

How to Normalize Data in R for Accurate Analysis

Authored by
stats writer

March 3, 2026

RECOMMENDED CITATION

stats writer (2026). *How to Normalize Data in R for Accurate Analysis*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=133648>

Data normalization in R is the process of standardizing data values to a common scale, allowing for more accurate analysis and comparison. This is achieved by transforming the data to a distribution that follows a predefined pattern, such as a normal distribution. In R, this can be done using various methods such as scaling, centering, and standardization. This ensures that the data is consistent and unbiased, making it easier to interpret and draw meaningful conclusions. Normalization is an essential step in data analysis and is widely used in various fields, including statistics, machine learning, and data science. With R, users have access to various built-in functions and packages that facilitate the normalization process, making it a powerful tool for handling and analyzing data.

Normalize Data in R

In most cases, when people talk about "normalizing" variables in a dataset, it means they'd like to scale the values such that the variable has a mean of 0 and a standard deviation of 1.

The most common reason to normalize variables is when you're conducting some type of multivariate analysis (i.e. you want to understand the relationship between several predictor variables and a response variable) and you want each variable to contribute equally to the analysis.

When variables are measured at different scales, they often do not contribute equally to the analysis. For example, if the values of one variable range from 0 to 100,000 and the values of another variable range from 0

to 100, the variable with the larger range will be given a larger weight in the analysis.

This is common when one variable measures something like salary (\$0 to \$100,000) and another variable measures something like age (0 to 100 years).

By normalizing the variables, we can be sure that each variable contributes equally to the analysis. Two common ways to normalize (or "scale") variables include:

Min-Max Normalization: $(X - \min(X)) / (\max(X) - \min(X))$

Z-Score Standardization: $(X - \mu) / \sigma$

Next, we'll show how to implement both of these techniques in R.

How to Normalize (or "Scale") Variables in R

For each of the following examples, we'll use the built-in R dataset *iris* to illustrate how to normalize or scale variables in R:

```
#view first six rows of iris dataset  
head(iris)
```

```
# Sepal.Length Sepal.Width Petal.Length Petal.Width  
Species
```

```
#1 5.1 3.5 1.4 0.2 setosa
```

```
#2 4.9 3.0 1.4 0.2 setosa
```

```
#3 4.7 3.2 1.3 0.2 setosa
```

```
#4 4.6 3.1 1.5 0.2 setosa
```

```
#5 5.0 3.6 1.4 0.2 setosa
```

```
#6 5.4 3.9 1.7 0.4 setosa
```

Min-Max Normalization

The formula for a min-max normalization is:

$$(X - \min(X)) / (\max(X) - \min(X))$$

For each value of a variable, we simply find how far that value is from the minimum value, then divide by the range.

To implement this in R, we can define a simple function and then use to apply that function to whichever columns in the iris dataset we would like:

```
#define Min-Max normalization function
```

```
min_max_norm <- function(x) {
```

```
(x - min(x)) / (max(x) - min(x))
```

```
}
```

```
#apply Min-Max normalization to first four columns in  
iris dataset
```

```
iris_norm <- as.data.frame(lapply(iris, min_max_norm))
```

```
#view first six rows of normalized iris dataset
```

```
head(iris_norm)
```

```
# Sepal.Length Sepal.Width Petal.Length Petal.Width
```

```
#1 0.22222222 0.6250000 0.06779661 0.04166667
```

```
#2 0.16666667 0.4166667 0.06779661 0.04166667
```

```
#3 0.11111111 0.5000000 0.05084746 0.04166667
```

```
#4 0.08333333 0.4583333 0.08474576 0.04166667
```

```
#5 0.19444444 0.6666667 0.06779661 0.04166667
```

```
#6 0.30555556 0.7916667 0.11864407 0.12500000
```

Notice that each of the columns now have values that range from 0 to 1. Also notice that the fifth column "Species" was dropped from this data frame. We can easily add it back by using the following code:

```
#add back Species column
```

```
iris_norm$Species <- iris$Species
```

```
#view first six rows of iris_norm  
head(iris_norm)
```

```
# Sepal.Length Sepal.Width Petal.Length Petal.Width  
Species  
#1 0.22222222 0.6250000 0.06779661 0.04166667 setosa  
#2 0.16666667 0.4166667 0.06779661 0.04166667 setosa  
#3 0.11111111 0.5000000 0.05084746 0.04166667 setosa  
#4 0.08333333 0.4583333 0.08474576 0.04166667 setosa  
#5 0.19444444 0.6666667 0.06779661 0.04166667 setosa  
#6 0.30555556 0.7916667 0.11864407 0.12500000 setosa
```

Z-Score Standardization

The drawback of the min-max normalization technique is that it brings the data values towards the mean. If we want to make sure that outliers get weighted more than other values, a z-score standardization is a better technique to implement.

The formula for a z-score standardization is:

$$(X - \mu) / \sigma$$

For each value of a variable, we simply subtract the mean value of the variable, then divide by the standard

deviation of the variable.

To implement this in R, we have a few different options:

1. Standardize one variable

If we simply want to standardize one variable in a dataset, such as *Sepal.Width* in the iris dataset, we can use the following code:

```
#standardize Sepal.Width
iris$Sepal.Width <- (iris$Sepal.Width -
mean(iris$Sepal.Width)) / sd(iris$Sepal.Width)

head(iris)

# Sepal.Length Sepal.Width Petal.Length Petal.Width
Species
#1 5.1 1.01560199 1.4 0.2 setosa
#2 4.9 -0.13153881 1.4 0.2 setosa
#3 4.7 0.32731751 1.3 0.2 setosa
#4 4.6 0.09788935 1.5 0.2 setosa
#5 5.0 1.24503015 1.4 0.2 setosa
#6 5.4 1.93331463 1.7 0.4 setosa
```

The values of *Sepal.Width* are now scaled such that the

mean is 0 and the standard deviation is 1. We can even verify this if we'd like:

```
#find mean of Sepal.Width
```

```
mean(iris$Sepal.Width)
```

```
# 2.034094e-16 #basically zero#find standard deviation of Sepal.Width
```

```
sd(iris$Sepal.Width)
```

```
# 1
```

2. Standardize several variables using the scale function

To standardize several variables, we can simply use the *scale* function. For example, the following code shows how to scale the first four columns of the iris dataset:

```
#standardize first four columns of iris dataset
```

```
iris_standardize <- as.data.frame(scale(iris))
```

```
#view first six rows of standardized datasethead(iris_standardize)
```

```
# Sepal.Length Sepal.Width Petal.Length Petal.Width
```

```
#1 -0.8976739 1.01560199 -1.335752 -1.311052
#2 -1.1392005 -0.13153881 -1.335752 -1.311052
#3 -1.3807271 0.32731751 -1.392399 -1.311052
#4 -1.5014904 0.09788935 -1.279104 -1.311052
#5 -1.0184372 1.24503015 -1.335752 -1.311052
#6 -0.5353840 1.93331463 -1.165809 -1.048667
```

Note that the *scale* function, by default, attempts to standardize every column in a data frame. Thus, we would get an error if we attempted to use `scale(iris)` because the *Species* column is not numeric and cannot be standardized:

```
scale(iris)
```

```
#Error in colMeans(x, na.rm = TRUE) : 'x' must be
numeric
```

However, it is possible to standardize only certain variables in a data frame while also keeping all other variables the same by using the *dplyr* package. For example, the following code standardizes the variables *Sepal.Width* and *Sepal.Length* while keeping all other variables the same:

```
#load dplyr package  
library(dplyr)  
  
#standardize Sepal.Width and Sepal.Length iris_new <-  
iris %>% mutate_each_(list(~scale(.) %>% as.vector),  
vars = c("Sepal.Width", "Sepal.Length"))  
  
#view first six rows of new data frame  
head(iris_new)  
  
# Sepal.Length Sepal.Width Petal.Length Petal.Width  
Species  
#1 -0.8976739 1.01560199 1.4 0.2 setosa  
#2 -1.1392005 -0.13153881 1.4 0.2 setosa  
#3 -1.3807271 0.32731751 1.3 0.2 setosa  
#4 -1.5014904 0.09788935 1.5 0.2 setosa  
#5 -1.0184372 1.24503015 1.4 0.2 setosa  
#6 -0.5353840 1.93331463 1.7 0.4 setosa
```

Notice that *Sepal.Length* and *Sepal.Width* are standardized such that both variables have a mean of 0 and a standard deviation of 1, while the other three variables in the data frame remain unchanged.