

How to Easily Remove Duplicate Rows in R

Authored by
stats writer

December 3, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Easily Remove Duplicate Rows in R*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=104609>

In the realm of data cleaning and preparation within the R programming environment, identifying and eliminating redundant observations is a critical step to ensure the integrity and accuracy of subsequent statistical analysis. Data duplication often arises during data merging, collection processes, or simple input errors. Handling these redundancies efficiently is paramount for reliable data science workflows.

Fortunately, R provides robust, built-in mechanisms for this task. The fundamental tool in Base R is the duplicated() function, which effectively flags rows that match previous entries exactly. This function returns a logical vector that is then leveraged for powerful subsetting operations on the original data structure.

This guide details two primary, authoritative methods for duplicate removal--using core Base R functionality and employing the powerful functions available in the dplyr package. We provide clear, practical examples illustrating how to apply these techniques not only to filter entire rows but also to identify unique records based solely on specific columns, allowing users to tailor the approach to their specific data cleaning requirements.

You can utilize one of the following two standard methodologies to meticulously remove duplicate rows from a data frame in R, depending on your preference for package usage and workflow integration:

Method 1: Utilize Base R Functions

```
#remove duplicate rows across entire data frame
```

```
df
```

```
#remove duplicate rows across specific columns of data frame
```

```
df[, ]
```

Method 2: Employ the dplyr Package (Tidyverse)

```
#remove duplicate rows across entire data frame
```

```
df %>%
```

```
distinct(.keep_all = TRUE)
```

```
#remove duplicate rows across specific columns of data frame
```

```
df %>%
```

```
distinct(var1, .keep_all = TRUE)
```

The subsequent sections delve into these methods with high detail and practical demonstrations,

starting with the setup of a sample data frame that contains identifiable duplicates.

Understanding Data Duplication in Data Analysis

Data quality is the foundation of reliable statistical inference and machine learning model training. The presence of duplicate records--where one or more rows are identical across all or a key subset of columns--can severely bias descriptive statistics, distort correlations, and lead to misleading predictive results. For instance, if a survey respondent is recorded twice, their responses are inappropriately double-counted, artificially inflating the weight of that specific observation.

Before attempting removal, it is critical to define what constitutes a duplicate. A duplicate might be an exact match across the entire row, meaning every single variable holds the same value. Alternatively, in scenarios where unique identifiers are missing, analysts may define a duplicate based on matching values within a critical subset of columns (e.g., matching name and date of birth, even if an unrelated timestamp column differs). Identifying the correct scope for duplication detection is the first analytical hurdle.

R provides flexible tools to handle both definitions. The Base R approach is highly efficient for simple tasks, while the `dplyr` approach offers enhanced readability and integration within the Tidyverse ecosystem, particularly when complex grouping and filtering steps are required immediately after duplicate removal. Choosing the right method depends on the complexity of your data cleaning pipeline and existing package dependencies.

Method 1: Leveraging Base R Functions for Duplicate Removal

The Base R method centers on the intrinsic `duplicated()` function. This function scans a vector or data frame and, for every element or row, checks if an identical value has occurred earlier in the sequence. If a row is the second, third, or subsequent instance of a match, `duplicated()` returns **TRUE**; otherwise, it returns **FALSE** for the first instance and for all unique rows.

The key to removal is understanding how to utilize this output. The `duplicated()` function generates a logical vector. Since we want to subset the data frame to keep only the unique (first occurrence) rows, we must logically negate the result using the exclamation point (!) operator. The expression `!duplicated(df)` returns **TRUE** for all original, unique rows and **FALSE** for every row flagged as a duplicate, effectively filtering the data frame to retain only the unique records.

Base R is also highly capable when we need to define uniqueness based on specific columns. By passing only a subset of the data frame (e.g., `df`) into the `duplicated()` function, R calculates uniqueness solely on those variables. When this result is negated and used for subsetting the original, full data frame, R retains the first row corresponding to each unique combination of the specified variables, discarding subsequent matches while preserving all other columns in the

retained row.

Method 2: Utilizing the Tidyverse Approach with dplyr

For users who adhere to the modern Tidyverse principles, the `dplyr` package offers a more syntactic and readable solution through the use of the `distinct()` function. `dplyr` emphasizes clarity and chaining operations using the pipe operator (`%>%`), making complex data manipulation sequences easier to write, read, and maintain compared to nested Base R bracket notation.

The `distinct()` function is specifically designed to eliminate duplicate rows. When called without any arguments specifying columns, it evaluates every column in the data frame and returns only rows that are entirely unique. This function seamlessly integrates into a Tidyverse workflow, allowing users to pipe data directly from preceding cleaning or transformation steps into the duplicate removal phase.

A key argument in `distinct()` is `.keep_all`. When filtering based on a subset of columns (e.g., `distinct(var1)`), the default behavior of `distinct()` is to return only the specified unique columns. However, in almost all practical data cleaning scenarios, we need to keep the entire row corresponding to the unique combination. Setting `.keep_all = TRUE` ensures that all original columns are retained alongside the unique combination identified by the function, preserving the full context of the unique record.

Setting Up the Demonstration Data Frame

To effectively illustrate both the Base R and `dplyr` methodologies, we must first establish a reproducible sample data frame containing known duplicate entries. This data frame represents a simple scenario where basketball players are grouped by team and position, allowing us to demonstrate both full-row and column-specific duplicate detection.

The data frame, named `df`, is defined below. Notice that row 1 and row 2 are identical (Team A, Guard), and row 5 and row 6 are identical (Team B, Center). Additionally, if we were only concerned with the 'team' column, we would find three duplicate occurrences of 'A' and two duplicate occurrences of 'B' after the first instance of each.

The following code snippet defines and displays the sample data used for all subsequent examples in this tutorial:

```
#define data frame  
df <- data.frame(team=c('A', 'A', 'A', 'B', 'B', 'B'),  
position=c('Guard', 'Guard', 'Forward', 'Guard', 'Center', 'Center'))
```

```
#view data frame
```

```
df

team position
1 A Guard
2 A Guard
3 A Forward
4 B Guard
5 B Center
6 B Center
```

Practical Application: Removing Duplicates Using Base R

We begin by applying the Base R technique to remove rows that are exact duplicates across all columns. This is achieved by passing the entire data frame `df` into the `duplicated()` function and then negating the resulting logical vector for subsetting.

As shown in the output below, R successfully identified and removed rows 2 (duplicate of 1) and 6 (duplicate of 5). The resulting data frame contains only the unique combinations of team and position. Crucially, the Base R method automatically preserves the first occurrence of any duplicate set.

The following R code demonstrates the removal of all fully duplicated rows from the sample data frame using the standard Base R approach:

```
#remove duplicate rows from data frame
df
```

```
team position
1 A Guard
3 A Forward
4 B Guard
5 B Center
```

Next, we demonstrate how to filter for uniqueness based only on a specified column--in this case, the 'team' column. We are instructing R to keep only the first row it encounters for each unique team identifier (A or B), regardless of the 'position' value.

To achieve this, we subset `df` to include only the 'team' column before passing it to `duplicated()`. The resulting logical vector is then applied to the full `df`, ensuring that all columns of the retained

unique rows are kept.

Observe that the output retains row 1 (the first instance of Team A) and row 4 (the first subsequent instance of Team B), discarding all other rows associated with those team names.

```
#remove rows where there are duplicates in the 'team' column  
df), ]
```

```
team position
```

```
1 A Guard
```

```
4 B Guard
```

Practical Application: Removing Duplicates Using dplyr

For those prioritizing code clarity and leverage the broader Tidyverse toolset, the `dplyr` package offers a streamlined solution. Before using its functions, we must load the package using the `library()` command. The core function here is `distinct()`, which is highly intuitive for filtering unique rows.

To remove duplicates across the entire data frame, we pipe `df` into `distinct()` and specify the argument `.keep_all = TRUE`. While `.keep_all = TRUE` is often redundant when no columns are specified (as `distinct()` defaults to analyzing all columns), it is good practice to include it for consistency, ensuring that the final output includes all original variables.

The output confirms that, like the Base R method, `distinct()` preserves the first observed unique rows (Team A Guard, Team A Forward, Team B Guard, Team B Center), eliminating rows 2 and 6. The `distinct()` approach is often preferred for its clear, verb-based syntax.

```
library(dplyr)
```

```
#remove duplicate rows from data frame
```

```
df %>%
```

```
distinct(.keep_all = TRUE)
```

```
team position
```

```
1 A Guard
```

```
2 A Forward
```

```
3 B Guard
```

```
4 B Center
```

It is important to understand the role of the `.keep_all` argument. When `distinct()` is used without specifying grouping variables, it assesses all columns for uniqueness. However, when we specify one or more columns as criteria for uniqueness, `.keep_all` is essential; it instructs `dplyr` to retain the complete row associated with the unique combination identified, preventing the resulting data frame from being truncated to only the grouping columns.

Finally, we apply `distinct()` to identify unique rows based only on the 'team' column, mirroring the second Base R example. We pass the column name 'team' directly into the `distinct()` function, followed by `.keep_all = TRUE` to ensure the 'position' column is also maintained in the output.

This approach yields the same logical outcome as the Base R method for column-specific filtering: it selects the first instance of Team A (Row 1) and the first instance of Team B (Row 4), discarding the rest. The Tidyverse syntax provides a clear linguistic structure: "take the data frame, then find rows distinct by team, and keep all columns."

library(dplyr)

```
#remove duplicate rows from data frame
df %>%
distinct(team, .keep_all = TRUE)
```

```
team position
1 A Guard
2 B Guard
```

Key Considerations and Best Practices

When dealing with duplicate rows, a critical, often unspoken rule is the handling of tie-breaking. Both the Base R `duplicated()` function (when used with negation) and the `distinct()` function adhere to the principle of "first occurrence wins." This means that among a set of identical rows, the row appearing first in the original data frame is the one that is preserved, and all subsequent matching rows are removed. If the order of your data frame is arbitrary, this behavior is acceptable. However, if your data frame is ordered by a variable like 'date' or 'timestamp', you may need to sort the data frame strategically **before** running the duplicate removal process to ensure the record you wish to keep (e.g., the most recent entry) is always the one preserved.

Another important practical consideration arises when working with numerical data, particularly those stored as floating-point numbers. Due to the way computers handle floating-point arithmetic, two numbers that appear identical might be stored with minute differences, causing them to fail an exact comparison test by `duplicated()` or `distinct()`. While our example used simple character strings, analysts dealing with precision issues in numerical data should consider rounding the

relevant columns to a consistent number of decimal places before attempting duplicate removal to ensure accurate matching.

In conclusion, both Base R and `dplyr` provide robust tools for data deduplication in R. Choosing between them often comes down to workflow preference. Base R offers minimal dependency overhead and high performance for basic tasks, relying on powerful indexing and logical vector manipulation. Conversely, `distinct()` provides superior readability, especially when chaining multiple data cleaning steps, making it the preferred method for many modern data scientists working within the Tidyverse framework. Mastering both ensures complete versatility in any data cleaning scenario.

The following tutorials explain how to perform other common functions in R: