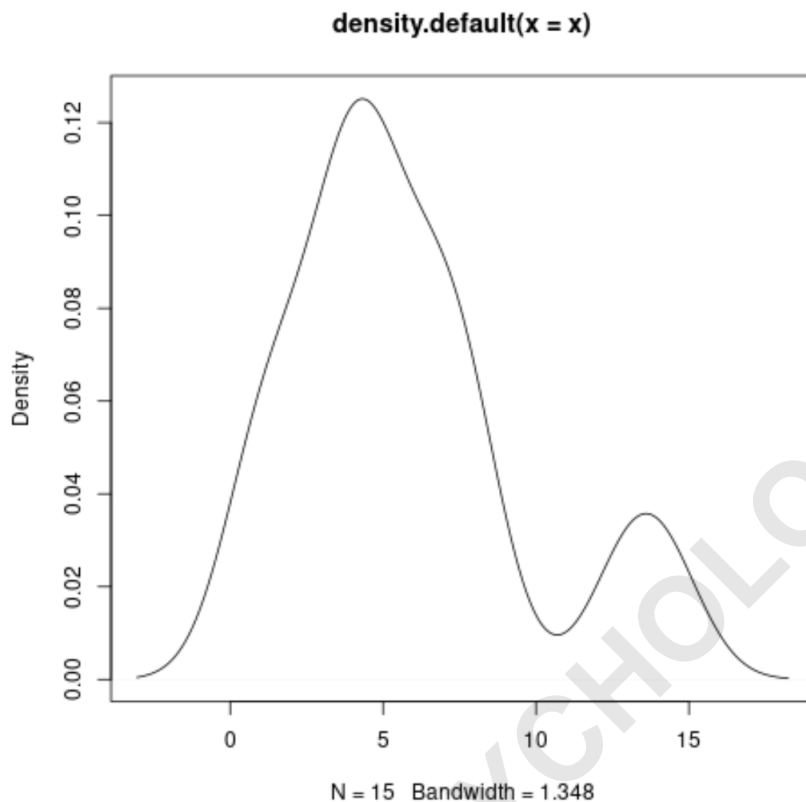
Complementing the histogram is the **density curve**, which provides a smoothed estimate of the probability density function. Unlike the blocky representation of the histogram, the density curve offers a continuous visualization of the distribution, making the detection of modality and skewness even smoother and cleaner. The area under the density curve sums to 1.0, representing the total probability.

To visualize the density curve in R, we first use the `density()` function to calculate the kernel density estimate, and then use the `plot()` function to display the result:

```
#produce density curve  
plot(density(x))
```

This curve provides a refined visual confirmation of the distribution shape. A steep rise and gentle,

extended tail on the right side indicates a strong positive skew, aligning perfectly with the boxplot and the mean-median relationship. Together, the histogram and density curve are vital tools for assessing the underlying structure of the data, which is paramount for selecting appropriate statistical modeling techniques in subsequent analyses.



## Synthesizing Univariate Findings for Robust Data Understanding

The completion of the numerical summaries, frequency tables, and visual charts provides a comprehensive understanding of our single variable,  $x$ . Each of these components contributes a unique layer of detail, ensuring that the analysis is robust. The **summary statistics** provided quantitative metrics for the typical value (mean 5.71, median 5) and the extent of variation (standard deviation 3.86). This initial numerical assessment immediately highlighted a tendency toward positive skewness.

The **frequency table** then substantiated these findings by showing the exact counts, revealing that the value 4 was the mode and that the data concentration lay predominantly below the mean. The visualizations--the **boxplot**, the **histogram**, and the **density curve**--visually confirmed the asymmetry and the high concentration of observations in the lower ranges, with a long tail extending toward the maximum value of 14.2.

In practice, this holistic approach to univariate analysis allows the analyst to make critical decisions

early on. For instance, based on the identified skewness and the notable difference between the mean and median, one might decide that non-parametric tests or data transformations (like log transformation) might be more appropriate than standard parametric tests (like t-tests) for future inferential analysis. This phase acts as a safeguard, ensuring that all assumptions are tested and that the raw data is fully understood before advanced modeling begins.

Each of these charts and numerical calculations gives us a unique, yet cohesive, way to characterize the distribution of values for our variable. For those interested in expanding their proficiency in statistical computing, you can find many more expert R tutorials focused on data science and statistical modeling on this site.

ARABPSYCHOLOGY.COM