

# How to Remove Columns with Identical Values in R Data Frames

Authored by  
**Mohammed loot**

November 13, 2025

## RECOMMENDED CITATION

Mohammed loot (2025). *How to Remove Columns with Identical Values in R Data Frames*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=24079>

## Introduction: The Necessity of Removing Constant Columns in Data Analysis

In the realm of statistical computing and data analysis using the [R programming language](#), working with large and complex [data frames](#) is standard practice. A common challenge encountered during the data preprocessing phase is identifying and eliminating columns that contain only a single, constant value across all rows. These columns, often referred to as zero-variance columns, provide no discriminatory information for modeling or analytical tasks. Their inclusion can unnecessarily increase computational overhead, complicate model interpretation, and occasionally cause errors in certain statistical procedures which rely on variance. Therefore, maintaining a clean and efficient dataset necessitates the removal of these redundant features.

The requirement to remove columns with zero variance arises frequently, especially when handling highly sparse data, resulting from wide merges, or after specific feature engineering steps. For instance, if a dataset is filtered to include only observations from a single geographic region, a 'Region' column might suddenly become constant. Keeping such a column is detrimental to efficiency. We must adopt robust and efficient methods to automate this cleaning process. Fortunately, [R](#) offers several powerful approaches to achieve this goal, primarily categorized into solutions utilizing base [R](#) functionality and those leveraging the performance-optimized Tidyverse ecosystem, particularly the [dplyr](#) package.

This comprehensive guide will explore the two principal methods available for purging constant columns from your R data structure. We will meticulously detail the underlying mechanics of each approach, compare their performance characteristics, and provide practical, runnable code examples using a shared dataset. Understanding both the base R and the Tidyverse solutions ensures that data scientists can select the optimal method based on their specific environment constraints and performance requirements. Here are the most common ways to efficiently remove these invariant columns:

### Method 1: Utilizing Base R for Column Filtering

The base [R](#) environment provides highly capable functions for manipulating data structures without relying on external packages. The most elegant and widely accepted technique for removing constant columns in base R involves a combination of the [Filter\(\)](#) function and logic related to uniqueness. The core principle is straightforward: a column is considered constant if the count of its unique values is exactly one. By applying a conditional function across all columns of the [data frame](#), we can effectively filter the structure to retain only the relevant variables.

The [Filter\(\)](#) function is designed to apply a logical test--provided as a function argument--to the elements of an object, returning only those elements for which the test evaluates to `TRUE`. When applied to a [data frame](#), [Filter\(\)](#) iterates column by column. The crucial logical test we supply is

`function(x) length(unique(x)) > 1`. Here, `unique(x)` extracts all non-duplicated values from the current column (represented by `x`). Subsequently, `length()` counts how many unique values were found. If this count is greater than one, the condition is met, signifying that the column contains varying data and should be kept. If the length equals one, the column is constant and automatically excluded from the resulting data frame. This method is highly transparent and relies solely on functions built into the core R distribution.

This particular method uses only functions from base R to remove any columns in the data frame that contain the same value in each row of the column. This method works by passing a function to the **Filter()** function, which filters the columns to only display the ones that have more than one unique value in the column. While the base R approach is universally accessible and requires no additional package installations, its performance characteristics must be considered for extremely large datasets. Base R functions, particularly those involving iterative application like `Filter()` combined with `unique()` and `length()`, can sometimes be less performant than vectorized or C++ backend operations offered by specialized packages.

**#remove any columns that contain all of the same value**

**Filter(function(x) length(unique(x))>1, df)**

## Method 2: Leveraging the Tidyverse with dplyr

For users who rely on the modern Tidyverse ecosystem, the [dplyr](#) package offers a highly optimized and syntactically clean alternative for data manipulation. The Tidyverse philosophy emphasizes readable, chained operations using the pipe operator (`%>%`), which often leads to code that is both faster to execute and easier to understand. The preferred [dplyr](#) approach utilizes the `select()` function in conjunction with specialized selection helpers to filter columns based on their contents.

The key to the [dplyr](#) solution is the combination of `select()` and `where()`, along with the distinct count function `n_distinct()`. The `select()` function is used for column selection, and `where()` allows us to apply a predicate (a function returning a single logical value) to each column. The predicate we use is `~n_distinct(.) > 1`. The tilde (`~`) denotes a shorthand anonymous function, where `.` represents the current column [vector](#) being evaluated. The `n_distinct()` function, highly optimized within [dplyr](#), calculates the number of unique values in the column. If this count is greater than one, the column is selected and retained. This method works by selecting only the columns in the data frame where the number of distinct values in the column is greater than 1.

A significant advantage of the [dplyr](#) method is its performance, particularly when dealing with large datasets. Functions within the Tidyverse are often implemented using C++ backends (via the `Rcpp` package), which dramatically accelerates operations compared to pure [R](#) implementations. This

performance gain, coupled with the highly intuitive syntax that fits seamlessly into existing Tidyverse workflows, makes it the preferred method for many professional data analysts. Note that for extremely large [data frames](#), the [dplyr](#) approach is likely to be faster. Both of these methods will produce the same result, but the speed difference can be substantial.

### **library(dplyr)**

```
#remove any columns that contain all of the same value  
df %>% select(where(~n_distinct(.) > 1))
```

## **Setting Up the Example Data Frame**

To demonstrate both methods practically, we will construct a sample [data frame](#) representing statistics for various basketball players. This dataset is intentionally designed to include one column--the 'steals' column--where every row holds the exact same value (2). This simulates a common scenario where a variable lacks variation, making it a prime candidate for removal during preprocessing. This setup ensures we can clearly observe which columns are retained and which are successfully eliminated by the filtering techniques described previously. The following examples show how to use each of these methods in practice with the data frame below.

The created data frame, named `df`, includes various numerical and categorical variables. Specifically, the columns 'team', 'position', 'points', and 'assists' all exhibit variability, meaning they contain multiple unique values. The 'steals' column, however, is invariant. Before applying the cleaning methods, viewing the data structure confirms the presence of this zero-variance column, providing a clear expectation for the output of our filtering procedures. Understanding the initial state of the data is crucial for validating the success of the applied methods.

By using this controlled example, we can visually confirm that both the base R [Filter\(\)](#) approach and the [dplyr](#) `select(where())` method yield identical, correct results, thereby confirming the functional equivalence of the two techniques in this context. While the results are identical, the choice between the two often comes down to performance considerations and personal preference regarding code style.

### **#create data frame**

```
df <- data.frame(team=c('A', 'A', 'A', 'A', 'B', 'B', 'B', 'B'),  
position=c('G', 'G', 'F', 'F', 'G', 'G', 'F', 'F'),  
points=c(22, 28, 31, 35, 34, 45, 28, 31),  
steals=c(2, 2, 2, 2, 2, 2, 2, 2),  
assists=c(8, 10, 12, 12, 8, 4, 3, 9))
```

```
#view data frame
```

```
df
```

```
team position points steals assists
```

```
1 A G 22 2 8
```

```
2 A G 28 2 10
```

```
3 A F 31 2 12
```

```
4 A F 35 2 12
```

```
5 B G 34 2 8
```

```
6 B G 45 2 4
```

```
7 B F 28 2 3
```

```
8 B F 31 2 9
```

## Example 1: Remove Columns with Same Value Using Base R

Applying the base R method provides a quick and robust way to clean the dataset without external dependencies. We pass our data frame `df` and the uniqueness check function into the `Filter()` function. As established, this function evaluates each column independently using the following syntax in base R:

```
#remove any columns that contain all of the same value
```

```
Filter(function(x) length(unique(x))>1, df)
```

```
team position points assists
```

```
1 A G 22 8
```

```
2 A G 28 10
```

```
3 A F 31 12
```

```
4 A F 35 12
```

```
5 B G 34 8
```

```
6 B G 45 4
```

```
7 B F 28 3
```

```
8 B F 31 9
```

For the 'team', 'position', 'points', and 'assists' columns, the number of unique values is clearly greater than one. However, when the function evaluates the 'steals' column, it finds only the value '2', resulting in `length(unique(steals))` being equal to 1. Since 1 is not greater than 1, the condition fails (returns `FALSE`), and the 'steals' column is excluded from the final result. Notice that this returns all columns from the data frame except for the **steals** column, which contained the value **2** in every single row of the column.

Every other column in the data frame had at least two unique values in the column, ensuring their

retention. The resulting [data frame](#) confirms that the operation was successful. The output preserves all rows and only the columns that exhibit variation. It is essential to remember that this operation does not modify the original data frame `df` in place; rather, it returns a new, cleaned data frame, which must be explicitly assigned to a variable if you wish to use it further.

## Example 2: Remove Columns with Same Value Using `dplyr`

Another way to remove columns that contain the same value in every single row is to use the robust and performance-optimized syntax from the **`dplyr`** package. Executing the Tidyverse approach requires loading the necessary package first. Once [`dplyr`](#) is loaded, the pipeline syntax provides a streamlined method for the same filtering operation:

### `library(dplyr)`

```
#remove any columns that contain all of the same value
df %>% select(where(~n_distinct(.) > 1))
```

```
team position points assists
```

```
1 A G 22 8
```

```
2 A G 28 10
```

```
3 A F 31 12
```

```
4 A F 35 12
```

```
5 B G 34 8
```

```
6 B G 45 4
```

```
7 B F 28 3
```

```
8 B F 31 9
```

Notice that this also returns all columns from the data frame except for the **`steals`** column, which contained the same value in every single row of the column. The output mirrors the result achieved by the base R method, successfully eliminating the 'steals' column. This equivalence reinforces the idea that both methods are functionally sound, providing the user with flexibility based on their comfort level with different R ecosystems.

In summary, while the base R method offers zero-dependency flexibility, the [`dplyr`](#) method provides superior performance and integrates seamlessly with other Tidyverse data manipulation tools. Choosing between them often depends on the scale of the operation and whether your project already incorporates the Tidyverse suite. For general-purpose scripting, both are excellent tools for ensuring your analytical [vectors](#) are free of constant values, thereby improving the quality and efficiency of subsequent analysis or modeling steps.

## Additional Resources for R Data Manipulation

Mastering the art of data cleaning and preparation is fundamental to effective data science in [R](#). The ability to quickly identify and manage zero-variance features is just one component of a larger skillset involving variable transformation, missing value imputation, and efficient subsetting. To further enhance your proficiency in R data manipulation, we recommend exploring tutorials covering related tasks.

We encourage readers to delve deeper into the documentation for both base R functions and the Tidyverse packages, particularly focusing on how to vectorize operations for maximum speed. Efficient column management is vital for maintaining scalable analytical workflows.

The following list provides links to helpful guides on common R data tasks:

Tutorials on optimizing R code performance.

Guides on conditional column selection based on data type.

In-depth coverage of the Tidyverse selection helpers.

By continually refining your data preparation techniques, you ensure that your analytical models are built upon the strongest possible foundation of clean, relevant data.