Ans!

## K - means clustering

Cost function for k-means

$$Cost(T) = \sum_{z \in T} \sum_{x \in C_z} \|x - z\|^2$$

it can also be written as

$$cost(C_1, \dots, C_k; Z_1, \dots, Z_k) = \sum_{j=1}^{k} \sum_{x \in C_j} \|x - z_j\|^2$$

$\Rightarrow$ For any specific $j = 1, \dots k$

$$cost(C_j; z_j) = \sum_{x \in C_j} \|x - z_j\|^2 \qquad -①$$

This represents the sum of squared distances from every point in 'C' to the centroid 'z'

we want to proove

$$Cost(C_j; z_j) = Cost(C_j; mean(C_j)) + |C_j| \|z_j - mean(C_j)\|^2 \qquad -②$$

where  mean$(C_j)$ is the mean (centroid) of points in set $C_j$

$$\Rightarrow \quad mean(C_j) = \frac{1}{|C_j|} \sum_{x \in C_j} x \qquad -③$$

where $|C_j|$ is the cardinality of $C_j$

Expanding eq ①,

$$Cost(C_j; z_j) = \sum_{x \in C_j} \|x - z_j\|^2$$

$$= \sum_{x \in C_j} \left( \|x\|^2 - 2x \cdot z_j + \|z_j\|^2 \right)$$

$$= \sum_{x \in C_j} \|x\|^2 - 2z_j \sum_{x \in C_j} x + \|z_j\|^2 \sum_{x \in C_j} 1$$

$$= \sum_{x \in C_j} \|x\|^2 - 2z_j |C_j| \, mean(C_j) + |C_j| \|z_j\|^2 \quad -④ \quad \left( \text{Using } ③ \right.$$

Expanding  Cost $(C_j; \text{mean}(C_j))$,

$$\text{Cost}\left(C_j ; \text{mean}(C_j)\right) = \sum_{x \in C_j} \| x - \text{mean}(C_j) \|^2 \qquad \left(\text{From } \textcircled{1}\right)$$

$$= \sum_{x \in C_j} \| x \|^2 - 2\,\text{mean}(C_j) \sum_{x \in C_j} x + \| \text{mean}(C_j) \|^2 \sum_{x \in C_j} 1$$

$$= \sum_{x \in C_j} \| x \|^2 - 2\,|C_j| \, \| \text{mean}(C_j) \|^2 + |C_j| \, \| \text{mean}(C_j) \|^2 \left(\text{Using } \textcircled{3}\right)$$

$$= \sum_{x \in C_j} \| x \|^2 - |C_j| \, \| \text{mean}(C_j) \|^2 \qquad\qquad - \textcircled{5}$$


Subtracting $\textcircled{4}$ & $\textcircled{5}$  $(\textcircled{4} - \textcircled{5})$

$\text{cost}(C_j; z_j) - \text{cost}(C_j; \text{mean}(C_j))$

$$= \left( \sum_{x \in C_j} \| x \|^2 - 2 z_j \, |C_j| \, \text{mean}(C_j) + |C_j| \, \| z_j \|^2 \right) -$$

$$\left( \sum_{x \in C_j} \| x \|^2 - |C_j| \, \| \text{mean}(C_j) \|^2 \right)$$

$$= |C_j| \, \| z_j \|^2 - 2 z_j |C_j| \, \text{mean}(C_j) - |C_j| \, \| \text{mean}(C_j) \|^2$$

$$= |C_j| \left( \| z_j \|^2 - 2 z_j \, \text{mean}(C_j) + \| \text{mean}(C_j) \|^2 \right)$$

$$= |C_j| \, \| z_j - \text{mean}(C_j) \|^2$$


$\Rightarrow$ The equation can be re-written as

$$\text{cost}(C_j; z_j) = \text{Cost}(C_j; \text{mean}(C_j)) + |C_j| \, \| z_j - \text{mean}(C_j) \|^2$$

is  same  as  eq $\textcircled{2}$

**Hence Proved**

<u>Ans 2</u>    k-means algorithm is defined as:

For

initial centers $z_1, \ldots, z_k \in \mathbb{R}^d$ & clusters $C_1, \ldots, C_k$
repeat until there is no change in cost

for each $j$: $C_j \leftarrow \{x \in S$ whose closest center is $z_j\}$
for each $j$: $z_j \leftarrow mean(C_j)$

We have to prove that the cost of K-means algorithm
is monotonically decreasing.

<u>Proof</u> :

The k-means algorithm has 2 main steps

1. <u>Assignment Step</u>
where each data point is assigned to its
nearest centroid

2. <u>Updation step</u>
The centroids for each cluster are calculated again
as the mean of all the points in the cluster

Let $C_j^{(t)}$ be the set of points in cluster $C_j$ at iteration '$t$'

with $z_j^{(t)}$ as it's centroid.

The cost function becomes

$$Cost^{(t)} = \sum_{j=1}^{k} \sum_{x \in C_j^{(t)}} \|x - z_j^{(t)}\|^2$$

For the function to be monotonically decreasing

$$Cost^{(t+1)} \leq Cost^{(t)}$$

<u>Step 1</u>: <u>Assignment Step</u>

At each iteration, each point $x$ is assigned to the nearest
centroid.

$\Rightarrow$ distance of points at iteration '$t$' $= \|x - z_j^{(t)}\|^2$
    & iteration '$t+1$' $= \|x - z_j^{(t+1)}\|^2$

Since points are reassigned to a closer centroid in iteration
'$t+1$' by definition

$$\Rightarrow \quad \| x - z_j^{(t+1)} \|^2 \leq \| x - z_j^{(t)} \|^2$$

Therefore, the cost associated with each point cannot increase i.e. we see that summing effect cannot increase total cost. which is the distance

It can also be thought as

$$\text{cost}\left( c_1^{(t+1)}, \ldots, c_k^{(t+1)} ; z_1^{(t)}, \ldots, z_k^{(t)} \right) \leq \text{cost}\left( c_1^{(t)}, \ldots c_k^{(t)} ; z_1^{(t)}, \ldots z_k^{(t)} \right)$$

## Step 2 : Updation Step

We calculate new centroids $z_j^{(t+1)}$ as the mean of all the points in cluster $c_j^{(t+1)}$

By definition of mean, it minimizes the sum of squared distances from the points in the cluster to its centroid (or mean)

Therefore, on calculating the centroids of the clusters at iteration 't+1' cannot increase the cost.

$$\Rightarrow \quad \sum_{x \in c_j^{(t+1)}} \| x - z_j^{(t+1)} \|^2 \leq \sum_{x \in c_j^{(t)}} \| x - z_j^{(t)} \|^2$$

It can also be thought as

$$\text{cost}\left( c_1^{(t+1)}, \ldots c_k^{(t+1)} ; z_1^{(t+1)}, \ldots z_k^{(t+1)} \right) \leq \text{cost}\left( c_1^{(t+1)}, \ldots, c_k^{(t+1)} ; z_1^{(t)}, \ldots, z_k^{(t)} \right)$$

On combining the 2 steps, we get

$$\text{cost}^{(t+1)} = \sum_{j=1}^{k} \sum_{x \in c_j^{(t+1)}} \| x - z_j^{(t+1)} \|^2$$

$$\leq \sum_{j=1}^{k} \sum_{x \in c_j^{(t+1)}} \| x - z_j^{(t)} \|^2$$

$$\leq \sum_{j=1}^{k} \sum_{x \in c_j^{(t)}} \| x - z_j^{(t)} \|^2$$

$$\leq \text{cost}^{(t)}$$

Hence    Proved