

# How to Easily Create Side-by-Side Boxplots in R

Authored by  
**stats writer**

December 4, 2025

## RECOMMENDED CITATION

stats writer (2025). *How to Easily Create Side-by-Side Boxplots in R*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=105233>

Creating comparative charts is a fundamental requirement in statistical analysis. When dealing with continuous variables grouped by categorical factors, generating side-by-side boxplots in the R programming language is both highly efficient and straightforward. This technique allows researchers to immediately discern differences in central tendency, spread, and skewness across multiple distributions.

The core methodology involves two main steps: first, preparing and grouping the data appropriately, typically using R's formula syntax; and second, invoking a specialized plotting function, such as `boxplot()` in base R or `geom_boxplot()` within the powerful `ggplot2` package. Furthermore, these visualizations can be extensively customized using various graphical parameters to highlight specific statistical differences and enhance readability for presentations and reports. Below, we provide comprehensive examples illustrating how to achieve these comparative plots, detailing the nuances of both visualization environments.

Side-by-side boxplots serve as critical tools for exploratory data analysis (EDA). They condense the five-number summary--minimum, first quartile (Q1), median, third quartile (Q3), and maximum--into a compact visual form, allowing users to quickly assess the similarities and critical differences between various statistical distributions. This rapid qualitative assessment is essential before moving on to more complex inferential testing.

This comprehensive tutorial outlines the procedural steps for generating side-by-side boxplots using two distinct environments within R: the native plotting functions available in **base R**, and the versatile grammar-of-graphics approach provided by the **ggplot2** package. For demonstration purposes, we will utilize the following sample data frame, which contains performance data grouped by specific teams:

```
#create data frame  
df <- data.frame(team=rep(c('A', 'B', 'C'), each=8),  
points=c(5, 5, 6, 6, 8, 9, 13, 15,  
11, 11, 12, 14, 15, 19, 22, 24,  
19, 23, 23, 23, 24, 26, 29, 33))
```

```
#view first 10 rows  
head(df, 10)
```

```
team points  
1 A 5  
2 A 5  
3 A 6  
4 A 6
```

5 A 8  
6 A 9  
7 A 13  
8 A 15  
9 B 11  
10 B 11

## The Importance of Side-by-Side Boxplots in Data Visualization

In the field of data analysis, effective visualization is paramount for conveying complex statistical summaries quickly. When researchers need to compare the distribution of a continuous variable (like 'points' in our example) across several discrete categories (like 'team'), side-by-side boxplots offer an unparalleled snapshot of comparative statistics. They allow viewers to instantly compare the median (the line inside the box), the interquartile range (IQR, the length of the box), and the presence of potential outliers (individual points outside the whiskers).

A significant advantage of this visual method is its robustness to dataset size and its ability to clearly highlight differences in data spread. For instance, if one team's boxplot is significantly longer than another's, it immediately signals higher variability or dispersion in their points scored. Conversely, if the medians are widely separated, it suggests a statistically meaningful difference in the average performance metrics of the groups. Understanding these visual cues helps analysts prioritize which groups require deeper statistical investigation.

While histograms and density plots show the shape of a single distribution, they become cluttered and difficult to interpret when overlaid for multiple groups. Boxplots solve this spatial challenge by arranging each group's summary vertically or horizontally on the same plane, maximizing clarity and minimizing visual noise. This makes them indispensable for comparative reporting, particularly when presenting findings to non-technical audiences.

## Setting Up the Sample Data in R

To successfully generate side-by-side boxplots in R, the input data must be correctly structured. R typically requires a **long format**, where one column holds the continuous variable (the values to be plotted), and a second column holds the categorical grouping variable (the factor by which the plot is split). In our example, the `points` column represents the continuous values, and the `team` column represents the categorical groups (A, B, and C).

The code snippet provided above uses the `data.frame()` function to construct a basic data frame named `df`. This structure is the workhorse of statistical computing in the R programming language, organizing data into rows and columns akin to a spreadsheet. The `rep()` function efficiently

creates the categorical groupings, ensuring each team label is repeated eight times, resulting in 24 observations in total. The `points` vector then provides the corresponding numerical values for each observation.

Verifying the structure with `head(df, 10)` confirms that the data is ready for plotting. The grouping variable (`team`) will dictate where the breaks occur on the axis, and the outcome variable (`points`) will determine the vertical or horizontal extent of each box. When using plotting functions in R, we often employ the formula interface (`y ~ x`), where `y` is the continuous variable (`points`) and `x` is the grouping variable (`team`).

## Creating Side-by-Side Boxplots Using Base R

The simplest and most direct way to generate these plots in R is by using the built-in `boxplot()` function, part of **base R** graphics. This function is highly efficient for quick visualizations and utilizes the standard formula syntax (`dependent variable ~ independent variable`) common throughout R's statistical modeling tools. When using the formula `df$points ~ df$team`, R automatically understands that it must create separate boxplots for the `points` data corresponding to each level defined in the `team` variable.

The following code snippet demonstrates the creation of standard **vertical** side-by-side boxplots. Notice the use of several arguments which are crucial for professional presentation:

`col`: Specifies the color fill for the boxes, improving aesthetic appeal.

`main`: Sets the title of the overall plot.

`xlab` and `ylab`: Provide clear labels for the x-axis and y-axis, respectively, ensuring the viewer understands what is being measured.

These arguments ensure that the resulting plot is not only statistically sound but also visually informative and ready for inclusion in a report.

**#create vertical side-by-side boxplots**

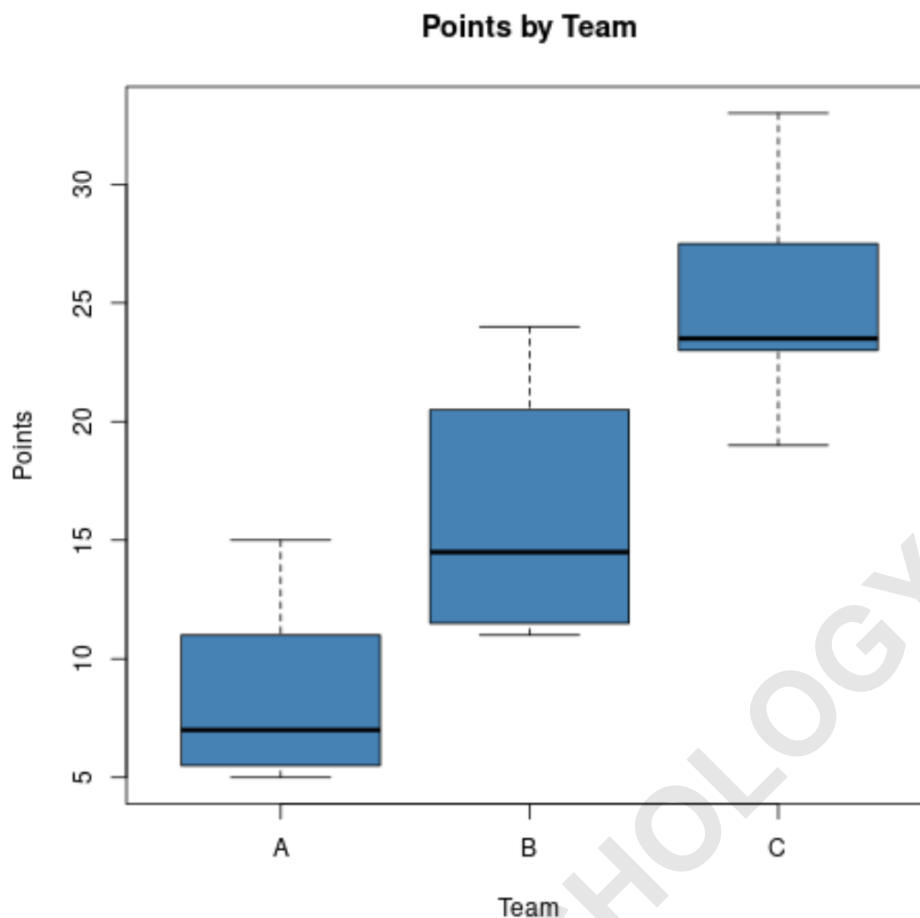
```
boxplot(df$points ~ df$team,
```

```
col='steelblue',
```

```
main='Points by Team',
```

```
xlab='Team',
```

```
ylab='Points')
```



Analyzing the resulting vertical plot reveals immediate insights into team performance: Team A displays the lowest median and the tightest distribution, while Team C exhibits the highest median and the greatest overall spread, suggesting more volatile or varied performance metrics.

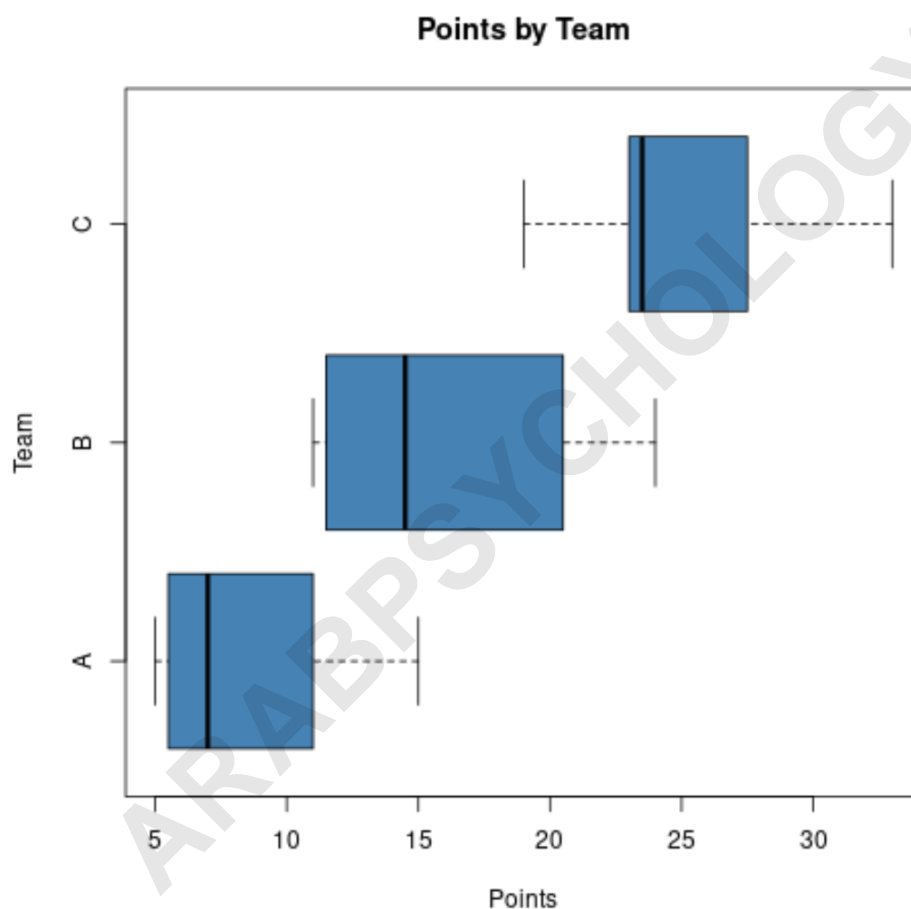
## Generating Horizontal Boxplots in Base R

While vertical boxplots are the default and perhaps most common display format, there are instances--especially when category names are long or numerous--where a **horizontal** orientation is preferable for improved readability. Base R facilitates this transformation through the use of a simple, logical argument within the `boxplot()` function: `horizontal=TRUE`. This parameter flips the axes, placing the grouping variable (Team) on the vertical axis and the continuous variable (Points) on the horizontal axis.

When switching to a horizontal orientation, it is absolutely crucial to adjust the axis labels to maintain clarity. The label that was previously assigned to `xlab` (Team) must now be conceptually linked to the vertical axis, and the label that was previously `ylab` (Points) must now describe the horizontal axis. Failing to update these labels results in a misleading and confusing chart.

The code below demonstrates this modification, ensuring that the chart title and color scheme remain consistent while providing the essential axis flip.

```
#create horizontal side-by-side boxplots  
boxplot(df$points ~ df$team,  
col='steelblue',  
main='Points by Team',  
xlab='Points',  
ylab='Team',  
horizontal=TRUE)
```



This horizontal representation confirms the same statistical findings but presents the information in a different visual hierarchy. This format is often easier for comparing the magnitude of the continuous variable across the different groups, as the length of the box (IQR) is directly comparable along the primary horizontal axis.

## Advanced Visualization with the ggplot2 Package

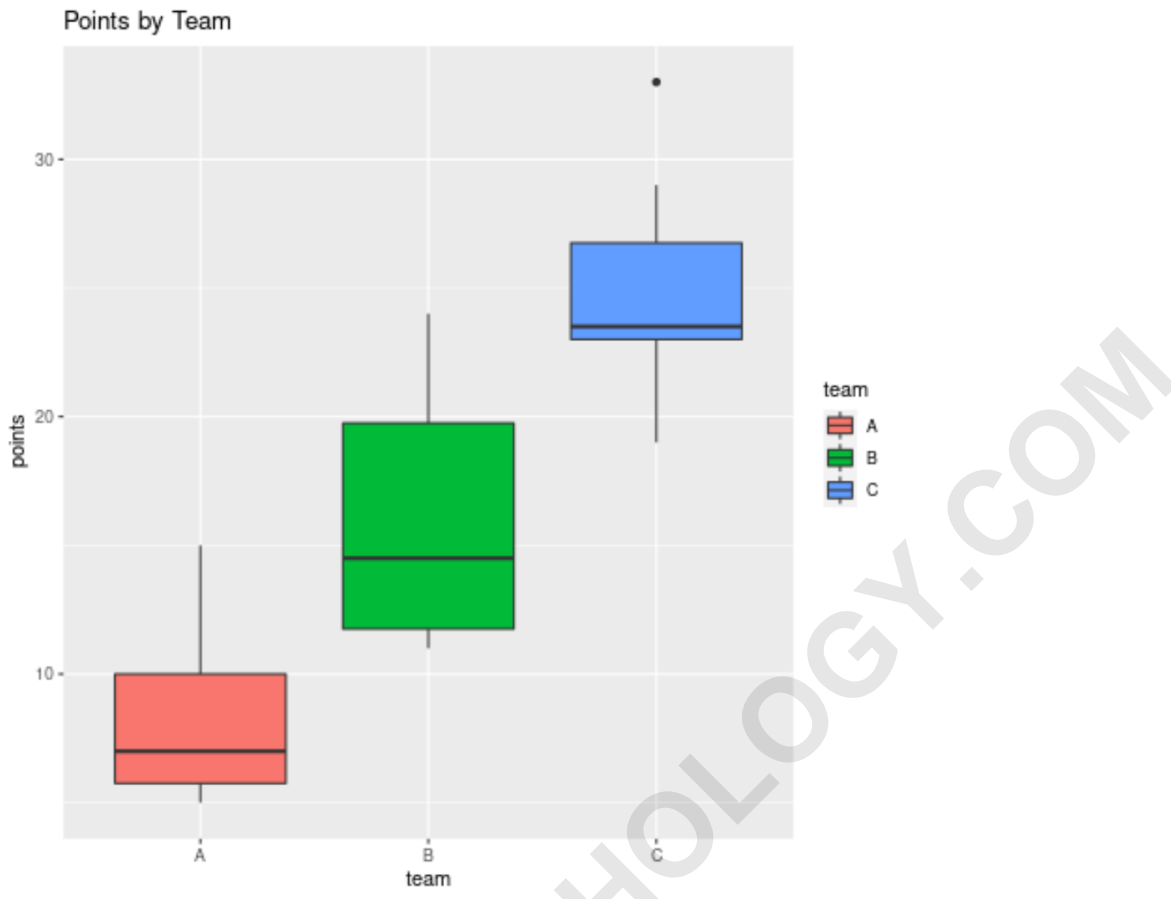
While base R offers rapid plotting capabilities, the **ggplot2** package, developed by Hadley Wickham, has become the industry standard for creating sophisticated, publication-quality graphics in R. **ggplot2** operates on the "grammar of graphics," building plots in layers: first defining the data, then mapping aesthetics (axes, color, size), and finally adding geometric objects (the boxplots themselves).

To use **ggplot2**, the package must first be loaded into the R session using the `library()` function. The core function is `ggplot()`, which takes the data frame (`df`) and defines the aesthetic mappings (`aes()`). For a boxplot, we map the categorical variable (`team`) to the x-axis and the continuous variable (`points`) to the y-axis. The actual boxes are drawn using the geometry function `geom_boxplot()`.

A major benefit of using **ggplot2** is the ease of aesthetic mapping. By assigning `fill=team` within the `aes()` call, we instruct the system to automatically color the boxes based on the team identity, enhancing the separation and visual appeal without requiring manual color specification like in base R. The plot title is added using the `ggtitle()` function.

### **library(ggplot2)**

```
#create vertical side-by-side boxplots
ggplot(df, aes(x=team, y=points, fill=team)) +
geom_boxplot() +
ggtitle('Points by Team')
```



The resulting plot offers a cleaner, more modern aesthetic compared to base R, and the automatic legend generation (not explicitly coded but included by default due to `fill=team`) clearly links the colors back to the group identities, greatly improving the quality of the [visualization](#).

### Creating Horizontal Boxplots with `ggplot2`'s `coord_flip()`

Just as base R offers the `horizontal=TRUE` argument, `ggplot2` provides a dedicated coordinate system function to achieve the same result: `coord_flip()`. This function is added as an additional layer to the existing plot structure. Unlike base R, where the user must manually swap `xlab` and `ylab`, `coord_flip()` handles the axis labels and orientation internally, simplifying the code required for transformation.

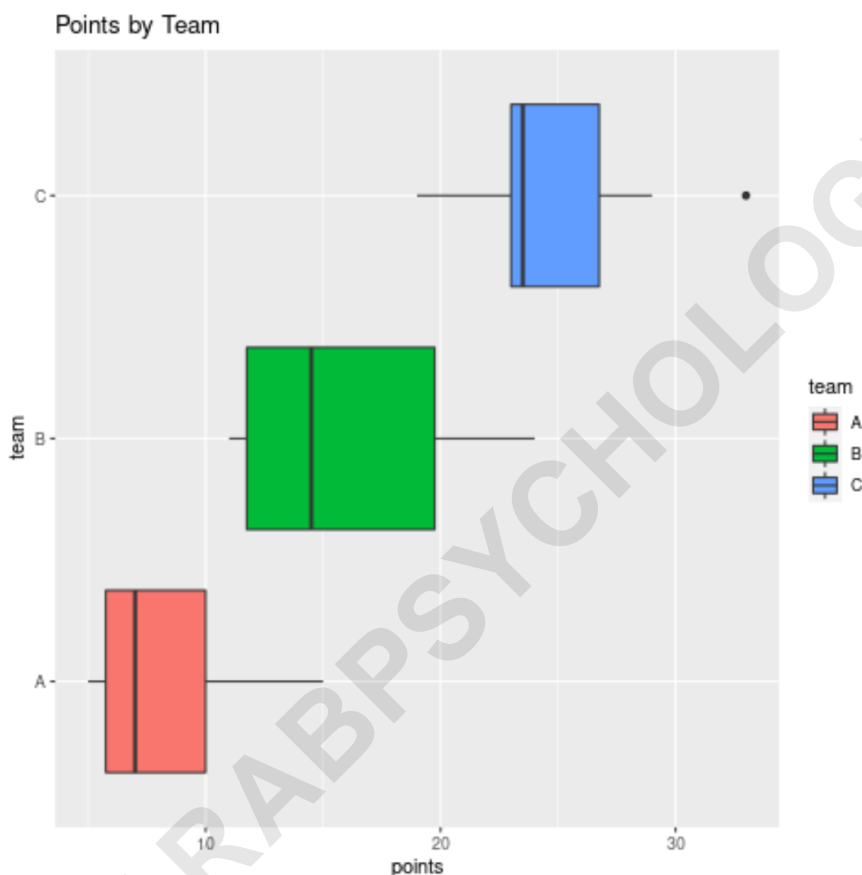
When `coord_flip()` is applied, the aesthetic mappings defined in `aes(x=team, y=points)` are internally transposed, meaning the categorical variable (`team`) shifts to the vertical axis, and the continuous variable (`points`) shifts to the horizontal axis. This preserves the interpretability of the graph while achieving the desired horizontal layout.

The following code block demonstrates how to incorporate this function into the existing `ggplot2`

structure:

```
library(ggplot2)
```

```
#create horizontal side-by-side boxplots  
ggplot(df, aes(x=team, y=points, fill=team)) +  
geom_boxplot() +  
coord_flip() +  
ggtitle('Points by Team')
```



## Conclusion and Interpretation

Generating side-by-side boxplots is a robust analytical technique for comparing distributions across groups. Whether you choose **base R** for rapid, scriptable graphics or **ggplot2** for sophisticated, customizable visuals, both environments within the R programming language offer clear pathways to visualize grouped data.

Based on our visual analysis of the "Points by Team" data, regardless of orientation, the key statistical findings are consistent:

**Central Tendency:** Team C exhibits a significantly higher median performance than Team B, which in turn performs better than Team A.

**Variability:** Team C has the largest interquartile range (IQR), indicating the highest spread in points. Team A has the smallest IQR, suggesting the most consistent performance among its members.

**Skewness and Outliers:** All teams appear relatively symmetrically distributed, and in this specific sample, no extreme outliers were plotted outside the whiskers, which typically represent 1.5 times the IQR above Q3 or below Q1.

Mastering both the base R `boxplot()` function and the multilayered approach of `ggplot2` ensures that data scientists have the flexibility to produce the appropriate visualization for any given analytical task, optimizing both speed of creation and quality of presentation.

ARABPSYCHOLOGY.COM