

Learning Data Grouping and Summarization with dplyr in R

Authored by
Mohammed loot

November 7, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning Data Grouping and Summarization with dplyr in R*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=12459>

[Data analysis](#) thrives on clarity, and achieving that often requires transforming vast tables of raw observations into concise, actionable reports. At the heart of this transformation lie two fundamental processes: **grouping** and **summarizing** data. Grouping allows us to segment a large dataset into meaningful subsets based on shared characteristics (e.g., all cars with four cylinders), while summarizing involves calculating aggregated statistics (like the mean or count) for each of those segments. Mastering these techniques is indispensable for uncovering meaningful patterns and deriving robust insights from complex real-world datasets.

Fortunately, the [dplyr](#) package offers a world-class solution for these tasks. As a core member of the [Tidyverse](#) ecosystem within the [R programming language](#), **dplyr** provides a consistent and highly readable set of functions--often called "verbs"--that streamline data manipulation. These verbs are specifically optimized for efficiency and speed, ensuring that even large-scale data aggregation becomes manageable and intuitive.

This expert tutorial guides you through the essential steps and advanced applications of **dplyr**'s primary aggregation functions: `group_by()` and `summarize()`. By the end of this guide, you will be proficient in creating powerful, segmented statistical summaries that form the backbone of quantitative reporting.

Setting Up the Environment: Installing and Loading the dplyr Package

Before any sophisticated data manipulation can begin, the necessary tools must be properly integrated into your R environment. While **dplyr** is frequently installed as part of the broader [Tidyverse](#) meta-package, understanding how to manage it individually ensures a clean and controlled setup. The installation step is generally a one-time requirement, executing the `install.packages()` command only when the package is not yet present on your machine.

In contrast, loading the package using the `library()` function is mandatory every time you begin a new R session where you plan to use **dplyr**'s functions. This command makes the package's functions available for use in your current workspace. The following code block details the standard procedure for both installation and loading, ensuring your environment is ready for data wrangling.

It is crucial to verify that **dplyr** is loaded successfully, as all subsequent operations rely on accessing its functional toolkit. Once loaded, you gain immediate access to the expressive and fast data manipulation verbs that define the package.

If dplyr is not yet installed on your system, execute the following line:

```
install.packages('dplyr')
```

```
# Load the dplyr package into the current session
```

```
library(dplyr)
```

Understanding the Data and Core Syntax

To provide a concrete, reproducible example, we will employ the well-known R built-in dataset, **mtcars**. This dataset is a standard resource for statistical demonstrations, containing information on 32 different car models across 11 variables. Key variables include fuel efficiency (`mpg`, measured in miles per gallon), the number of engine cylinders (`cyl`), gross horsepower (`hp`), and vehicle weight (`wt`). The goal of our analysis will be to summarize fuel efficiency based on cylinder count.

Before proceeding with any aggregation, it is standard practice to inspect the dataset's dimensions and structure. This confirms the data's integrity and helps analysts understand the categorical variables available for grouping. Below, we confirm the size of the dataset and preview its initial rows to familiarize ourselves with its columnar structure.

Obtain the total number of rows (observations) and columns (variables) of mtcars

```
dim(mtcars)
```

```
32 11
```

View the first six rows of the mtcars dataset to understand its structure

```
head(mtcars)
```

```
mpg cyl disp hp drat wt  qsec vs am gear carb
Mazda RX4 21.0 6 160 110 3.90 2.620 16.46 0 1 4 4
Mazda RX4 Wag 21.0 6 160 110 3.90 2.875 17.02 0 1 4 4
Datsun 710 22.8 4 108 93 3.85 2.320 18.61 1 1 4 1
Hornet 4 Drive 21.4 6 258 110 3.08 3.215 19.44 1 0 3 1
Hornet Sportabout 18.7 8 360 175 3.15 3.440 17.02 0 0 3 2
Valiant 18.1 6 225 105 2.76 3.460 20.22 1 0 3 1
```

The core functional principle of **dplyr** is its seamless use of the [pipe operator](#) (`%>%`). This operator dramatically improves code readability by allowing data operations to be chained together sequentially. It takes the output of the preceding function (or verb) and automatically passes it as the primary input to the next function. For complex grouping and summarization tasks, the workflow is typically structured around three key steps, forming a clear data pipeline:

The initial **data frame** serves as the starting input.

The `group_by()` function segments the data frame based on the specified categorical column(s).

The `summarize()` function calculates and returns the desired aggregate statistics for each newly formed group.

This chained approach yields the following standard syntax, which is the foundational blueprint for nearly all data aggregation tasks in **dplyr**:

```
data %>%  
group_by(col_name) %>%  
summarize(summary_name = summary_function)
```

Note: While both `summarize()` (U.S. spelling) and `summarise()` (U.K. spelling) are functionally identical in **dplyr**, this tutorial will consistently use the U.S. spelling.

Statistical Summarization I: Measures of Central Tendency

One of the most frequently requested tasks in quantitative analysis is determining the typical or central value of a distribution. These [measures of central tendency](#)--primarily the mean and the median--provide a singular reference point that aims to characterize the center of the data spread. When these measures are calculated across different groups, they immediately highlight how the typical value shifts based on the categorical variable.

For our **mtcars** example, we calculate both the arithmetic average (mean) and the middle value (median) of miles per gallon (`mpg`), grouped by the number of cylinders (`cyl`). This direct comparison instantly quantifies the fuel efficiency trade-offs between, for example, a 4-cylinder vehicle and an 8-cylinder vehicle.

A crucial technical detail is the inclusion of the argument `na.rm = TRUE` within summary functions like `mean()` or `median()`. This instruction ensures that the function safely handles any [missing values](#) (represented by NA) by systematically removing them before calculating the aggregate statistic. Failing to include this step, especially in real-world data, would often result in the function returning NA for the entire group, rendering the analysis inconclusive.

Find the mean mpg by cylinder group

```
mtcars %>%  
group_by(cyl) %>%  
summarize(mean_mpg = mean(mpg, na.rm = TRUE))
```

```
# A tibble: 3 x 2
```

```
cyl mean_mpg
```

```
1 4 26.7
```

```
2 6 19.7
```

```
3 8 15.1
```

```
# Find the median mpg by cylinder group
mtcars %>%
group_by(cyl) %>%
summarize(median_mpg = median(mpg, na.rm = TRUE))

# A tibble: 3 x 2
  cyl median_mpg
  <dbl> <dbl>
1     4     26
2     6    19.7
3     8    15.2
```

Analyzing the output reveals the subtle yet important distinction between the two measures. The **mean** (arithmetic average) is highly susceptible to the influence of extreme values (outliers) or skewed data distributions. Conversely, the **median**, which represents the 50th percentile, is considered a more robust measure of central tendency because it is less affected by these data anomalies. The difference between the mean (26.7) and median (26.0) for 4-cylinder cars suggests a slightly non-symmetrical distribution in that group.

Statistical Summarization II: Assessing Data Dispersion

While knowing the center of the data is essential, it provides no information about the internal consistency or variability within that group. To fully characterize a distribution, analysts must also rely on [measures of dispersion](#) (or spread). High dispersion indicates that the observations are widely scattered around the mean, implying low consistency. Low dispersion, conversely, suggests that the values are tightly clustered, indicating high reliability in the average measure.

dplyr facilitates the calculation of several key measures of spread. These include the **standard deviation** (`sd()`), which measures the average distance from the mean; the **interquartile range** (`IQR()`), which captures the spread of the middle 50% of the data (Q3 minus Q1); and the **median absolute deviation** (`mad()`), which is a highly resistant and robust measure of spread derived from the median. Calculating these metrics by cylinder group provides deeper insight into the consistency of fuel efficiency across engine types.

Standard Deviation (SD): This is the most common measure, indicating the typical amount of variation from the group mean. It is sensitive to outliers.

Interquartile Range (IQR): Because it excludes the top and bottom 25% of the data, the IQR is a more robust measure than the standard deviation when dealing with potential outliers.

Median Absolute Deviation (MAD): This is the most robust measure of spread, recommended for

data analysis where extreme values could significantly distort the SD.

By running the following code, we calculate all three dispersion measures simultaneously within a single, efficient `summarize()` operation, demonstrating the flexibility of the **dplyr** framework:

```
# Find standard deviation (sd), IQR, and mad by cylinder group
```

```
mtcars %>%
```

```
group_by(cyl) %>%
```

```
summarize(sd_mpg = sd(mpg, na.rm = TRUE),
```

```
iqr_mpg = IQR(mpg, na.rm = TRUE),
```

```
mad_mpg = mad(mpg, na.rm = TRUE))
```

```
# A tibble: 3 x 4
```

```
cyl sd_mpg iqr_mpg mad_mpg
```

```
1 4 4.51 7.60 6.52
```

```
2 6 1.45 2.35 1.93
```

```
3 8 2.56 1.85 1.56
```

The results clearly indicate that 6-cylinder cars have the smallest values across all three dispersion metrics (SD, IQR, and MAD). This strongly suggests that the fuel efficiency (`mpg`) of 6-cylinder cars is the most homogeneous and consistent group within the **mtcars** dataset.

Advanced Summarization Techniques: Counting and Quantiles

Beyond traditional statistical moments like central tendency and dispersion, effective data summarization often requires specialized functions for frequency checks and defining specific distribution thresholds. **dplyr** provides powerful tools for these advanced needs.

Finding Counts by Group

A fundamental first step in any grouped analysis is assessing the sample size of each group. If a group contains too few observations, any statistics derived from it may be unreliable. The `n()` function is indispensable here, providing a simple, accurate count of the rows belonging to each defined group. Furthermore, `n_distinct()` allows us to count the number of unique values for a specified column within those groups, useful for assessing variety.

Using **mtcars**, calculating the row count with `n()` reveals the frequency distribution of cars across cylinder types. Simultaneously, `n_distinct(mpg)` shows how much variation in fuel efficiency exists within each category, based on the number of unique `mpg` recordings:

Find total row count and unique row count of mpg by cylinder

```
mtcars %>%
group_by(cyl) %>%
summarize(count_mpg = n(),
u_count_mpg = n_distinct(mpg))
```

A tibble: 3 x 3

```
cyl count_mpg u_count_mpg
```

```
1 4 11 9
```

```
2 6 7 6
```

```
3 8 14 12
```

The count shows that 8-cylinder cars form the largest segment in this dataset (14 observations), while 6-cylinder cars are the smallest (7 observations). This contextual count is vital when interpreting the previously calculated means and medians.

Calculating Percentiles by Group

Percentiles and [quantiles](#) are essential for establishing performance benchmarks and understanding data thresholds. They tell us the specific value below which a certain percentage of observations fall. For instance, determining the 90th percentile reveals the performance level achieved by the top 10% of vehicles within each cylinder group.

The `quantile()` function is used for this purpose. It requires specifying the column of interest and the `probs` argument, which defines the desired quantile (expressed as a decimal between 0 and 1). Here, we calculate the 90th percentile (`probs = .9`) of `mpg` to identify the fuel efficiency benchmark for high performance in each segment:

Find the 90th percentile of mpg for each cylinder group

```
mtcars %>%
group_by(cyl) %>%
summarize(quant90 = quantile(mpg, probs = .9))
```

A tibble: 3 x 2

```
cyl quant90
```

```
1 4 32.4
```

```
2 6 21.2
```

```
3 8 18.3
```

The output provides immediate, targeted insights: 90% of all 4-cylinder cars in the dataset achieve 32.4 mpg or less, establishing a clear high-performance benchmark for that category.

Additional Resources for Data Wrangling

The utility of the [dplyr](#) package extends far beyond the paired use of `group_by()` and `summarize()`. It is intentionally designed as a comprehensive, grammar-based toolkit for nearly all aspects of data transformation and preparation.

For comprehensive documentation, detailed functional descriptions, and excellent visualization cheat sheets related to **dplyr** and the broader [Tidyverse](#), analysts are strongly encouraged to consult the official package website [here](#).

To construct robust, multi-step data processing pipelines, analysts typically combine aggregation steps with other essential **dplyr** verbs. For instance, data is often pre-processed by [filtering data frame rows](#) (using `filter()`) to isolate specific subsets, followed by grouping and summarizing, and finally concluded by [arranging rows in certain orders](#) (using `arrange()`) for final presentation. Mastering the chaining of these functions allows complex data manipulation processes to be executed with unparalleled clarity and computational efficiency.