

Large Scale Data Clustering Using Various-Widths Clustering Approach

Mrs. Harshal Agashe¹, Prof. Satish Banait²

¹M.E Student, Computer Engineering, KKWIEER Nashik, Maharashtra, India,

²Assistant Professor, Computer Engineering, KKWIEER Nashik, Maharashtra, India,

ABSTRACT

To perform a clustering widely used and most powerful technique is used i.e k-nearest neighbor. This approach required large computational cost for high dimensional datasets. The proposed work focuses on k-NN is based on various clustering widths on large scale data. We are proposing modified kNN approach with MapReduce parallel computing algorithm and clusters grouping with goal of improving the performance in terms of clustering time, pre-processing costs and querying cost while working with high dimensional data. First we are presenting the kNN method using various width clustering to efficiently extract the kNNs for input query object from the dataset. The given dataset is clustered using global width then each cluster that satisfies its predefined criteria i.e threshold value is recursively clustered using their local width. To prune unlikely clusters triangle inequality was used earlier, but we designed tree based approach in which centers of clusters grouped into the tree based index to maximize the more clusters pruning. To reduce the processing time and clustering time, we designed parallel computing algorithm based on MapReduce.

Keyword : - Clustering, k-Nearest Neighbor, Tree Index, large scale data, MapReduce.

1. INTRODUCTION

To perform a clustering widely used and most powerful technique is used i.e k-nearest neighbor. It plays an most used in unsupervised learning, also to measuring an quality of cluster this is mainly used. In clustering quality indexes has been proposed in many years and different indexes are used in different area also. Clustering is useful for grouping similarities also an decision making in the machine learning which including data mining, document information retrieval, image segmentation. Clustering is task of to find

homogeneous groups of the studied objects. Many researchers is interested to developed an clustering using various types of algorithms. The clustering main issue is we don't have any kind of knowledge, information of data. Using hidden pattern the quality of clustering will be measure.

In many research domains K-nearest neighbor (kNN) is widely used for information retrieval and classification process. The input set of objects P and test query object Q, the kNN query extracts the k similar objects to Q from the set P. Since from last decade, the kNN problem was extensively studied by various researchers. There are number of methods introduced in order to compute extract or approximate outcomes based on requirements of applications and end users. The existing approximate based methods aiming to achieve more efficiency at the quality accuracy cost. The exact methods are very costly, but they produce the correct results. Therefore it is required to have exact technique over high dimensional data rates with minimum cost. Recently exact kNN based method proposed. This is novel kNN approach based on different widths in the cluster using that we can

improve the efficiency. So that to improve this method by adding parallel computation framework and tree based cluster groupings. In single clustering[12] are used with following properties.

Large scale of dataset is task to perform an clustering in limited time and producing an better output also easy to understand it so this is focused on the preprocessing the data i.e construction of index in different distribution of data they are divide into two parts those 1.Tree based indexes.2.Flat based indexes .Tree based indexes is used binary partition method for construct a tree for given dataset. The k-NN is use flat indexes approach and the triangle inequality for efficient to prune the node of tree. Tree based indexes have technique to partition a data to recursive fashion for building the a tree of data .Using Various-width clustering approach and find k-NN search [4]using triangle inequality to efficiently find the query object and computing cost to partitioning the clusters.This process involved following operation those are cluster width learning, partitioning and merging And then construction of cluster into tree like index [8].To reduce the processing time and clustering time we designed parallel computing algorithm

In the below sections we are going to discuss about related work done for the proposed research area. We refer some existing research paper for completing this task. It is given as follow:

2. RELATED WORK

A.Gupta [1] presented a new the method an extraction anomalies of Nuclear Power Plant time sequence created an Data with the help the Fixed Widths Clustering Algorithm. Basically the time data will be recorded successive point in the time. To find their anomalies as well as correlation and pattern of time series. Causes will be found their corrective action are taken. Using an dynamic method is to decide which cluster width will be used for clustering the data. The algorithm for fixed-width clustering is based on the outline in . Anomaly detection using fixed width clustering is a three stage process, (1) normalization, (2) cluster formation ,and (3) cluster labeling. Anomaly detection using brute force method on SAX is an $O(n^2)$ complexity algorithm. We have used it for first dataset because, this algorithm is very accurate and it allows dimensionality/ numerosity reduction of the original dataset. Hence we can reduce the size of original dataset before processing it.

A Ezugwu [2] presented a new approach Performances Characterization of heterogeneous distribution of cluster resources in evaluation technique to the measurement based is characterise into specific performance and their contribution of individual resources of clusters configurations .Also identify the initial parameters are considered at the stage of selection and their allocation an computational node of an all application for execution.For selection of resources parameter(Proceesor speed, associated with the bandwidth and size of memory in each processor) and also their interactions. To map out strategies necessary for determining and improving the performance of clusters with varying resource specification requirements. At the end of this research, we would have demonstrated the need for resource characterization and how it could be used to determine and predict resource quality and also improve performance for user application on a high scale.

A .Almalawi [3] proposed a new method an data driven technique clustering to detect an attack on SCADA system .It is method of data acquisition in supervisory control systems having small salient part to control the critical infrastructure, for e.g power plants, various energy grids, various water distribution systems in that it is automatically identify an normal and critical states of an given system.Also it extract proximity which is based detection rules which will help to identifying states for monitoring purposes. a novel data-driven clustering approach that removes the need for domain experts and the purely “normal” SCADA data to build the detection models.This approach is based on the assumptions that “normal states”, that are represented by a combination of the status and values, of multivariate process parameters in a SCADA system can be clustered into finite groups of dense clusters, and critical states in the n -dimensional space will take the form of noise data, also called outliers.

K.Hajebi Hong Zhang[4] was introduced an new method an fast approximate an nearest neighbor by using k nearest neighbor graph.The algorithm construct a graph is offline phase using its nearest neighbor.when query to a new point then it will perform an hill climbing to start with randomly sample to the node in given graph and lazy learner which they doesn't learn anything from training sample data and used for just classification. We build a k-nearest neighbor (k-NN) graph and perform a greedy search on the graph to find the closest node to the query. introduce the Graph Nearest Neighbor Search algorithm (GNNS) and analyze its performance. compare the GNNS algorithm with the KD-tree and LSH methods on a real-world dataset as well as a synthetically generated dataset.

N.Madicar [5] proposed a new method the time series clustering which are parameter free Subsequences with various width clustering technique. Sub-sequence of Time Series (STS) is a clustering of the subsequences in single time of an series which are their subparts or subsequences in that a single time series. It also the similar group of pattern in time series are combine or the cluster centroid is represent an data of every group of cluster. Algorithm can produce the results very similar to that of the STS clustering algorithm in . In addition, our algorithm could perform better in some cases since the widths of the clusters are allowed to vary. This means that the clustering results can contain the clusters with different widths, by removing the existing requirement that every subsequence to be clustered must have one same length.

G.Karypis [6] was introduced a new approach CHAMELEON : Hierarchical clustering algorithm using an dynamic modeling. In this algorithm breakdown if incorrect choice of parameter of data is clustered or if they not able to find the property of clusters. In clusters data is consist off diverse shapes, density, and sizes. Using this approach we determine the smiler property of two cluster based on dynamic model. An two clusters will combine their inter-connectivity and close to each other and various type of data and their similarity matrix will be constructed. In the clustering process, two clusters are merged only if the inter-connectivity and closeness (proximity) between two clusters are comparable to the internal inter-connectivity of the clusters and closeness of items within the clusters. The merging process using the dynamic model presented in this paper facilitates discovery of natural and homogeneous clusters.

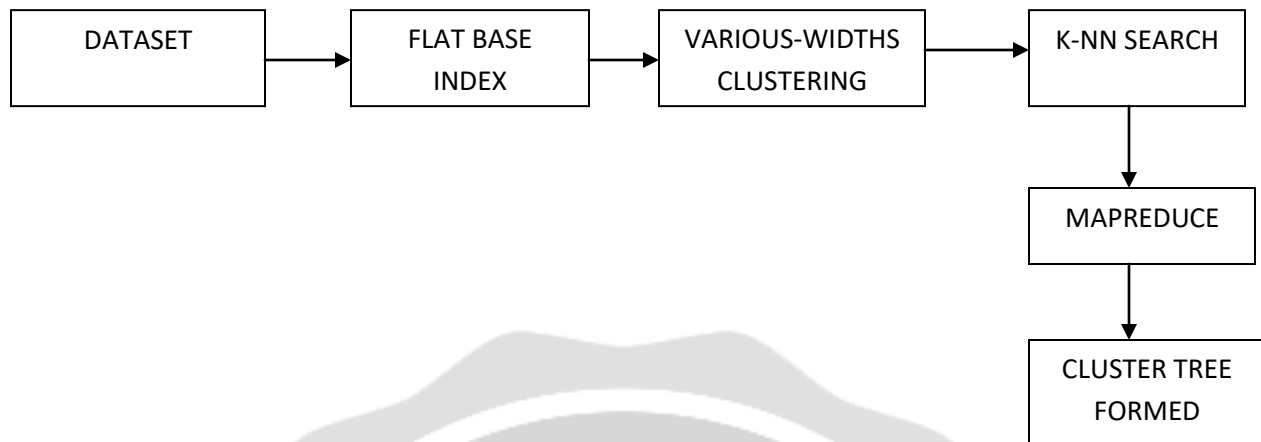
X.Wang [7] present a new method An fast extract k nearest neighbor in high dimensional data with the help of k-means and triangle inequality. This known as kMkNN (k- Means for k-Nearest Neighbor search) using we can accelerate an finding the nearest neighbors. Preprocess of k-means using an metric trees, kd-trees, or ball-tree, kMkNN. Using triangle inequality we reduce the computing distance calculation.

T.Chiang [12] presented a new method for k-NN using ranking based on multilabel classification. Ranking based model is to know to neighbor's labels which are more comfort candidate using weighted KNN-based strategy after that assign higher weights in which candidates having more vote among all the member. The weight are form using generalized searching pattern technique. So possible to improved the accuracy of multilabel data. The approach exploits a ranking model to learn which neighbor's labels are more trustable candidates for a weighted KNN-based strategy, and then assigns higher weights to those candidates when making weighted-voting decisions. The weights can then be determined by using a generalized pattern search technique.

S.guha[13] was proposed a new method ROCK: A Robust Clustering Algorithm are used their Categorical attributes. In that clustering algorithm to study data having an Boolean and categorical attributes. Using an distances between points for clustering which are not appropriate for Boolean and categorical attributes so we use of links to find the similarities among their data points. This method is non-metric similarities measures in the relevant in situation. So that we exhibits the good scalability of clustering properties.

J.Maillo[17] was proposed a new method A MapReduce based k-Nearest Neighbor Approach for Big Data Classification. This model allows simultaneously classification of large amounts of data. In map phase to determine k-NN from different splits into the data. After that reduce the number of stage will be compute the definitive neighbors which they are getting in the map phase. This model allows the k-NN classifier to scale to datasets arbitrary size, just simply adding more computing nodes if necessary. This model allows us to simultaneously classify large amounts of unseen cases (test examples) against a big (training) dataset. To do so, the map phase will determine the k-nearest neighbors in different splits of the data. Afterwards, the reduce stage will compute the definitive neighbors from the list obtained in the map phase. The designed model allows the k- Nearest neighbor classifier to scale to datasets of arbitrary size, just by simply adding more computing nodes if necessary.

3. SYSTEM ARCHITECTURE



A. kNNWC: This is an existing method of k-nearest neighbor approach which is based on various-widths clustering. In that using K-NN for searching nearest objects using triangle inequality approach we required low computing cost for partitioning the clusters of various distribution of data.

B. Flat based index: Flat indexes cluster the objects in the data set into a number of clusters. Based on centers and radii of the clusters, and also the triangle inequality which is applying into prune the clusters which they cannot contain the results. The performance of flat indexes depends on the quality of clusters.

C. Fix-Widths Clustering: Fixed width clustering will create the number of clusters having fixed radius (width) w . Here the width w is a parameter should be defined by the user what the user wants the width of cluster.

D. Various-Widths Clustering: Clustering will perform using various width approach for large scale data by performing the cluster width learning i.e. creation of cluster by using threshold value which is given up to largest size of cluster will be formed after they are divided to number of clusters with the help of width suits cluster.

E. k-Nearest Neighbours Search: All instances correspond to points in an n -dimensional Euclidean space. To find nearest neighbor in that cluster.

F. Tree Prune: In tree pruning removing some part of tree after tree has been built. In that various tree pruning methods are used for construction of tree.

G. Map-Reduce: MapReduce is a programming model which is associated with the implementation for processing and it generates large data sets having parallel or distributed algorithm of a cluster. MapReduce program can be composed of a Map() and procedure (method) which performs filtering and sorting and a Reduce() method that performs a summary operation. The "MapReduce System" is also known as "infrastructure" or "framework" also the processing about marshalling about distributed servers so that running the various tasks at the same time in parallel fashion, managing that all communications and data transformation between the various parts of the system and it will provide for redundancy and fault tolerance.

3. RESULT

1. Reduction Rate of Distance Computations: Finding k-NN in distance computations the performance is calculated using the following formula

$$RD = 1 - (m \% (q * n))$$

2.Speedup Rate: To calculate the overall performance of k- NN in clustering algorithm and report to the speedup rate. Let t_1 is the total running time will required for k-NN and t_2 be the running time for algorithm W. Then the speed-up rate of overall algorithm will be

$$W = t_1/t_2$$

Table 9.3: Accuracy of dataset having different k values

Dataset	K values	Existing method	Proposed method
KDD	10	86.95	90.45
KDD	50	65.12	70.45
KDD	100	45.39	49.99
KDD	200	40.15	45.02
SpamBase	10	87.45	89.29
SpamBase	50	87.45	89.29
SpamBase	100	34.19	40.27
SpamBase	200	32.25	39.01
Shuttle	10	87.56	89.78
Shuttle	50	63.56	68.47
Shuttle	100	55.23	59.49
Shuttle	200	53.33	54.22
Waveform	10	82.91	87.23
Waveform	50	63.19	71.95
Waveform	100	56.33	64.66
Waveform	200	52.2	58.98
Waveform-+noise	10	88.92	92.73
Waveform-+noise	50	73.26	79.96
Waveform-+noise	100	59.28	65.27
Waveform-+noise	200	55.23	62.33



Table 9.4: Construction Time

Dataset	K values	Existing method	Proposed method
KDD	10	5123	4082
KDD	50	6346	4938
KDD	100	7189	5925
KDD	200	7882	6844
SpamBase	10	4289	3294
SpamBase	50	5238	4388
SpamBase	100	6343	5126
SpamBase	200	7147	5321
Shuttle	10	4073	3576
Shuttle	50	5892	4085
Shuttle	100	6372	5388
Shuttle	200	7325	5925
Waveform	10	4698	3397
Waveform	50	5238	4489
Waveform	100	6343	5524
Waveform	200	7146	6549
Waveform-+noise	10	4073	3277
Waveform-+noise	50	5092	4488
Waveform-+noise	100	6723	5984
Waveform-+noise	200	7998	6897

5. CONCLUSIONS

Various clustering widths on large scale data to find k-NN by applying MapReduce parallel computing algorithm and grouping of clusters to improve the performance. The given dataset is clustered using global width then each cluster will satisfies its predefined criteria i.e threshold value is recursively clustered using their local width. Tree based approach in which centers of clusters grouped into the tree based index to maximise the more clusters pruning. To reduce the processing time and clustering time using parallel computing algorithm based on MapReduce to adequately prune additional clusters for various distribution to clusters will be form. In this research, we aims at producing compact and well separated clusters from high dimensional data by using grouping of center (radii) of cluster into tree like structure to adequately prune additional clusters and learning the optimum clustering radius for various distribution to minimize the overlapping between the clusters.

5. ACKNOWLEDGEMENT

I would like to thanks to Prof. Dr. S. S. Sane, professor and Head of Department of Computer Engineering and Prof. Dr.K. N. Nandurkar, Principal, K.K.W.I.E.E.R., Nashik for their kind support and their suggestion. We would also expressed special thanks to all staff members of computer engineering department and our colleagues who knowing and unknowing will help us to complete this work. not mandatory.

6. REFERENCES

- [1] Aditya Gupta, Durga Toshniwal "Extracting Anomalies from Time Sequences Derived from Nuclear Power Plant Data by Using Fixed Width Clustering Algorithm" 2013 International Conference on Advances in Computing, Communications and Informatics (ICACCI)
- [2] Absalom E. Ezugwu, Marc E. Frincu, Sahalu B. Junaidu "Performance Characterization of Heterogeneous Distributed Commodity Cluster Resources" 2014 IEEE.
- [3] Abdulmohsen Almalawi, Adil Fahad, Zahir Tari, Abdullah Alamri, Rayed AlGhamdi, and Albert Y. Zomaya, Fellow "An Efficient Data-Driven Clustering Technique to Detect Attacks in SCADA Systems", IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY, VOL. 11, NO. 5, MAY 2016
- [4] Kiana Hajebi and Yasin Abbasi-Yadkori and Hossein Shahbazi and Hong Zhang "Fast Approximate Nearest-Neighbor Search with k-Nearest Neighbor Graph". hajebi, abbasiya, shahbazi, hzhangg@ualberta.ca
- [5] Navin Madicar Haemwaan Sivaraks Sura Rodpongpun Chotirat Ann Ratanamahatana, "Parameter-Free Subsequences Time Series Clustering with Various-width Clusters", 2013 5th International Conference on Knowledge and Smart Technology (KST)
- [6] George Karypis Eui-Hong (Sam) Han Vipin Kumar "CHAMELEON: A Hierarchical Clustering Algorithm Using Dynamic Modeling", To Appear in the IEEE Computer
- [7] Xueyi Wang "A Fast Exact k-Nearest Neighbors Algorithm for High Dimensional Search Using k- Means Clustering and Triangle Inequality", Proceedings of International Joint Conference on Neural Networks, San Jose, California, USA, July 31 August 5, 2011
- [8] Dantong Yu and Aidong Zhang "ClusterTree: Integration of Cluster Representation and Nearest- Neighbor Search for Large Data Sets with High Dimensions", IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 15, NO. 5, SEPTEMBER/OCTOBER 2003
- [9] Adrien Gaidon, Zaid Harchaoui, Cordelia Schmid "Recognizing activities with cluster-trees of tracklets", BMVC, Sep 2012, Guildford, United Kingdom. 2012.
- [10] QING-BAO LIU, SU DENG, CHANG-HUI LU, BO WANG, YONGFENG ZHOU "RELATIVE DENSITY BASED K-NEAREST NEIGHBORS CLUSTERING ALGORITHM" Proceedings of the Second International Conference on Machine Learning and Cybernetics, Wan, 2-5 November 2003
- [11] Chanop Silpa-Anan Richard Hartley "Optimised KD-trees for fast image descriptor matching", 2008 IEEE Hosein Alizadeh, Behrouz Minaei-Bidgoli and Saeed K. Amirgholipo
- [12] "A New Method for Improving the Performance of K Nearest Neighbor using Clustering Technique"
- [13] Sudipto Guha, Rajeev Rastogi, Kyuseok Shim "ROCK: A Robust Clustering Algorithm for Categorical Attributes"
- [14] D.A. White, R. Jain, Similarity Indexing with the SS-Tree, Proc. 12th Intl. Conf. Data Eng., pp. 516-523, Feb. 1996.
- [15] R. Kurniawati, J.S. Jin, and J.A. Shepherd, The SS+- Tree: An Improved Index Structure for Similarity Searches in a High-Dimensional Feature Space, Proc. SPIE Conf. Storage and Retrieval for Image and Video Databases, pp. 13-24, Feb. 1997.
- [16] K. Chakrabarti and S. Mehrotra, The Hybrid Tree: An Index Structure for High Dimensional Feature Spaces, Proc. 16th Intl Conf. Data Eng., pp. 440-447, Feb. 2000.
- [17] J. Maillo, Isaac Triguero "A MapReduce-based k- Nearest Neighbor Approach for Big Data Classification", 2015 IEEE, BigDataSE/ISPA