

How to Calculate Euclidean Distance in R Using the dist() Function

Authored by
stats writer

December 4, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Calculate Euclidean Distance in R Using the dist() Function*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=105357>

The `dist()` function in R is a fundamental utility in statistical computing, serving the critical purpose of calculating proximity or dissimilarity between observations within a dataset. Specifically, it computes a distance matrix, which quantifies the separation between every row pair (representing individual data points or vectors) in a given matrix or data frame. Understanding these distances is paramount in numerous analytical applications, including clustering algorithms like K-means or hierarchical clustering, where the grouping of similar items is based directly on their computed proximity.

The ability to measure dissimilarity is central to many multivariate analysis techniques. The `dist()` function provides a versatile interface to calculate these measures using a variety of metrics, making it highly adaptable for different types of data and analytical goals. When dealing with numerical data, the choice of the distance metric--whether it is the familiar Euclidean distance or less common options like Canberra--can drastically influence the outcome of subsequent analyses, such as pattern recognition or anomaly detection.

Understanding the `dist()` Function in R

The **`dist()`** function in R is designed to handle the crucial step of generating a distance matrix. This matrix is essential for visualizing relationships between data points, identifying outliers, and serving as input for algorithms that require a measure of proximity between data points. The function operates primarily on the rows of a matrix or data frame, treating each row as a single observation or vector in a high-dimensional space.

By default, `dist()` employs the standard Euclidean distance, but its real power lies in the flexibility offered by its `method` parameter. This allows researchers and analysts to rigorously test which definition of "distance" best captures the underlying structure of their specific dataset, which is a key consideration in complex statistical modeling.

The Core Syntax and Available Methods

The `dist()` function employs a straightforward syntax, yet it offers powerful flexibility through its `method` argument, allowing users to select the most appropriate measure for their data structure. The basic structure requires an input matrix or data frame `x` and an optional specification of the distance calculation method.

The basic structure is:

```
dist(x, method="euclidean")
```

Here is a breakdown of the key parameters:

x: This parameter specifies the input data structure, which must be the name of the matrix or data frame containing the observations (rows) and variables (columns).

method: This crucial parameter determines the mathematical formula used to calculate the distance. The default is "**euclidean**", which is suitable for many standard applications. However, R supports several alternative metrics, including "**maximum**", which is the Chebyshev distance; "**manhattan**", often called the L1 norm; "**canberra**", a weighted version of Manhattan distance; "**binary**", which implements the Jaccard dissimilarity index for binary data; and "**minkowski**", a general framework encompassing Euclidean and Manhattan metrics.

Preparing the Sample Dataset

To provide clear and practical illustrations of how these various distance metrics function, we will utilize a small, defined dataset. This matrix, consisting of four vectors (a, b, c, and d), will serve as the consistent input for all subsequent examples. Each row represents an observation in a four-dimensional space, allowing us to compute the dissimilarity between these observations across the chosen metrics.

The code below defines the four individual vectors and subsequently combines them into a single matrix named `mat`, which we will use for calculating the distance matrix:

```
#define four vectors
```

```
a <- c(2, 4, 4, 6)
```

```
b <- c(5, 5, 7, 8)
```

```
c <- c(9, 9, 9, 8)
```

```
d <- c(1, 2, 3, 3)
```

```
#row bind four vectors into matrix
```

```
mat <- rbind(a, b, c, d)
```

```
#view matrix
```

```
mat
```

```
a 2 4 4 6
```

```
b 5 5 7 8
```

```
c 9 9 9 8
```

```
d 1 2 3 3
```

Example 1: Calculating Euclidean Distance

The **Euclidean distance**, often referred to as the standard straight-line distance or L2 norm, is the most common metric used in geometric and statistical applications. It represents the shortest path

between two points in Euclidean space. In the context of the **dist()** function, this metric is applied by default when the `method` argument is omitted, making it the most frequently used measure of dissimilarity, especially when variables are continuous and similarly scaled.

Mathematically, the Euclidean distance (D) between two vectors, A and B, is calculated by summing the squared differences of their corresponding elements, and then taking the square root of that sum. This calculation is crucial because it accounts for differences across all dimensions equally, providing a holistic measure of separation. The formal definition is given by:

$$\text{Euclidean distance} = \sqrt{\sum(A_i - B_i)^2}$$

The following R code computes the distance matrix, showing the Euclidean distance between every pair of rows in our sample matrix `mat`:

```
#calculate Euclidean distance between each row in matrix  
dist(mat)
```

```
a b c  
b 4.795832  
c 10.148892 6.000000  
d 3.872983 8.124038 13.190906
```

The output is presented in a compact, lower-triangular matrix format. Interpreting this matrix allows us to quantify the similarity or dissimilarity between our observations:

The distance between row a and row b is **4.795832**.

The distance between row a and row c is **10.148892**.

The distance between row a and row d is **3.872983**. This indicates that observations 'a' and 'd' are the closest pair in this dataset according to the Euclidean metric.

The distance between row b and row c is **6.000000**.

The distance between row b and row d is **8.124038**.

The distance between row c and row d is **13.190906**. This demonstrates that observations 'c' and 'd' are the most distant pair.

Example 2: Calculating Maximum Distance

The **Maximum distance** method, also known as Chebyshev distance or L-infinity norm, defines the distance between two vectors as the largest absolute difference among any of their corresponding dimensions. Unlike the Euclidean distance, which averages differences across all dimensions, the Maximum distance is sensitive only to the single dimension where the disparity is greatest. This makes it particularly useful in scenarios where the largest deviation in any single

variable is the most critical measure of dissimilarity, such as quality control where the maximum deviation from a target specification is crucial.

For two vectors A and B, the Maximum distance is calculated simply as: $\max(|A_i - B_i|)$. This metric tends to be less influenced by high dimensionality compared to Euclidean distance, but it heavily weights differences caused by outliers if they occur in corresponding dimensions.

The following code utilizes the `method="maximum"` argument to compute the Maximum distance matrix between the rows of `mat`:

```
#calculate Maximum distance between each row in matrix
```

```
dist(mat, method="maximum")
```

```
a b c  
b 3  
c 7 4  
d 3 5 8
```

For instance, analyzing the distance between 'a' (2, 4, 4, 6) and 'c' (9, 9, 9, 8): the absolute differences are $|2-9|=7$, $|4-9|=5$, $|4-9|=5$, and $|6-8|=2$. The maximum of these differences is 7, which corresponds precisely to the value shown in the resulting distance matrix.

Example 3: Calculating Canberra Distance

The **Canberra distance** is a specialized distance metric that is particularly effective when dealing with data that are subject to scale dependencies or variations near the origin (zero). This metric calculates the sum of fractional differences between paired elements, where the absolute difference is divided by the sum of the absolute values of the elements themselves. Because it normalizes differences by the magnitude of the data points, it is highly sensitive to small changes when both values are close to zero.

The formula for Canberra distance between two vectors A and B is given by:

$$\text{Canberra distance} = \sum |A_i - B_i| / (|A_i| + |B_i|)$$

This metric is often preferred in applications like ecological studies or time series analysis where relative changes are more important than absolute magnitude. Using the `method="canberra"` parameter allows R to apply this specific calculation:

```
#calculate Canberra distance between each row in matrix
```

```
dist(mat, method="canberra")
```

```
a b c
b 0.9552670
c 1.5484515 0.6964286
d 1.1428571 1.9497835 2.3909091
```

The results show that row 'b' and row 'c' are the closest points according to the Canberra metric (0.696). This finding differs from the Euclidean analysis, underscoring the importance of selecting the appropriate distance measure based on the underlying data characteristics and analytical goals.

Example 4: Calculating Binary Distance

The **Binary distance** (or Jaccard distance) is specifically designed for vectors containing binary (0/1) data. If the input matrix contains continuous data, as ours does, R treats any non-zero value as '1' and zero values as '0' before calculating the distance. This metric measures dissimilarity as the proportion of non-shared characteristics (elements where one vector is 1 and the other is 0) relative to the total number of dimensions where at least one vector is 1.

More formally, the Binary distance between two binary vectors, A and B, is calculated as the proportion of elements that are different divided by the total number of non-zero elements present in either vector. It is exceptionally useful for comparing sets or assessing similarity in presence/absence data, such as market basket analysis or biological feature comparison.

Given that our sample data contains only positive, non-zero integer values, the `dist()` function interprets every element as '1' for the purpose of binary comparison. Consequently, the vectors are considered identical in terms of shared presence across all dimensions.

#calculate Binary distance between each row in matrix

```
dist(mat, method="binary")
```

```
a b c
b 0
c 0 0
d 0 0 0
```

Since every element in every vector (a, b, c, d) is non-zero, R treats them as identical binary vectors (1, 1, 1, 1). Thus, the binary distance between all pairs is zero, indicating perfect similarity based on shared non-zero presence. Had our dataset included zeros, the results would reflect the degree of feature overlap between the rows.

Example 5: Calculating Minkowski Distance

The **Minkowski distance** is a powerful generalization that encompasses both the Euclidean distance and the Manhattan distance (City Block distance). It introduces a parameter, p , which dictates the power to which the absolute differences are raised. By adjusting p , the analyst can control the sensitivity of the metric to large coordinate differences, allowing for fine-tuning of the distance measurement to suit the distribution of the data.

The Minkowski distance (D) between two vectors, A and B , is calculated using the following general formula:

$$\text{Minkowski distance} = (\sum |a_i - b_i|^p)^{1/p}$$

In this formula, i represents the i th element in each vector, and the integer p is the parameter defining the norm. This flexibility is what makes it valuable: when $p = 1$, the Minkowski distance is equivalent to the Manhattan distance (L1 norm); when $p = 2$, it is equivalent to the Euclidean distance (L2 norm).

The following R code demonstrates how to compute the distance matrix using the Minkowski metric, specifying a value of $p=3$. This choice means the differences are cubed before being summed, which places a higher emphasis on coordinates that have large deviations compared to the standard Euclidean approach.

#calculate Minkowski distance between each row in matrix

```
dist(mat, method="minkowski", p=3)
```

```
a b c
b 3.979057
c 8.439010 5.142563
d 3.332222 6.542133 10.614765
```

Interpreting and Applying Distance Metrics

Selecting the appropriate distance metric is perhaps the most critical step when using the `dist()` function for statistical analysis. The choice fundamentally dictates how similarity is quantified and subsequently affects the outcome of algorithms that rely on these measurements, such as clustering, manifold learning, and dimensionality reduction. For instance, if data contains many variables with small but meaningful differences, Manhattan distance might be preferred because it is less sensitive to extreme outliers than Euclidean distance.

In practical data science, researchers often experiment with several metrics to determine which

one produces the most meaningful groupings or insights specific to the domain knowledge. The use of the **Minkowski distance** allows analysts to systematically explore the effects of different norms (i.e., different values of p) without changing the underlying computational function. Ultimately, the power of the R `dist()` function lies in providing this variety of robust tools seamlessly within a single command structure.

Summary of `dist()` Utility

The `dist()` function is an essential component of the R statistical environment, offering a robust and flexible method for quantifying dissimilarity between multivariate observations. By allowing users to easily switch between standard metrics like Euclidean and specialized metrics like Canberra or Minkowski, it facilitates advanced exploration of data structure and relationships. Mastering this function is key for anyone involved in statistical modeling, pattern recognition, or machine learning tasks that require a precise measure of object proximity and the computation of a reliable distance matrix.

#define four vectors

```
a <- c(2, 4, 4, 6)
```

```
b <- c(5, 5, 7, 8)
```

```
c <- c(9, 9, 9, 8)
```

```
d <- c(1, 2, 3, 3)
```

```
#row bind four vectors into matrix
```

```
mat <- rbind(a, b, c, d)
```

```
#view matrix
```

```
mat
```

```
a 2 4 4 6
```

```
b 5 5 7 8
```

```
c 9 9 9 8
```

```
d 1 2 3 3
```

```
#calculate Euclidean distance between each row in matrix
```

```
dist(mat)
```

```
a b c
```

```
b 4.795832
```

```
c 10.148892 6.000000
```

d 3.872983 8.124038 13.190906

#calculate Maximum distance between each row in matrix
dist(mat, method="maximum")

a b c
b 3
c 7 4
d 3 5 8

#calculate Canberra distance between each row in matrix
dist(mat, method="canberra")

a b c
b 0.9552670
c 1.5484515 0.6964286
d 1.1428571 1.9497835 2.3909091

#calculate Binary distance between each row in matrix
dist(mat, method="binary")

a b c
b 0
c 0 0
d 0 0 0

#calculate Minkowski distance between each row in matrix
dist(mat, method="minkowski", p=3)

a b c
b 3.979057
c 8.439010 5.142563

d 3.332222 6.542133 10.614765

ARABPSYCHOLOGY.COM