



## 9 Distance Measures in Data Science

Distance measures in data science refer to algorithms that quantify the similarity or dissimilarity between two or more objects. These algorithms are commonly used in a wide range of data science applications, including clustering, classification, recommendation systems, and more.

The choice of distance measure can have a significant impact on the performance of a data science model. It is important to carefully consider which distance measure is most appropriate for a given problem, as different distance measures may be more or less suitable depending on the characteristics of the data.

In this article, we will explore nine different distance measures that are commonly used in data science. We will discuss the definition, formula, and pros and cons of each distance measure, and provide examples to illustrate how they can be applied. By the end of this article, you should have a solid understanding of the different distance measures available and how to choose the right one for your data science problem.

### Euclidean Distance

Euclidean distance, also known as L2-Norm, is a measure of the straight-line distance between two points in Euclidean space. It is calculated as the square root of the sum of the squares of the differences between the coordinates of the points.

The formula for Euclidean distance between two points  $p$  and  $q$  is as follows:

$$d(p, q) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2}$$

where  $p$  and  $q$  are the coordinates of the two points, and  $n$  is the number of dimensions.

For example, suppose we have two points in two-dimensional space,  $p(1, 2)$  and  $q(4, 6)$ . The Euclidean distance between these two points can be calculated as follows:

$$d(p, q) = \sqrt{(4 - 1)^2 + (6 - 2)^2} = \sqrt{9 + 16} = \sqrt{25} = 5$$

Euclidean distance is a commonly used distance measure because it is easy to understand and compute. It is also well-suited for continuous variables and data with a Euclidean structure, such as images.

## Manhattan Distance

Manhattan distance, also known as L1-Norm or taxicab norm, is a measure of the distance between two points in a grid-like structure, such as a city block. It is calculated as the sum of the absolute differences between the coordinates of the points.

The formula for the Manhattan distance between two points  $p$  and  $q$  is as follows:

$$d(p, q) = |q_1 - p_1| + |q_2 - p_2| + \dots + |q_n - p_n|$$

where  $p$  and  $q$  are the coordinates of the two points, and  $n$  is the number of dimensions.

For example, suppose we have two points in two-dimensional space,  $p(1, 2)$  and  $q(4, 6)$ . The Manhattan distance between these two points can be calculated as follows:

$$d(p, q) = |4 - 1| + |6 - 2| = 3 + 4 = 7$$

Manhattan distance is a popular choice for data with a grid-like structure, such as text data or image data. It is also less sensitive to outliers than Euclidean distance and may be more appropriate for data with skewed distributions.

## Cosine Similarity

Cosine similarity is a measure of similarity between two non-zero vectors of an inner product space that measures the cosine of the angle between them. It is commonly used in data science to compare the similarity of documents, such as articles or reviews, based on the vector space model of document representation.

The formula for cosine similarity between two vectors  $p$  and  $q$  is as follows:

$$\cos(p, q) = (p \cdot q) / (||p|| \cdot ||q||)$$

where  $p$  and  $q$  are the vectors,  $*$  represents the dot product, and  $||p||$  and  $||q||$  represent the magnitudes of the vectors.

For example, suppose we have two vectors  $p$  and  $q$  represented as follows:

$$p = [1, 2, 3]$$

$$q = [4, 5, 6]$$

The cosine similarity between these two vectors can be calculated as follows:

$$\cos(p, q) = (1 * 4 + 2 * 5 + 3 * 6) / (\sqrt{1^2 + 2^2 + 3^2} * \sqrt{4^2 + 5^2 + 6^2}) = 32 / (\sqrt{14} * \sqrt{77}) = 32 / (7.81 * 8.77) = 0.84$$

Cosine similarity ranges from -1 to 1, where 1 indicates that the vectors are identical, 0 indicates that the vectors are orthogonal (perpendicular) and have no similarity, and -1 indicates that the vectors are opposed and have maximum dissimilarity.

Cosine similarity is a popular choice for comparing the similarity of text data, as it is insensitive to the magnitude of the vectors and only considers the orientation of the vectors. It is also efficient to compute and does not require the vectors to be normalized.

## Jaccard Index

The Jaccard index, also known as the Jaccard coefficient, is a measure of the similarity between two sets. It is calculated as the size of the intersection of the sets divided by the size of the union of the sets.

The formula for the Jaccard index between two sets  $A$  and  $B$  is as follows:

$$J(A, B) = |A \cap B| / |A \cup B|$$

where  $|A \cap B|$  is the number of elements that are common to both sets  $A$  and  $B$ , and  $|A \cup B|$  is the total number of elements in both sets.

For example, suppose we have two sets  $A$  and  $B$  represented as follows:

$$A = \{1, 2, 3, 4\}$$

$$B = \{3, 4, 5, 6\}$$

The Jaccard index between these two sets can be calculated as follows:

$$J(A, B) = |\{3, 4\}| / |\{1, 2, 3, 4, 5, 6\}| = 2 / 6 = 1/3$$

The Jaccard index ranges from 0 to 1, where 1 indicates that the sets are identical and 0 indicates that the sets have no elements in common.

The Jaccard index is a popular choice for comparing the similarity of categorical data, as it only considers the presence or absence of elements in the sets and is insensitive to the order or magnitude of the elements. It is also efficient to compute and does not require the sets to be normalized.

## Hamming Distance

Hamming distance is a measure of the difference between two strings of equal length. It is calculated as the number of positions at which the corresponding symbols are different.

The formula for Hamming distance between two strings  $s$  and  $t$  is as follows:

$$d(s, t) = \sum (s_i \neq t_i \text{ for } s_i, t_i \text{ in } \text{zip}(s, t))$$

where  $s$  and  $t$  are the strings, and  $\text{zip}$  is a function that returns an iterator of tuples, where the  $i$ -th tuple contains the  $i$ -th element from each of the input iterables.

For example, suppose we have two strings  $s$  and  $t$  represented as follows:

$s = \text{"abcdef"}$

$t = \text{"abcxyz"}$

The Hamming distance between these two strings can be calculated as follows:

$$d(s, t) = \sum (s_i \neq t_i \text{ for } s_i, t_i \text{ in } \text{zip}(s, t)) = \sum (\text{True}, \text{True}, \text{True}, \text{False}, \text{False}, \text{False}) = 3$$

The Hamming distance is a popular choice for comparing the difference between strings, such as DNA sequences or error-correcting codes. It is also efficient to compute and does not require the strings to be normalized.

# Minkowski Distance

Minkowski distance is a generalized form of the Euclidean distance and the Manhattan distance. It is a measure of the distance between two points in a Euclidean space and is defined as the sum of the absolute differences of their coordinates raised to the power of  $p$  and then taking the  $p$ th root of the result.

The formula for the Minkowski distance between two points  $x$  and  $y$  in an  $n$ -dimensional space is as follows:

$$d(x, y) = (\sum |x_i - y_i|^p)^{1/p}$$

where  $x$  and  $y$  are the points,  $x_i$  and  $y_i$  are the  $i$ -th coordinates of the points  $x$  and  $y$ , respectively, and  $p$  is a positive integer parameter called the Minkowski exponent.

When  $p = 1$ , the Minkowski distance reduces to the Manhattan distance, and when  $p = 2$ , it reduces to the Euclidean distance. For other values of  $p$ , the Minkowski distance is referred to as the generalized Minkowski distance.

Suppose we have two points  $x$  and  $y$  in a two-dimensional space represented as follows:

$$x = (3, 4)$$

$$y = (6, 8)$$

We can calculate the Minkowski distance between these two points using the following formula:

$$d(x, y) = (\sum |x_i - y_i|^p)^{1/p}$$

where  $p$  is a positive integer parameter called the Minkowski exponent.

For example, if we set  $p = 1$ , the Minkowski distance reduces to the Manhattan distance, which is calculated as follows:

$$d(x, y) = (|3 - 6| + |4 - 8|) = (3 + 4) = 7$$

If we set  $p = 2$ , the Minkowski distance reduces to the Euclidean distance, which is calculated as follows:

$$d(x, y) = \sqrt{(3 - 6)^2 + (4 - 8)^2} = \sqrt{9 + 16} = \sqrt{25} = 5$$

The Minkowski distance is a useful measure of distance in many applications, including data clustering, pattern recognition, and machine learning. It is also efficient to compute and is not sensitive to the scale of the coordinates.

## Chebyshev Distance

Chebyshev Distance, also known as the Chessboard Distance or Tchebychev Distance, is a measure of distance between two points in a multidimensional space. It is defined as the maximum of the absolute differences between the coordinates of the two points. This distance measure is often used in cases where the shape of the data is not known and the distance measure should not be affected by the scale of the variables. It has a variety of applications, including image processing, pattern recognition, and machine learning.

To calculate the Chebyshev distance between two points  $x$  and  $y$ , with coordinates  $(x_1, x_2, \dots, x_n)$  and  $(y_1, y_2, \dots, y_n)$ , respectively, we use the following formula:

$$d(x, y) = \max(|x_1 - y_1|, |x_2 - y_2|, \dots, |x_n - y_n|)$$

For instance, let's consider two points in a 2D space with coordinates  $(2, 3)$  and  $(5, 7)$ . The Chebyshev distance between these two points is:

$$d((2, 3), (5, 7)) = \max(|2 - 5|, |3 - 7|) = \max(3, 4) = 4$$

The Chebyshev distance is a metric, meaning that it satisfies the following properties:

$$d(x, y) \geq 0 \text{ (non-negativity)}$$

$$d(x, y) = 0 \text{ if and only if } x = y \text{ (identity of indiscernibles)}$$

$$d(x, y) = d(y, x) \text{ (symmetry)}$$

$$d(x, z) \leq d(x, y) + d(y, z) \text{ (triangle inequality)}$$

## Haversine Distance

Haversine Distance, also known as Great Circle Distance, is a measure of the distance between two points on the surface of a sphere. It is commonly used to calculate the distance between two points on the Earth's surface, such as the distance between two cities.

The formula for Haversine Distance between two points x and y, with coordinates (latitude1, longitude1) and (latitude2, longitude2), respectively, is as follows:

$$d(x, y) = 2 * R * \text{asin}(\sqrt{\sin^2((\text{latitude2} - \text{latitude1})/2) + \cos(\text{latitude1}) * \cos(\text{latitude2}) * \sin^2((\text{longitude2} - \text{longitude1})/2)})$$

where R is the radius of the sphere (e.g., 6371 km for the Earth), and asin, sin, and cos are the inverse sine, sine, and cosine functions, respectively.

For example, let's consider two points on the Earth's surface with coordinates (40.7128° N, 74.0060° W) and (35.6895° N, 139.6917° E). The Haversine Distance between these two points is:

$$d((40.7128^\circ \text{ N}, 74.0060^\circ \text{ W}), (35.6895^\circ \text{ N}, 139.6917^\circ \text{ E})) = 2 * 6371 \text{ km} * \text{asin}(\sqrt{\sin^2((35.6895^\circ - 40.7128^\circ)/2) + \cos(40.7128^\circ) * \cos(35.6895^\circ) * \sin^2((139.6917^\circ - 74.0060^\circ)/2)}) = 10850 \text{ km}$$

## Sørensen-Dice Index

Sørensen-Dice Index, also known as Sørensen Index or Dice's Coefficient, is a measure of the similarity between two sets. It is a widely used measure in various fields such as information retrieval, data mining, and natural language processing.

The Sørensen-Dice Index is calculated using the following formula:

$$SDI(A, B) = 2 * |A \cap B| / (|A| + |B|)$$

where A and B are the two sets, |A| and |B| are the number of elements in each set, and  $A \cap B$  is the intersection of the two sets (the elements that are common to both sets).

To better understand the Sørensen-Dice Index, let's consider an example. Suppose set A contains the elements {apple, banana, cherry, dragonfruit} and set B contains the elements {apple, cherry, lemon, orange}. The Sørensen-Dice Index of these two sets can be calculated as follows:

$$SDI(\{\text{apple, banana, cherry, dragonfruit}\}, \{\text{apple, cherry, lemon, orange}\}) = 2 * |\{\text{apple, cherry}\}| / (4 + 4) = 2 * 2 / 8 = 0.5$$

This means that the Sørensen-Dice Index of these two sets is 0.5, or 50%. This tells us that there is a 50% overlap between the elements in the two sets.

The Sørensen-Dice Index ranges from 0 to 1, where 0 indicates that the sets have no common elements and 1 indicates that the sets are identical. It is a useful measure when comparing the similarity of categorical data, such as the presence or absence of certain keywords in a document.

One important property of the Sørensen-Dice Index is that it is symmetric, meaning that the similarity between two sets is the same regardless of the order of the sets. This is in contrast to measures such as Jaccard Index, which is not symmetric. Another advantage of the Sørensen-Dice Index is that it is easy to interpret and understand. It gives a clear and intuitive sense of the overlap between two sets and is therefore widely used in various applications.

## Conclusion

here are several distance measures that are commonly used in data science to compare the similarity or dissimilarity between two or more data points. These measures include Euclidean distance, Manhattan distance, Minkowski distance, Mahalanobis distance, Hamming distance, Levenshtein distance, Chebyshev distance, Haversine distance, and Sørensen-Dice index. Each measure has its own strengths and limitations, and it is important to choose the appropriate measure based on the nature and characteristics of the data being compared.

If you are looking to take your data science skills to the next level and learn more about these and other advanced techniques, consider enrolling in Skillslash's [Advanced Data Science and AI program](#). This comprehensive program covers a wide range of topics including machine learning, deep learning, natural language processing, and more. You will gain the knowledge and skills you need to succeed in today's competitive data science job market and make a meaningful impact in your career. Don't miss this opportunity to take your data science career to new heights. Enroll today!

Overall, **Skillslash** also has in store, exclusive courses like [Data Science Course In jaipur](#), [Best system design Course](#) to ensure aspirants of each domain have a great learning journey and a secure future in these fields. To find out how you can make a career in the IT and tech field with Skillslash, contact the student support team to know more about the course and institute.



# 9 DISTANCE MEASURES IN DATA SCIENCE

[www.skillslash.com](http://www.skillslash.com)

