

Learning to Display Percentages on Histograms Using ggplot2

Authored by
Mohammed loot

November 3, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning to Display Percentages on Histograms Using ggplot2*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=9143>

The Challenge of Displaying Relative Frequency in ggplot2

Histograms are fundamental tools in [R programming language](#) for visualizing the distribution of data. By default, the popular [ggplot2 package](#) calculates and displays the absolute counts (or frequencies) of observations falling into specific bins or categories on the y-axis. While this is useful for understanding raw magnitude, it often fails to provide the necessary context for comparison, especially when dealing with samples of varying sizes or when the total population size is the primary focus. To create more insightful and comparative visualizations, data analysts frequently need to convert these raw counts into relative frequencies, expressed as percentages. This conversion allows stakeholders to immediately grasp the proportional breakdown of the data, enhancing the clarity and interpretability of the [histogram](#).

The process of transforming absolute counts into percentages within the [ggplot2 package](#) requires leveraging its powerful internal statistical transformations. Specifically, we must instruct the geometric object (`geom_bar`` or `geom_histogram``) to calculate the proportion of each bar relative to the total number of observations in the dataset. This is achieved through a combination of aesthetic mapping functions and internal variables that [ggplot2](#) generates during the plot creation process. Successfully implementing this technique transforms a simple frequency chart into a sophisticated visualization of proportional distribution.

The solution presented here simplifies this complex statistical transformation, providing a clear and replicable methodology for achieving percentage display on the y-axis. It primarily relies on the utilization of the `..count..`` variable, which represents the height of each bar, and wrapping this calculation within the `sum()`` function to establish the denominator (the total count). Furthermore, the necessary formatting--converting decimal proportions (e.g., 0.25) into recognizable percentage strings (e.g., 25%)--is handled efficiently by the complementary [scales package](#), ensuring a polished and professional output.

The Core Syntax: Understanding the ggplot2 Mechanics

To effectively transition from raw counts to proportional data, we must manipulate the aesthetic mapping (`aes``) within the `geom_bar`` layer. When creating a bar chart or histogram, [ggplot2](#) automatically calculates the count for each category. We intercept this internal calculation using the variable `..count..``. The formula `(..count..)/sum(..count..)`` is the core statistical engine that performs the desired conversion, dividing the count of the specific category by the total count across all categories. This results in a proportion, where the y-axis values range from 0 to 1.

Once the proportional calculation is set up, the next critical step involves proper labeling of the y-axis. Since the calculated values are fractions, they must be formatted as percentages for visual clarity. This is where the specialized [scales package](#) becomes indispensable. By calling `scale_y_continuous(labels=percent)``, we utilize the `percent`` function from the [scales package](#),

which automatically handles the multiplication by 100 and appends the percentage symbol (%) to the axis labels. This two-part approach--the statistical transformation in `geom_bar` and the formatting in `scale_y_continuous`--forms the backbone of displaying relative frequencies in [ggplot2](#).

The following basic syntax encapsulates the necessary libraries and functions required to execute this transformation. Note that this structure is primarily designed for visualizing the distribution of a [categorical variable](#), where the input data represents distinct groups rather than continuous numerical bins. For continuous data, slight adjustments involving `geom_histogram` might be necessary, but the proportional calculation logic remains fundamentally the same. Understanding and internalizing this template is the first step toward advanced customization.

You can use the following basic syntax to display percentages on the y-axis of a [histogram](#) in [ggplot2](#):

```
library(ggplot2)
```

```
library(scales)
```

```
#create histogram with percentages  
ggplot(data, aes(x = factor(team))) +  
geom_bar(aes(y = (..count..)/sum(..count..))) +  
scale_y_continuous(labels=percent)
```

The following examples show how to use this syntax in practice.

Example 1: Implementing the Basic Percentage Histogram

This first example demonstrates the application of the core syntax to a simple, manually defined dataset. We define a data frame containing two variables: a categorical variable called 'team' (A, B, or C) and a numerical variable called 'points'. Since we are interested in visualizing the distribution of the categorical variable 'team'--specifically, the percentage contribution of each team to the total dataset--we use `geom_bar`. The use of `geom_bar` is appropriate here because it calculates counts for discrete categories, which is distinct from `geom_histogram`, which bins continuous data.

Within the `ggplot()` call, we map the 'team' variable to the x-axis, ensuring it is treated as a factor using `aes(x = factor(team))`. The crucial transformation occurs within `geom_bar`, where we explicitly map the y-aesthetic to our proportional calculation: `aes(y = (..count..)/sum(..count..))`. This tells [ggplot2](#) to calculate the relative frequency for each team. Finally, `scale_y_continuous(labels=percent)` ensures the resulting proportions (e.g., 0.36) are rendered as human-readable percentages on the vertical axis (e.g., 36%). This approach provides an

immediate visual representation of the team distribution without requiring mental calculation of fractions.

This setup is highly effective for exploratory data analysis, allowing analysts to quickly determine if categories are balanced or if one group disproportionately dominates the dataset. By displaying the proportions directly, the resulting visualization is self-contained and requires minimal external explanation regarding the scale of the y-axis, making it suitable for direct inclusion in reports and presentations.

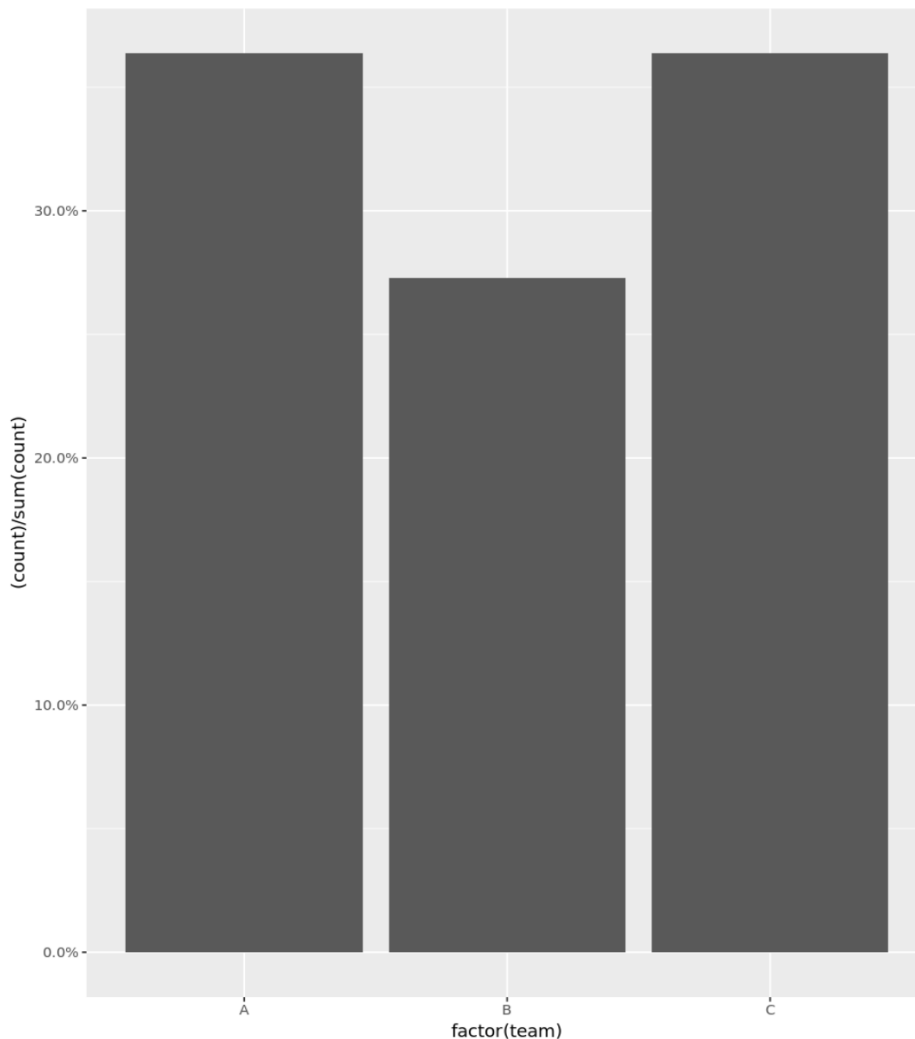
Example 1: Basic Histogram with Percentages

The following code shows how to create a [histogram](#) for categorical variables with percentages displayed on the y-axis:

```
library(ggplot2)
library(scales)

#define data frame
data <- data.frame(team = c('A', 'A', 'A', 'A', 'B', 'B', 'B', 'C', 'C', 'C', 'C'),
  points = c(77, 79, 93, 85, 89, 99, 90, 80, 68, 91, 92))

#create histogram with percentages
ggplot(data, aes(x = factor(team))) +
  geom_bar(aes(y = (..count..)/sum(..count..))) +
  scale_y_continuous(labels=percent)
```



Example 2: Refining Output - Removing Decimal Places

While the initial percentage display is functional, it often includes several decimal places (e.g., 36.3636%), which can clutter the y-axis and reduce readability. In many reporting scenarios, displaying only whole numbers or rounding to a specific level of precision is preferable. The [scales package](#) provides fine-grained control over label formatting through the `percent_format()` function, which is a highly customizable alternative to the basic `percent` shortcut used previously.

The key argument within `percent_format()` is `accuracy`. By setting `accuracy = 1L` (where `1L` indicates an integer or whole number accuracy), we instruct the function to round the calculated percentages to the nearest whole number before displaying them on the axis. This significantly cleans up the visualization, ensuring that the proportions are displayed concisely (e.g., 36%, 27%). It is important to remember that this rounding occurs only for display purposes; the underlying calculation of the bar heights remains based on the precise proportional value.

Implementing this refinement replaces the simple `labels=percent` call with the more detailed `labels = scales::percent_format(accuracy = 1L)`. Using the double colon notation (`scales::`) explicitly calls the function from the package, which is good practice for clarity, although it is not strictly necessary if the package has already been loaded via `library(scales)`. This customization step is essential for creating publication-quality graphics where precision and visual simplicity are prioritized.

Example 2: Histogram with Percentages (Remove Decimal Places)

You can use the **accuracy** argument to only display whole numbers as percentages on the y-axis as well:

```
library(ggplot2)
```

```
library(scales)
```

```
#define data frame
```

```
data <- data.frame(team = c('A', 'A', 'A', 'A', 'B', 'B', 'B', 'C', 'C', 'C', 'C'),
```

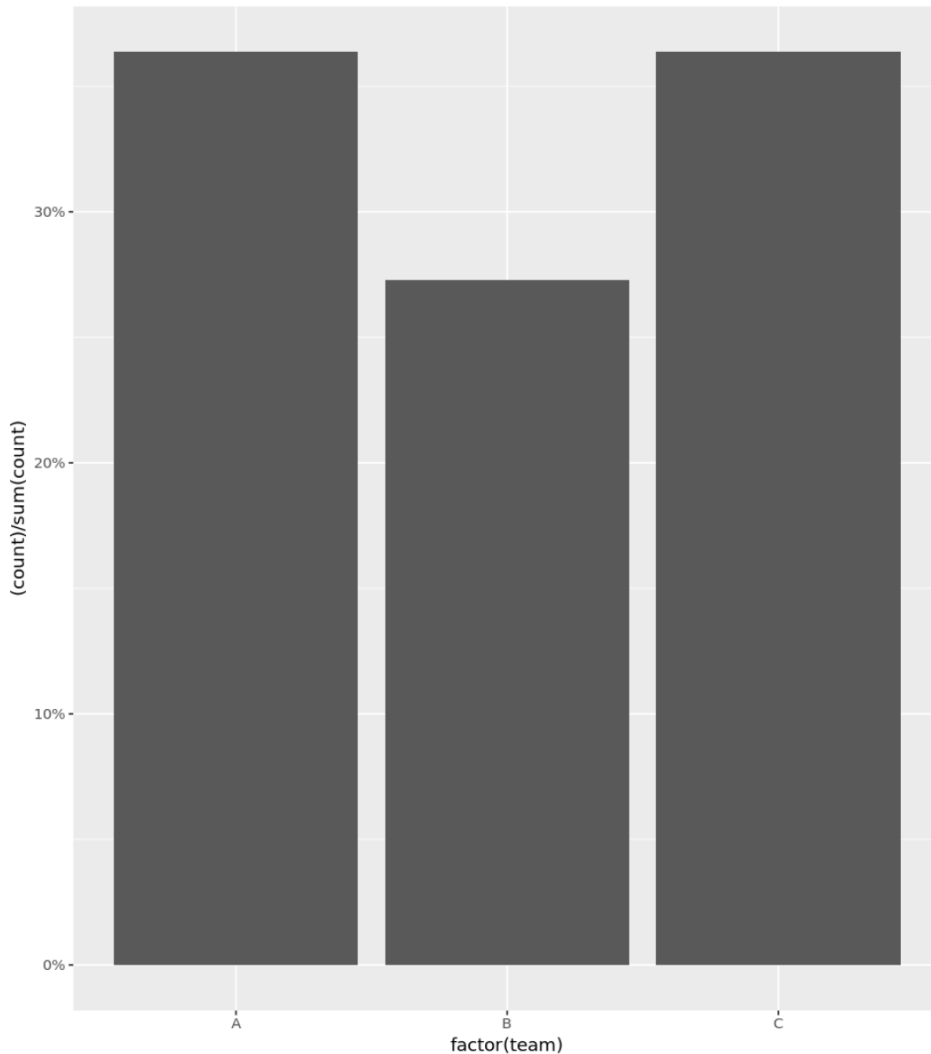
```
points = c(77, 79, 93, 85, 89, 99, 90, 80, 68, 91, 92))
```

```
#create histogram with percentages
```

```
ggplot(data, aes(x = factor(team))) +
```

```
geom_bar(aes(y = (..count..)/sum(..count..))) +
```

```
scale_y_continuous(labels = scales::percent_format(accuracy = 1L))
```



Example 3: Customizing Aesthetics for Professional Visualizations

Beyond the essential data transformation, effective data visualization demands attention to aesthetics. A standard [ggplot2](#) plot often uses default gray backgrounds and generic labels, which may not meet professional presentation standards. Customizing the plot involves adjusting colors, adding meaningful titles, relabeling axes, and selecting an appropriate theme. These enhancements increase the plot's visual appeal and ensure that the message derived from the data is communicated clearly and effectively.

In this advanced example, we integrate several customization layers. We change the bar color using the `fill` argument within `geom_bar`, setting it to `'lightblue'` for a cleaner look. More importantly, we use the `labs()` function to define a descriptive title and informative axis labels. Instead of the default `'factor(team)'` and `'proportion'`, we use clear labels like `'Team'` and `'Percent of Total'`. Finally, applying `theme_minimal()` removes the standard gray background and grid lines,

resulting in a cleaner, minimalist aesthetic that focuses the viewer's attention purely on the data distribution.

This level of customization demonstrates the flexibility of [ggplot2](#). By combining the statistical power of the percentage calculation with thoughtful aesthetic choices, the resulting figure is not merely a data output but a polished, narrative-driven artifact suitable for high-stakes reporting or academic publication. It confirms that the underlying proportional data is accurately represented while optimizing the visual experience for the audience.

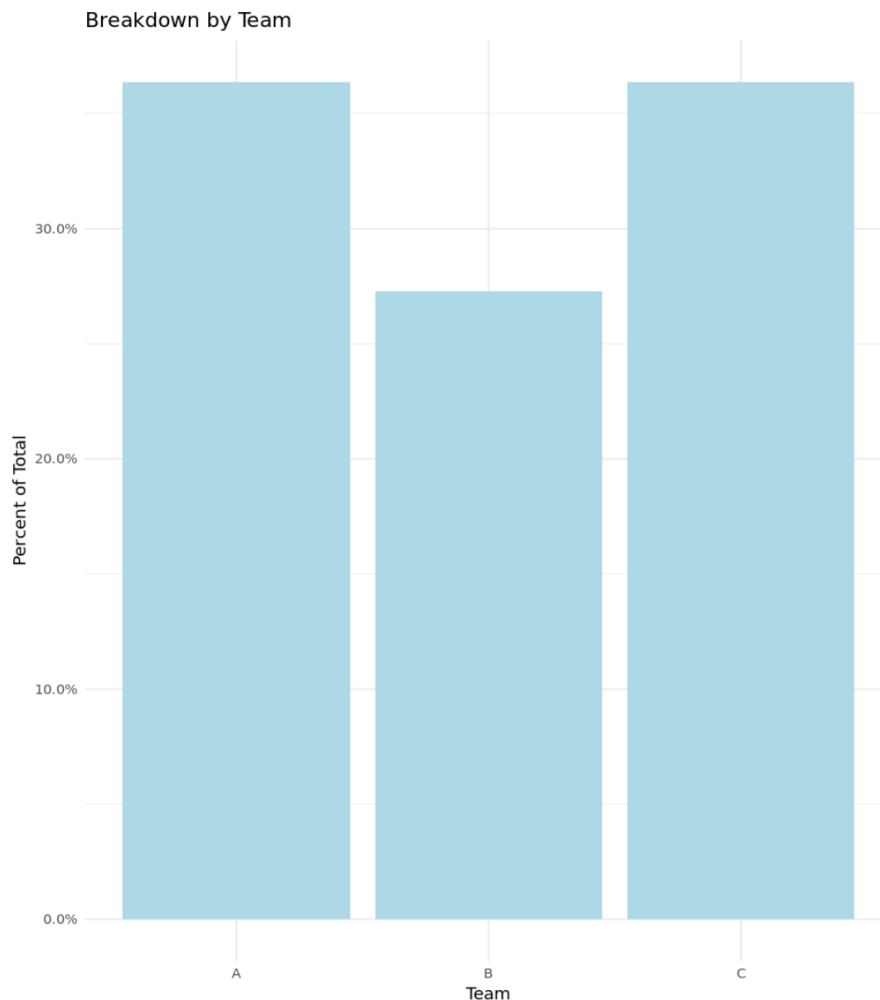
Example 3: Custom [Histogram](#) with Percentages

The following code shows how to create a [histogram](#) with percentages shown on the y-axis and a custom title, axis labels, and colors:

```
library(ggplot2)
library(scales)

#define data frame
data <- data.frame(team = c('A', 'A', 'A', 'A', 'B', 'B', 'B', 'C', 'C', 'C', 'C'),
  points = c(77, 79, 93, 85, 89, 99, 90, 80, 68, 91, 92))

#create histogram with percentages and custom aesthetics
ggplot(data, aes(x = factor(team))) +
  geom_bar(aes(y = (..count..)/sum(..count..)), fill = 'lightblue') +
  scale_y_continuous(labels=percent) +
  labs(title = 'Breakdown by Team', x = 'Team', y = 'Percent of Total') +
  theme_minimal()
```



Why Relative Frequencies are Essential for Data Comparison

The decision to display relative frequencies (percentages) rather than absolute counts is often driven by the need for contextual comparison. Absolute counts are inherently tied to the sample size. If one category has a count of 50 and another has a count of 10, this might seem like a large disparity. However, if the total sample size changes, or if these counts are drawn from subpopulations of vastly different sizes, the raw count comparison becomes misleading. Percentages, on the other hand, normalize the data by expressing each count as a proportion of the whole, providing a standardized measure that is independent of the overall sample magnitude.

Consider a scenario where you are comparing survey responses across two geographical regions. Region A had 1,000 respondents, and 200 answered 'Yes' (20%). Region B had 100 respondents, and 30 answered 'Yes' (30%). If we only plotted the raw counts (200 vs. 30), Region A would appear overwhelmingly dominant. By converting these to percentages, the visualization correctly highlights that, proportionally, Region B actually has a higher frequency of 'Yes' responses. Therefore, using percentages is crucial when the goal is to assess the intensity or distribution

pattern within the data, rather than the sheer volume of observations.

Furthermore, percentages simplify interpretation for a general audience. Most people intuitively understand a percentage breakdown, whereas interpreting statistical count variables often requires understanding the underlying dataset size. By employing the techniques demonstrated using [ggplot2](#) and the [scales package](#) in [R programming language](#), data professionals can ensure their visualizations are both statistically accurate and highly accessible, bridging the gap between complex data analysis and effective communication.

Additional Resources and Summary

Mastering the display of relative frequencies is a critical skill for anyone using the [ggplot2 package](#) for data visualization. The method is straightforward: calculate the proportion using the internal `..count..` variable within `geom_bar` and apply the necessary percentage formatting using the functions available in the [scales package](#). This technique moves beyond simple frequency plots to offer true insight into the distribution of categorical data.

For those interested in further expanding their visualization capabilities in [R programming language](#), exploring additional features of the `percent_format()` function--such as customizing the percentage sign or handling missing values--can lead to even more nuanced and tailored graphics.

Related: