

# KMeans Clustering Implementation

## Technology/Concept:

- Distributed Machine Learning: Apache Spark (MLlib) can be used for large-scale data processing and machine learning tasks like KMeans clustering, enabling efficient handling of large datasets across distributed clusters.
- Hadoop: For large datasets, Hadoop with MapReduce can be used to parallelize the KMeans computation.
- Data Storage: Large input datasets can be stored in HDFS, Amazon S3, or other distributed storage systems for efficient access and processing.

## Steps:

1. First step to read the info files.
2. Then parse the info.
3. Initialize the K clusters.
4. Now perform iteration until converge.
5. Get nearest points mapper.
6. Combined the mapper results using reducer.
7. Update the centroids.
8. Check the space between points and centroid.
9. After converge results as follows.

## Output:

```
In [2]: 1 # Reading the text file
        2 text = sc.textFile("C:/Users/ysrivastava/BDA/Assignment/KMeans/points.txt")
        3 text.collect()

Out[2]: ['3.023\t5.138',
        '3.075\t4.989',
        '2.321\t5.350',
        '3.328\t4.944',
        '3.195\t5.186',
        '3.034\t4.656',
        '3.484\t4.877',
        '3.023\t4.951',
        '2.707\t4.871',
        '2.894\t5.031',
        '2.880\t5.207',
        '2.780\t5.196',
        '2.941\t4.951',
        '3.213\t5.123',
        '2.972\t4.941',
        '3.138\t5.460',
        '3.247\t5.185',
        '3.602\t5.405',
        '2.857\t4.606',
        '2.800\t4.801']

In [6]: 1 # Passing the data points to remove the tabs
        2 def parse(points):
        3     var = points.split('\t')
        4     result = []
        5     for x in var:
        6         value = float(x)
        7         result.append(value)
        8     return np.array(result)

In [7]: 1 #Calling the parsing function
        2 arr = text.map(lambda m: parse(m)).cache()
        3 arr.collect()

Out[7]: [array([3.023, 5.138]),
        array([3.075, 4.989]),
        array([2.321, 5.35 ]),
        array([3.328, 4.944]),
        array([3.195, 5.186]),
        array([3.034, 4.656]),
        array([3.484, 4.877]),
        array([3.023, 4.951]),
        array([2.707, 4.871]),
        array([2.894, 5.031]),
        array([2.88, 5.207]),
        array([2.78, 5.196]),
        array([2.941, 4.951]),
        array([3.213, 5.123]),
        array([2.972, 4.941]),
        array([3.138, 5.46]),
        array([3.247, 5.185]),
        array([3.602, 5.405]),
        array([2.857, 4.606]),
        array([2.8, 4.801])]
```

```

In [18]: 1 Dist = float(0.1)
          2 kp = [np.array([float(2) , float(5)]), np.array([float(6), float(2)]), np.array([float(1), float(1)])]
          3 dist2 = 1.0
          4 type(kp)
          5 print(kp)

[array([2., 5.]), array([6., 2.]), array([1., 1.])]

In [21]: 1 nearest = arr.map(lambda p: process_map(p, kp))
          2 print(nearest.take(5))

[(0, (array([3.023, 5.138]), 1)), (0, (array([3.075, 4.989]), 1)), (0, (array([2.321, 5.35 ]), 1)), (0, (array([3.328, 4.944]), 1)), (0, (array([3.195, 5.186]), 1))]

In [23]: 1 Status = nearest.reduceByKey(lambda pointer1 , pointer2: pairs(pointer1 , pointer2))
          2 print(Status.take(4))

[(0, (array([1496.991, 2502.065]), 500)), (2, (array([1016.494, 501.689]), 500)), (1, (array([2496.329, 995.104]), 500))]

In [24]: 1 new_points = Status.map(lambda map_object: mapper(map_object)).collect()
          2 print(new_points)

[(0, array([2.993982, 5.00413 ])), (2, array([2.032988, 1.003378])), (1, array([4.992658, 1.990208]))]

In [25]: 1 updated_distance = all_points_dist(new_points)
          2 print(updated_distance)

51.92141868048006

In [27]: 1 def change_points(new_points, kp):
          2     for (cluster_number, point) in new_points:
          3         kp[cluster_number] = point
          4     return kp
          5
          6 change_points(new_points, kp)
          7 print(kp)

[array([2.993982, 5.00413 ]), array([4.992658, 1.990208]), array([2.032988, 1.003378])]

```

```

In [30]: 1 while dist2 > Dist:
          2     print(dist2)
          3     near = arr.map(lambda p: process_map(p, kp))
          4     point_Stats = nearest.reduceByKey(lambda pointer1 , pointer2: pairs(pointer1 , pointer2))
          5     new_Points = point_Stats.map(lambda map_object: mapper(map_object)).collect()
          6     print(new_Points)
          7     dist2 = sum(np.sum((kp[iK] - p) ** 2) for (iK, p) in new_Points)
          8
          9     for (iK, p) in new_Points:
         10         kp[iK] = p
         11
         12 results(kp)

1.0
[(0, array([2.993982, 5.00413 ])), (2, array([2.032988, 1.003378])), (1, array([4.992658, 1.990208]))]
2.99398#5.00413
4.99266#1.99021
2.03299#1.00338

```