



CS506 Lecture

<input checked="" type="checkbox"/> Override filters	<input type="checkbox"/>
☰ Room	KCB 104

Notes Distance and Dissimilarity:

Data in the matrix:

$$\begin{array}{c} \text{n data points} \left\{ \begin{pmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1m} \\ \vdots & \ddots & \vdots & & \vdots \\ x_{i1} & \dots & x_{ij} & \dots & x_{im} \\ \vdots & & \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nj} & \dots & x_{nm} \end{pmatrix} \right. \\ \underbrace{\hspace{10em}}_{\text{m features}} \end{array}$$

Over here n represents the data points but m represents as a sort of features almost like X and Y coordinates.

Dissimilarity function:

- to uncover interesting structures from our data we need a way to compare data points

- A dissimilarity function is a function that takes two objects (data points) and returns a large value if these objects are dissimilar (or I guess not similar)
- in a sense, it is a way to compare data points
- distance is one we will use

Distance

- properties
 - $d(i,j) = 0$ if and if $i = j$
 - $d(i,j) = d(j,i)$
 - $d(i,j) \leq d(i,k) + d(k,j)$: Triangle inequality
- Minkowski Distance
 - $x = (x_1, \dots, x_d) \wedge y = (y_1, \dots, y_d)$ for some d dimensional real space $p \geq 1$
 - p is a Parameter
 - d is the dimension

$$L_p(x, y) = \left(\sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

- Euclidean Distance:
 - The distance between two points is the length of the line between them
- Manhattan Distance:

- The distance between two points is the length of going up and across from two points (**NOT DIAGONAL**)

Jaccard Similarity:

- Example: union of all the words used in two documents in an 2d matrix or array or vectors
 - How would you would u compare the distance, you can use the Minkowski distance
 - If you have two documents with only 2 words differing the Manhattan distance would be 2 which means u lose a lot of info. Thus this is jaccard Similarity

$$JSim(x, y) = \frac{|x \cap y|}{|x \cup y|}$$

- Similarity

$$JDist(x, y) = 1 - \frac{|x \cap y|}{|x \cup y|}$$

Cosine Similarity:

- The cosine of the angle of the two vectors
 - Two proportional vectors: 1
 - Two orthogonal vectors: 0
 - Two opposite vectors : -1
- To get a corresponding dissimilarity function we can do
 - $1 - \text{cos_similarity}(x,y)$
- Why would you want angle instead
 - Say you have two things talking about the same topic, but one uses more words and is longer over all
 - Thus distance would show them farther apart and imply they are very different
 - Thus angle will show that they are still quite similar even though one is longer and the other is shorter