

# Learning Data Summarization in R with the ``summarize()`` Function

Authored by  
**Mohammed loot**

November 12, 2025

## RECOMMENDED CITATION

Mohammed loot (2025). *Learning Data Summarization in R with the ``summarize()`` Function*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=23851>

The core competency of modern [data science](#) hinges upon the ability to efficiently distill vast quantities of raw data into manageable, actionable insights. Data summarization is not merely an optional step; it is the fundamental process that underpins effective [Exploratory Data Analysis \(EDA\)](#) and prepares datasets for advanced applications like machine learning. By calculating metrics such as central tendency, dispersion, and specific quantiles, practitioners gain immediate validation of assumptions, identify potential outliers, and understand the inherent shape of their data distributions. This initial statistical overview is crucial for framing subsequent analytical questions and ensuring the robustness of any predictive modeling effort.

Within the statistical computing environment provided by the [R programming language](#), the most systematic and powerful approach to data manipulation and summarization is achieved through the tools housed in the [dplyr package](#). This package, which forms a cornerstone of the larger [Tidyverse](#) ecosystem, provides a coherent grammar for data analysis. Specifically, the function **summarize()** (or its fully equivalent alias, **summarise()**) is the primary engine for reducing a data frame down to a single row of computed metrics. This function converts observational data into immediate statistical information, marking a critical transition from raw data to analytical findings.

The philosophy behind **dplyr** is predicated on readability and composability, allowing analysts to chain data transformation steps together using the distinctive [pipe operator \(`%>%`\)](#). This syntactic approach ensures that complex workflows remain clear and intuitive, describing a sequence of operations that the data undergoes. While **summarize()** is the focus here, it often works in concert with other **dplyr** verbs, most notably **group\_by()**, to perform targeted analyses across different subsets of the data.

**Important Note:** The functions **summarize()** and **summarise()** are mathematically and functionally identical within the **dplyr** environment. The choice between the two spellings rests entirely with the user's preference for American or British English dialect, ensuring maximum flexibility regardless of geographical location.

To provide practical, hands-on demonstrations of the **summarize()** function, this comprehensive tutorial will utilize the well-known built-in R dataset, **mtcars**. This dataset contains 32 observations detailing the performance and design characteristics of various automobile models. Before diving into the statistical summaries, it is essential to inspect the structure of the data we will be manipulating:

```
#view first six rows of mtcars
```

```
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
```

## Setting Up the R Environment and Loading the dplyr Package

To effectively utilize external functionalities within an [R programming language](#) session, the necessary package must first be installed onto the local system and subsequently loaded into the current working environment. While the **dplyr** package is arguably the most utilized tool for data wrangling, it is not loaded by default when R starts. Therefore, making its functions, including **summarize()** and the highly versatile pipe operator (`%>%`), available requires a specific procedure.

The initial installation step, executed via the `install.packages('dplyr')` command, is a one-time requirement per machine setup. Once the package files are safely housed on your computer, you only need to execute the `library(dplyr)` command at the beginning of any R session where **dplyr** functions are needed. This action loads the package into memory, granting immediate access to its entire suite of data manipulation verbs. Neglecting this crucial step will result in R returning errors when attempting to call functions like **summarize()**, as the environment will not recognize the commands.

It is standard practice for professional R scripts to place all necessary `library()` calls at the very top, ensuring that the necessary dependencies are met before any data processing begins. This systematic approach ensures reproducibility and clarity, making it easy for others (or your future self) to understand the package requirements for the analysis being performed.

**#install dplyr package (run this only once)**

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

**#load dplyr package (run this every session)**

```
library(dplyr)
```

## Fundamental Summarization: Calculating Metrics for a Single Variable

The simplest and most direct application of the **summarize()** function involves calculating a single aggregate metric across a single column of data. This operation is most cleanly executed using the [pipe operator](#) (`%>%`), which systematically directs the entire input data frame (in this case, **mtcars**) into the first argument of the **summarize()** function. Within the function call, the user defines a new column name--such as `mean_mpg`--and assigns it the result of an appropriate aggregation function,

like `mean(mpg)`.

Consider the need to determine the overall average fuel efficiency for all vehicles in the **mtcars** dataset. We calculate the mean miles per gallon (**mpg**) using the following piped command. Crucially, notice the inclusion of the argument `na.rm = TRUE` within the `mean()` function. This essential safety measure instructs R to safely remove or ignore any missing values (represented by `NA`) during the calculation. If missing values were present and this argument was omitted, the calculation would typically return `NA` as the result, which severely compromises the reliability of the statistical summary.

**#calculate mean of mpg column**

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

```
mean_mpg  
1 20.09062
```

The resulting output confirms that the average fuel efficiency across the observed car models is approximately **20.09** MPG. However, the true utility of **summarize()** extends far beyond simple averages. By substituting the aggregation function, analysts can rapidly calculate a wide array of [descriptive statistics](#), including measures of location (median) and spread (standard deviation, minimum, maximum). For instance, calculating the median fuel efficiency is often a robust alternative to the mean, particularly when the data distribution is known to be skewed or potentially affected by extreme outliers, as the median is less sensitive to these anomalies.

To calculate the median fuel efficiency, we simply replace the `mean()` function with `median()`, maintaining the same clean syntax:

**#calculate median of mpg column**

```
mtcars %>%  
summarize(median_mpg = median(mpg, na.rm = TRUE))
```

```
median_mpg  
1 19.2
```

This result indicates that half of the vehicles in the dataset achieve 19.2 MPG or less, and half achieve 19.2 MPG or more. Furthermore, **summarize()** integrates perfectly with R's built-in `quantile()` function, which is critical for identifying specific percentiles. Percentiles help analysts understand data spread and define performance thresholds. Below, we determine the 90th percentile for the **mpg** variable:

```
#find 90th percentile of mpg column
mtcars %>%
summarize(quant90 = quantile(mpg, probs = .9))
```

```
quant90
1 30.09
```

The outcome, **30.09**, establishes that only 10% of the cars in this specific collection demonstrate a fuel efficiency exceeding this value, thereby defining a clear high-performance boundary within the dataset.

## Efficiency Through Parallelization: Summarizing Multiple Variables

While calculating a single metric is fundamental, real-world data analysis almost always demands obtaining summary statistics for several variables concurrently. The **summarize()** function within [dplyr package](#) is engineered for this exact purpose, allowing multiple summary calculations to be executed within a single function call, separated simply by commas. This capability dramatically enhances code readability, reduces processing overhead, and is vastly more efficient than running individual aggregation commands for each variable of interest.

To illustrate this powerful feature, we will calculate the 90th percentile for three distinct performance metrics in one operation: **mpg** (miles per gallon), **qsec** (quarter mile time), and **disp** (engine displacement). This comparative approach is exceptionally valuable when an analyst needs to understand the distribution tails of several related, yet distinct, continuous variables simultaneously.

```
#find 90th percentile of multiple columns
mtcars %>%
summarize(quant90mpg = quantile(mpg, probs = .9),
quant90qsec = quantile(qsec, probs = .9),
quant90disp = quantile(disp, probs = .9)))
```

```
quant90mpg quant90qsec quant90disp
1 30.09 19.99 396
```

The resulting output is a tidy, single-row data frame containing three newly created columns, each clearly labeled according to the names defined in the **summarize()** function. This immediate, tabular summary condenses multivariate information effectively. We can interpret these high-end thresholds as follows: the 90th percentile for **mpg** is **30.09**; for **qsec**, it is **19.99**; and for **disp**, it is **396**. This structured output ensures that comparative metrics are presented cleanly side-by-side.

This structure effectively showcases the capacity of **summarize()** to handle numerous calculations simultaneously. Analysts are free to incorporate as many variables and corresponding summary functions as their analysis requires, provided the chosen aggregation function (e.g., `mean()`, `sd()`, or `quantile()`) is mathematically appropriate for the data type of the column being processed. This efficiency is paramount when dealing with large datasets where manual, variable-by-variable summarization would be tedious and error-prone.

## Advanced Analysis: Grouped Summarization using `group_by()`

While overall summary statistics provide a baseline understanding, the most insightful data analysis often requires understanding how metrics behave within distinct subcategories of the data. This necessity introduces the most potent combination in the [dplyr package](#): chaining the **group\_by()** function immediately before the **summarize()** function. This pairing implements the classical "split-apply-combine" paradigm, transforming a global calculation into a targeted, categorical comparison.

The **group\_by()** verb logically segments the entire data frame into discrete subsets, based on every unique value found within the specified categorical column(s). Once the data is grouped, any subsequent **dplyr** operation, such as **summarize()**, is applied independently to each of these subsets. For instance, we might want to calculate the average fuel efficiency (**mpg**) not for the dataset as a whole, but separately for vehicles categorized by their number of cylinders (**cyl**).

```
#find mean of mpg grouped by cyl
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
```

The output of this operation is no longer a single-row summary, but a concise table detailing the mean **mpg** for each unique cylinder configuration (4, 6, and 8). This result immediately reveals the structural differences in fuel efficiency: 4-cylinder cars average 26.7 MPG, 6-cylinder cars average 19.7 MPG, and 8-cylinder cars average 15.1 MPG. This categorized insight is significantly more profound than the overall mean of 20.09 MPG calculated previously, as it quantifies how the categorical variable (engine size) drives variation in the continuous variable (fuel efficiency).

This grouped summarization technique is indispensable for comparative analysis and hypothesis testing in [data science](#). It allows researchers to quickly and accurately assess the influence of different categorical factors on continuous outcomes. The seamless integration of `group_by()` and `summarize()` is widely considered the most powerful analytical combination within the [dplyr package](#), enabling the production of clear, structured, and highly informative analytical results with minimal code.

## Expanding Your Data Analysis Toolkit

Achieving mastery over the `summarize()` function represents a fundamental milestone in becoming proficient with R for data analysis. It provides the essential capability to condense complex data frames into readily interpretable, foundational statistics. When combined with `group_by()`, this functionality extends to support intricate segmentation and comparative analysis, offering a flexible and highly readable approach to statistical reporting.

As you continue advancing your skills in data analysis, it is highly beneficial to explore other complementary functions provided by the [Tidyverse](#) and `dplyr`. Functions such as `count()` are excellent for quickly generating frequency tables, while `mutate()` is essential for creating new, transformed variables based on existing data columns before performing any summarization. By integrating these tools, analysts can build robust and reproducible data pipelines that handle everything from raw data ingress to sophisticated statistical reporting.

The following resources offer further guidance on related R data manipulation tasks and the broader ecosystem of data science tools: