

How to Calculate Bray-Curtis Dissimilarity in R with dist()

Authored by
stats writer

December 1, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Calculate Bray-Curtis Dissimilarity in R with dist()*.
PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=102945>

The Bray-Curtis dissimilarity index is a fundamental metric used extensively in ecological studies and biological research. It provides a robust measure of the compositional difference, or dissimilarity, between two distinct sites or samples. Unlike Euclidean distance, the Bray-Curtis index is specifically designed for count data, such as the abundance of different species, making it highly sensitive to changes in community structure.

Mathematically, this metric quantifies the difference between two datasets by comparing the sum of the absolute differences in species counts across the sites, relative to the total abundance found across both sites combined. This normalization process ensures that the result is bounded and easily interpretable. While the manual calculation provides deep insight into the underlying mechanism, modern statistical analysis often relies on dedicated software environments. In the statistical programming language, R, calculating this index is streamlined, typically involving specialized functions like the base R dist() function or more specialized packages like vegan.

Understanding how to both conceptualize and calculate the Bray-Curtis dissimilarity is essential for researchers aiming to compare biodiversity, track environmental impacts, or analyze community shifts over time or space. This comprehensive guide will detail the formula, demonstrate a manual calculation, and walk through the required steps for efficient computation using R.

Understanding Bray-Curtis Dissimilarity (Definition and Context)

The Bray-Curtis dissimilarity index, sometimes referred to as the Sorensen quantitative index, serves as a cornerstone method for measuring community dissimilarity between two different sites, often denoted as Site i and Site j . This measurement is crucial in fields such as microbial ecology, plant science, and zoology, where researchers need to quantify how different two environments are solely in terms of the biological composition found within them. The method is particularly robust because it accounts for both the presence/absence of species and their relative abundance, providing a quantitative measure rather than just a qualitative comparison.

One of the primary advantages of utilizing the Bray-Curtis dissimilarity is its reliance on abundance data. It is a distance measure derived from the differences in raw counts, making it non-Euclidean and sensitive to the total number of individuals sampled. This contrasts sharply with binary measures (like Jaccard distance) which only consider the mere presence or absence of a species. By incorporating abundance, Bray-Curtis provides a richer representation of community dynamics, highlighting sites that might share the same set of species but differ dramatically in the dominance or rarity of those species.

The measure essentially quantifies the proportion of shared observations relative to the total number of observations, transformed into a dissimilarity value. Its widespread acceptance across ecological studies stems from its intuitive interpretation and its effectiveness in dealing with datasets characterized by many zeros, a common occurrence in biological surveys where many

species are rare or locally absent. It is vital to recognize that while it is often called a distance metric, technically, it only satisfies the properties of a semi-metric, as it does not meet the triangle inequality requirement, a crucial distinction in advanced spatial statistics.

The Mathematical Formula Behind Bray-Curtis

The calculation of Bray-Curtis dissimilarity is performed using a straightforward algebraic formula derived from comparing the minimum counts of shared elements against the total counts observed. It is crucial to understand the components of this formula before attempting any computational implementation. The formula yields a value that reflects the lack of similarity between Site i and Site j .

The formula is presented as follows:

$$BC_{ij} = 1 - (2 * C_{ij}) / (S_i + S_j)$$

Where the variables represent specific ecological counts:

C_{ij}: This represents the sum of the lesser counts for each species found in both sites. For any given species, you take the minimum count between Site i and Site j , and then sum these minimum values across all species recorded in the study. This component essentially captures the shared abundance between the two communities.

S_i: This is the total number of specimens (total abundance) counted across all species at site i . This is calculated simply by summing the counts of all species observed in that specific site.

S_j: Analogously, this is the total number of specimens (total abundance) counted across all species at site j .

Alternatively, the dissimilarity can be expressed by summing the absolute differences in species abundance and normalizing by the total abundance:

$$BC_{ij} = (\text{Sum } |X_{ik} - X_{jk}|) / (\text{Sum } X_{ik} + \text{Sum } X_{jk})$$

Where X_{ik} is the abundance of species k at site i , and X_{jk} is the abundance of species k at site j . Note that the two formulations are mathematically equivalent. The first formulation emphasizes the shared component (C_{ij}), while the second emphasizes the difference component (the numerator containing the absolute differences). Computational implementations, particularly in R, often use the difference formulation for efficiency.

Interpreting the Bray-Curtis Index

A significant strength of the Bray-Curtis metric lies in its standardized range, which always spans between 0 and 1. This normalized range makes interpretation straightforward and allows for direct

comparison across different studies or datasets, regardless of the initial magnitude of species counts. Understanding the endpoints of this scale is essential for drawing meaningful conclusions from the results of ecological studies.

The range endpoints are interpreted as follows:

Value of 0 indicates that the two sites have zero dissimilarity. In practical terms, this means the sites are identical in composition and abundance. They share the exact same set of species, and the relative number of specimens for each species is also precisely the same. A value close to zero suggests a very high degree of similarity between the compared communities.

Value of 1 signifies complete dissimilarity between the two sites. This outcome occurs when the sites share no species whatsoever, meaning the set of species found at Site i is mutually exclusive from the set of species found at Site j . A value approaching one indicates that the communities are entirely different in their species composition.

Intermediate values, such as 0.5, suggest that half of the total abundance observed across both sites represents differences in species composition or abundance skew. Researchers typically use these values to cluster sites into distinct groups, visualizing the relationships using methods like Non-metric Multidimensional Scaling (NMDS) or Principal Coordinates Analysis (PCoA). The Bray-Curtis metric is the most common input for these ordination techniques when analyzing species composition data.

A Practical Example of Manual Calculation

To solidify the understanding of the formula, let us walk through a detailed example of calculating the Bray-Curtis dissimilarity manually. Suppose a marine biologist collects data on the abundance of five different macro-invertebrate species (A, B, C, D, and E) across two distinct sampling locations, Site 1 and Site 2. The raw count data collected illustrates the quantitative nature of the comparison.

The resulting species count data is summarized in the following matrix, where the rows represent the sites and the columns represent the specific species:

		Count of Species				
		A	B	C	D	E
Site 1		4	0	2	7	8
Site 2		3	6	0	4	11

To apply the Bray-Curtis formula, we must first calculate the three core components: C_{ij} , S_i , and S_j .

This preparatory step involves summing the total abundance for each site and determining the minimum shared counts for each species:

Calculation of **C_{ij}** (Sum of minimum shared abundances):

Species A: $\min(4, 3) = 3$

Species B: $\min(0, 6) = 0$

Species C: $\min(2, 0) = 0$

Species D: $\min(7, 4) = 4$

Species E: $\min(8, 11) = 8$

$C_{ij} = 3 + 0 + 0 + 4 + 8 = 15$

Calculation of **S_i** (Total abundance at Site 1): $4 + 0 + 2 + 7 + 8 = 21$

Calculation of **S_j** (Total abundance at Site 2): $3 + 6 + 0 + 4 + 11 = 24$

This intermediate calculation process can often be visualized in a table format for clarity, showcasing the contribution of each species to the overall similarity and dissimilarity between the two sites. The image below helps conceptualize how the minimum values are derived from the raw data:

		Count of Species				
		A	B	C	D	E
Site 1		4	0	2	7	8
Site 2		3	6	0	4	11

$$C_{ij} = 3 + 0 + 0 + 4 + 8 = 15$$

$$S_i = 4 + 0 + 2 + 7 + 8 = 21$$

$$S_j = 3 + 6 + 0 + 4 + 11 = 24$$

We can now plug these derived numerical values back into the standardized Bray-Curtis dissimilarity formula:

$$BC_{ij} = 1 - (2 * C_{ij}) / (S_i + S_j)$$

$$BC_{ij} = 1 - (2 * 15) / (21 + 24)$$

$$BC_{ij} = 1 - (30 / 45)$$

$$BC_{ij} = 1 - 0.6666\dots$$

$$BC_{ij} \approx 0.33$$

The resulting Bray-Curtis dissimilarity value between these two sampling locations is approximately **0.33**. This outcome indicates a moderate level of difference in community structure. Specifically, roughly one-third of the total abundance contributes to the compositional differences, suggesting that while the sites share many elements, Site 2 possesses a higher overall abundance and a different distribution profile for Species B and E compared to Site 1.

Setting Up Data for R Implementation

While manual calculation is excellent for small datasets and conceptual understanding, computational analysis of large ecological matrices requires specialized statistical software. The R environment is the standard platform for processing community data and calculating distance metrics efficiently. Before calculating the Bray-Curtis index in R, the data must be correctly formatted into a data frame or matrix structure.

The fundamental requirement for community ecology data in R is the sites-by-species matrix format. In this structure, each row must represent a distinct sample or site (the observations), and each column must represent a single, unique species or variable (the attributes). The cell values within the matrix must contain the abundance counts for that specific species at that specific site. Failure to transpose data correctly will lead to erroneous results, as standard distance functions in R calculate distances between rows (sites).

Using the data from our previous example, we initiate the process by defining the data frame within the R console:

```
#create data frame
```

```
df <- data.frame(A=c(4, 3),  
B=c(0, 6),  
C=c(2, 0),  
D=c(7, 4),  
E=c(8, 11))
```

```
#view data frame
```

```
df
```

```
A B C D E  
1 4 0 2 7 8  
2 3 6 0 4 11
```

In this created data frame, Row 1 corresponds to Site 1 (or Site *i*) with its associated species counts, and Row 2 corresponds to Site 2 (or Site *j*). The columns A through E correctly represent the five different species observed. This structure is now ready for the application of distance

measurement functions.

Calculating Bray-Curtis Dissimilarity in R (Using Base R Logic)

While the specialized ecological package `vegan` is typically used for complex analyses, it is possible and instructive to calculate the Bray-Curtis index using core base R functions, especially if one is working with a simplified dataset or wishes to verify the underlying algebraic process. This approach relies on applying the absolute difference formulation of the index directly to the data frame.

The calculation requires two key steps: first, calculating the numerator (the sum of absolute differences between sites for each species), and second, calculating the denominator (the sum of total abundance across both sites). The resulting ratio provides the dissimilarity index. This specific manual code structure is often used to demonstrate the equivalence between the mathematical formula and the computational implementation:

#calculate Bray-Curtis dissimilarity

```
sum(apply(df, 2, function(x) abs(max(x)-min(x)))) / sum(rowSums(df))
```

```
0.3333333
```

Let's break down this base R code snippet to confirm its alignment with the formula:

`apply(df, 2, function(x) abs(max(x)-min(x)))`: This command iterates over the columns (the '2' argument) of the data frame `df`. For each species column (vector `x`), it calculates the absolute difference between the counts in Site 1 and Site 2 (which correspond to the max and min values within that column, since there are only two rows/sites). This yields the differences: $|4-3|=1$, $|0-6|=6$, $|2-0|=2$, $|7-4|=3$, $|8-11|=3$.

`sum(...)`: This outer function sums these absolute differences ($1 + 6 + 2 + 3 + 3 = 15$). This represents the numerator ($\sum |X_{ik} - X_{jk}|$).

`sum(rowSums(df))`: This calculates the total abundance across all sites. `rowSums(df)` yields `10`, and summing this gives 45. This is the denominator ($S_i + S_j$).

The final division $15 / 45$ yields 0.3333333.

The computational result, 0.3333333, confirms the value of **0.33** that we previously calculated by hand using the shared abundance formulation. It is important to remember that this specific base R structure is optimized for calculating the distance between exactly two rows (two sites). For calculating a full distance matrix involving many sites, specialized functions are much more efficient and reliable.

Alternative R Packages for Ecological Dissimilarity

For large-scale ecological studies involving dozens or hundreds of sites, the base R approach demonstrated above becomes cumbersome and prone to error. The standard professional method involves using the `vegan` package, which is the cornerstone package for multivariate analysis in ecology. The `vegan` package uses the general purpose dist() function implicitly or explicitly, providing tailored methods for calculating Bray-Curtis and dozens of other ecological distance measures.

When the `vegan` package is loaded, the dist() function (or preferably, the dedicated `vegdist()` function) can be instructed to use the Bray-Curtis method. This calculates the full distance matrix, where every site is compared against every other site. Using `vegdist()` is considered best practice:

```
# Install and load the vegan package (if not already installed)
```

```
# install.packages("vegan")
```

```
library(vegan)
```

```
# Calculate the full distance matrix using vegdist
```

```
bray_dist_matrix <- vegdist(df, method="bray")
```

```
# View the result (distance between Site 1 and Site 2)
```

```
bray_dist_matrix
```

```
# 1
```

```
# 2 0.3333333
```

Alternatively, the base dist() function can sometimes be adapted, although it relies on specific package implementations or specific arguments that might vary depending on the environment. The `vegdist()` function guarantees the correct implementation of the Bray-Curtis formula as defined in the ecological literature. Furthermore, `vegdist()` is optimized to handle large matrices efficiently, making it the preferred choice for multivariate statistical tasks such as ordination (NMDS) or cluster analysis.

Considerations and Limitations

While the Bray-Curtis dissimilarity index is powerful, it is not without its considerations and limitations, particularly when dealing with zero values and standardization. One key feature to note is that Bray-Curtis is inherently sensitive to differences in total community size (or sequencing depth, in genomics). If Site 1 has a total count of 100 individuals and Site 2 has 1000, this disparity in effort or productivity will influence the resulting dissimilarity, even if the relative proportions of

species are somewhat similar. Therefore, it is often recommended to standardize the input data (e.g., transform counts to relative proportions or perform rarefaction) before calculating Bray-Curtis, especially when comparing samples with vastly different total abundances.

Another crucial distinction is that the Bray-Curtis index is not truly metric; as mentioned earlier, it fails the triangle inequality. This is important in advanced spatial analysis, as certain geometric procedures rely on fully metric distances. However, for practical purposes in descriptive ecological studies and standard ordination methods, this semi-metric property rarely causes problems and is generally accepted.

Finally, the Bray-Curtis index is sensitive to zeros. If a species is absent from both sites, it contributes zero to both the numerator (absolute difference) and the denominator (total abundance), thus having no influence on the dissimilarity result. This contrasts with distance metrics like Euclidean distance, which can inflate the perceived difference due to shared zeros. This 'double-zero' exclusion makes Bray-Curtis highly suitable for community data, where shared absences are common and typically uninformative for community structure comparison.

Conclusion: Mastering Ecological Distance Metrics

The dist() function, when properly employed either directly or via specialized ecological packages like `vegan`, provides a reliable computational approach for quantifying community differences. The Bray-Curtis index stands out as the preferred measure in quantitative ecology due to its intuitive interpretation, its sensitivity to species abundance, and its robust handling of sparse count data typical of biodiversity surveys.

By mastering both the theoretical foundation--understanding the roles of shared abundance and total counts--and the practical implementation in R, researchers can effectively characterize the structural differences between biological communities. Whether performing manual validation or running large-scale multivariate analyses, the Bray-Curtis dissimilarity index remains an indispensable tool for comparative distance measurement in environmental and biological sciences.