

International Journal of Computer Science and Mobile Computing



A Monthly Journal of Computer Science and Information Technology

ISSN 2320-088X

IMPACT FACTOR: 6.017

IJCSMC, Vol. 6, Issue. 8, August 2017, pg.84 – 88

Map Reduce Design for Data Clustering

Priyanka A

Computer Science & VTU, India

priya.bhat.a@gmail.com

Abstract— The data is of no use unless or until the useful information is retrieved from it. It is impractical and inefficient to use traditional database techniques since the data is generated in practical scenario is huge in volume and it is exchanged in faster way than ever before. That is why a new open source framework called Hadoop comes in to existence. Hadoop is used to process huge amount of data called big data. To extract useful information and to derive some meaningful conclusions from massive amount of data, datamining techniques can be used. Among many datamining techniques, clustering is most popular one. Clustering is a process of binding data members with similar characteristics in to one group whereas dissimilar data members are distributed into different groups. In this paper K-Means clustering algorithm is implemented using map reduce programing model. The key idea behind the KMeans implementation is design of mapper and reducer routines. The technique presented here asks to enter the number of clusters that is K and large data file. It finds centroid for each cluster and distance from each item to centroid of the cluster is computed. The data is assigned to the group with which the distance of the data is least.

Keywords— Hadoop, Data Mining, Big Data, K-Means, Clustering

I. INTRODUCTION

In today's world, the production of data which is coming quickly from various sources creates difficulty in their collection, storage, analysis, management and processing. The generated information becomes useless unless or otherwise they are analysed correctly. Big data provides a set of techniques to handle these challenges in order to grab the opportunities that they offer.

Data mining is one such technique used to discover useful information from huge data set. But existing data mining techniques cannot be applied directly on huge volume of data because it increases the complexity of processing. Clustering is one such approach in data mining used to analyse large volume of data generated by real world applications. Clustering is a process of partitioning data in such a way that data with similar characteristics grouped into single cluster and data with dissimilar characteristics belong to different clusters. In traditional method this process becomes tedious as the data increases in quantity. There are different types of clustering algorithm used to partition data in to clusters.

A. Clustering Algorithms

1) *Hierarchical Algorithms:* In this algorithm data is ordered depending on the method of nearness (proximity).The proximity is obtained by the middle nodes. A dataset is represented by dendrogram, where every information represented by leaf nodes. As the hierarchy continues the initial group (cluster) is progressively divided in to several groups.

2) *Partitioning Algorithms:* These algorithms divide datasets into a different number of partitions, in which every partition is considered as a single cluster. The following requirements must be satisfied by these clusters. Each cluster should contain minimum of one data object. Each data object must belong to only one cluster.

3) *Density-Based Algorithms*: In this algorithm, based on the region of density and connectivity the data items are grouped into clusters. They are closely linked with point –nearest neighbours.

4) *Grid-Based Algorithms*: In this data object space is divided in to grids. This algorithm scans the input dataset only once to find statistical value for grids. Based on accumulated grid data grid based clustering is applied.

5) *Model-Based Algorithms*: In this method clustering is done based on some already defined mathematical model. This approach is depending on an assumption that data is generated by a combination of probability distribution. This method automatically determines the number of clusters based on statistical data.

In this paper, KMeans clustering algorithm is implemented for large data set. The KMeans clustering algorithm comes under the category of partition based algorithm. The K-Means algorithm is an unsupervised algorithm used to group dataset with similar characteristic into one cluster and data with dissimilar properties belong to other cluster. This algorithm is used to extract some meaningful information from the available dataset. It allocates n objects to k partitions in such a way that every data item belong to nearby partition.

The main requirement for efficient clustering method is

- Intra-clustering similarity should be maximum
- Inter clustering similarity should be minimum

B. About Hadoop

To implement KMeans clustering on big data an open source java framework called Hadoop is used. Hadoop framework works based on Map Reduce programming model, which consists of mapper and reducer task.

•The Mapper Task: This task is the one which takes input data and changes it into a set of data, where every single component is broken down into tuples of the form key-value pairs.

•The Reducer Task: The output of the mapper task is given as an input to this task. It combines the data tuples into a tuples of smaller set. This task is executed after the execution of mapper task.

II. RELATED WORK

Purnawansyah and Haviluddin [3] in this paper they proposes KMeans clustering implementation in network traffic activities. Here, k-means was used to identify bandwidth used by the users per day. The algorithm was applied on 456 data set and grouping is done based on low, medium and high bandwidth usage. This algorithm was implemented using MATLAB2013 tool. The information analysed in this paper was information exchange over a period of 30 minutes per day for three units (Rectorat, science, forestry) in five months. The data was normalized before sending through network between 0.05 and 1. The normalized value was calculated for the data using original, minimum, maximum of the raw data.

QingHe, XinJinChangyingDu, FuzhenZhuang and Zhongzhi Shi [4] tells about Clustering in extreme learning machine feature space in which ELM K-Means and ELM NMF uses ELM feature on Kmeans to solve the clustering problem. It was a partition based algorithm which takes dataset from document corpus as well as from UCI machine learning repository. If the number of nodes was more than 300 than this algorithm gives very good accuracy else optimal performance cannot be obtained.

R Madhuri, M Ramakrishna Murty, JVR Murthy, PVGD Prasad Reddy and Suresh C Satapathy [5] says Cluster analysis on different datasets using K-modes and K-prototype algorithms. This was a kind of partitioning based algorithm and takes data set from Mixed Numeric and Categorical data .The quality of cluster generated by applying this method is not very good.

III. EXISING SYSTEM

Clustering is an important technique in data mining to retrieve a meaningful information. There are different types of clustering algorithms existing which are applied on different types of datasets. These techniques uses traditional database management tools to store the data. The traditional storage systems can store only some limited amount of data. As the data size increases data access and management becomes a tedious task. When the huge amount of data is given to these algorithm, the time it takes to process these data sets also increases and cluster quality reduces. So we need to have a tool to process massive volume of data.

IV. PROPOSED SYSTEM

The graphical representation of the proposed system is shown below.

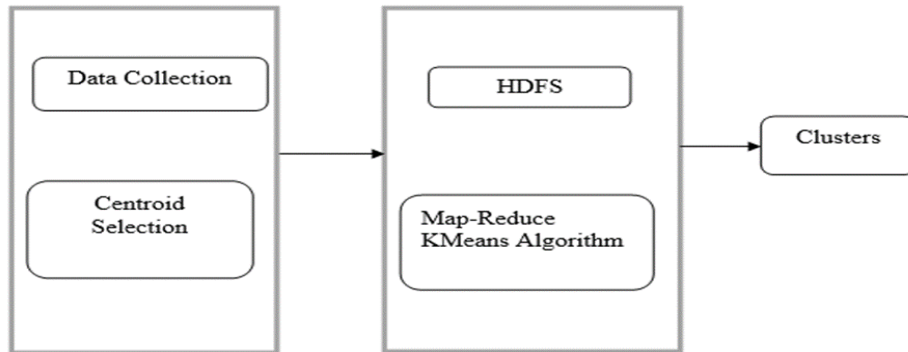


Fig. 1 Model of Proposed System

Proposed system uses an open source framework called Hadoop. It consists of HDFS and map reduce programming model. Since HDFS uses data replication concept, the loss of data is impossible. This file system is capable of managing large amount of data, data access becomes easier. In this project KMeans clustering algorithm is implemented using map reduce construct. The mapper function computes Euclidean distance between cluster centroid and data point and the reducer function recalculates the cluster centroid. So computation time it takes to group the data with similar characteristic in to a same cluster is reduced and cluster quality is improved.

V. IMPLEMENTATION

In this section map reduce methods are implemented using Hadoop framework to cluster huge amount of data set using KMeans clustering algorithm.

A. Algorithm for map reduce KMeans

- Read the data set
- Compute the vector value for each data point and store it in a vector file
- Input number of clusters
- Select the cluster centroids randomly and store it on a cluster file
- Run the Hadoop components
- Copy the vector and cluster file onto HDFS file system
- Run MapReduce K-Means algorithm
- Copy the cluster centroids obtained in the last iteration into local file system
- Based on cluster centroids assign the data into different cluster

B. Mapper Function

This method takes input as a list of <key, value >pair where key is the location of the data object and value is its content.

This method produces output in the form of <key, value>pair in which key is index of the cluster and value is the data object belonging to that cluster. This function performs following task.

- Calculate the distance between data object and each cluster centroids
- Assign the data object to its nearest cluster centroid
- Repeat above two steps until all data objects are processed

C. Reducer Function

It takes input as a list of < key, value > pair. The output produced by this function is of the form <key, value > pair where key represents the index of the cluster and value specifies the centroid for the new cluster. This function performs following tasks.

- Computes sum of data objects in each cluster
- Computes number of data objects in each cluster
- Computes the mean of data objects belonging to cluster. This is called as new cluster centroid
- Repeat the above steps until the termination condition reaches.

The dataset is collected from UCI Machine Learning Dataset Repository. This dataset consists 20000 text documents from 20 different newsgroups.

VI. RESULTS

Before executing Kmeans algorithm, it is essential to check Hadoop components are running successfully. This is done by executing jps command. In case if any component is not started successfully, then necessary steps has to be taken for successful running of individual Hadoop component.

```
3602 NameNode
4516 Jps
4087 ResourceManager
3720 DataNode
3929 SecondaryNameNode
4206 NodeManager
```

Fig. 2 Check all Components of Hadoop are running

The output of map reduce program is stored in a default file called part-r-00000. This file is present in HDFS. Here, part-r-00000 contains centroids of each cluster after 3 iterations.

```
key: cluster1; value: 10; len: 10000; 0:0.0935823033117506; 1:0.19960329275427124; 2:0.06617479520849329; 3:0.05346452883637152;
4:0.05267916124306217; 5:0.042303564202445965; 6:0.034029583376327145; 7:0.02660908311845816; 9:0.021031675463567785; 10:0.0618766724024157
11:0.01750112344827458; 12:0.03410228327471742; 13:0.02585567658316511; 14:0.06763244174741027; 15:0.030200863403487346;
16:0.015274723043685465; 17:0.022079345834488423; 18:0.03900233240847896; 19:0.01949958068144829; 20:0.008168717445797792;
21:0.055770158015045125; 22:0.008593127061553628; 23:0.011981899596627541; 24:0.024937594967648448; 25:0.006927086247048727;
26:0.050402992439369444; 27:0.00583125744951417; 28:0.02683542794339059; 29:0.007041596321460505; 30:0.0062623818215725964;
31:0.0034287465477656114; 32:0.0033827463191707335; 33:0.027766598391349734; 34:0.028090626608201236; 35:0.02836487839255726;
36:0.027730688673797962; 37:0.026871552816280033; 38:0.028174450271513177; 39:0.02657875919191581; 40:0.02657875919191581;
41:0.010482539708936908; 42:0.027174403627042654; 43:0.004768559397574071; 44:0.008962833798851759; 45:0.023356075863491595;
46:0.003618167075451892; 47:0.008785778652242757; 48:0.023310190984727825; 49:0.006282208156596274; 50:0.0037264225351381734;
51:0.0023738227607800283; 52:0.004103818564066353; 53:0.009536112521309318; 54:0.0033322634972872294; 55:0.0026419626195341643;
56:0.00664486380653439; 57:0.006032715873399032; 58:0.005723874963389806; 59:0.01487766873256681; 60:0.005097065620224921;
61:0.0013650516505624506; 62:0.017067173993653446; 63:0.0021330135880663725; 64:0.0173275818711065; 65:0.009612877821194;
66:6.281907435523367E-4; 67:0.0030230754390547033; 68:0.021353178189579165; 69:0.0604712695338728; 70:0.006548062463132712;
71:0.010323602966588833; 72:0.0037038919574375385; 73:0.00814255977946581; 74:8.988441554112201E-4; 75:0.00374442226547774;
76:0.008683930836655935; 77:0.010344894606406358; 78:0.0028720328564132724; 79:0.008683930836655935; 80:0.004712627542059863;
81:0.008683930836655935; 82:0.014531232039444933; 83:0.0102259126459103; 84:0.008439560442309516; 85:0.0015493959041958949;
86:0.0037335824786941803; 87:0.003262199837864346; 88:0.011156292029074547; 89:7.747010227958201E-4; 90:0.018196160489014516;
91:0.0015174417734737237; 92:0.02128359708397813; 93:0.0010613646427263441; 94:0.0020031830338299032; 95:0.004292023127784116;
96:0.0019964295906154403; 97:4.781893114085104E-4; 98:0.003702624491215896; 99:0.0175002275949392; 100:0.0020538211943303317;
101:9.993509039144994E-4; 102:0.0035649136080189523; 103:0.0022073548411977706; 104:4.4570714505355017E-4; 105:0.003886872079293705;
106:0.00438553739042519; 107:0.01022353495443138; 108:0.0019964295906154403; 109:8.687058092438253E-4; 110:6.188335947208321E-4;
111:0.009161821647912622; 112:0.18950416229653771; 113:0.19325440059458276; 114:2.513826043237800E-4; 115:0.0015502763536108602;
116:0.0010389397623468309; 117:0.004166692797183029; 118:0.0011907977602909764; 119:0.003045472145857076; 120:8.716718773147899E-4;
121:0.03015055575627733; 122:0.007264730026113076; 123:0.012477592961530578; 124:0.0018784567407591221; 125:0.0035743408250219696;
126:0.0012257987208393304; 127:0.009789588052395298; 128:0.0034501942534816915; 129:0.003406354329375045; 130:0.002287606450186719;
131:0.00536716557567568; 132:3.696744058474283E-4; 133:0.002957212152575642; 134:0.012390678763227145; 135:0.009120640780260035;
136:0.004593783512343064; 137:0.012884603302795464; 138:0.02382423327475366; 139:0.0012369264217678274; 140:8.050722515060109E-4;
```

Fig. 3 Part-r-00000 file

Here text files are grouped based on final centroids obtained as a result of the execution of map reduce programs. The output is 20 clusters with number of text documents from each subdirectory assigned to each cluster.

VII. CONCLUSIONS

The amount of data exchange in these days requires high quantity of data processing. This project implements K-Means Algorithm in a single node cluster using Hadoop framework to group given text documents based on certain patterns. This application provides an efficient method and robust system to group data with similar characteristics, decreases the problems in handling huge volume of data and also reduces the execution time.

REFERENCES

- [1] P. P. Anchalia, A. K. Koundinya and S. N. K., "MapReduce Design of K-Means clustering Algorithm", International Conference on Information Science and Applications (ICISA), Suwon, pp. 1-5,2013.
- [2] Adil Fahad,Najlala Alshatri,Zahir Tari, Abdullah Alamri, "A Survey of clustering algorithms for big data", IEEE, Vol 2, pp 267-279,September 2014.
- [3] Purnawansyah and Haviluddin,"K-Means clustering implementation in network traffic activities", IEEE, International Conference on Computational Intelligence and Cybernetics, Makassar, pp. 51-54, 2016.
- [4] Qing He, Xin Jin, Changying Du, Fuzhen Zhuang, and Zhongzhi Shi., "Clustering in extreme learning machine feature space.Neurocomputing",pp 88-95, 2014.
- [5] R Madhuri, M Ramakrishna Murty, JVR Murthy, PVGD Prasad Reddy, and Suresh C Satapathy, "Cluster analysis on different data sets using k-modes and k-prototype algorithms", In ICT and Critical Infrastructure: Proceedings of the 48th Annual Convention of Computer Society of India-Vol II, pp 137-144, Springer, 2014.
- [6] Apache Hadoop: <http://hadoop.apache.org/>
- [7] www.bogotobogo.com/Hadoop/BigData_hadoop
- [8] www.michael-noll.com/tutorials