

Learning dplyr: How to Ungroup Data After Using group_by()

Authored by
Mohammed loot

October 27, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning dplyr: How to Ungroup Data After Using group_by()*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=4467>

Introduction: Managing State in Data Manipulation with ungroup()

Data manipulation is a cornerstone of modern analysis, and few tools are as central to this process in the [R](#) ecosystem as the [dplyr](#) package. The power of [dplyr](#) lies in its ability to perform highly efficient, readable operations on datasets. However, when analysts use functions like [group_by\(\)](#) to calculate metrics based on categorical variables, they fundamentally change the state of the resulting [data frame](#). This change is often referred to as "grouping" the data, meaning subsequent operations are executed group-wise rather than on the entire dataset. While grouping is essential for aggregate statistics, leaving a [data frame](#) in this grouped state can lead to significant and often subtle errors in later analytical steps, making the use of the [ungroup\(\)](#) function absolutely critical for maintaining data integrity and predicting function behavior.

The primary purpose of the [ungroup\(\)](#) function in [dplyr](#) is straightforward: it removes the grouping attribute that was applied by a preceding [group_by\(\)](#) operation. When a dataset is grouped, it carries metadata indicating which variables define the groups. If this metadata remains attached, any further functions that are sensitive to grouping (such as window functions, aggregations, or even some joining operations) will still operate within those previously defined boundaries. This behavior is typically desirable immediately after grouping, but it becomes a silent trap when moving to subsequent, unrelated analytical tasks where the entire dataset should be treated as a single entity again. Therefore, disciplined data wrangling mandates the explicit use of [ungroup\(\)](#) immediately after the group-dependent calculations are completed, thereby restoring the dataset to its default, ungrouped state.

Understanding when and why to use [ungroup\(\)](#) involves recognizing the lifecycle of a grouped operation. The analyst first defines the groups, then performs a calculation (usually aggregation via [summarize\(\)](#) or mutation via [mutate\(\)](#)), and finally, they must clean up the state change using [ungroup\(\)](#). This principle ensures that the output of one section of the data pipeline does not inadvertently influence the results of the next, unrelated section. Neglecting this step is a common error for those new to the [dplyr](#) framework, as the results might look correct for the current step, but the residual grouping structure can silently corrupt later calculations. The following detailed examples will illustrate precisely how and why this function is essential for robust data pipelines in [R](#).

Understanding Grouped Data Frames in R

When working within the [dplyr](#) environment, particularly when employing the [group_by\(\)](#) function, we are not just visually sorting the data; we are applying an intrinsic property to the tibble or [data frame](#) object. This property is crucial because it dictates the scope of subsequent operations. If we group a dataset by a variable like "team," [dplyr](#) internally stores this grouping information. When a function like [summarize\(\)](#) is called, it respects these groups, calculating summary statistics

(means, counts, standard deviations) separately for each unique value within the grouping variable. This approach provides immense efficiency and clarity when performing aggregate analysis, replacing cumbersome loops or split-apply-combine methods previously necessary in base [R](#).

However, the persistence of grouping attributes is where careful management becomes necessary. In modern [dplyr](#), the resulting output of a grouping operation, even after aggregation, often retains some form of grouping structure unless explicitly told otherwise. For instance, if you group by two variables (Team and Quarter) and then summarize, the resulting summarized table might still be grouped by the first variable (Team). If you then pass this partially grouped result to another function designed for whole-dataset operations, the function might still behave as if it needs to respect the residual grouping, leading to unexpected results or errors. This statefulness is a powerful feature but requires the analyst to be constantly aware of the current state of their [data frame](#), which is where the clarity provided by [ungroup\(\)](#) becomes indispensable.

To demonstrate the initial setup, we begin with a simple [data frame](#) detailing sports statistics. We will first create this data structure in [R](#) before proceeding to group and summarize it. This foundational dataset allows us to clearly observe the impact of grouping operations and the subsequent necessity of removing those attributes using [ungroup\(\)](#). The data represents individual player performance metrics across two distinct teams, 'A' and 'B', which will serve as our grouping variable throughout the examples.

Example 1: Basic Application of ungroup() after summarize()

In the most typical scenario, [ungroup\(\)](#) is used immediately following a function like [summarize\(\)](#), which reduces the number of rows in the dataset by calculating aggregate statistics for each group. We first initialize our sample [data frame](#), which contains observations for 'team', 'points', and 'assists'.

```
#create data frame
```

```
df <- data.frame(team=c('A', 'A', 'A', 'B', 'B', 'B'),  
points=c(14, 18, 22, 26, 36, 34),  
assists=c(5, 4, 4, 8, 7, 3))
```

```
#view data frame
```

```
df
```

```
team points assists
```

```
1 A 14 5
```

```
2 A 18 4
```

```
3 A 22 4
```

4 B 26 8

5 B 36 7

6 B 34 3

Now, we proceed to calculate the mean number of **points**, grouped by **team**, using a standard [dplyr](#) pipeline. Notice the inclusion of [ungroup\(\)](#) at the end of the chain. This specific placement ensures that once the aggregation is complete, the resulting summarized table is immediately stripped of any grouping attributes, making it safe for any subsequent calculations that should treat the results as independent rows. If we omitted [ungroup\(\)](#) here, and later tried to calculate a grand total or rank these two teams against each other using a window function, the grouping structure might interfere, potentially causing the function to fail or calculate the ranking within the group (which, in a two-row summarized table, would be meaningless).

library(dplyr)

```
#calculate mean of points, grouped by team
df_new <- df %>%
group_by(team) %>%
summarize(mean_points = mean(points)) %>%
ungroup()
```

```
#view results
```

```
df_new
```

```
# A tibble: 2 x 2
```

```
team mean_points
```

```
1 A 18
```

```
2 B 32
```

By executing this syntax, we successfully calculated the desired aggregated metric. Critically, because we used [summarize\(\)](#), we reduced the dataset to two rows--one for each team--and we explicitly ensured that this two-row result is now treated as a standard, ungrouped dataset thanks to the terminal [ungroup\(\)](#) call. However, this method of aggregation inherently results in the loss of all non-aggregated columns, such as the **assists** column, a common trade-off when using [summarize\(\)](#). Addressing the need to retain original row details while calculating group statistics requires a different approach, which still relies heavily on the appropriate use of the [ungroup\(\)](#) function.

The Role of mutate() in Grouped Operations

A frequent requirement in data analysis is to calculate a group mean or total and then append that calculated value back onto the original dataset, ensuring that every row retains its individual information while simultaneously carrying the group-level metric. When using [summarize\(\)](#), as seen above, all individual-level variables are dropped because the function collapses the data. To avoid this loss and retain auxiliary columns like **assists**, we leverage the [mutate\(\)](#) function instead of [summarize\(\)](#). The [mutate\(\)](#) function creates or modifies columns, and when used within a [group_by\(\)](#) chain, it calculates the new value group-wise but applies that resulting value back to every row within that group, thus preserving the original row structure.

When [mutate\(\)](#) is used after [group_by\(\)](#), it calculates the mean **points** for Team A and assigns '18' to every row belonging to Team A, and calculates the mean **points** for Team B and assigns '32' to every row belonging to Team B. This powerful combination allows us to perform group-level calculations without losing the granularity of the original data. The output dataset maintains its original row count (six rows) but now includes the new column, **mean_points**, which contains the aggregate data calculated according to the grouping variable. This technique is fundamentally important for tasks such as calculating the difference between an individual score and the group average, or identifying outliers within their respective categories.

Just like with [summarize\(\)](#), when using [mutate\(\)](#) after [group_by\(\)](#), the resulting [data frame](#) remains grouped by the specified variables. This residual grouping status is precisely why [ungroup\(\)](#) is necessary. If we were to perform a calculation immediately after the [mutate\(\)](#) step--for example, trying to calculate a grand total of points across all teams using another [summarize\(\)](#) call--the operation would still respect the 'team' grouping, yielding two separate totals (one for A and one for B) instead of the single overall total intended. By diligently adding [ungroup\(\)](#), we ensure that the dataset is ready for subsequent, whole-dataset operations, eliminating any ambiguity regarding the intended scope of future calculations.

Example 2: Retaining Columns using mutate() and ungroup()

To successfully retain the **assists** column while calculating the team mean, we replace [summarize\(\)](#) with [mutate\(\)](#), and critically, we still utilize [ungroup\(\)](#) to reset the data state immediately following the group-dependent calculation. This sequence represents the best practice for adding group statistics to row-level data.

library(dplyr)

```
#calculate mean of points, grouped by team, and add as a new column
df_new <- df %>%
group_by(team) %>%
```

```
mutate(mean_points = mean(points)) %>%
ungroup()

#view results
df_new

# A tibble: 6 x 4
  team points assists mean_points
1 A     14     5      18
2 A     18     4      18
3 A     22     4      18
4 B     26     8      32
5 B     36     7      32
6 B     34     3      32
```

This resulting [data frame](#), `df_new`, successfully retains the **assists** column and has added a new column, **mean_points**, which correctly shows the mean points value for each team across all six rows. This is precisely the desired outcome when needing to combine individual observation data with group-level summaries. The integrity of the original dataset is preserved while the necessary analytical context is added. The use of [mutate\(\)](#) ensures the row count remains unchanged, but the grouping attribute applied by [group_by\(\)](#) is still present after the mutation step if [ungroup\(\)](#) is omitted.

The inclusion of [ungroup\(\)](#) in this pipeline is not merely a formality; it is a critical safeguard. Since we used [ungroup\(\)](#), we can now perform any subsequent calculations on this modified [data frame](#) without worrying that those calculations will be silently restricted by the 'team' grouping. For example, if we wanted to calculate the overall correlation between **points** and **assists** in this new table, we would expect a single correlation coefficient derived from all six rows. If the data remained grouped, the correlation calculation might attempt to run separately within Team A and Team B, potentially yielding a warning or, worse, an incorrect result that is silently aggregated based on the residual grouping state.

This emphasizes the importance of managing state explicitly in data science workflows. When chaining many operations together, the failure to reset the data frame state using [ungroup\(\)](#) creates technical debt. This debt manifests as hidden constraints that can break pipelines later on when new steps are introduced, requiring tedious debugging to determine why a function is behaving unexpectedly. The simple, explicit act of calling [ungroup\(\)](#) guarantees that the output of this segment of the code is a clean, standard tibble, ready for the next analytical phase.

The Necessity of Explicit State Management

The core principle behind using `ungroup()` is based on the concept of explicit state management. In `dplyr`, grouping is a state change that persists until explicitly reversed. If we didn't use the `ungroup()` function, the rows of the `data frame` would still carry the grouping metadata associated with the 'team' variable. While this metadata is invisible when simply viewing the table, it becomes highly influential when subsequent functions interact with the data frame. For instance, any subsequent call to a window function (like `row_number()` or `lag()`) would operate relative to the groups, rather than relative to the entire dataset. This could have unintended and analytically unsound consequences, particularly if the next step involves calculations that are meant to span the entire population, such as calculating overall percentiles or global rankings.

Consider a scenario where the next step is to prepare the data for visualization, perhaps by calculating an overall standardization score across all observations. If the data remains grouped, the standardization might be calculated separately for Team A and Team B, resulting in two distinct standardization distributions which would be inappropriate if the goal was to compare all players equally against the global mean and standard deviation. Furthermore, if the grouped data is exported and later read back into a different environment or system that doesn't fully respect the `dplyr` grouping attributes, it might lead to unpredictable behavior or even data corruption warnings, highlighting the fragility of relying on persistent grouping attributes across different stages of a workflow.

Therefore, adopting the practice of appending `ungroup()` immediately after the group-dependent operation (whether it is `summarize()` or `mutate()`) serves as a necessary safety measure. It acts as a clear delimiter, separating the 'grouped calculation phase' from the 'ungrouped subsequent analysis phase'. This discipline not only prevents errors but also significantly improves code readability and maintainability, making the pipeline easier to audit, debug, and share with colleagues. The explicit nature of `ungroup()` ensures that the analyst is always in full control of how their data is being processed at every step of the `R` script.

Conclusion: Best Practices for Data Wrangling

The `ungroup()` function is a deceptively simple yet profoundly important component of the `dplyr` toolkit. It acts as the necessary counterweight to `group_by()`, allowing analysts to perform powerful group-wise calculations while ensuring their resulting dataset is free from hidden constraints that could compromise future steps. Whether reducing the dataset using `summarize()` or augmenting it using `mutate()`, the best practice dictates that the grouping attribute must be consciously removed once the group-specific task is complete. This attention to state management is what distinguishes robust, reliable data pipelines from those prone to difficult-to-trace errors.

By consistently incorporating [ungroup\(\)](#) into your data manipulation workflows, you ensure clarity, prevent unexpected aggregation behavior, and guarantee that your data is always ready for global operations. This commitment to cleaning up the data state is a hallmark of professional data analysis in [R](#). Always remember the mantra: group, calculate, and then immediately ungroup.

Additional Resources

For those interested in mastering advanced data wrangling techniques in [R](#), the following tutorials explain how to perform other common tasks within the [dplyr](#) framework: