

Implementing K-Means Clustering Algorithm Using MapReduce Paradigm

Botcha Chandrasekhara Rao¹, Medara Rambabu²

¹M.Tech Scholar, Department of Computer Science & Engineering, Pydah Kaushik College of Engineering, Gambheeram Village, Anandapuram, Mandal-531163. Visakhapatnam, AP, India

²HOD & Associate Professor, Department of Computer Science & Engineering, Pydah Kaushik College of Engineering, Gambheeram Village, Anandapuram, Mandal-531163. Visakhapatnam, AP, India

Abstract: Clustering is a useful data mining technique which groups' data points such that the points within a single group have similar characteristics, while the points in different groups are dissimilar. Partitioning algorithm methods such as k-means algorithm is one kind of widely used clustering algorithms. As there is an increasing trend of applications to deal with vast amounts of data, clustering such big data is a challenging problem. Recently, partitioning clustering algorithms on a large cluster of commodity machines using the MapReduce framework have received a lot of attention. Traditional way of clustering text documents is Vector space model, in which tf-idf is used for k-means algorithm with supportive similarity measure. This project exhibits an approach to cluster text documents in which results obtained by executing map reduce k-means algorithm on single node cluster show that the performance of the algorithm increases as the text corpus increases.

Keywords: Vector space model, map reduce, text, clustering, map reduce k-means, Hadoop

1. Introduction

Big Data[1] is the term applied to data sets whose size is beyond the ability of the commonly used software tools to capture, manage, and process within a tolerable elapsed time. Big data can be analysed with the software tools commonly used as part of advanced analytics disciplines such as predictive analytics, data mining, text analytics and statistical analysis. The semi-structured and unstructured data may not fit well in traditional data warehouses based on relational databases. Hadoop is a platform that provides both distributed storage and computational capabilities. It is an open source software project that enables the distributed processing of large data sets across clusters of commodity servers. It is designed to scale up from a single server to thousands of machines, with a very high degree of fault tolerance. Hadoop is a distributed master-slave architecture that consists of Hadoop distributed file system (HDFS) for storage and Map-Reduce for computational capabilities.

Mathematical models like Boolean model, probabilistic model, vector space model are proposed to use in information retrieval systems. Vector space model (VSM) [2] is most popular and widely used model.

2. Literature Survey

A. Vector Space Model

In VSM [3], A term document matrix X of size $t \times d$ is created, where t is the terms (words) of the documents and d represents documents and find the frequency of the terms in each and every document. This term frequency cannot be able to find its

importance in the whole corpus rather cannot compare with other documents. Term frequency - Inverse document frequency (TF-IDF) overcomes the limitations of term frequency and calculated term weighting to find how important a word is to a document in a collection or corpus.

The weight vector for document d is $V_d = [w_{1,d}, w_{2,d}, w_{3,d}, \dots, w_{N,d}]^T$, where

$$w_{t,d} = (1 + \log \text{tf}_{t,d}) \times \log \frac{N}{\text{df}_t}$$

$\text{tf}_{t,d}$ is the term frequency of term t in document d (a local parameter)

For any two documents d_j and d_k , their similarity is

$$\text{sim}(d_j, d_k) = \frac{\vec{d}_j \cdot \vec{d}_k}{\|\vec{d}_j\| \|\vec{d}_k\|} = \frac{\sum_{i=1}^n w_{i,j} w_{i,k}}{\sqrt{\sum_{i=1}^n w_{i,j}^2} \sqrt{\sum_{i=1}^n w_{i,k}^2}}$$

where w_i is the weight of the term (Here tf-idf value) n is the number of dimensions of document vector

B. K-Means Clustering

K-Means is one of the simplest unsupervised learning algorithms that solve the well known clustering problem. The procedure follows a simple and easy way to classify a given data set through a certain number of clusters (assume k clusters). The number of clusters should match the data. An incorrect choice of the number of clusters will invalidate the whole process. An empirical way to find the best number of clusters is to try K-means clustering with different number of clusters and measure the resulting sum of squares.

The basic K-means Algorithm is as follows:

- Step 1: Select K points as initial centroids
- Step 2: repeat
- Step 3: Form K clusters by assigning each point to its Closest centroid
- Step 4: Recompute the centroid of each cluster
- Step 5: Until centroids do not change.

The key step of Basic K-means algorithm is selection of proper initial centroids. Initial clusters are formed by random initialization of centroids.

C. Hadoop

Hadoop is a platform that provides both distributed storage and computational capabilities. It is an open source software project that enables the distributed processing of large data sets across clusters of commodity servers. It is designed to scale up from a single server to thousands of machines, with a very high degree of fault tolerance. Hadoop[7] is a distributed master-slave architecture that consists of Hadoop distributed file system (HDFS) for storage and Map-Reduce for computational capabilities. Hadoop can handle all types of data from disparate systems: structured, unstructured, log files, pictures, audio files, communications records, email – just about anything you can think of, regardless of its native format. Even when different types of data have been stored in unrelated systems, you can dump it all into your Hadoop cluster with no prior need for a schema.

Hadoop Distributed File System (HDFS)[5] is a file system that spans all the nodes in a Hadoop cluster for data storage. The HDFS splits large data files into chunks that are managed by different nodes in the cluster. Each chunk is replicated across several nodes to address single node outage or fencing scenarios.

D. MapReduce Programming

MapReduce[6] runs as a series of jobs, with each job essentially a separate Java application that goes out into the data and starts pulling out information as needed. Based on the MapReduce design, records are processed in isolation via tasks called Mappers. The output from the Mapper tasks is further processed by a second set of tasks, the Reducers, where the results from the different Mapper tasks are merged together. Using MapReduce instead of a query gives data seekers a lot of power and flexibility, but also adds a lot of complexity.

The *Map* and *Reduce* functions of *MapReduce* are both defined with respect to data structured in (key, value) pairs. *Map* takes one pair of data with a type in one data domain, and returns a list of pairs in a different domain:

$$\text{Map}(k1, v1) \rightarrow \text{list}(k2, v2)$$

The *Map* function is applied in parallel to every pair in the input dataset. This produces a list of pairs for each call. After that, the MapReduce framework collects all pairs with the same key from all lists and groups them together, creating one group for each key.

The *Reduce* function is then applied in parallel to each group, which in turn produces a collection of values in the same domain:

$$\text{Reduce}(k2, \text{list}(v2)) \rightarrow \text{list}(v3)$$

Each *Reduce* call typically produces either one value $v3$ or an empty return, though one call is allowed to return more than

one value. The returns of all calls are collected as the desired result list.

Thus the MapReduce framework transforms a list of (key, value) pairs into a list of values.

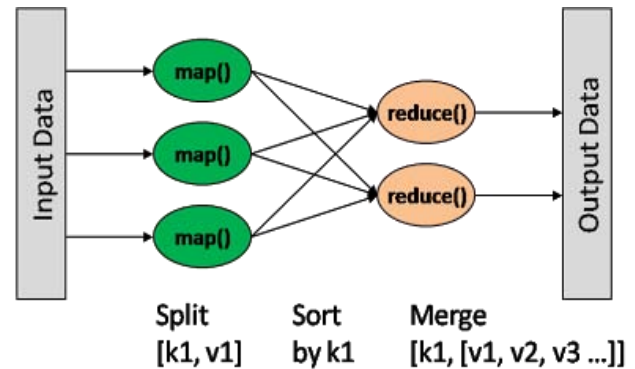


Figure 1: MapReduce Tasks

3. Mapreduce K-Means Clustering

The mapreduce k-means clustering approach for processing big text corpus [4] can be done by the following steps:

- 1) Give the Sequence file from directory of text documents as input.
- 2) Tokenize and generate TF-IDF vector for each document from sequence file.
- 3) Apply Map-Reduce K-Means algorithm to form k clusters.

A. Sequence File from Directory of Text Documents

Map reduce programming is coined to process huge data sets in parallel and distributed environment. Suppose, we select the input data from a document set, where the text files in the directory are small in size. Since HDFS and Mapreduce are optimized for large files, convert the small text files into larger file i.e., SequenceFile format. SequenceFile is a hadoop class, which allows us to write document data in terms of binary <key, value> pairs, where key is a Text with unique document id and value is Text content within the document in UTF-8 format. SequenceFile packs the small files and process whole file as a record. Since the SequenceFile is in binary format, we could not able to read the content directly but faster for read /write operations.

B. Creating TF-IDF Vectors

The sequence file from the previous step is fed as Input to create vectors. The TF-IDF vectors are calculated in Mapreduce by the following steps:

Step 1: Tokenization: The input fed to map function is in format of <key, value> pairs, where key is the document name and value as document content. The outcome of reduce function is also <key, value> pair where key is document name and value are tokens (words) present in that document.

Ex: Key: /acq1.txt: Value: [macandrews, forbes, holdings, bids, revlon, mcandrews, forbes, holdings, inc,said, offer, dlrs, per, share, all, revlon, group]

Step 2: Dictionary file: This step assign unique number to each token in all documents. The input format for the map function is <document name, wordlist> and the output of reduce function is <Word, uniqueid>.

Ex: Key: accounts: Value: 152.

Step 3: Frequency count: The number of times the word appears globally in all documents is calculated in this step. The input to this map function is <docname, words> and output format is <word id, 1>. The value of output of the map function is accumulated and find the sum in reduce function. The output format of reduce function is <word id, count>.
 Ex: Key: 50: Value: 2

Step 4 : Calculate term frequency : The map reduce function in this step takes input as <docname, wordlist > and counts the number of times each word or term t_i occurs in that document d_j . The outcome of this step is in the format of <docname, { t_i : count}>.

Ex: In the below example, acq1.txt is document name and the values are in the format of list of (wordid:count)
 Key:/acq1.txt: Value: {3258:1.0, 3257:1.0, 157:2.0 ...}

Step 5 : Calculate tf-idf value : The output of step 4 is taken as input to map function of this step and calculates the weight vector as $tf \times tf\text{-idf}$ value of each term t_i in each document d_j as specified in equation(1). The output format of the result is <docname, { t_1 : tf-idf1, t_2 :tf-idf2 ... t_i :tf-idfi}>.
 For example: The output format of this step is as follows:
 Key: acq1.txt: Value: {3258:0.12728, 3257:0.12728, 462:0.08060 ...}

C. MapReduce K-Means Algorithm

The implementation of map reduce k-means accepts two input files. One input file contains the documents with each term and its tf-idf values, and the second is k initial centroids file. The set of k initial centroids are selected randomly. In every iteration, the map reduce framework splits the input data into M splits and then processed in parallel as shown in Figure 2.

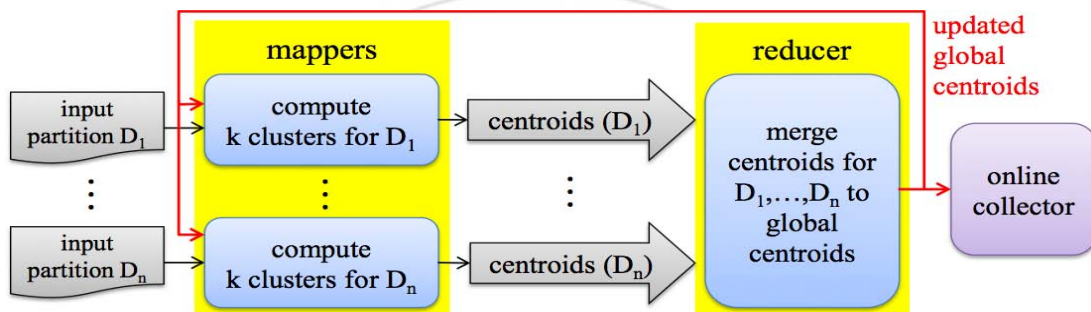


Figure 2: MapReduce Framework for K-Means Algorithm

The map function read the document in the format of <docname, { t_1 :tf-idf1, t_2 :tf-idf2, ... t_i :tf-idfi}> along with randomly selected set of k initial centroids. The map function determines the document set which are closer to the centroids by calculating cosine similarity measure and emits records containing all documents data with kcentroids in the format as <k-center, {docname, { t_1 :tfidf1, t_2 :tf-idf2, ... t_i : tf-idfi}>}. The reducer function receives the output of map function k-centroids along with closest documents bound to it. and calculates new k-centroid. The mapper and reducer algorithm of k-means is explained as below:

Algorithm for Mapper

Input: A set of objects $X = \{x_1, x_2... x_n\}$, A Set of initial Centroids $C = \{c_1, c_2, ,c_k\}$
 Output: An output list which contains pairs of (C_i, x_j) where $1 \leq i \leq n$ and $1 \leq j \leq k$

Procedure

```

M1 ← {x1, x2... xm}
current_centroids ← C
Distance (p, q) = √Σ(p_i - q_i)² (where p_i (or q_i) is the coordinate of p (or q) in dimension i)
for all xi ∈ M1 such that 1 ≤ i ≤ m do
    bestCentroid ← null
    minDist ← ∞
    for all c ∈ current_centroids do
        dist ← distance (xi, c)
        if (bestCentroid = null || dist < minDist)
            then
                minDist ← dist
    
```

```

        bestCentroid ← c
    end if
end for
emit (bestCentroid, xi)
i += 1
end for
return Outputlist
    
```

Algorithm for Reducer

Input: (Key, Value), where key = bestCentroid and Value = Objects assigned to the lpgri; 1\} x centroid by the mapper

Output: (Key, Value), where key = oldCentroid and value = newBestCentroid which is the new centroid value calculated for that bestCentroid

Procedure

```

Outputlist ← outputlist from mappers
∅ ← {}
newCentroidList ← null
for all β outputlist do
    centroid ← β.key
    object ← β.value
    [centroid] ← object
end for
for all centroid ∈ ∅ do
    newCentroid, sumofObjects,
    sumofObjects ← null
    for all object ∈ ∅ [centroid] do
    
```



```

sumofObjects += object
numofObjects += 1
    end for
newCentroid ← (sumofObjects +
numofObjects)
emit (centroid, newCentroid)
end for
end
    
```

The outcome of the k-means map reduce algorithm is the cluster points along with bounded documents as <key, value> pairs, where key is the cluster id and value contains in the form of vector: weight. The weight indicates the probability of vector be a point in that cluster. For Example: Key: 92: Value: 1.0: [32:0.127,79:0.114, 97:0.114, 157:0.148 ...].

The final output of the program will be the cluster name, file name: number of text documents that belong to that cluster.

4. Experimental Results

This section presents the results obtained by executing map reduce k-means clustering algorithm on cluster of machines. The experimentation with 20_newsgroups dataset is explained in detailed in below sections.

A. Environment setup

The experimentation is conducted in single node cluster. The single node is configured with I7processor, 4 GB memory, 64GB hard disk along with JDK 1.7.0 and hadoop 2.4.1 version. The Operating system used is Ubuntu 14.04LTS.

B. Dataset description

The 20 Newsgroups data set is a collection of approximately 20,000 newsgroup documents, partitioned (nearly) evenly across 20 different newsgroups. It was originally collected by Ken Lang, probably for his [Newsweeder: Learning to filter netnews](#) paper, though he does not explicitly mention this collection. The 20 newsgroups collection has become a popular data set for experiments in text applications of machine learning techniques, such as text classification and text clustering. The data is organized into 20 different newsgroups, each corresponding to a different topic. Some of the newsgroups are very closely related to each other (e.g. [comp.sys.ibm.pc.hardware](#) / [comp.sys.mac.hardware](#)), while others are highly unrelated (e.g. [misc.forsale](#) / [soc.religion.christian](#)).

C. Data pre-processing

The data available here are in .tar.gz bundles. You will need tar and unzip to open them. Each subdirectory in the bundle represents a newsgroup; each file in a subdirectory is the text of some newsgroup document that was posted to that newsgroup. Create a sequence file from 20_newsgroups text documents. Sequence file is passed as input to find tf-idf value of each term in every text document. The tf-idf file is fed as input to map reduce k-means algorithm for form k number of clusters.

D. Results

In a single node cluster, the map reduce k-means algorithm is executed with 20_newsgroups dataset with selected number of text files of 20,000 from 20 different topics. The algorithm

tested in single node cluster with 20_newsgroups database and the performance of the map reduce algorithm increases efficiently by increasing the number of input files and increasing number of nodes in the cluster. The runtime execution in terms of seconds can be given as the CPU time spent for each iteration.

E. Output

```

Combine Input records=0
Combine output records=0
Reduce input groups=0
Reduce shuffle bytes=3594132
Reduce input records=0
Reduce output records=0
Spilled Records=0
Shuffled Maps=1
Failed Shuffles=0
...
GC time elapsed (ms)=1749
CPU time spent (ms)=19630
...
    
```

Time taken for iteration 0

```

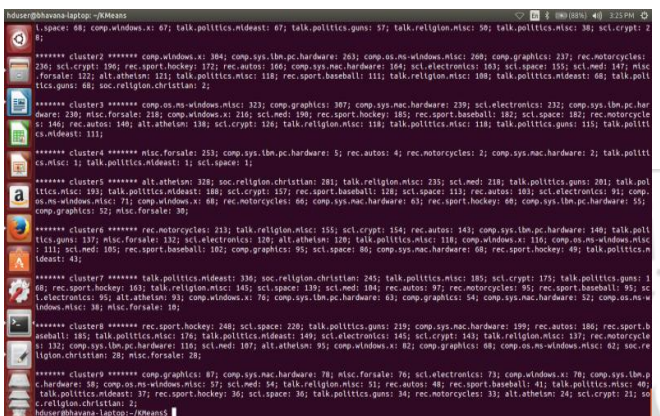
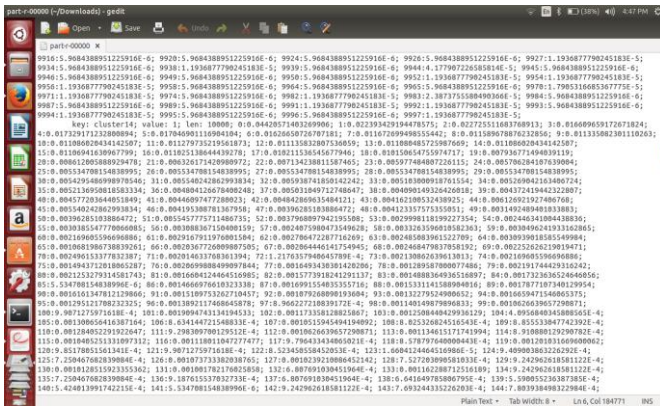
Total time spent by all map tasks (ms)=20518
Total time spent by all reduce tasks (ms)=7908
Total vcore-seconds taken by all map tasks=20518
Total vcore-seconds taken by all reduce tasks=7908
Total megabyte-seconds taken by all map tasks=31018432
Total megabyte-seconds taken by all reduce tasks=8997744
...
GC time elapsed (ms)=1624
CPU time spent (ms)=19730
...
    
```

Time taken for iteration 1

```

FILE: Number of bytes written=7375679
FILE: Number of read operations=0
FILE: Number of write operations=0
...
GC time elapsed (ms)=1711
CPU time spent (ms)=19640
...
    
```

Time taken for iteration 2



Clusters

5. Applications and Future Work

Information retrieval techniques are widely popular in most of the search engines to efficiently organize and retrieve information systems. Most of the data in internet is in the format of unstructured and semi structured. Currently clustering techniques are used to organize and group the similar data objects to retrieve search results faster. Traditional way of clustering text documents is Vector space model, in which tf-idf is used for k-means algorithm with supportive similarity measure. As the data is enormously increasing day by day, elastic resources are required to store and compute. Hadoop framework supports to store and compute big data in parallel and distributed platform with the help of HDFS and Map reduce.

This paper exhibits an approach to cluster text documents using vector space model. The results obtained by executing map reduce k-means algorithm on single node cluster shows that the performance of the algorithm increases as the text corpus increases. In the current approach of map reduce k-means algorithm, the initial centroids are selected randomly. In future, Canopy clustering algorithm can be used as initial step to get initial centroids and then fed this as input to map reduce k-means algorithm. In future, we can try to cluster text corpus by map reduced text clustering based on frequent itemsets.

References

[1] Hadoop The definitive guide.

[2] G. Salton, A. Wong, C. S. Yang. A vector space model for automatic indexing. Communications of the ACM, version.18 n.11, pages.613-620, Nov. 1975.

[3] George Tsatsaronis and Vicky Panagiotopoulou. A generalized vector space model for text retrieval based on semantic relatedness. Proceedings of the EACL 2009 student research workshop, pages 70-78, April 2009.

[4] J Ditttrich, JA Quiané-Ruiz. Efficient big data processing in Hadoop MapReduce. Proceedings of the VLDB Endowment, 2012 - dl.acm.org, Volume 5 Issue 12, August 2012, Volume 5 Issue 12, August 2012.

[5] Sanjay, G., G. Howard, and L. Shun-Tak, The Google file system, in Proceedings of the nineteenth ACM symposium on Operating systems principles. ACM: Bolton Landing, NY, USA 2003,

[6] Jeffrey, D. and G. Sanjay,. MapReduce: simplified dataprocessing on large clusters. Commun. ACM, 51(1): pages 107-113 2008.

[7] Apache Lucene Hadoop[EB/OL]. <http://hadoop.apache.org/>.

Author Profile



Botcha Chandrasekhara rao is currently pursuing his 2 years M.Tech in CSE at Pydah Kaushik College of Engineering, Gambheeram Village, Anandapuram Mandalam, Visakhapatnam. His area of interest includes Cloud Computing.



Medara Rambabu is currently working as HOD & Associate Professor in Department of Computer Science & Engineering at Pydah Kaushik College of Engineering, Gambheeram Village, Anandapuram Mandalam, Visakhapatnam. He completed his M.Tech from JNTUK in 2010, and he completed his B.Tech from AU, in 2004. He has more than 6 years of teaching experience in various engineering colleges. His research interest includes Computer Organization, Data Communication Systems, Computer Networks and DM & DW.