

# How to Extract Text Before a Space in R

Authored by  
**stats writer**

November 21, 2025

## RECOMMENDED CITATION

stats writer (2025). *How to Extract Text Before a Space in R*. PSYCHOLOGICAL SCALES.  
Retrieved from <https://scales.arabpsychology.com/?p=98792>

## Introduction: Mastering String Extraction in R

Extracting specific components from text strings is a fundamental task in data cleaning and analysis, particularly when working within the **R** programming environment. Data often arrives in formats where numerical values or essential identifiers are concatenated with descriptive text, requiring precise methods to isolate the desired element. When the goal is to retrieve the initial segment of a string, delimited by the very first space or whitespace character, **R** offers several robust and efficient functions.

While general functions like `strsplit()` and `substr()` are powerful tools for general string manipulation, they often require multiple steps or complex indexing to achieve the simple task of extracting content before the first space. Fortunately, specialized functions exist in both **Base R** and external packages like `stringr` that streamline this process dramatically. Choosing the right tool depends on factors such as performance requirements, code readability, and familiarity with **regular expressions**.

This guide details the two most effective strategies for extracting a string segment preceding the first whitespace: the powerful substitution capabilities of **Base R's `gsub()` function** combined with **regex**, and the streamlined, readable approach provided by the `word()` function from the **stringr package**. We will explore the mechanics behind both methods and demonstrate their practical application using a structured **data frame** example.

## Core Methods for Extracting Strings Before a Whitespace

The challenge of isolating the first word or initial value in a string is common across many domains. In **R**, this task can be approached using specialized string processing functions. The core principle involves identifying the delimiter (the space) and instructing the function to either remove everything after that delimiter or extract only the elements before it.

The following two methods are the most commonly employed techniques, offering distinct advantages regarding syntax and dependency management:

You can use the following methods to extract a string before a whitespace in **R**:

### Method 1: Leveraging Base R with `gsub()`

This approach utilizes the **Base R** function `gsub()`, which performs global substitutions based on **regular expressions**. By defining a pattern that matches the first space and everything that follows it until the end of the string, we can replace that entire matched sequence with an empty string (`" "`), effectively retaining only the prefix we desire.

```
gsub( ".*$", "", my_string)
```

The **regex pattern** `.*$` is critical here: it searches for a literal space (), followed by any character (`.`) repeated zero or more times (`*`), extending until the end of the string (`$`). By replacing this entire pattern with nothing, the function efficiently isolates the content before the first occurrence of a space.

## Method 2: Utilizing the `stringr` Package's `word()` Function

For those who prefer the readable syntax provided by the **Tidyverse** ecosystem, the **stringr package** offers a highly specialized function called `word()`. This function is designed specifically for tokenization, allowing users to extract specific words or sequences from a string based on delimiters (usually whitespace).

The `word()` function simplifies the logic significantly. Instead of requiring a complex **regular expression**, we simply tell it to extract the first word.

```
library(stringr)
```

```
word(my_string, 1)
```

By specifying the index `1` as the second argument, we instruct the **word() function** to return the substring that occurs before the first defined boundary, which defaults to whitespace. This approach is generally considered cleaner and less error-prone for simple extraction tasks compared to managing **regex** within **Base R**.

Both of these examples effectively extract the string before the first space in the hypothetical string called `my_string`. We will now apply these concepts to a real-world **data frame** structure.

## Setting Up the Sample Data Frame

To properly illustrate how these extraction methods work in a typical data analysis scenario, we will use a small sample **data frame**. This structure mimics common datasets where a column contains mixed numerical and textual data (e.g., measurements followed by units). Our objective is to isolate the numerical distance value from the unit of measurement ("miles").

The following code snippet demonstrates the creation and structure of our example **data frame**:

```
#create data frame
df <- data.frame(athlete=c('A', 'B', 'C', 'D'),
```

```
distance=c('23.2 miles', '14 miles', '5 miles', '9.3 miles'))
```

```
#view data frame
```

```
df
```

```
athlete distance
```

```
1 A 23.2 miles
```

```
2 B 14 miles
```

```
3 C 5 miles
```

```
4 D 9.3 miles
```

In this setup, the `distance` column is a character vector where the numeric value is separated from the unit string ("miles") by a single space. Our goal is to generate a new column, `distance_amount`, containing only the numerical component (e.g., "23.2", "14", "5", "9.3").

This transformation is crucial because standard mathematical operations cannot be performed directly on character vectors containing text. By isolating the numerical part, we prepare the data for subsequent analytical steps, such as conversion to a numeric data type, aggregation, or visualization.

### Example 1: Base R Extraction Using `gsub()`

The `gsub()` function, part of **Base R**, offers a powerful, dependency-free method for complex string manipulations. When targeting the content before the first space, we must precisely define the pattern that needs to be removed.

The core logic relies on the substitution pattern `.*$`. This pattern identifies the first space in the string, along with every character following that space until the end marker (`$`). By replacing this entire sequence with an empty string (`" "`), we effectively truncate the string immediately before the unit text begins.

The following code shows how to apply this method to extract the string before the space in each entry of the `distance` column of the **data frame**:

```
#create new column that extracts string before space in distance column
```

```
df$distance_amount <- gsub(" .*$", "", df$distance)
```

```
#view updated data frame
```

```
df
```

```
athlete distance distance_amount
```

```
1 A 23.2 miles 23.2
2 B 14 miles 14
3 C 5 miles 5
4 D 9.3 miles 9.3
```

As observed in the result, the newly created column `distance_amount` successfully contains only the numerical value, detached from the descriptive unit. This method is highly efficient and relies solely on functions built into the core **R** environment, making it an excellent choice when avoiding external package dependencies is a priority.

One important consideration when using `gsub()` is the potential for complexity if the delimiter is not a simple space, or if the string structure is variable. However, for a fixed structure like `VALUE SPACE UNIT`, the `.*$` **regular expression** provides a reliable and concise solution.

**Related:**

## Example 2: String Extraction Using `stringr::word()`

The **stringr package** is part of the **Tidyverse** suite and is designed to provide a cohesive, consistent, and user-friendly set of functions for string manipulation. The `word()` function is arguably the most straightforward solution for extracting tokens based on word position.

Unlike `gsub()`, which removes unwanted parts, **word()** explicitly extracts the desired part. By setting the index argument to `1`, we specify that we want the first element before the default delimiter (which is any sequence of whitespace). This enhances code clarity, especially for those less familiar with intricate **regular expressions**.

The following code demonstrates the use of the **word() function** from the **stringr** package to achieve the exact same extraction result:

```
library(stringr)
```

```
#create new column that extracts string before space in distance column
df$distance_amount <- word(df$distance, 1)
```

```
#view updated data frame
df
```

```
athlete distance distance_amount
1 A 23.2 miles 23.2
2 B 14 miles 14
```

3 C 5 miles 5

4 D 9.3 miles 9.3

The output is identical to the `gsub()` example, confirming the efficacy of this method. The primary difference lies in the syntax and required dependencies. Using **stringr** requires loading the package, but in return, the resulting code `word(df$distance, 1)` is highly descriptive of its intent: "get the first word from the distance column."

Notice that the new column called **distance\_amount** contains the string before the space in the strings in the **distance** column of the **data frame**.

## Comparison and Best Practices

When deciding between **Base R's** `gsub()` and the **stringr package's** `word()` function, several factors should influence the choice:

**Dependency Management:** The `gsub()` method requires zero external packages, relying solely on **Base R**. If minimizing dependencies is crucial for deployment or reproducibility, `gsub()` is the preferred choice. The **stringr package** is lightweight but must be installed and loaded.

**Readability and Intent:** The `word(x, 1)` syntax clearly conveys the intent of extracting the first token, making the code easier to maintain and understand for analysts unfamiliar with **regular expressions**. The `gsub()` method requires understanding the complex interaction of the **regex pattern** `.*$`.

**Performance:** For large datasets, both methods are typically efficient, but **Base R functions** are often marginally faster as they are highly optimized C implementations. However, the performance difference is usually negligible unless processing millions of strings.

Ultimately, the choice often comes down to ecosystem preference. If the user already works within the **Tidyverse** (which includes `dplyr`, `ggplot2`, etc.), using `stringr::word()` maintains consistency. If the environment is strictly **Base R**, `gsub()` is the logical choice.

## Alternative Base R Method: The `strsplit()` Function

While `gsub()` provides a one-step substitution, another common **Base R** tool is `strsplit()`. This function splits a string vector into a list of character vectors, using a specified delimiter (the space, in this case).

Using `strsplit()` is slightly more verbose because the result is a nested list structure that requires indexing and unlisting to return the desired vector. However, it is an essential function to

understand for more complex splitting scenarios.

To extract the string before the space using `strsplit()`, you would perform the following steps:

Apply `strsplit()` to the string vector using a space (" ") as the separator.

Use `sapply()` or a similar indexing method to extract the first element (index ) from each resulting sub-vector within the list.

#### # Example of strsplit approach

```
split_list <- strsplit(df$distance, " ")  
df$distance_amount_strsplit <- sapply(split_list, "[", 1)
```

While this approach is functionally correct, it involves an intermediate list structure and an additional step (`sapply`) to flatten the result back into a vector suitable for column assignment, making `gsub()` or `stringr::word()` generally more compact and readable for this specific task.

## Summary of String Extraction Techniques

Effective string extraction in **R** relies on the ability to isolate relevant tokens based on delimiters. For the common task of extracting a string before the first space, both **Base R** and the **stringr package** provide powerful and distinct solutions.

Use `gsub()` with the **regex pattern** `.*$` for a dependency-free, high-performance solution that leverages the power of regular expression substitution.

Use `stringr::word(x, 1)` for a concise, highly readable approach that simplifies the task by abstracting away the need for explicit regular expression management.

By understanding these techniques, data professionals can efficiently clean and prepare character data in **R**, ensuring that numerical values embedded within text fields are correctly isolated for quantitative analysis.