# K-Means Clustering

# what is K-Means?

- given $K > 0$ (where $K$ is the number of clusters)
- a set of $N$ $d$-dimensional objects
  - Clustering is the process of grouping a set of $N$ $d$-dimensional (2-dimensional, 3-dimensional, etc.) objects into $K$ clusters of similar objects.
  - Objects should be similar to one another within the same cluster and dissimilar to those in other clusters.
- K-Means is a distance-based clustering algorithm.
- K-Means clustering (aka unsupervised learning) is a data mining algorithm to cluster, classify, or group your $N$ objects based on their attributes or features into $K$ number of groups (so-called clusters).

# For example

- It can be used to find a group of consumers with common behaviors.

- To cluster documents based on the similarity of their contents

- Selection of $K$ is specific to the application or problem domain.

# K-Means algorithm

1. Partition *N* objects into *K* nonempty subsets.

2. Compute the centroids of the clusters in the current partition (the centroid is the center, or mean point, of the cluster). For all $i$ = 1, 2, ..., *K*, compute $\boldsymbol{\mu}i$ as:

$$\mu_i = \frac{1}{|c_i|} \sum_{j \in c_i} x_j, \ \forall i$$

3. Assign each object to the cluster with the nearest centroid. This is simply attributing the closest cluster to each data point:

$$c_i = \{j : d(x_j, \mu_i) \leq d(x_j, \mu_l), l \neq i, j = 1, ..., n\}$$

where *d*(*a*, *b*) is the distance function for two points: *a* and *b*.

4. Stop when there are no more new assignments. Otherwise, go back to step 2. Basically, we repeat steps 2 and 3 until convergence.

# K-Means Distance Function

- Euclidean distance, Manhattan distance, Inner product space, Maximum norm, Your own custom function (any metric you define over the *d*-dimensional space)

Let:

$$X = (X_1, X_2, ..., X_d)$$

$$Y = (Y_1, Y_2, ..., Y_d)$$

then:

$$distance(X, Y) = \sqrt{(X_1 - Y_1)^2 + (X_2 - Y_2)^2 + \ldots + (X_d - Y_d)^2}$$

The Euclidean distance function has some interesting properties:

- distance(i,j) ≥ 0
- distance(i,i) = 0
- distance(i,j) = distance(j,i)
- distance(i,j) ≤ distance(i,k) + distance(k,j)

# MapReduce Solution for K-Means Clustering

*Example 12-1. K-Means clustering algorithm*

```
 1 // k = number of desired clusters
 2 // delta = acceptable error for convergence
 3 // data = input data
 4 kmeans(k, delta, data) {
 5     // initialize the cluster centroids
 6     initial_centroids = pick(k, data);
 7
 8     // this is how we broadcast centers to mappers
 9     writeToHDFS(initial_centroids);
10
11     // iterate as long as necessary
12     current_centroids = initial_centroids;
13     while (true) {
14         // theMapReduceJob() does 2 tasks:
15         //     1. uses current_centroids in map()
16         //     2. reduce() creates new_centroids and writes it to HDFS
17         theMapReduceJob();
18         new_centroids = readFromHDFS();
19         if change(new_centroids, current_centroids) <= delta {
20             // we are done, terminate loop-iteration
21             break;
22         }
23         else {
24             current_centroids = new_centroids;
25         }
26     }
27
28     result = readFromHDFS();
29     return result;
30 }
```

*Example 12-2. K-Means clustering algorithm: change() method*

```
1 change(new_centroids, current_centroids) {
2     new_distance = [sum of squared distance in the new_centroids];
3     current_distance = [sum of squared distance in the current_centroids];
4     changed = absoulteValue(new_distance - current_distance);
5     return changed;
6 }
```

# MapReduce Solution: map()

- Read the cluster centroids into memory from a SequenceFile1 (note that we may also use Redis or memcached for persisting cluster centroids). In Hadoop's implementation, this will be done in the setup() method of the mapper class.

- Iterate over each cluster centroid for each input key-value pair. In Hadoop's implementation, for the map() function, the key is generated by Hadoop and ignored (not used).

- Compute the Euclidean distances and save the nearest center with the lowest distance to the input point (as a $d$-dimensional vector).

- Write the key-value pair to be consumed by reducers, where the key is the nearest cluster center to the input point (and the value is a $d$-dimensional vector). Both the key and the value are of the Vector data type.

# MapReduce Solution: map()

*Example 12-3. K-Means MapReduce: map() function*

```java
1  public class KmeansMapper ... {
2
3      private List<Vector> centers = null;
4
5      private List<Vector> readCentersFromSequenceFile() {
6          // read cluster centroids from a SequenceFile,
7          // which is a set of key-value pairs
8          ...
9      }
10
11     // called once at the beginning of the map task
12     public void setup(Context context) {
13         this.centers = readCentersFromSequenceFile();
14     }
17      * @param key is MapReduce generated, ignored here
18      * @param value is the d-dimensional Vector (V1, V2, ..., Vd)
19      */
20     map(Object key, Vector value) {
21         Vector nearest = null;
22         double nearestDistance = Double.MAX_VALUE;
23         for (Vector center : centers) {
24             double distance = EuclideanDistance.calculateDistance(center, value);
25             if (nearest == null) {
26                 nearest = center;
27                 nearestDistance = distance;
28             }
29             else {
30                 if (nearestDistance > distance) {
31                     nearest = center;
32                     nearestDistance = distance;
33                 }
34             }
35         }
36
37         // prepare key-value for reducers
38         emit(nearest, value);
39     }
```

# MapReduce Solution: combine()

```
1 /**
2  * @param key is the Centroid
3  * @param values is a list of Vectors
4  */
5 combine(Vector key, Iterable<Vector> values) {
6   // all dimensions in sum Vector are initialized to 0.0
7   Vector sum = new Vector();
8   for (Vector value : values) {
9     // note that value.length = d,
10    // where d is the number of dimensions for input objects
11    for (int i = 0; i < value.length; i++) {
12      sum[i] += value[i];
13    }
14  }
15
16  emit(key, sum);
17 }
```

# MapReduce Solution: reduce()

- The main task of the reduce() function is to recenter.

- Each reducer iterates over each value vector and calculates the mean vector. Once you have found the mean, this is the new center; your final step is to save it into a persistent store (such as a SequenceFile).

```
1  /**
2   * @param key is the Centroid
3   * @param values is a list of Vectors|
4   */
5  reduce(Vector key, Iterable<Vector> values) {
6      // all dimensions in newCenter are initialized to 0.0
7      Vector newCenter = new Vector();
8      int count = 0;
9      for (Vector value : values) {
10         count++;
11         for (int i = 0; i < value.length; i++) {
12             newCenter[i] += value[i];
13         }
14     }
15
16     for (int i = 0; i < key.length; i++) {
17         // set new mean for each dimension
18         newCenter[i] = newCenter[i] / count;
19     }
20
21     emit(key.ID, newCenter);
22 }
```