

Mining of Massive Dataset

Hao Jiteng, Zhou Lizhi, Yang Fangzhou

June 18, 2012

Abstract

K-means is a simple yet useful clustering algorithm. It's underline natural implies that it could be parallelized or distributized. Some implementation of kmeans employs platforms such as Hadoop and CUDA to boost the process of mining of massive dataset. In our project we implement k-means algorithm on Apache Hadoop Project. We ran our algorithm on our tiny cluster. Evaluation has been done to measure our algorithm.

1 Introduction

1.1 Hadoop

Following is the introduction of Hadoop on its project homepage [1].

The ApacheTM HadoopTM project develops open-source software for reliable, scalable, distributed computing.

The Apache Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using a simple programming model. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage. Rather than rely on hardware to deliver high-availability, the library itself is designed to detect and handle failures at the application layer, so delivering a highly-available service on top of a cluster of computers, each of which may be prone to failures.

We mainly use the Hadoop MapReduce framework rather than HDFS.

1.2 k-means

There are many materials that introduce k-means algorithms. MacQueen, J. et al. proposed k-means algorithm in [4]. In [3] the author introduced a simple k-means MapReduce algorithm.

The dataset is relatively large comparing to the memory size. Thus we need a mechanism to deal with the incompatibility of memory. MapReduce is a solution to this problem. The detail of our algorithm is documented in the following section.

1.3 Dataset

2 Clustering Algorithm

In this section we discuss the k-means clustering algorithm used in our implementation. Note that to find proper initial clusters, which may lead to fewer iterations and good result, we also implemented a Canopy clustering MapReduce algorithm.

2.1 Canopy Clustering

2.2 k-means Clustering

Mahout Project [2] is a data mining framework under ApacheTM Foundation. It contains a k-means implementation. The blog [5] gives a very detailed view of the algorithm. Our algorithm mainly based on the idea of Mahout Project. Actually, this algorithm is very similar to the BFR algorithm in our lectures.

The algorithm is built up by the following MapReduce phases,

1. The k-means iteration, output every cluster centroid if converged.
2. Assign every point to a known cluster and output the result.

The detail of the first phase is described as follows

2.2.1 Iterations

1. Mapper<LongWritable, Text>→ <LongWritable,KmeansCluster>. It reads the input file content as value. The key is the value offset in the file. Then it convert the value to a VectorDoubleWritable, which is used to represent the feature vector. It finds the nearest cluster to the vector, and output cluster id as key, a new cluster containing only point as value.
2. Combiner<LongWritable, KmeansCluster>→ <LongWritable, KmeansCluster>. It reads the output from Mapper and combine those tuples who have the same cluster id(meaning that they are assigned to the same cluster) using KmeansCluster.omitCluster(). This function reduced the network transmission flow because the actual meaningful information need to be communicated between different nodes are only the N, SUM and SUMSQ of clusters, which is described in the lecture slides of BFR algorithm.
3. Reducer<LongWritable, KmeansCluster>→ <LongWritable, KmeansCluster>. It reads the cluster id as key and the KmeansCluster as value. It adds the N, SUM and SUMSQ. The result leads to the combination of

clusters. Finally the reducer outputs the result clusters of this single iteration. These clusters are input of next iteration. During reducer, if the movement of one cluster is less than a threshold, it's said to be "Converged". If a cluster is converged, a counter in context will increase by one.

4. If, in the driver, the counter equals the number of clusters, meaning that all clusters are converged, this phase is finished.

2.2.2 Assign Point to Clusters

After the iterations, the clusters are stable. Next we are going to assign every point to its nearest cluster. This phase only requires mapper to do all the works, since this procedure is highly parallelized. Each two points are independent of each other.

References

- [1] Apache Foundation. Hadoop project. <http://hadoop.apache.org/>.
- [2] Apache Foundation. Mahout project. <http://mahout.apache.org/>.
- [3] Ricky Ho. K-means clustering in map reduce. <http://horicky.blogspot.com/2011/04/k-means-clustering-in-map-reduce.html>.
- [4] J. MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, page 14. California, USA, 1967.
- [5] Leo Zhang. Learning mahout: K-means clustering. <http://www.cnblogs.com/vivounicorn/archive/2011/10/08/2201986.html>.