

1.

It is important to define or select similarity measures in data analysis. However, there is no commonly accepted subjective similarity measure. Results can vary depending on the similarity measures used. Nonetheless, seemingly different similarity measures may be equivalent after some transformation.

Suppose we have the following 2-D data set:

	$A_1$	$A_2$
$x_1$	1.5	1.7
$x_2$	2	1.9
$x_3$	1.6	1.8
$x_4$	1.2	1.5
$x_5$	1.5	1.0

(a) Consider the data as 2-D data points. Given a new data point,  $x = (1.4, 1.6)$  as a query, rank the database points based on similarity with the query using Euclidean distance, Manhattan distance, supremum distance, and cosine similarity.

2. Consider the following Term-Frequency vector.

Document Vector or Term-Frequency Vector

<b>Document</b>	<b>team</b>	<b>coach</b>	<b>hockey</b>	<b>baseball</b>	<b>soccer</b>	<b>penalty</b>	<b>score</b>	<b>win</b>	<b>loss</b>	<b>season</b>
Document1	5	0	3	0	2	0	0	2	0	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	0
Document4	0	1	0	0	1	2	2	0	3	0

1) Find Cosine similarity matrix of the given documents.

2) Also, find dissimilarity matrix by using

- Euclidean distance
- Manhattan distance
- Minkowski distance ( $h=3$ )
- Supremum distance

So finally, you are getting one similarity matrix and 4 dissimilarity matrix for the given four documents.

### 3. Dissimilarity between binary attributes:

Find the distance between objects (patients) by considering only on asymmetric attributes. (Gender is symmetric attribute.)

Relational Table Where Patients Are Described by Binary Attributes

<i>name</i>	<i>gender</i>	<i>fever</i>	<i>cough</i>	<i>test-1</i>	<i>test-2</i>	<i>test-3</i>	<i>test-4</i>
Jack	M	Y	N	P	N	N	N
Jim	M	Y	Y	N	N	N	N
Mary	F	Y	N	P	N	P	N

For the given data, find  $d(\text{Jack}, \text{Jim})$ ,  $d(\text{Jack}, \text{Mary})$ ,  $d(\text{Jim}, \text{Mary})$ .

Also, ensure that your program should work correctly for any other possible size of data.