

When to Use stat="identity" in ggplot2 Plots

Authored by
Mohammed looti

March 23, 2026

RECOMMENDED CITATION

Mohammed looti (2026). *When to Use stat="identity" in ggplot2 Plots*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=3314>

Crafting effective data visualizations is an essential skill in modern data analysis, and [ggplot2](#) stands as the preeminent package in [R](#) for generating intricate and publication-quality graphics. Among its many capabilities, the [geom_bar\(\)](#) function is frequently utilized to construct [bar charts](#). However, this function often presents a challenge to new users because it operates in two fundamentally different modes, depending on how the underlying data is structured and how the `stat` argument is employed. Understanding this core dichotomy is critical for accurately representing your analytical findings.

The primary source of confusion stems from [geom_bar\(\)](#)'s default behavior, which automatically performs a statistical transformation on the data. This article serves as a comprehensive guide to navigating these two operational modes, focusing intently on the scenario where specifying [stat="identity"](#) is not just optional, but absolutely indispensable for plotting pre-aggregated summary data.

Method One: Leveraging Default Frequency Counting

When a user calls [geom_bar\(\)](#) and only maps a variable to the [x aesthetic](#), the function automatically engages its default statistical transformation: [stat_count\(\)](#). This powerful default behavior is specifically engineered to visualize the [frequency distribution](#) of a single categorical variable, streamlining the process of exploratory data analysis.

In this mode, [ggplot2](#) performs an internal calculation, counting the number of occurrences of each unique value present in the variable mapped to the [x-axis](#). The resulting bar heights on the y-axis then represent these calculated frequencies. This automatic aggregation is highly convenient for quickly assessing how often different categories appear within your dataset, allowing for a direct visual summary without requiring any pre-processing or summarization of the raw input data.

The syntax for leveraging this default behavior is remarkably simple, requiring only the dataset and the categorical variable of interest:

```
ggplot(df, aes(x)) +  
geom_bar()
```

Here, the visualization library processes the raw observations provided in the [x variable](#), grouping them and calculating their respective [counts](#). These [counts](#) are then plotted as the vertical heights of the bars. This approach is perfect when dealing with raw observation data, such as visualizing the number of transactions per month or the distribution of survey responses.

Method Two: Plotting Pre-Aggregated Data with `stat="identity"`

When your data has already been summarized--meaning you have pre-calculated metrics like totals, averages, or percentages--you must instruct `ggplot2` to bypass its default counting mechanism. This is achieved by explicitly setting the `stat` argument to `"identity"`. The term `"identity"` indicates that the bar heights should directly reflect the values provided in the data, without any further aggregation or transformation.

The crucial difference when using `stat="identity"` is the requirement to map two aesthetics: an `x variable` (the categories) and a `y variable` (the bar heights). The `y variable` must contain the precise numerical values--whether they represent a `sum`, mean, or other summary statistic--that you intend to visualize for each category.

The necessary syntax for this explicit control is as follows:

```
ggplot(df, aes(x, y)) +  
geom_bar(stat="identity")
```

By enforcing `stat="identity"`, you are directing the plot to treat the values in the `y variable` as the literal measure of bar height. This method is essential when working with a `data frame` that already contains summarized data, ensuring that the visual representation accurately matches your pre-calculated metrics.

The Role of Statistical Transformations in `ggplot2`

To fully appreciate the distinction between the two modes of `geom_bar()`, it is necessary to understand the concept of `statistical transformations` (or stats) within the `ggplot2` grammar of graphics. Every geometric object (geom) in the package is paired with a default stat, which defines how raw data is processed or summarized before being mapped to visual properties.

For most geoms, such as `geom_point()`, the default stat is `stat_identity()`, meaning the geom plots the data points exactly as provided in the dataset. `geom_bar()`, however, is a notable exception, defaulting to `stat_count()`. This default transformation automatically groups the data by the `x variable`, calculates the number of observations in each group, and maps these `counts` to the `y aesthetic`.

The moment you specify `stat="identity"`, you are overriding this automatic grouping and counting process. You are effectively telling `ggplot2` that the data provided in the `y variable` is already the final, aggregated metric that should be plotted directly. This transformation is crucial when visualizing metrics like pre-calculated averages, totals, or other `descriptive statistics` that do not require further computation by `ggplot2` itself.

Practical Demonstration: Understanding the Data Frame

To concretely illustrate the distinct outcomes produced by these two methods, we will utilize a simple [data frame](#) created in [R](#). This dataset models points scored by individual players across three different basketball teams, providing both raw observation data and the potential for summary analysis.

#create data frame

```
df <- data.frame(team=rep(c('A', 'B', 'C'), each=4),  
points=c(3, 5, 5, 6, 5, 7, 7, 8, 9, 9, 9, 8))
```

#view data frame

```
df
```

```
team points
```

```
1 A 3
```

```
2 A 5
```

```
3 A 5
```

```
4 A 6
```

```
5 B 5
```

```
6 B 7
```

```
7 B 7
```

```
8 B 8
```

```
9 C 9
```

```
10 C 9
```

```
11 C 9
```

```
12 C 8
```

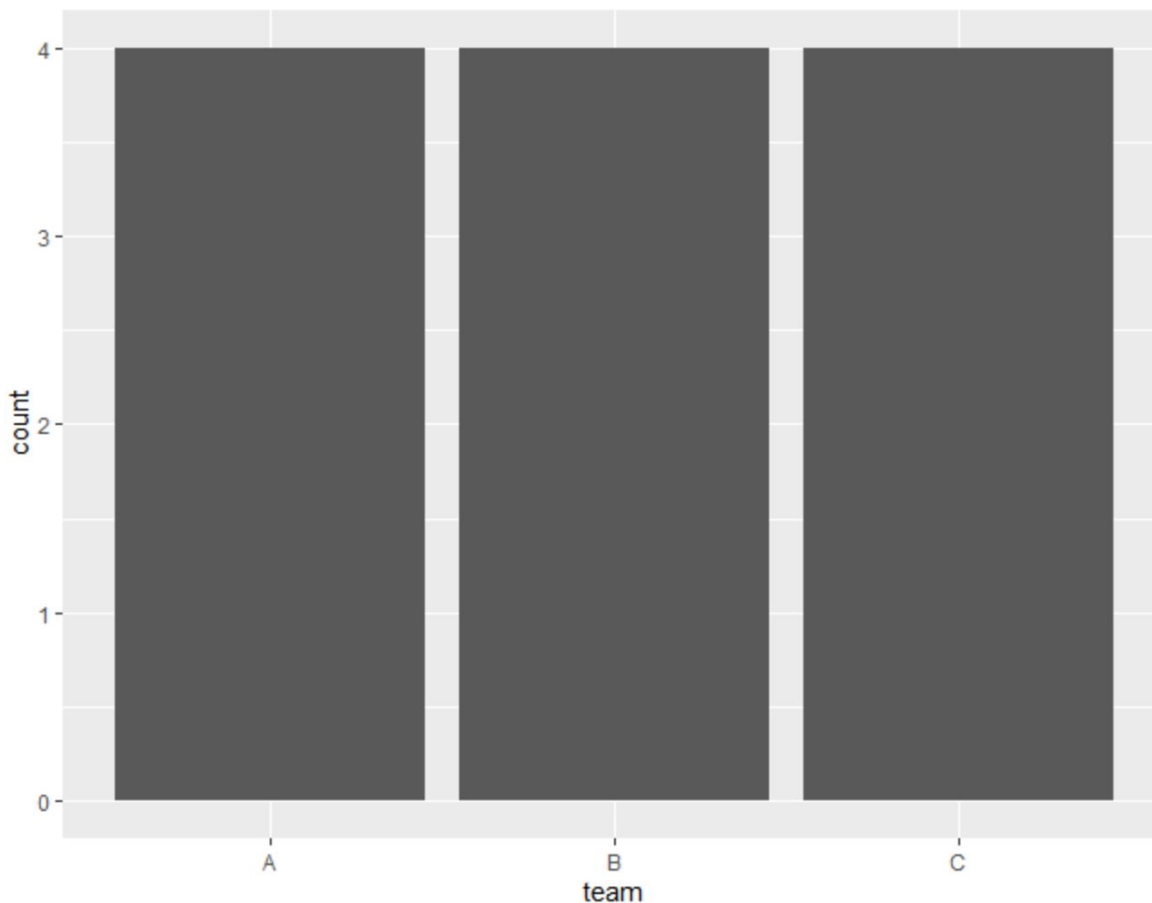
The [data frame](#), named `df`, contains 12 observations. The `team` column is a [categorical variable](#) (A, B, C), and the `points` column is a numerical measure of individual player scores. Critically, this structure allows us to perform two different visualizations: first, counting how many observations belong to each team (frequency), and second, calculating the total [sum](#) of points for each team (identity).

Case Studies: Visualizing Counts vs. Identity

We begin by creating a [bar chart](#) that visualizes the distribution of teams using the default [geom_bar\(\)](#). Since we only map `team` to the [x aesthetic](#), [ggplot2](#) automatically counts the instances of each team entry in the raw data.

library(ggplot2)

```
#create bar chart to visualize occurrence of each unique value in team column  
ggplot(df, aes(team)) +  
geom_bar()
```

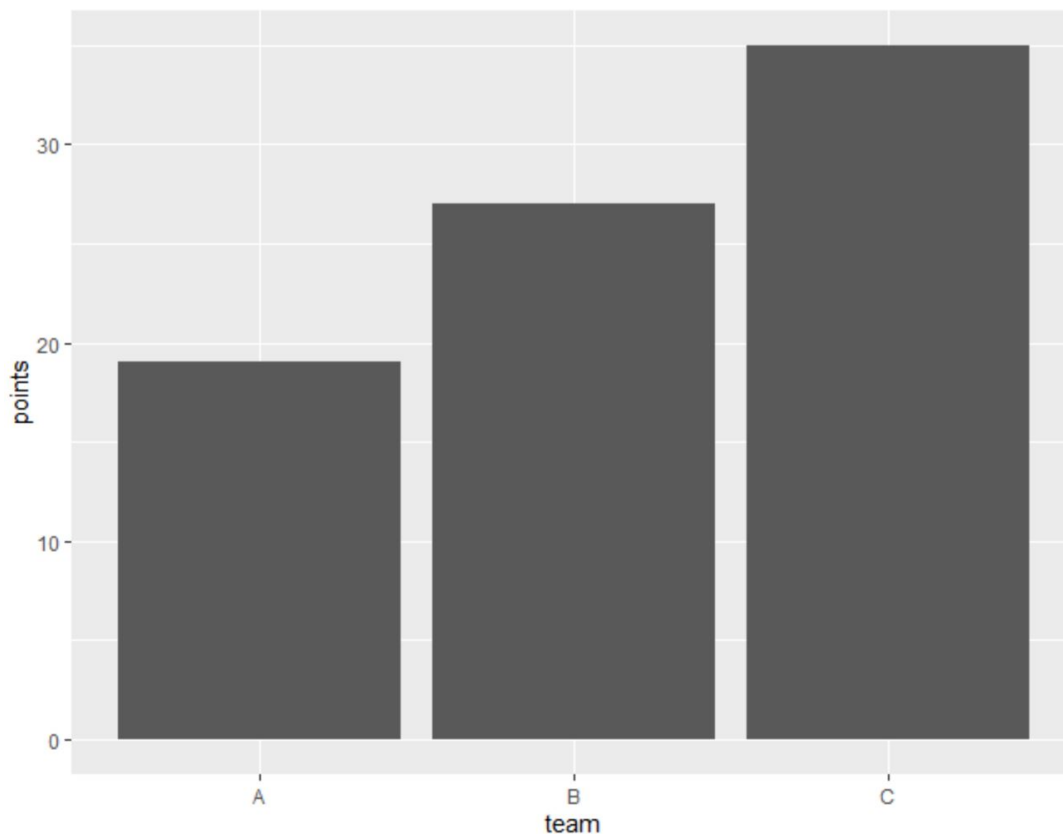


In this initial output, the [x-axis](#) shows the unique team categories (A, B, C). The [y-axis](#) reflects the [count](#) of observations for each team. As we structured the [data frame](#) such that each team appears exactly four times, all bars possess a uniform height of 4. This clearly demonstrates the function of the default [stat_count\(\)](#) in summarizing raw data.

Next, we pivot to plotting the pre-calculated measure--the total points scored by each team--using [stat="identity"](#). Here, we must explicitly map both `team` to the [x aesthetic](#) and `points` to the [y aesthetic](#). Because we are using [stat="identity"](#), [geom_bar\(\)](#) is instructed to internally sum the points for each category and use those totals directly as the bar heights.

library(ggplot2)

```
#create bar chart to visualize sum of points, grouped by team  
ggplot(df, aes(team, points)) +  
geom_bar(stat="identity")
```



The resulting [bar chart](#) now displays the [sum](#) of points for each team along the [y-axis](#). Based on the source data: Team A totals 19 points (3+5+5+6), Team B totals 27 points (5+7+7+8), and Team C totals 35 points (9+9+9+8). The bar heights precisely reflect these [sums](#), confirming that [stat="identity"](#) successfully mapped the aggregated values to the visual height, fulfilling the goal of plotting a pre-calculated metric.

Choosing the Right Approach and Critical Requirements

The decision between using the default [geom_bar\(\)](#) (which uses [stat_count\(\)](#)) and specifying [stat="identity"](#) is determined entirely by the format of your input data and the analytical goal of the visualization. Selecting the appropriate method is paramount for generating accurate and meaningful [bar charts](#).

Opt for the **default** [geom_bar\(\)](#) when your [data frame](#) contains raw, unsummarized observations of a [categorical variable](#), and your objective is to visually represent the frequency or [count](#) of

each category. This is the fastest way to visualize distributions and determine if categories are balanced or skewed.

Conversely, you must use `geom_bar(stat="identity")` when your **data frame** already holds the calculated summary values (e.g., **sums**, averages, medians) that you want to plot as bar heights. The critical requirement here is the explicit mapping of both the **x variable** (categories) and the **y variable** (the numerical measure). Failure to supply a **y variable** when using `stat="identity"` will result in an error, as **ggplot2** will not know which values to use for the visual height. If you need to perform more complex aggregations than simple counts or sums, consider utilizing other advanced **statistical transformations** provided by **ggplot2**, such as `stat_summary()`.

Conclusion

The effective utilization of `geom_bar()` in **ggplot2** hinges on a clear understanding of the `stat` argument. The default mechanism automatically calculates and plots the **counts** of unique **categorical data**, making it ideal for visualizing raw distributions. Conversely, setting `stat="identity"` is the necessary command for plotting data that has already been summarized, requiring the user to supply the numerical values directly via the **y aesthetic**. By correctly applying these two distinct methods, you ensure that your visualizations are faithful representations of your data structure and analytical intent, providing clear and powerful visual insights.

Additional Resources

The following tutorials explain how to perform other common tasks in **ggplot2**: