

# Data Mining

## Assignment 2

-Suraj Prathik Kumar (2016101)

Assumptions for the Code - Input: Number of Clusters(k) and Number of data Points (Including the Points)

Question 1: Code Submitted : File Name - Question 1.py

Question 2: Code Submitted : File Name - Question 2.py

Question 3: Code Submitted : File Name - Question 3.py

Question 4:

a) Yes, the algorithm is able to find the right clusters given in the Question

```
Question 2 x
/Users/surajprathikkumar/Desktop/Dataminingass2/bin/python "/Users/surajprathikkumar/Desktop/Dataminingass2/Question 2.py"
Value of K - 3
No. of DataPoints - 15
Point - 1
Point - 2
Point - 3
Point - 4
Point - 5
Point - 8
Point - 9
Point - 10
Point - 11
Point - 12
Point - 24
Point - 28
Point - 32
Point - 36
Point - 40
initial [[1.0], [11.0], [28.0]]
cluster [[[1.0], [2.0], [3.0], [4.0], [5.0]], [[8.0], [9.0], [10.0], [11.0], [12.0]], [[24.0], [28.0], [32.0], [36.0], [40.0]]]
Seedafter [array([3.]), array([10.]), array([32.])]
cluster [[[1.0], [2.0], [3.0], [4.0], [5.0]], [[8.0], [9.0], [10.0], [11.0], [12.0]], [[24.0], [28.0], [32.0], [36.0], [40.0]]]
Seedafter [array([3.]), array([10.]), array([32.])]
Done
Cluster :-
[[1.0], [2.0], [3.0], [4.0], [5.0]]
[[8.0], [9.0], [10.0], [11.0], [12.0]]
[[24.0], [28.0], [32.0], [36.0], [40.0]]
Process finished with exit code 0
```

b) Yes, the algorithm is able to find the right clusters given in the Question

```
Question 2 x
/Users/surajprathikkumar/Desktop/Dataminingass2/bin/python "/Users/surajprathikkumar/Desktop/Dataminingass2/Question 2.py"
Value of K - 3
No. of DataPoints - 15
Point - 1
Point - 2
Point - 3
Point - 4
Point - 5
Point - 8
Point - 9
Point - 10
Point - 11
Point - 12
Point - 24
Point - 28
Point - 32
Point - 36
Point - 40
initial [[1.0], [2.0], [3.0]]
cluster [[[1.0]], [[2.0]], [[3.0], [4.0], [5.0], [8.0], [9.0], [10.0], [11.0], [12.0], [24.0], [28.0], [32.0], [36.0], [40.0]]]
Seedafter [array([1.]), array([2.]), array([17.07692308])]
cluster [[[1.0]], [[2.0], [3.0], [4.0], [5.0], [8.0], [9.0]], [[10.0], [11.0], [12.0], [24.0], [28.0], [32.0], [36.0], [40.0]]]
Seedafter [array([1.]), array([5.16666667]), array([24.125])]
cluster [[[1.0], [2.0], [3.0]], [[4.0], [5.0], [8.0], [9.0], [10.0], [11.0], [12.0]], [[24.0], [28.0], [32.0], [36.0], [40.0]]]
Seedafter [array([2.]), array([8.42857143]), array([32.])]
cluster [[[1.0], [2.0], [3.0], [4.0], [5.0]], [[8.0], [9.0], [10.0], [11.0], [12.0]], [[24.0], [28.0], [32.0], [36.0], [40.0]]]
Seedafter [array([3.]), array([10.]), array([32.])]
cluster [[[1.0], [2.0], [3.0], [4.0], [5.0]], [[8.0], [9.0], [10.0], [11.0], [12.0]], [[24.0], [28.0], [32.0], [36.0], [40.0]]]
Seedafter [array([3.]), array([10.]), array([32.])]
Done
Cluster :-
[[1.0], [2.0], [3.0], [4.0], [5.0]]
[[8.0], [9.0], [10.0], [11.0], [12.0]]
[[24.0], [28.0], [32.0], [36.0], [40.0]]
Process finished with exit code 0
```

- c) In First and Second part, Initial Seeds are different {1, 11, 28} and {1, 2, 3} respectively. The final Cluster Formed are same but the number of iteration are less in First part as compared to Second part. Choice of seed play in important factor in determining Faster Convergence. The seeds should not be very close to each other and should fall in the range of similar datapoints.  
Example: Here there are 3 range [1-5], [8-11] and [24-40]  
Selecting seed from each of the Range will give us faster Convergence and take Less iterations than selecting all the seeds from same range.