

Learning to Import Data with the R scan() Function: A Practical Guide

Authored by
Mohammed loot

November 11, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning to Import Data with the R scan() Function: A Practical Guide*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=17075>

The capacity to efficiently import external data is an essential cornerstone of any analytical or statistical programming environment. Within the **R** language, one of the foundational input/output utilities available for reading raw data from a file into a session is the **`scan()`** function. This tool proves exceptionally valuable when researchers or developers must process simple, highly structured files, prioritizing execution speed and precise control over the resulting data type assignment--often as a prerequisite to building more complex data frame objects.

While contemporary **R** development often integrates specialized packages, such as `readr` or `data.table`, for high-volume, streamlined data ingestion, gaining a deep understanding of **`scan()`** provides critical insight into the underlying mechanisms by which **R** handles raw file input and parsing. The function operates sequentially, reading data element by element, line by line. This method grants the user meticulous control over the intended output structure, whether the goal is a simple, homogenous **vector** or a partitioned **list** of vectors corresponding to columns.

Understanding the Purpose and Flexibility of `scan()`

The **`scan()`** function operates primarily as a low-level workhorse, specifically engineered for reading sequential streams of numerical or character data that are separated by whitespace or a user-defined delimiter. Crucially, its behavior differs significantly from functions tailored for rectangular, tabular data (like `read.csv` or `read.table`). By default, **`scan()`** aggregates all parsed input elements into a single, cohesive **vector**. To deviate from this default and handle input as multiple fields or columns, the user must explicitly define the expected output structure utilizing the fundamental `what` argument.

This inherent flexibility renders **`scan()`** highly effective and efficient when the task involves importing large volumes of homogeneous data. For instance, if a file contains only a long, uninterrupted series of temperature readings, stock prices, or physical measurements, **`scan()`** excels at rapid ingestion. However, the function's true power is unlocked when processing structured data by carefully defining the specific data type expected for each column. This guidance ensures data integrity is maintained throughout the parsing process, instructing **R** how to interpret ambiguous text strings.

The function is robust enough to handle various data sources, capable of reading directly from a named file path on the disk or dynamically processing raw data input streamed from the standard input (i.e., the console). This fundamental, direct approach highlights **`scan()`**'s utility in situations requiring rapid scripting, prototyping, or when complex **R** data frame creation is not the immediate analytical requirement.

Core Syntax and Essential Parameters of `scan()`

To leverage the full power of `scan()`, developers must master its core arguments, particularly `file` and `what`. The foundational syntax structure that guides the function's behavior is deceptively simple, yet it allows for powerful customization:

```
scan(file = "", what = double(), ...)
```

The two primary parameters govern the data source and the format that **R** is instructed to anticipate for the data elements:

file: This argument serves to specify the absolute or relative path and name of the source file from which the data is to be retrieved. If this argument is intentionally left as an empty string (`""`), the function defaults to reading from the standard input, making it useful for interactive data entry directly within the **R** console.

what: This is arguably the most critical argument, as it explicitly defines the data type or types expected in the input stream. If the user provides a single atomic type (for instance, `double()` for numerical data or `character()` for text strings), `scan()` will read all available data and combine it into a single output **vector** of that specified type. Conversely, if a **list** of empty or typed structures (e.g., `list(character(), double(), double())`) is supplied, the function interprets the input file as structured with corresponding columns and returns the output as a **list**, where each element represents one parsed column from the source file.

The ellipses (`...`) denote a variety of secondary, but often necessary, arguments documented in the comprehensive `scan()` manual. These include `sep`, which is vital for specifying the field delimiter when the data uses characters other than standard whitespace (e.g., commas or tabs); `skip`, used to ignore non-data lines such as header rows; and `nmax`, which allows the user to limit the total number of data items read, useful for testing large files.

Practical Demonstration: Importing Delimited Data

To fully appreciate the practical capabilities of the `scan()` function, we will examine a highly common data task: reading structured data from a simple delimited file. Specifically, we will use a **CSV** (Comma-Separated Values) file. While higher-level functions are generally preferred for **CSV** processing, `scan()` offers a precise, low-level method for handling this input format.

Imagine a scenario where we have a file named `data.csv` containing aggregated performance data for several hypothetical teams. The file is located at the path `C:\Users\Bob\Desktop\data.csv` and holds the following structure:

```
team, points, assists
```

```
'A', 78, 12
```

```
'B', 85, 20
```

'C', 93, 23

'D', 90, 8

'E', 91, 14

Since this source file contains three distinct columns--one character column (team identifier) and two numeric columns (points and assists)--we must explicitly instruct [R](#) to read three separate fields. We achieve this by providing the `what` argument as a [list](#) containing three placeholder elements (in this case, empty strings `" "`). Using a [list](#) informs [scan\(\)](#) to expect multiple columns. While using empty strings allows [R](#) to attempt type inference, defining the types explicitly (e.g., `list(character(), double(), double())`) is generally the safest practice for robust scripts.

Beyond specifying the column structure, we must also configure [scan\(\)](#) to ignore the header line (using `skip=1`) and correctly identify the comma as the primary field delimiter (using `sep=","`). The following [R](#) code demonstrates the precise implementation required to read this structured data and capture it within an [R](#) object:

Read data.csv into a list, skipping the header line

```
data <- scan("C:UsersBobDesktopdata.csv", what = list("", "", ""), sep = ",", skip = 1)
```

View the resulting list structure

```
data
```

```
]
```

```
"team" "A" "B" "C" "D" "E"
```

```
]
```

```
"points" "78" "85" "93" "90" "91"
```

```
]
```

```
"assists" "12" "20" "23" "8" "14"
```

Analyzing the Output Structure and Data Transformation

Upon successful execution of the code above, the [scan\(\)](#) function correctly parses the delimited file contents. Because a [list](#) was supplied to the `what` argument, the resulting output object, named `data`, is consequently structured as a [list](#). Within this object, each individual element of the [list](#) corresponds directly to one field or column imported from the original [CSV](#) file.

Examining the structure, we observe that the first element (labeled `]`) contains all values sourced from the first column (Team), the second element (`]`) holds the values from the second column

(Points), and so forth. This output format is fundamentally a collection of individual [vector](#)s, where each [vector](#) is stored as a distinct element nested within the broader [list](#) structure. It is important to note that since we used `list("", "", "")` with `skip=1`, the values in the resulting vectors represent the data rows from the file, excluding the initial header row.

If the ultimate objective is to transform this raw input into a standard [R](#) data frame, the resulting [list](#) can be seamlessly converted using the built-in `as.data.frame()` function. Nevertheless, the direct output provided by [scan\(\)](#) is valuable as it grants immediate access to the raw column [vectors](#), which is often advantageous for specific statistical or programmatic functions that mandate [vector](#) input rather than a data frame structure.

Verifying Data Integrity and Class Type

A necessary and often critical step following any data import procedure in [R](#) is the verification of the resulting object's structure and class. Errors in parsing or incorrect specification of the `what` argument can result in unexpected data classes, which can subsequently cause failures in downstream analytical processes or statistical model fitting.

In the preceding example, our expectation was that the output would be a [list](#), given that we provided a [list](#) of desired data types to the `what` argument. We can definitively confirm this structural classification using the standard `class()` function, which reports the fundamental class of any [R](#) object:

```
# View class of data object
class(data)
```

```
"list"
```

The resulting output, `"list"`, successfully confirms that the [scan\(\)](#) function has read the raw data from the file and correctly organized it into the requested multiple-column structure. Had we instead specified only a single data type (for example, `what = character()`), the resulting object's class would have been `"character"`, indicating a single, concatenated [vector](#) containing all data items from the file, irrespective of their column separation.

Contextualizing Data Import in R

While [scan\(\)](#) remains indispensable for scenarios demanding precise, low-level control over data parsing, data analysts regularly utilize a broader spectrum of functions tailored to specific file formats and data complexities. Understanding where [scan\(\)](#) fits into the overall data ingestion ecosystem is key to effective workflow management. The following approaches detail alternative

methods for importing various file types into [R](#), providing a wider context for robust data preparation techniques:

For reading files where column data occupies specific, fixed character positions, the function `read.fwf()` is the tool of choice.

When importing standard rectangular files with complex delimiters and quoting rules, the workhorse functions `read.table()` or `read.csv()` offer higher-level abstraction than **`scan()`**.

For proprietary formats, such as Microsoft Excel files, specialized packages like `readxl` are necessary to handle binary file structures.