

Homework - 1

1) vectors $x = (1, 1, 0, 1, 0, 1)$
 $y = (1, 1, 1, 0, 0, 1)$

a) Manhattan distance $\doteq \text{mod}(x_1 - x_2)$

$$\begin{aligned} &\Rightarrow |1-1| + |1-1| + |0-1| + |1-0| + |0-0| + \\ &= |0| + |0| + |-1| + |1| + |0| + |0| \\ &= 0 + 0 + 1 + 1 + 0 + 0 \\ &= 2 \end{aligned}$$

b) Euclidean distance $\doteq \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

$$\begin{aligned} Ed &= \sqrt{|1-1|^2 + |1-1|^2 + |0-1|^2 + |1-0|^2 + |0-0|^2 + |1-1|^2} \\ &= \sqrt{0^2 + 0^2 + (-1)^2 + (1)^2 + 0^2 + 0^2} \\ &= \sqrt{1^2 + 1^2} = \sqrt{2} \end{aligned}$$

c) supremum distance $\div \max(|x_i - y_i|)$

$$= \max(|1-1|, |1-1|, |0-1|, |1-0|, |0-0|, |1-1|)$$

$$= \max(0, 0, 1, 1, 0, 0)$$

$$= 1$$

d) cosine distance = 1 - cosine similarity

$$\text{cosine similarity} = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

$$\mathbf{x} \cdot \mathbf{y} = 1^*1 + 1^*1 + 0^*1 + 1^*0 + 0^*0 + 1^*1$$

$$\mathbf{x} \cdot \mathbf{y} = 1 + 1 + 0 + 0 + 0 + 1$$

$$\mathbf{x} \cdot \mathbf{y} = 3$$

$$\|\mathbf{x}\| = \sqrt{1^2 + 1^2 + 0^2 + 1^2 + 0^2 + 1^2} = \sqrt{4} = 2$$

$$\|\mathbf{y}\| = \sqrt{1^2 + 1^2 + 1^2 + 0^2 + 0^2 + 1^2} = \sqrt{4} = 2$$

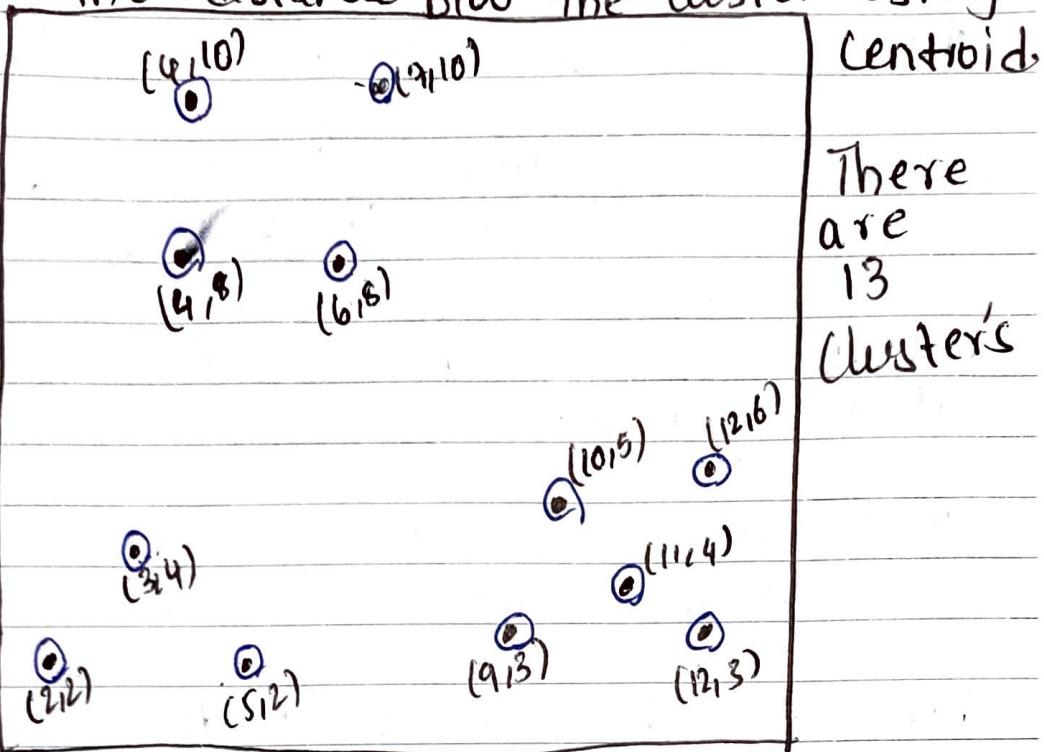
$$\text{cosine similarity} = \frac{3}{4}$$

$$\text{cosine distance} = 1 - \frac{3}{4} = 1/4 = 0.25,$$

* At first each point is its own cluster (1)

Find the distance b/w the cluster using

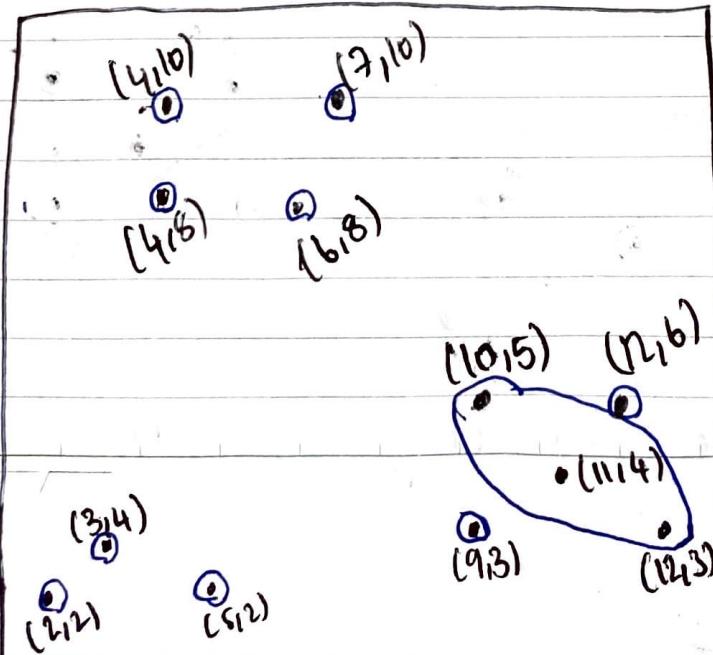
30)



consider the first iteration we can see that the nearest points (or) the distance between $(11,4)$ $(10,5) = \sqrt{1^2 + 1^2} = \sqrt{2}$

$$(11,4) (12,3) = \sqrt{1^2 + 1^2} = \sqrt{2} = 1.41$$

Hence combine $(10,5)$ $(11,4)$ $(12,3)$ as cluster & consider it as 'A'



Here there are 10 cluster's

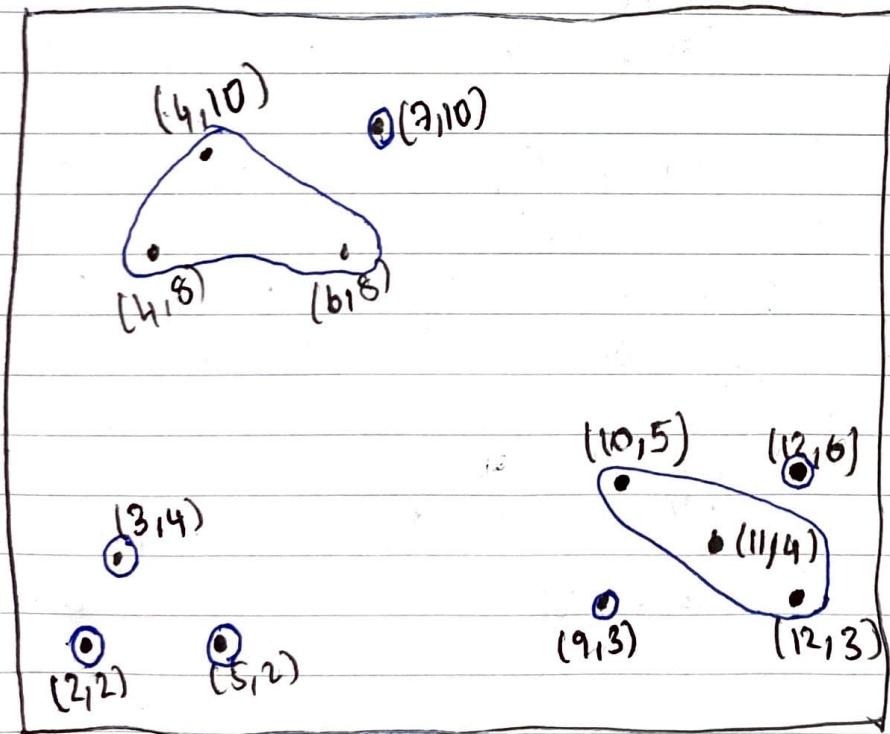
The centroid of the cluster is $(11,4)$

$$\text{distance blw } (11,4) - (9,3) \\ (11,4) - (12,6) \} = 2.24$$

But we can see here that the distance blw

$$(4,10) (4,8) = \sqrt{0^2 + 2^2} = 2 \\ (4,8) (6,8) = \sqrt{2^2 + 0^2} = 2$$

These are the two closest distances found Hence make $(4,10), (4,8), (6,8)$ as a cluster.



Here there are only 8 cluster's

(2) Consider the cluster $(4,8)$ $(6,8)$ $(4,10)$ as B

Find the centroid cluster B and that is $\left(\frac{4+6+4}{3}, \frac{8+8+10}{3} \right) = (4.66, 8.66)$

Now we have to find the distance b/w the centroid of B & $(7,10)$

$$(4.66, 8.66), (7,10) = \sqrt{(7-4.66)^2 + (10-8.66)^2} \\ \approx 2.69$$

But this greater than the distance b/w centroid A & $(9,3)$, $(12,6) \Rightarrow 2.24$

calculate distance b/w $(2,2)$ $(3,4)$
 $(2,2)$ $(5,2)$
 $(3,4)$ $(5,2)$

$$(2,2), (3,4) = \sqrt{1^2 + 2^2} \approx 2.24$$

$$(2,2), (5,2) = \sqrt{3^2 + 0} = 3$$

$$(3,4), (5,2) = \sqrt{2^2 + 2^2} = \sqrt{8} \approx 2.83$$

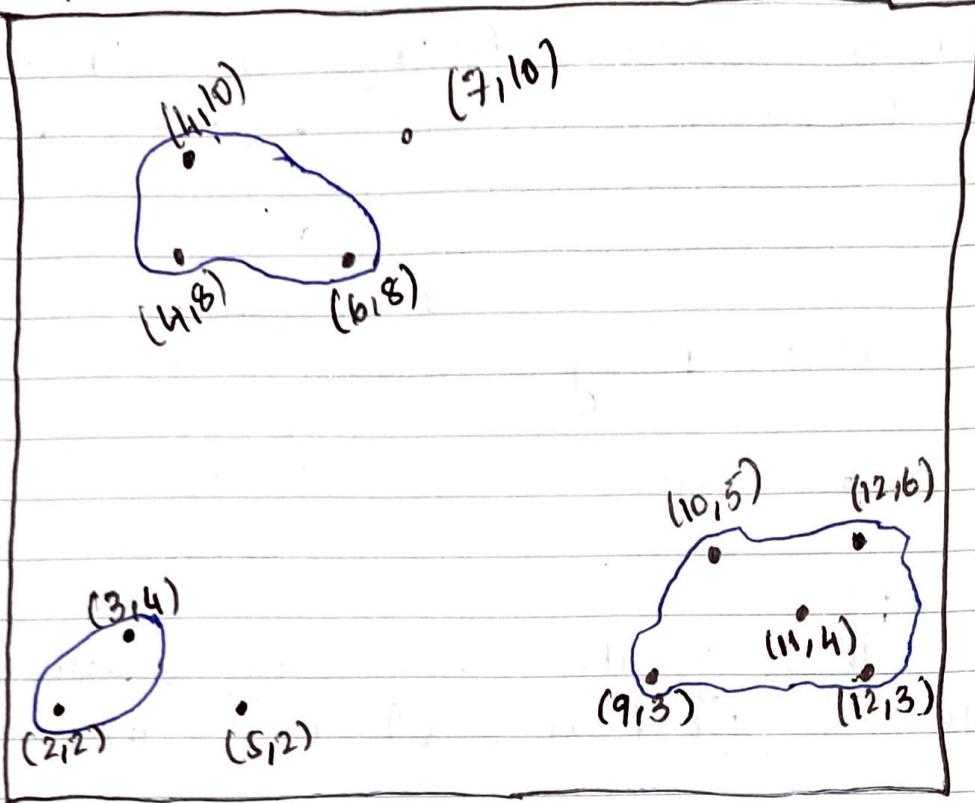
Now we can see that there 2 smallest distances that are ~~2.24~~ between

$$(2,2), (3,4) \quad (11,4), (9,3) \quad \{ 2.24 \\ (11,4), (9,3) \quad \{ (11,4), (12,6) \quad 2.24$$

\therefore cluster both these $(2,2)$, $(3,4)$

& Add $(9,3)$ & $(12,6)$ to the cluster
'A'

Here there are only 5 clusters



let the cluster $(2, 2) \& (3, 4)$ be 'C'

we know the centroid of $(2, 2) \& (3, 4) \Rightarrow (4.66, 8.66)$

The centroid of A is $\left(\frac{10+11+9+12+11}{5}, \frac{5+6+4+3+3}{5} \right)$

$$= (10.8, 4.2)$$

The centroid of C is $\left(\frac{5}{2}, \frac{6}{2} \right)$
 $(2.5, 3)$

Now calculate the distance b/w

$$(2.5, 3), (5, 2) = \sqrt{6.25 + 1} \\ = \sqrt{7.5} = 2.73$$

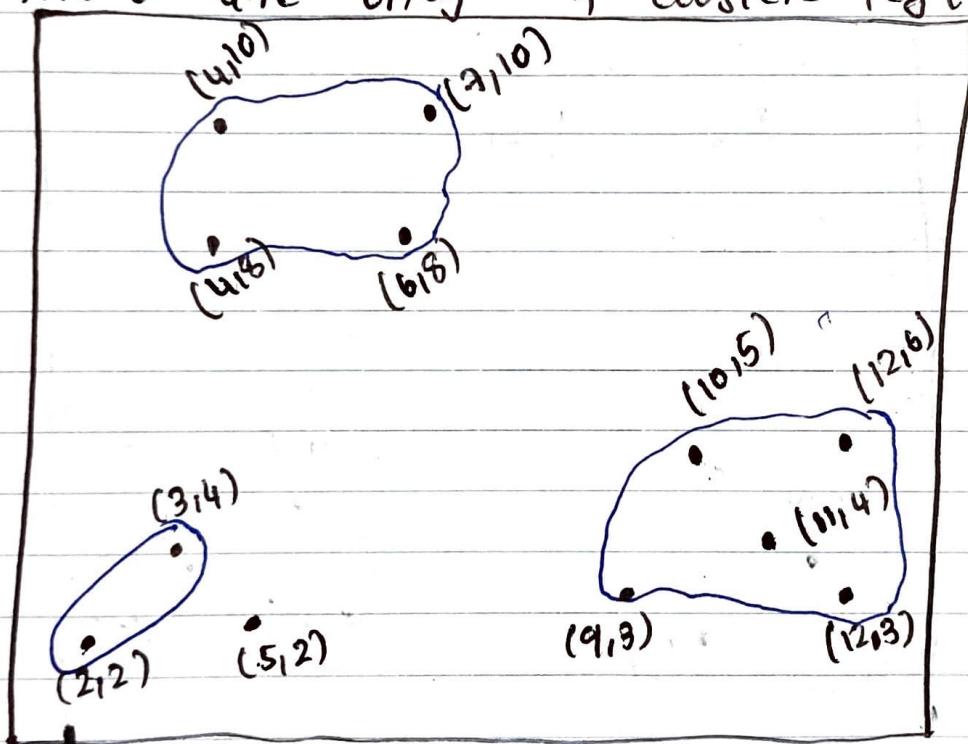
(3)

The distance b/w $(4.66, 8.66)$ $(7, 10)$ is
 $= 2.69$

Here we can see that

distance between cluster B & $(7, 10)$ is near Hence add $(7, 10)$ with the cluster 'B'

There are only 4 clusters left



The centroid of A is $(10.8, 4.2)$

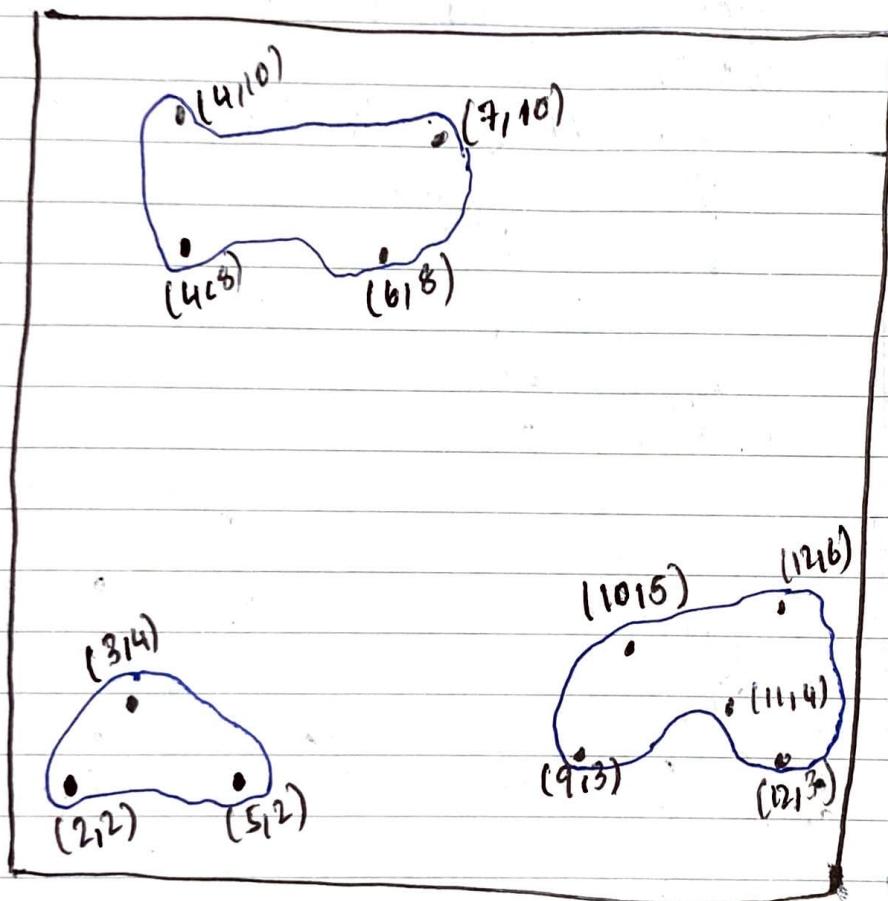
The centroid of B is $\left(\frac{4+4+7+16}{4}, \frac{10+10+8+8}{4} \right)$

$$= (5.25, 9)$$

Now calculate the distance b/w
the centroid of C & (5, 2)

$$(2.5, 3) (5, 2) = \sqrt{7.5} = 2.73$$

Now Add (5, 2) to the cluster 'C'



This is the 5th iteration & there
are 3 clusters A, B, C

~~WTF~~

48)

The first 5-ten shingles

- * Considering white-spaces as underscore the shingles must be set i.e there should not be any shingles that are duplicate
- * The given sentences can be written as

"The-most-effective way to represent documents as sets for the purpose of identifying lexically-similar documents is to construct from the document the set of short strings that appear within it."

The Shingles are

- 1) "The-m"
- 2) "he-mo"
- 3) "e-mos"
- 4) "-most"
- 5) "most-"
- 6) "ost-e"
- 7) "st-ef"
- 8) "t-eff"
- 9) "-effe"
- 10) "effec"

The main usage of shingling is to convert document into sets.