

Learning Quantiles by Group with R: A Step-by-Step Guide

Authored by
Mohammed loot

November 5, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning Quantiles by Group with R: A Step-by-Step Guide*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=10274>

The Significance of Quantiles in Data Analysis

In the expansive domain of [descriptive statistics](#), **quantiles** serve as fundamental measures for understanding data distribution. They function by dividing a ranked dataset into continuous intervals, ensuring that each interval contains an equal proportion of data points. Unlike simple summary statistics such as the mean or standard deviation, quantiles offer a robust view of the spread, shape, and central tendency of the data.

The most widely utilized quantiles are **quartiles**, which partition the data into four segments (representing the 25th, 50th, and 75th percentiles). However, the definition of quantiles is highly flexible, extending to percentiles (100 divisions), deciles (10 divisions), or any custom fractional division required by the analysis. Calculating these critical threshold values is indispensable for practical applications, including determining compensation structures, evaluating performance variability across different units, or accurately flagging potential statistical outliers.

When analyzing real-world, complex datasets, a single, overall quantile calculation is often insufficient. Data is rarely homogeneous; it typically comprises distinct subgroups (e.g., product lines, geographical regions, or departments) whose distributions differ significantly. Aggregating the data before calculation often masks critical trends specific to these subgroups. This complexity necessitates calculating [quantiles](#) conditional on specific categorical variables, a task optimally handled by the powerful data manipulation capabilities available in [R](#).

Establishing the R Environment and Dependencies

To execute efficient, grouped statistical calculations within [R](#), we rely almost exclusively on the renowned **tidyverse** collection of packages. Specifically, the [dplyr](#) package is the cornerstone of this workflow, providing a grammar for data manipulation that is both highly readable and intuitive. The [dplyr](#) philosophy encourages a sequential workflow, leveraging the pipe operator (`%>%`) to clearly chain operations from data input to final summary output.

The core strategy for grouped quantile calculation involves a precise two-step process. First, the `group_by()` function is used to segment the original dataset based on the chosen categorical variable. Second, the `summarize()` function is applied to the grouped data, performing the desired calculation--in this case, invoking R's built-in [quantile function](#)--independently on each defined subgroup. This ensures that the resulting quantile boundaries are calculated autonomously for every unique level found within the grouping variable.

Before any analysis can commence, the necessary dependencies must be available in the current R session. If the [dplyr](#) package has not been installed on your system, the command `install.packages("dplyr")` must be run first. Following installation, the package must be loaded using the `library()` function, preparing the environment for data manipulation tasks.

The Core Methodology: Grouping and Summarization in R

The standard syntax for calculating quantiles by group is remarkably straightforward, combining the data frame, the grouping directive, and the summary operation into a cohesive chain. This structure relies on three key elements from [dplyr](#) and the powerful base R [quantile function](#). This methodology is highly flexible, allowing analysts to specify multiple quantiles simultaneously within a single summary call.

The first step in defining the calculation involves creating a vector, conventionally named `q`, which stores the probability values (ranging from 0 to 1) corresponding to the desired quantiles. These values are then passed directly to the mandatory `probs` argument within the [quantile function](#). This general structure provides a robust template suitable for virtually any data analysis task requiring grouped statistics or customized distributional insights.

The following syntax block demonstrates the calculation of the 25th, 50th, and 75th [quantiles](#) for a continuous variable (`numeric_variable`), segmented by the levels of a factor or character variable (`grouping_variable`):

library(dplyr)

```
#define quantiles of interest
```

```
q = c(.25, .5, .75)
```

```
#calculate quantiles by grouping variable
```

```
df %>%
```

```
group_by(grouping_variable) %>%
```

```
summarize(quant25 = quantile(numeric_variable, probs = q),
```

```
quant50 = quantile(numeric_variable, probs = q),
```

```
quant75 = quantile(numeric_variable, probs = q))
```

The output generated by this operation is a concise, tidy data frame. Each row in the resulting frame corresponds to a unique subgroup defined by the `grouping_variable`, and the columns display the specific calculated quantiles. The following detailed example illustrates the application of this powerful structure using a practical dataset.

Practical Example 1: Calculating Standard Quartiles by Group

To solidify the concept of grouped quantile calculation, let us construct a simple, relatable dataset designed to track the historical performance (measured in "wins") of three distinct, hypothetical sports teams: A, B, and C. This scenario perfectly highlights why aggregated statistics fail: the teams likely have vastly different performance histories, and pooling their data would severely

distort the true distribution of wins for any single team.

The R code below first generates the sample data frame and then applies the three-part methodology described previously. We specifically target the standard quartiles (25%, 50%, and 75%) as they provide a foundational, immediately understandable summary of each team's win distribution.

By executing `group_by(team)`, we explicitly instruct [R](#) to partition the data before the `summarize()` function is executed. Notably, the 50th quantile, commonly known as the **median**, is a crucial metric here. Since the median is less susceptible to the influence of extreme values (outliers) than the arithmetic mean, it provides a more robust and reliable measure of the typical performance center point for meaningful inter-team comparison.

library(dplyr)

```
#create data
df <- data.frame(team=c('A', 'A', 'A', 'A', 'A', 'A', 'A', 'A',
'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
'C', 'C', 'C', 'C', 'C', 'C', 'C', 'C'),
wins=c(2, 4, 4, 5, 7, 9, 13, 13, 15, 15, 14, 13,
11, 9, 9, 8, 8, 16, 19, 21, 24, 20, 19, 18))

#view first six rows of data
head(df)

team wins
1 A 2
2 A 4
3 A 4
4 A 5
5 A 7
6 A 9

#define quantiles of interest
q = c(.25, .5, .75)

#calculate quantiles by grouping variable
df %>%
group_by(team) %>%
summarize(quant25 = quantile(wins, probs = q),
quant50 = quantile(wins, probs = q),
quant75 = quantile(wins, probs = q))
```

```
team quant25 quant50 quant75
1 A 4 6 10
2 B 9 12 14.2
3 C 17.5 19 20.2
```

The results immediately reveal profound differences in performance variability and central tendency among the teams. For example, the median number of wins (`quant50`) for Team C (19) is dramatically higher than that of Team A (6). Moreover, a powerful insight emerges when comparing the lower quartile: the 25th percentile for Team C (17.5 wins) exceeds the 75th percentile for Team A (10 wins), clearly demonstrating the significant performance gap and the necessity of this grouped analysis.

Extending the Analysis with Custom Quantile Definitions

While quartiles offer an excellent generalized summary, many analytical requirements demand a finer granularity to truly map the data distribution--perhaps requiring deciles (10% increments) or even vigintiles (5% increments). Fortunately, this level of flexibility is trivially achieved in [R](#) by simply redefining the probability vector `q`.

For instance, if the analytical goal is to divide the data into five segments (quintiles), the probability vector is defined as `q = c(.2, .4, .6, .8)`. This capability is critical for establishing specific thresholds, such as identifying the top 20% of customer spenders or defining the bottom quintile of organizational performance for targeted intervention.

Crucially, the underlying procedural structure involving `group_by()` and `summarize()` remains unchanged, ensuring the method is easily adaptable across various analytical challenges. The following code demonstrates how to calculate the 20th, 40th, 60th, and 80th [quantiles](#) for each team's wins, providing a more detailed view of their performance spread:

#define quantiles of interest

```
q = c(.2, .4, .6, .8)
```

```
#calculate quantiles by grouping variable
```

```
df %>%
```

```
group_by(team) %>%
```

```
summarize(quant20 = quantile(wins, probs = q),
```

```
quant40 = quantile(wins, probs = q),
```

```
quant60 = quantile(wins, probs = q),
```

```
quant80 = quantile(wins, probs = q))
```

```
team quant20 quant40 quant60 quant80
```

```
1 A 4 4.8 7.4 11.4
2 B 9 10.6 13.2 14.6
3 C 16.8 18.8 19.2 20.6
```

The resulting granular output uncovers distributional characteristics missed by the standard quartiles. For Team A, the difference between the 60th percentile (7.4 wins) and the 80th percentile (11.4 wins) is substantial, indicating that their top 20% of outcomes represent a significant leap in performance compared to their typical results.

Focusing on Specific Percentiles: Identifying Extreme Thresholds

In various specialized analytical contexts, the focus often narrows to a single, extreme percentile boundary. This is particularly relevant when assessing risk (e.g., the lowest 5% of asset values), quality control (e.g., failure rates above the 99th percentile), or high-end performance analysis (e.g., the threshold for the top 10%).

When calculating a singular percentile, the complexity of defining the vector `q` is eliminated. Instead, the desired probability value can be passed directly to the `probs` argument within the [quantile function](#). This approach simplifies the syntax significantly while preserving the necessary grouped structure provided by [dplyr](#).

The following example calculates the 90th percentile of wins for each team. This value represents the minimum number of wins required to be categorized within that team's top 10% of historical performance outcomes:

```
#calculate 90th percentile of wins by team
df %>%
group_by(team) %>%
summarize(quant90 = quantile(wins, probs = 0.9))
```

```
team quant90
```

```
1 A 13
2 B 15
3 C 21.9
```

These 90th percentile results confirm the stark variability in peak performance across the groups. Team C requires 21.9 wins to meet its 90th percentile, confirming its exceptional top-end capability, whereas Team A achieves its equivalent high-end benchmark at 13 wins. Such insights are crucial for setting realistic, group-specific targets.

Interpreting and Applying Grouped Quantile Results

The proficiency in calculating grouped [quantiles](#) is a core competence in modern data science, providing insights that transcend the limitations of simple averages or standard deviations. When communicating these results, it is essential to emphasize that quantiles are based solely on the rank ordering of data points, which grants them high robustness against issues such as data skewness and the presence of extreme outliers.

Key practical applications derived from this type of grouped analysis include:

Benchmarking: Comparing the 50th percentile (median) across groups to establish typical, robust performance levels, as demonstrated by the comparison between Team A's median and Team C's significantly higher median.

Identifying Dispersion and Consistency: The calculation of the difference between the 75th and 25th [quantiles](#), known as the Interquartile Range (IQR), provides a measure of data spread that is inherently resistant to outliers. A smaller IQR suggests greater consistency or lower variability within that specific subgroup.

Policy and Threshold Setting: In corporate, financial, or governmental contexts, quantiles are routinely utilized to define critical boundaries. Examples include establishing the income level corresponding to a specific poverty line (a low percentile) or setting qualification requirements for bonuses (requiring performance above an 80th percentile threshold).

By mastering this technique within [R](#), analysts ensure that their [descriptive statistics](#) are precisely tailored to the specific distributional characteristics of their subgroups, thereby yielding more accurate, contextually relevant, and actionable conclusions.