

How to Classify Iris Species in R: A Step-by-Step Guide

Authored by
stats writer

December 4, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Classify Iris Species in R: A Step-by-Step Guide*.

PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=104904>

The Iris Dataset, readily available in R, stands as one of the most famous and fundamental datasets for illustrating fundamental concepts in machine learning and statistics. It is particularly renowned for its use in analyzing classification problems.

This dataset, originally compiled by botanist Edgar Anderson and popularized by Ronald Fisher in 1936, comprises 150 observations, 50 from each of the three species of Iris flowers: **Setosa**, **Versicolor**, and **Virginica**. Each observation includes four key morphological measurements: sepal length, sepal width, petal length, and petal width, all recorded in centimeters. Its simplicity and clarity make it an ideal starting point for data scientists new to R programming or statistical analysis, offering a perfect environment to practice tasks like classification, clustering, and exploratory data analysis (EDA).

The **iris** dataset is a highly structured, built-in resource within R, providing measurements on four distinct attributes for 150 samples across three species. This tutorial will guide you through the essential steps of exploring, summarizing, and visualizing this dataset using core R functions, establishing a strong foundation for handling more complex data structures.

Accessing the Built-in Iris Dataset in R

A significant advantage of the Iris Dataset is that it comes pre-loaded with the standard R installation, eliminating the need for external package installation or file imports. To formally load the dataset into your active R session, you use the fundamental `data()` command, specifying the name of the dataset. This action makes the data immediately accessible as a data frame named `iris`.

data(iris)

Once loaded, the first step in any data exploration is inspecting the structure. We use the powerful `head()` function to display the initial six observations, which confirms successful loading and provides an immediate sense of the data structure, including variable names and data types for the first few rows.

View the initial six records of the iris data frame

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
```

As shown in the output, the data contains four numeric variables (measurements in centimeters) and one categorical variable, `Species`, which identifies the flower type. The first six rows confirm that all these entries belong to the **setosa** species.

Statistical Summary and Descriptive Analysis

To gain a rapid understanding of the distribution and central tendencies within the data, R's built-in `summary()` function is invaluable. Applied to the `iris` data frame, it automatically detects the variable types and calculates appropriate descriptive statistics for each column, providing a comprehensive five-number summary plus the mean for quantitative variables, and frequency counts for factors (categorical variables).

```
# Generate descriptive statistics for all variables
summary(iris)
```

```
Sepal.Length Sepal.Width Petal.Length Petal.Width
Min. :4.300 Min. :2.000 Min. :1.000 Min. :0.100
1st Qu.:5.100 1st Qu.:2.800 1st Qu.:1.600 1st Qu.:0.300
Median :5.800 Median :3.000 Median :4.350 Median :1.300
Mean :5.843 Mean :3.057 Mean :3.758 Mean :1.199
3rd Qu.:6.400 3rd Qu.:3.300 3rd Qu.:5.100 3rd Qu.:1.800
Max. :7.900 Max. :4.400 Max. :6.900 Max. :2.500
Species
setosa :50
versicolor:50
virginica :50
```

The `summary()` output is crucial for quickly assessing data quality, identifying potential outliers, and understanding the spread of the features. For the four numeric columns (Sepal Length, Sepal Width, Petal Length, and Petal Width), we observe the following statistical measures, forming the basis of our exploratory data analysis:

Min: Represents the smallest recorded measurement in that feature.

1st Qu (First Quartile): The value below which 25% of the data falls (the 25th percentile).

Median: The middle value of the sorted dataset (the 50th percentile).

Mean: The arithmetic average of all measurements.

3rd Qu (Third Quartile): The value below which 75% of the data falls (the 75th percentile).

Max: Represents the largest recorded measurement in that feature.

For the categorical variable, `species`, the summary provides a frequency table, confirming the balanced nature of the Iris Dataset: there are exactly 50 observations for each of the three distinct species, making it an ideal dataset for training balanced classification problems models.

setosa: 50 observations.

versicolor: 50 observations.

virginica: 50 observations.

Inspecting Dataset Dimensions and Variable Names

Beyond the statistical summary, understanding the physical size of the dataset is critical for resource management and modeling strategy. The `dim()` function in R provides the dimensions of the data frame, outputting the number of rows (observations) and the number of columns (variables) in that specific order.

Determine the number of rows (observations) and columns (features)

```
dim(iris)
```

```
150 5
```

The output confirms that the `iris` dataset consists of 150 total observations and 5 variables. Knowing the exact structure of your dataset ensures that subsequent operations, such as splitting the data for training and testing in machine learning tasks, are correctly implemented.

Finally, explicitly retrieving the column names is important for programmatic access and ensuring consistency in script development. The `names()` function returns a character vector containing the exact labels for all variables within the Iris Dataset.

Display the variable (column) names

```
names(iris)
```

```
"Sepal.Length" "Sepal.Width" "Petal.Length" "Petal.Width" "Species"
```

Exploratory Data Visualization Using Base R Graphics

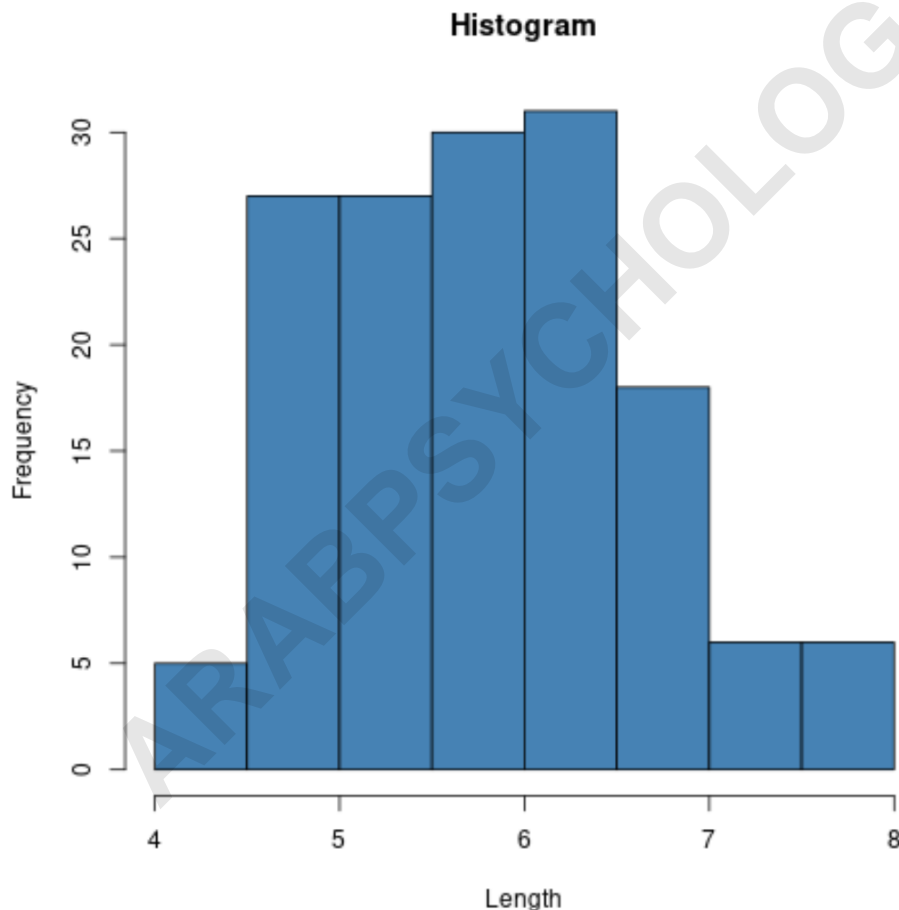
Statistical summaries are powerful, but graphical representations offer immediate insights into data distribution and relationships that numbers often obscure. R's base plotting functions allow for quick and effective visualization of the `iris` data, forming the backbone of effective exploratory data analysis (EDA). We will demonstrate three key plot types: the histogram, the scatterplot, and

the `boxplot`.

The `hist()` function is used to plot the frequency distribution of a single numeric variable. By examining the histogram of `Sepal.Length`, we can visually assess its spread, central tendency, and identify if the distribution is approximately normal or skewed.

Generate a histogram for the Sepal Length variable

```
hist(iris$Sepal.Length,  
col='steelblue',  
main='Distribution of Sepal Length',  
xlab='Sepal Length (cm)',  
ylab='Frequency')
```



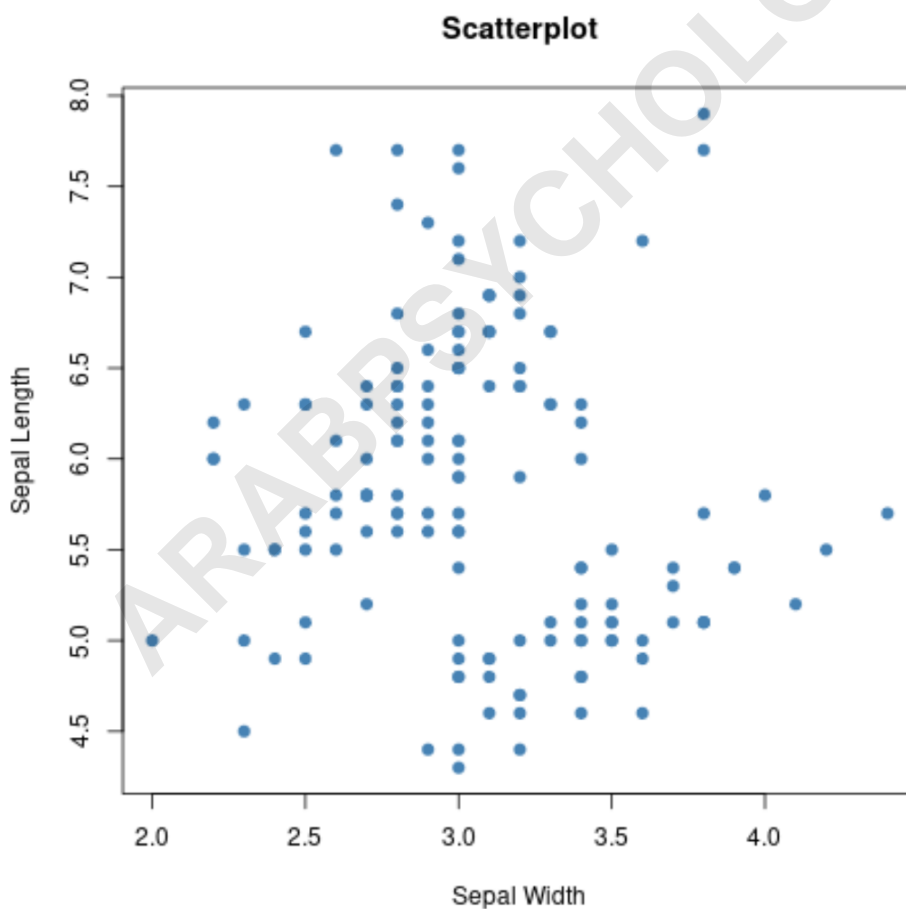
This visualization clearly shows a bimodal or potentially trimodal distribution in sepal length, suggesting that the different species contained within the dataset exhibit distinct measurement characteristics. This observation immediately validates the need for species-specific analysis, a common pattern in successful machine learning approaches to classification problems.

Analyzing Feature Correlation with Scatterplots

To analyze the relationship between two continuous variables, the scatterplot is the standard visualization tool. Using the `plot()` function in `R`, we can map one feature against another, allowing us to detect correlation patterns, clusters, or potential outliers that might influence modeling. For instance, we can plot `Sepal.Width` against `Sepal.Length` to see how these two measurements relate across the entire dataset.

Create a scatterplot showing the relationship between sepal width and length

```
plot(iris$Sepal.Width, iris$Sepal.Length,  
col='steelblue',  
main='Sepal Dimensions Scatterplot',  
xlab='Sepal Width (cm)',  
ylab='Sepal Length (cm)',  
pch=19)
```



The resulting scatterplot shows a cloud of points, hinting at a subtle positive correlation, though the presence of distinct clusters (which would be more obvious if points were colored by `Species`)

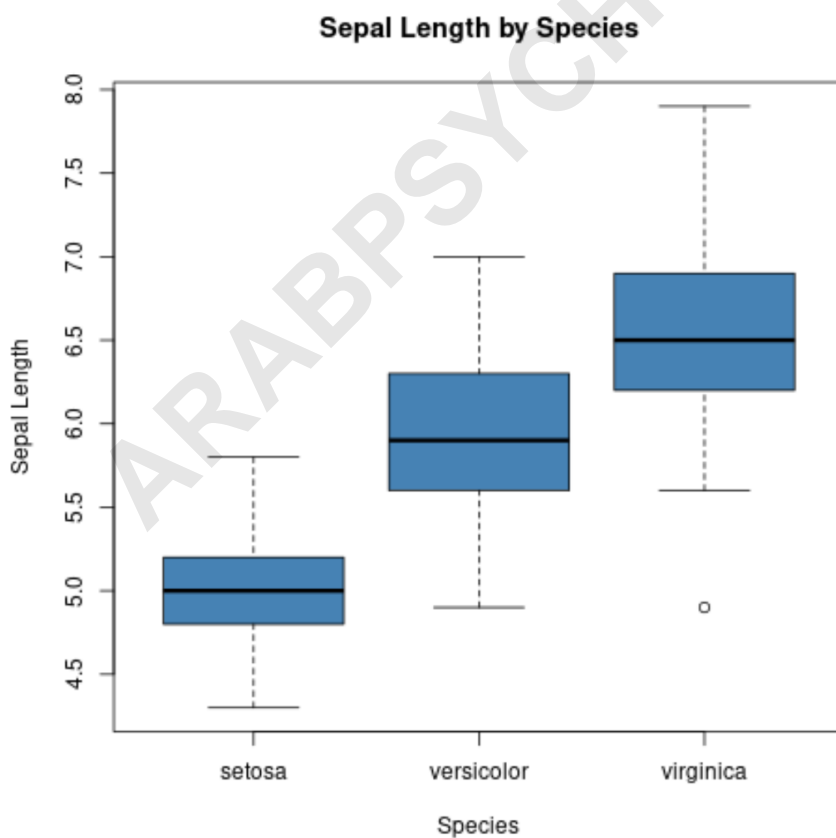
complicates a simple linear interpretation. In real-world data frame analysis, plotting all pairwise combinations is often the first step to identify the most predictive features.

Comparing Feature Distributions Across Species with Boxplots

When dealing with classification problems, it is essential to visualize how a numeric feature is distributed across different categorical groups. The `boxplot()` function is perfectly suited for this, allowing us to compare the median, quartiles, and range of a variable (like `Sepal.Length`) stratified by the grouping factor (`Species`).

Generate a boxplot comparing Sepal Length across the three species

```
boxplot(Sepal.Length~Species,  
data=iris,  
main='Sepal Length Distribution by Species',  
xlab='Iris Species',  
ylab='Sepal Length (cm)',  
col='steelblue',  
border='black')
```



The resulting [boxplot](#) clearly illustrates the distinct statistical differences between the species. The x-axis separates the three flower types, while the y-axis displays the distribution of sepal length for each group. We can immediately deduce that the **virginica** species generally exhibits the largest sepal length measurements, followed by **versicolor**, while **setosa** flowers have the shortest sepals. This high degree of separability between the groups is precisely why the [Iris Dataset](#) remains a perfect benchmark for testing statistical models and introducing concepts in [machine learning](#).

Conclusion and Next Steps in Data Analysis

Mastering the fundamentals of loading, summarizing, and visualizing the [Iris Dataset](#) using base R functions is a foundational skill for any aspiring data scientist. Through simple commands like `head()`, `summary()`, `hist()`, and `boxplot()`, we can quickly derive powerful insights into data structure, feature distributions, and inter-group differences.

The techniques demonstrated in this guide--statistical summarization, viewing feature distributions via the [histogram](#), and comparing groups using the [boxplot](#)--are universally applicable to any new dataset encountered in statistical research or [machine learning](#) projects.

To further advance your skills in data processing and analysis within R, consider exploring advanced topics such as data manipulation using the `dplyr` package, or building predictive models to formally classify the three species of Iris flowers.