

ASSIGNMENT - I

Discussed with :

Pushkar Mayumdar (18110132)
Harsh Patel (18110062)

1. We can show that $d(u, y) = \min |x_i - y_i|$ does not satisfy the metric property:

$$d(u, y) + d(y, z) = d(z, u)$$

Suppose $u = (0, 0)$
 $y = (10, 0)$
 $z = (10, 10)$

then, $d(u, y) = 0$
 $d(y, z) = 0$
 $d(u, z) = 10$

Hence, for these values of u, y, z :
 $d(u, y) + d(y, z) < d(u, z)$

Hence, $d(u, y) = \min |x_i - y_i|$ is NOT a metric

2.

$$\text{We have, cost}(C) = \sum_i \frac{1}{|C_i|} \sum_{x,y \in C_i} \|x - y\|_2^2$$

We show, $\text{cost}(C) = K$ Means Cost

Let \bar{R} be the mean of all points in the cluster,

$$\Rightarrow \text{cost}(C) = \sum_i \frac{1}{|C_i|} \sum_{x,y \in C_i} \|x - \bar{R} + \bar{R} - y\|^2$$

$$= \sum_i \frac{1}{2|C_i|} \sum_u \sum_y \| (x - \bar{R}) - (y - \bar{R}) \|^2$$

$$\text{cost}(C) = \sum_i \frac{1}{2|C_i|} \sum_u \sum_y \left(\|x - \bar{R}\|^2 + \|y - \bar{R}\|^2 - 2(x - \bar{R})^T(y - \bar{R}) \right)$$

Solving each summation separately :

$$\rightarrow \sum_u \sum_y \|x - \bar{R}\|^2 = \sum_u |C_i| \|x - \bar{R}\|^2$$

$$\rightarrow \sum_u \sum_y \|y - \bar{R}\|^2 = \sum_y \sum_u \|y - \bar{R}\|^2 = \sum_y |C_i| \|y - \bar{R}\|^2$$

$$\rightarrow \sum_u \sum_y 2(x - \bar{R})^T(y - \bar{R})$$

$$= 2 \sum_u \sum_y x^T y - x^T \bar{R} - \bar{R}^T y + \bar{R}^T \bar{R}$$

$$= 2 \left(|C_i|^2 \bar{R}^T \bar{R} - |C_i|^2 \bar{R}^T \bar{R} - |C_i|^2 \bar{R}^T \bar{R} + |C_i|^2 \bar{R}^T \bar{R} \right)$$

$$= 0$$

Hence,

$$\text{Cost}(C) = \sum_i \frac{1}{z|C_i|} \left(\sum_u |c_i| \|x - k\|^2 + \sum_y |c_i| \|y - k\|^2 \right)$$

$$\text{Cost}(C) = \sum_i \frac{1}{z|C_i|} \left(\cancel{2 \sum_u |c_i| \|u - k\|^2} \right)$$

$$(\text{As } \sum_u |c_i| \|u - k\|^2 = \sum_y |c_i| \|y - k\|^2)$$

$$\boxed{\text{Cost}(C) = \sum_i \sum_u \|u - k\|^2}$$

Hence, $\text{Cost}(C)$ = Cost of K Means

Algorithm for optimizing Cost(C)

- We can use any of the K Means objective optimization algorithms such as Lloyd or KMeans++.

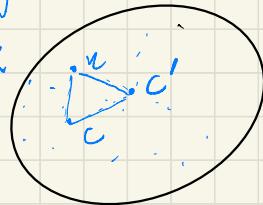
3. Let c be the optimal centre when we allow any arbitrary point to be the centre.

Let c' be the optimal centre when we allow only the data points to be the centre.

Let K be the partition created by c centre.

Now,

Partition K



$\forall x \in K$

$$d(c', x) \leq d(c, x) + d(c, c') \quad \text{--- (1)}$$

(as L1 norm follows triangle inequality)

We can choose c' to be the data point in K , st

$$d(c, c') \leq d(c, x) \quad \text{--- (2)}$$

Using (1) and (2)

$$d(c', x) \leq d(c, x) + d(c, c')$$

$$\leq d(c, x) + d(c, x)$$

$$d(c', x) \leq 2 d(c, x) \quad \forall x \in K$$

Cost of partition K using C' as centre =

$$\sum_{x \in K} d(C', x) \leq 2 \sum_{x \in K} d(C, x)$$

= 2 Optimal cost of clustering

→ Similarly we can do this for all cluster partitions

Hence, we proved that there is two ratio between optimal value when we require all cluster centers to be data points or allow arbitrary points to be centers.

Lloyd's algorithm for Euclidean k-median

The only change that we see from Lloyd's for K-Means is the criteria for recalculating centers. Here, we take median of the partition for recalculating center.

Lloyd for K-Median

- Choose k points arbitrary from data points
- repeat
 - Assign points to nearest centers
 - recalculate centers by taking median of all points. Then choose the data point closest to the median as center.
- until no (or small #) points change cluster
or when cluster centers dont shift much.

Let C be number of iteration

Time Complexity : $O(CKNd)$

where R = no. of clusters

N = no. of samples

d = dimension

To Show :- Cluster cost decreases with each iteration

- The first step of an iteration, we recalculate center by putting it as the median.

As we know, median gives the minimum cost when using L1 norm, this step will decrease cost.

- In the second step, we reassign points to their nearest centers, hence reducing cost of the whole clustering.

Hence, clustering cost always reduces with each iteration.

4 . <https://github.com/shivam15s/CS-328-HW-1/blob/main/q4.ipynb>

5. LINK FOR YOUTUBE VIDEO :

<https://youtu.be/KaRMekgtFCY>

Challenges :

1. Figuring out the attributes which can lead to good inferences
2. Handling the NaN values.
Used bfill, ffill for filling