

Measures of Distance in Data Mining

Clustering consists of grouping certain objects that are similar to each other, it can be used to decide if two items are similar or dissimilar in their properties.

In a **Data Mining** sense, the similarity measure is a distance with dimensions describing object features. That means if the distance among two data points is **small** then there is a **high** degree of similarity among the objects and vice versa. The similarity is **subjective** and depends heavily on the context and application. For example, similarity among vegetables can be determined from their taste, size, colour etc.

Most clustering approaches use distance measures to assess the similarities or differences between a pair of objects, the most popular distance measures used are:

1. Euclidean Distance:

Euclidean distance is considered the traditional metric for problems with geometry. It can be simply explained as the **ordinary distance** between two points. It is one of the most used algorithms in the cluster analysis. One of the algorithms that use this formula would be **K-mean**. Mathematically it computes the **root of squared differences** between the coordinates between two objects.

$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \cdots + (q_n - p_n)^2}$$
$$= \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

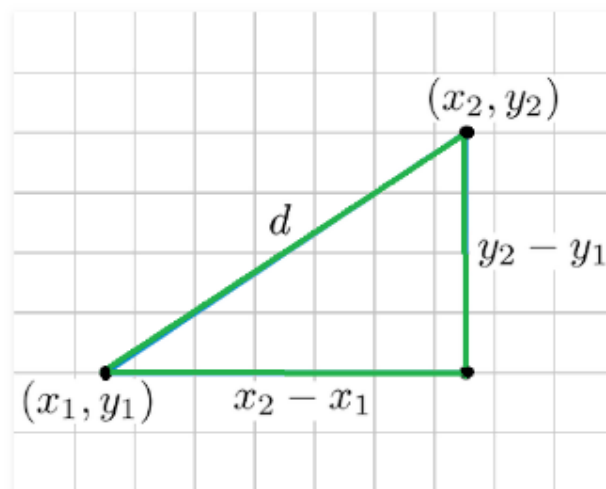


Figure – Euclidean Distance

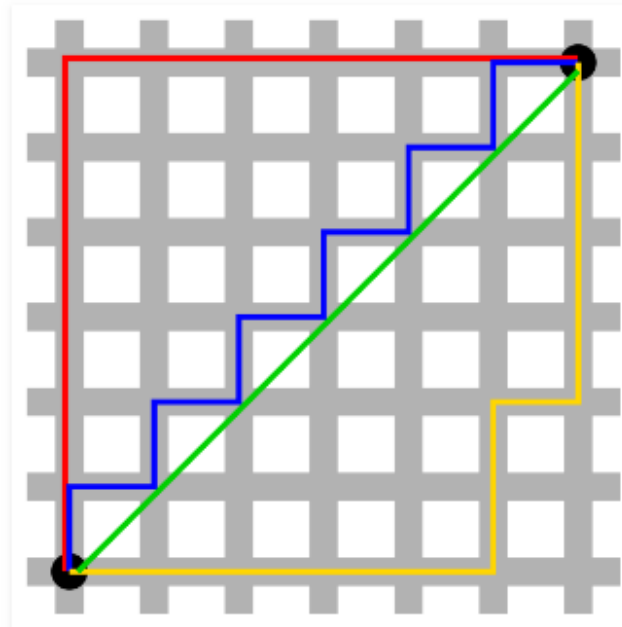
2. Manhattan Distance:

This determines the absolute difference among the pair of the coordinates.

Suppose we have two points P and Q to determine the distance between these points we simply have to calculate the perpendicular distance of the points from X-Axis and Y-Axis.

In a plane with P at coordinate (x_1, y_1) and Q at (x_2, y_2) .

Manhattan distance between P and Q = $|x_1 - x_2| + |y_1 - y_2|$



Here the total distance of the **Red** line gives the Manhattan distance between both the points.

3. Jaccard Index:

The Jaccard distance measures the similarity of the two data set items as the **intersection** of those items divided by the **union** of the data items.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

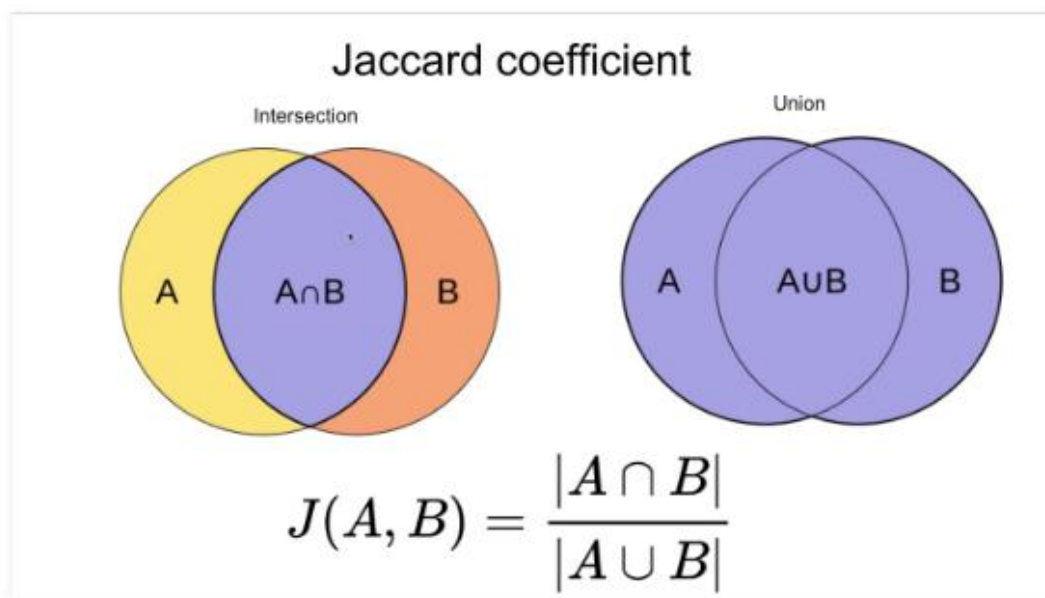


Figure – Jaccard Index

4. Minkowski distance:

It is the **generalized** form of the Euclidean and Manhattan Distance Measure. In an **N-dimensional space**, a point is represented as,

$$(x_1, x_2, \dots, x_N)$$

Consider two points P1 and P2:

$$P1: (x_1, x_2, \dots, x_N)$$

$$P2: (y_1, y_2, \dots, y_N)$$

Then, the Minkowski distance between P1 and P2 is given as:

$$\sqrt[p]{(x_1 - y_1)^p + (x_2 - y_2)^p + \dots + (x_N - y_N)^p}$$

- When **p = 2**, Minkowski distance is same as the **Euclidean** distance.
- When **p = 1**, Minkowski distance is same as the **Manhattan** distance.

5. Cosine Index:

Cosine distance measure for clustering determines the **cosine** of the angle between two vectors given by the following formula.

$$\text{sim}(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

Here (**theta**) gives the angle between two vectors and A, B are n-dimensional vectors.

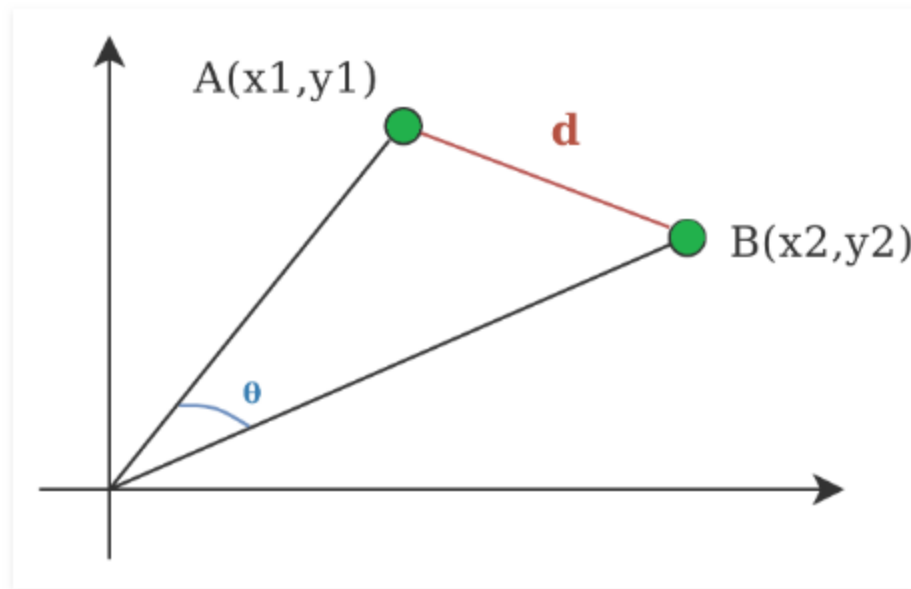


Figure – Cosine Distance

KNN CLASSIFIER

Eager vs Lazy learners

- Eager learner
 - Generalized model from training data set is constructed
 - Using the model the class of test data set is predicted
 - Example – decision tree
- Lazy learner
 - Training dataset is stored
 - On querying similarity between test data and training set records is calculated to predict the class of test data
 - Example – k-Nearest Neighbour

k-NN

- Non-parametric method used for classification
- Prediction for test data is done on the basis of its neighbour
- k is an integer (small), if k=1, k is assigned to the class of single nearest neighbour

Example

Name	Acid Durability	Strength	Class
Type-1	7	7	Bad
Type-2	7	4	Bad
Type-3	3	4	Good
Type-4	1	4	Good

Test-Data ➔ acid durability=3, and strength=7, class=?

- Calculated using distance measure like Euclidean

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

$$= \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

Name	Acid Durability	Strength	Class	Distance
Type-1	7	7	Bad	Sqrt((7-3) ² +(7-7) ²)=4
Type-2	7	4	Bad	5
Type-3	3	4	Good	3
Type-4	1	4	Good	3.6

Rank these attributes

Name	Acid Durability	Strength	Class	Distance	Rank
Type-1	7	7	Bad	4	3
Type-2	7	4	Bad	5	4
Type-3	3	4	Good	3	1
Type-4	1	4	Good	3.6	2

k=1

Name	Acid Dura bility	Strength	Class	Distance	Rank
Type-1	7	7	Bad	4	3
Type-2	7	4	Bad	5	4
Type-3	3	4	Good	3	1
Type-4	1	4	Good	3.6	2

Based on immediate neighbour, Good

k=2

Name	Acid Dura bility	Strength	Class	Distance	Rank
Type-1	7	7	Bad	4	3
Type-2	7	4	Bad	5	4
Type-3	3	4	Good	3	1
Type-4	1	4	Good	3.6	2

Based on two neighbours, Good

k=3

Name	Acid Durability	Strength	Class	Distance	Rank
Type-1	7	7	Bad	4	3
Type-2	7	4	Bad	5	4
Type-3	3	4	Good	3	1
Type-4	1	4	Good	3.6	2

Based on three neighbours, 2 Goods and 1 bad, majority ➔ Good