# Histograms in R with ggplot2

2024-07-27

**i** This work was developed using resources that are available under a Creative Commons Attribution 4.0 International License, made available on the SOLES Open Educational Resources repository by the School of Life and Environmental Sciences, The University of Sydney.





## About

The histogram is a bar plot that shows the frequency of (often) continuous values in a dataset. It helps identify patterns, outliers, and the shape of the distribution (skewness, kurtosis).

i Assumed knowledge

- You know how to install and load packages in R.
- You know how to import data into R.
- You recognise data frames and vectors.

#### Data structure

The data should be in a **long format** (also known as tidy data), where each row is an observation and each column is a variable (Figure 1). If your data is not already structured this way, reshape it manually in a spreadsheet program or in R using the pivot\_longer() function from the tidyr package.

$\mathbf{Sex}$	BW
F	2.15
М	2.55
$\mathbf{F}$	2.95
$\mathbf{F}$	2.70
$\mathbf{M}$	2.20
$\mathbf{F}$	1.85
Μ	2.55
Μ	2.60

Figure 1: Data should be in long format (left) where each row is an observation and each column is a variable. This is the preferred format for most statistical software. Wide format (right) is also common, but may require additional steps to analyse or visualise in some instances.

#### Data

For this cheatsheet we will use the entire possums dataset used in BIOL2022 labs.

#### Import data

```
library(readx1)
possums <- read_excel("possums.xlsx", sheet = 2)</pre>
```

## Plot

Use the different plots below to explore the use of histograms in R. Note that histograms only need one variable to be plotted, therefore we pick any one of the several continuous variables in the dataset in the **aes()** function.

## Version 1

There bare minimum code to create a histogram in R.

```
library(ggplot2)
ggplot(possums, aes(x = BW)) +
geom_histogram()
```

- (1) Load the ggplot2 package with library(ggplot2).
- (2) Create a canvas using ggplot() and specify the variable to be plotted on the x-axis with aes(x = BW).

(1)

2

(3)

(3) Add a histogram layer using geom\_histogram().

`stat\_bin()` using `bins = 30`. Pick better value with `binwidth`.



## Version 2

Add colours, labels, adjust bin width, and change the theme.

- (1) Use colour and fill arguments to change bar colors.
- (2) Adjust bin width with binwidth to control the detail of the histogram. A smaller value will increase the number of bars based on the range of the data.
- (3) Add axis labels using labs().
- (4) Use theme\_minimal() for a standardized appearance.



# Version 3

Plot both a histogram and a density plot at the same time.

```
рЗ
```



## Version 4

Compare a histogram with a standardised normal distribution.

```
library(ggplot2)
ggplot(possums, aes(x = AactiveTBLUP)) +
geom_histogram(aes(y = after_stat(density)),
            fill = "skyblue",
            color = "black",
            binwidth = 0.3) +
stat_function(fun = dnorm,
            args = list(mean = mean(possums$AactiveTBLUP),
                sd = sd(possums$AactiveTBLUP)),
            color = "red", linewidth = 1) +
theme_minimal()
```



#### Export

Use the ggsave() function to save the plot as an image file.

- The filename argument specifies the name of the file.
- The plot argument specifies the plot to be saved. In this case, the plot is stored in the object p3 (from Version 3).
- The width and height arguments specify the dimensions of the plot in inches.

ggsave(filename = "histogram.pdf", plot = p3, width = 7, height = 5)

The plot will be saved in the working directory, unless you specify a different path in the filename argument.