

# Calculate Hamming Distance in Python (With Examples)?

Authored by  
**stats writer**

December 15, 2025

## RECOMMENDED CITATION

stats writer (2025). *Calculate Hamming Distance in Python (With Examples)?*.  
PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=107546>

The concept of the Hamming distance is fundamental in fields ranging from information theory and coding to bioinformatics and data analysis. Coined by Richard Hamming, this metric provides a simple yet powerful way to quantify the difference between two sequences of equal length. Specifically, the **Hamming distance** measures the minimum number of substitutions required to change one sequence into the other, operating under the strict condition that both sequences--whether they are binary strings, numerical arrays, or character lists--must possess the exact same length.

In the context of data science and computation using Python, calculating this distance is essential for tasks like error detection in telecommunications or comparing DNA sequences. While manual calculation is straightforward for small data sets, leveraging optimized libraries becomes necessary when dealing with large-scale data. For robust and efficient calculation of the Hamming distance in Python, we predominantly rely on the specialized functions provided by the SciPy ecosystem, specifically within the `scipy.spatial.distance` module.

Before diving into the implementation, it is crucial to understand the definition precisely: the Hamming distance between two vectors is merely the summation of corresponding elements where differences exist. If two elements in the same position are not identical, they contribute 1 to the total distance count. This tutorial will explore how to apply this calculation effectively in Python, addressing the specific nuances of the SciPy implementation, which are often overlooked by newcomers.

To illustrate this principle, consider the comparison of two simple vectors,  $\mathbf{x}$  and  $\mathbf{y}$ , where we identify the corresponding positions that contain differing elements:

$\mathbf{x} =$

$\mathbf{y} =$

By comparing these element by element (index 0: 1=1; index 1: 2=2; index 2: 3!=5; index 3: 4!=7), we observe two mismatches. Therefore, the absolute Hamming distance between the two vectors would be **2**, representing the total number of non-matching elements across all corresponding positions.

## Implementing Hamming Distance using SciPy

While the theoretical calculation of the Hamming distance is based on counting absolute differences, the standard implementation within the SciPy library requires careful handling. SciPy is a core library in the Python scientific stack, providing advanced mathematical algorithms. The function we use is `scipy.spatial.distance.hamming()`, which is designed to efficiently compute the disparity between two input arrays or sequences.

A critical distinction of the SciPy implementation is that the `hamming()` function does not return the raw count of differences, but rather the **proportion** of elements that differ between the two arrays. This result is a float value between 0.0 and 1.0. For instance, if an array has 10 elements and 3 are mismatched, the function returns 0.3. This standardized proportional output is often useful for applications requiring normalized metrics, but if the absolute Hamming distance (the count) is needed, a conversion step is mandatory.

To retrieve the true, absolute Hamming distance count, we must multiply the proportional result returned by the SciPy function by the total length of the input array. Since the Hamming distance requires equal-length inputs, multiplying by the length of either array (`len(array1)` or `len(array2)`) will yield the correct integer count. Understanding this conversion is key to accurate results when using this library in Python projects.

## Syntax and Conversion for Absolute Distance

The base syntax for calling the function is straightforward, requiring two input sequences, `array1` and `array2`, which are typically passed as NumPy arrays or standard Python lists. However, since we are usually interested in the absolute count of dissimilarities, the multiplication factor must be integrated into the final calculation expression.

The core function call looks like this:

```
scipy.spatial.distance.hamming(array1, array2)
```

As established, this function returns the **percentage** (proportion) of corresponding elements that differ between the two arrays.

To obtain the desired absolute distance, we must append the multiplication step. This ensures that the output reflects the actual number of mismatched positions, suitable for direct interpretation in contexts like error counting.

Thus, to obtain the absolute Hamming distance (the count), we simply multiply the proportional output by the length of one of the arrays:

```
scipy.spatial.distance.hamming(array1, array2) * len(array1)
```

The following examples demonstrate this conversion process across different data types--binary, numerical, and categorical (string) arrays--showcasing the function's versatility within the Python environment.

## Calculation Example 1: Binary Data Comparison

The most traditional application of the Hamming distance involves comparing binary sequences, often found in digital communications or computer memory storage where data is represented by bits (0s and 1s). The distance in this context directly corresponds to the number of bit errors or positions where the bits flip state between the two sequences. This use case highlights the metric's efficacy as an error detection mechanism.

For this example, we define two binary arrays,  $\bar{x}$  and  $\bar{y}$ , both having a length of six elements. We first import the necessary `hamming` function from the `scipy.spatial.distance` module, a standard requirement when utilizing components of the SciPy library in Python scripts. We then execute the function and multiply the result by the length of  $\bar{x}$  (which is 6) to obtain the absolute count of differences.

The following code demonstrates the setup and calculation:

```
from scipy.spatial.distance import hamming

#define arrays
x =
y =

#calculate Hamming distance between the two arrays
hamming(x, y) * len(x)

2.0
```

Upon inspection, we compare the arrays:

Position 1:  $\bar{x}$  (1) vs  $\bar{y}$  (0) -> Mismatch (Count: 1)

Position 5:  $\bar{x}$  (1) vs  $\bar{y}$  (0) -> Mismatch (Count: 2)

All other elements are identical (0 and 0, 1 and 1, 1 and 1, 0 and 0).

The proportional result returned by `hamming(x, y)` would be  $2/6$ , or approximately 0.3333. Multiplying this by 6 gives 2.0. Therefore, the absolute Hamming distance between the two binary arrays is confirmed to be **2**.

## Calculation Example 2: Handling General Numerical Arrays

The utility of the Hamming distance extends beyond binary data. It can be applied effectively to any set of sequences where position-wise comparison of discrete values is meaningful. When dealing with numerical arrays containing integers, the calculation remains identical to the binary case: we

only care if the element at position  $i$  in array A is strictly equal to the element at position  $i$  in array B. The magnitude of the difference (e.g., 5 vs 10) is irrelevant; only the state of being different (a mismatch) contributes 1 to the distance count.

It is crucial to differentiate the Hamming distance from other metrics, such as the Euclidean distance or Manhattan distance, which measure the numerical magnitude of differences between points in space. The Hamming distance is purely a measure of categorical disparity--the number of positions where two values fail the equality test. This makes it suitable for comparing data records, IDs, or features that are ordinal or discrete.

In this second example, we define two numerical arrays,  $x$  and  $y$ , each containing five integer elements. We proceed with the same method of importing the function and calculating the absolute distance by multiplying the proportion by the array length:

```
from scipy.spatial.distance import hamming
```

```
#define arrays
```

```
x =
```

```
y =
```

```
#calculate Hamming distance between the two arrays
```

```
hamming(x, y) * len(x)
```

```
3.0
```

Analyzing the result, a distance of 3.0 confirms that three positions contain mismatched values:

```
Index 0 (7 vs 7) -> Match
```

```
Index 1 (12 vs 12) -> Match
```

```
Index 2 (14 vs 16) -> Mismatch (Count: 1)
```

```
Index 3 (19 vs 26) -> Mismatch (Count: 2)
```

```
Index 4 (22 vs 27) -> Mismatch (Count: 3)
```

This result illustrates the straightforward application of the `hamming` function in Python for general numeric data comparison, emphasizing that the distance ignores the size of the numerical gap between the differing elements.

### Calculation Example 3: Comparing String or Categorical Sequences

The Hamming distance is often applied to sequences composed of characters or other categorical labels, such as those found in DNA sequencing (A, T, C, G) or comparing textual representations. When `scipy.spatial.distance.hamming` processes arrays containing strings, it performs

element-wise comparison based on equality, treating each string element as a discrete category. This application is particularly common in bioinformatics, where comparing genetic sequences of the same length is critical for evolutionary analysis or identifying mutations.

For the Hamming distance to be valid for strings, it is usually applied to sequences that have been broken down into constituent elements, such as lists of characters or words, ensuring that the length requirement is met. If comparing two full strings like "KITTEN" and "SITTIN", one must compare them character by character. In Python, treating the strings as lists of characters allows SciPy to correctly handle the comparison.

In this example, we compare two arrays,  $x$  and  $y$ , which are composed of four character strings. Note that the arrays are identical except for the final element, which is the sole point of difference:

```
from scipy.spatial.distance import hamming
```

```
#define arrays
```

```
x =
```

```
y =
```

```
#calculate Hamming distance between the two arrays
```

```
hamming(x, y) * len(x)
```

```
1.0
```

The calculated absolute Hamming distance is **1**. This signifies that only one substitution is necessary to transform array  $x$  into array  $y$  (or vice versa). Specifically, the character 'd' at index 3 in  $x$  must be changed to 'r' to match  $y$ . This confirms the flexibility of the SciPy `hamming` function in handling various discrete data types, provided they are structured as comparable sequences of equal length.

## Limitations and Prerequisites for Hamming Distance

While the Hamming distance is highly useful, its application is strictly constrained by two major prerequisites. First and foremost, the two sequences being compared must be of **exactly equal length**. This constraint is integral to the definition, as the distance is determined by counting mismatches at corresponding positions. If the sequences differ in length, the comparison breaks down, and alternative distance metrics, such as the Levenshtein distance (which accounts for insertions and deletions), must be employed instead.

Secondly, the Hamming distance is only suitable for comparing categorical or discrete variables. It measures difference based on identity (A equals B or A does not equal B). If the data represents

continuous numerical values where the magnitude of the difference matters (e.g., comparing temperatures or spatial coordinates), metrics like Euclidean distance are far more appropriate. Applying the Hamming distance to continuous data essentially reduces a complex difference in magnitude to a binary mismatch state, potentially discarding critical information about the underlying relationship between the data points.

Furthermore, users must be mindful of data type handling in Python, especially when arrays contain mixed or complex objects. SciPy's `hamming` function relies on the underlying comparison operators defined for the elements. If attempting to compare strings where capitalization matters (e.g., 'A' versus 'a'), these will correctly register as a mismatch. Proper preprocessing, such as normalization (converting all characters to lowercase), might be required depending on whether the comparison should be case-sensitive or case-insensitive.

## Manual Calculation and Efficiency Considerations

For binary sequences, the Hamming distance can also be calculated manually using bitwise operations, which can be highly efficient, especially when working directly with NumPy arrays. The core idea is utilizing the **Exclusive OR (XOR)** operation. When two bits are XORed, the result is 1 if they are different (a mismatch) and 0 if they are the same (a match). By applying XOR across the two sequences and then summing the resulting array, we instantly obtain the absolute Hamming distance.

In Python, if the input sequences `x` and `y` are NumPy arrays of integers or booleans, the manual calculation involves `np.sum(x != y)` or `np.count_nonzero(x - y)`. While such manual methods offer excellent performance for optimized data structures like NumPy arrays, the `scipy.spatial.distance.hamming` function provides a robust, standardized interface that handles type checking and formatting concerns seamlessly, making it generally preferred for heterogeneous data types (such as the string arrays shown in Example 3).

When choosing between a specialized SciPy function and a manual NumPy operation, the decision often hinges on performance requirements and maintainability. For massive datasets, highly optimized NumPy operations might provide a slight edge, but for standard applications in data analysis where clear code and reliability across different discrete data types are paramount, the SciPy implementation, despite the extra conversion step (multiplying by length), offers superior flexibility and adherence to standardized scientific computing practices.

## Applications of Hamming Distance in Data Science

The application scope of the Hamming distance is broad and impactful across various technical disciplines. In **telecommunications**, it is fundamental to error detection and correction. By appending check bits to data words, a receiver can calculate the Hamming distance between the

received data and known valid code words. A distance of zero means no error, and a small, non-zero distance (often 1 or 2) indicates detectable and potentially correctable errors.

In **bioinformatics**, sequences like DNA or RNA are often compared using this metric. For instance, comparing two gene sequences of the same length allows researchers to quantify genetic divergence or identify the number of nucleotide substitutions that have occurred between two related species or individuals. A lower Hamming distance suggests closer genetic similarity, provided that only substitutions (and not insertions or deletions) are considered relevant.

Furthermore, in **machine learning and data retrieval**, Hamming distance is used for feature comparison, especially when dealing with encoded categorical features or binary hashes (LSH - Locality-Sensitive Hashing). It is employed to quickly estimate the similarity between two hashed documents or images. Its simplicity and computational speed make it an excellent choice for preprocessing steps where rapid dissimilarity quantification is necessary before invoking more complex algorithms.