# A STATISTICAL SURVEY ON THE BIG DATA MINING WITH THE RESULTS

R. Karunia Krishna Priya[1]   Dr.G. Vinodhini[2]   Dr.K. Sakthivel[3]

1. Research Scholar, Department of Computer Science and Engineering, Annamalai University, Annamalai Nagar, Tamil Nadu, India

2. Assistant Professor, Department of Computer Science and Engineering, Annamalai University, Annamalai Nagar, Tamil Nadu, India

3. Professor, Department of Computer Science and Engineering, K.S.Rangasamy College of Technology, Tiruchengode, Tamil Nadu, India

**Abstract:** Big Data Mining is one of the emerging technology which deals with the storage of huge data's in the form of the repositories. The Big data is used in the various fields for the disseminated usage. The induction of various data mining algorithms in the field has made the mining things more efficient in the repositories. Statistically surveying all the possible algorithms which produces the mining operation more feasible is analyzed in this paper and concluded with the best optional algorithm which makes the mining operation more efficient. The algorithms like K-means Clustering, Apriori Algorithm, Support Vector Machine Algorithm, k-NN Algorithm and Page Rank algorithms are taken into consideration which makes the big data mining operation more efficient. The different results are obtained from the algorithms with respect to the substantial and structural mining. The Statistical Result survey is produced and concluded in this paper.

*Index Terms: Substantial Mining, Structure Mining, Mining Operation Results, Performance Analysis.*

## I.INTRODUCTION

**B**ig data alludes to a procedure that is utilized when customary information mining and taking care of systems can't reveal the experiences and importance of the hidden information. Information that is unstructured or time delicate or just expansive can't be prepared by social database motors. This sort of information requires an alternate preparing approach called huge information, which utilizes gigantic parallelism on promptly accessible equipment.

Simply, huge information mirrors the changing scene we live in. The more things change, the more the progressions are caught and recorded as information. Take climate for instance. For a climate forecaster, the measure of information gathered far and wide about nearby conditions is significant. Consistently, it would bode well that neighborhood conditions direct territorial impacts and provincial impacts manage worldwide impacts, yet it could well be the a different way. Somehow, this climate information mirrors the traits of enormous information, where continuous handling is required for a huge measure of information, and where the huge number of sources of info can be machine created, individual perceptions or outside powers like sun spots.

Preparing data like this outlines why huge information has turned out to be so critical: Most information gathered now is unstructured and requires distinctive capacity and preparing than that found in customary social databases.

Accessible computational power is soaring, which means there are more chances to process enormous information.

The Internet has democratized information, relentlessly expanding the information accessible while additionally creating increasingly crude information. Information in its crude frame has no esteem. Information should be handled so as to be of important. Be that as it may, in this lies the innate issue of enormous information. Is preparing information from local question configuration to a usable knowledge worth the monstrous capital cost of doing as such? Or on the other hand is there just an

excess of information with obscure esteems to legitimize the bet of preparing it with enormous information devices? The vast majority of us would concur that having the capacity to anticipate the climate would have esteem, the inquiry is whether that esteem could exceed the expenses of crunching all the constant information into a climate report that could be depended on.

This paper surveys the existing mining operation methodologies with the materialized search with the traditional and structural approach. Through the statistical survey we can conclude the suggestable situations of the algorithms which will be the best suite for the mining operations in the Big Data Environment.

## II.RELATED STUDY

**K-MEANS** is one of the least complex unsupervised learning calculations that take care of the outstanding bunching issue. The strategy takes after a straightforward and simple approach to group a given informational index through a specific number of bunches (accept k groups) settled apriori. The primary thought is to characterize k focuses, one for each bunch. These focuses ought to be put cleverly due to various area causes distinctive outcome. Along these lines, the better decision is to put them however much as could reasonably be expected far from each other. The following stage is to take each guide having a place toward a given informational collection and partner it to the closest focus. At the point when no point is pending, the initial step is finished and an early gathering age is finished. Now we have to re-figure k new centroids as barycenter of the groups coming about because of the past advance. After we have these k new centroids, another coupling must be done between similar informational collection focuses and the closest new focus. A circle has been produced. Because of this circle we may see that the k focuses change their area well ordered until the point that no more changes are done or at the end of the day focuses don't move any more.

**APRIORI** is a calculation that is utilized for visit itemset mining and affiliation manage learning general value-based databases. The calculation is continued by the distinguishing proof of the individual things that are visit in the database and after that extending them to bigger itemsets as long as adequately those thing sets show up frequently enough in the database. These regular itemsets that are dictated by Apriori can be utilized for the assurance of affiliation rules which at that point feature general patterns.

**SUPPORT VECTOR MACHINE** With regards to machine learning, bolster vector machines that are otherwise called help vector systems are fundamentally administered learning models that accompany related learning calculations which at that point examine information that are utilized for the investigation of relapse and grouping. A SVM display is made that is a portrayal of the cases as focuses in space, that are additionally mapped with the goal that the cases of the different classes are then isolated by a reasonable hole that is should be as wide as could be expected under the circumstances.

**K-NEAREST NEIGHBORS** (k-NN) is a sort of apathetic learning or case based learning and is considered as a non-parametric strategy that is utilized for grouping and regression.\In both the specified cases, the information comprises of the k nearest preparing cases in the component space and the yield relies upon whether the calculation is being utilized for characterization or relapse. This kNN Algorithm is considered and is likewise among the least difficult of all machine learning calculations.

**PAGERANK** (PR) that was named after Larry Page who is one of the authors of Google is a calculation that is utilized by Google Search to rank the sites in their web index comes about. PageRank, that is the primary calculation that was utilized by the organization isn't the main calculation that is being utilized by Google to arrange web search tool comes about, yet it is the best-known method for estimating the significance of site pages.

## III. STATISTICAL SURVEY ON DIFFERENT MINING ALGORITHMS ON BIG DATA

In this section, we discuss various mining algorithms with respect to the big data environment. Each Algorithms are selected based on the importance of the systematic mining probabilities.

### K-MEANS ALGORITHM ON BIG DATA

The k-means calculation is outstanding for its productivity in grouping huge informational collections. Be that as it may, working just on numeric esteems forbids it from being utilized to group true information containing downright esteems. In this paper we display two calculations which stretch out the k-means calculation to all out spaces and areas with blended numeric and downright esteems. The k-modes calculation

utilizes a straightforward coordinating disparity measure to manage downright questions, replaces the methods for bunches with modes, and utilizations a recurrence based technique to refresh modes in the grouping procedure to limit the bunching cost work.

With these expansions the k-modes calculation empowers the grouping of straight out information in a manner like k-means. The k-models calculation, through the meaning of a consolidated disparity measure, additionally coordinates the k-means and k-modes calculations to take into consideration grouping objects depicted by blended numeric and straight out characteristics. We utilize the outstanding soybean sickness and credit endorsement informational indexes to exhibit the bunching execution of the two calculations.

$$J(V) = \sum_{i=1}^{c} \sum_{j=1}^{c_i} \left( \left\| x_i - v_j \right\| \right)^2$$

Equation 1: K-Means Clustering Algorithm Strategy for Mining

The most appealing property of the k-implies calculation in information mining is its effectiveness in bunching extensive informational indexes. Nonetheless, that it just takes a shot at numeric information restricts its utilization in numerous information mining applications in light of the inclusion of straight out information. The k-modes and k-models calculations have expelled this constraint and expanded the k-implies worldview into a bland information dividing activity for information mining.

The grouping execution of the two calculations has been assessed utilizing two genuine world informational indexes. The acceptable outcomes have exhibited the viability of the two calculations in finding structures in information. The adaptability tests have demonstrated that the two calculations are effective when bunching substantial complex informational collections as far as both the number of records and the quantity of bunches. These properties are vital to information mining. By and large, the k-modes calculation is quicker than the k-means and k-models calculation since it needs less cycles to focalize.
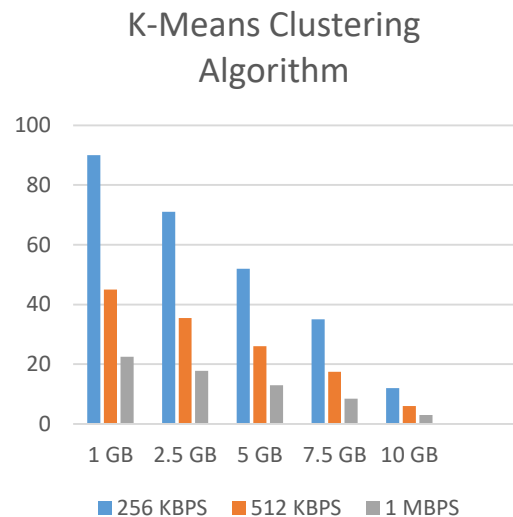
This paper has concentrated on the specialized issues of stretching out the k-implies calculation to group information with downright esteems. In spite of the fact that we have exhibited that the two new calculations function admirably on two known informational collections, we need to recognize this chiefly come about

because of us from the earlier learning to the informational indexes.

| Size (GB) | Time % based mining (256kbps) | Time % based mining (512kbps) | Time % based mining (1mbps) |
|---|---|---|---|
| 1 | 90 | 45 | 22.5 |
| 2.5 | 71 | 35.5 | 17.75 |
| 5 | 52 | 26 | 13 |
| 7.5 | 35 | 17.5 | 8.5 |
| 10 | 12 | 6 | 3 |

Table 1: The K-Means Time Sequence Representation of the mining operation with respect to the size in Repositories

The Underlying Graph representation which describes the graphical stands of the mining operation with respect to its input.



Graph 1: Graphical Representation of the K-means clustering algorithm with statistical data

## APRIORI ALGORITHM ON BIG DATA

Apriori is a calculation that is utilized for visit item set mining and affiliation manage learning general value-based databases. The calculation is continued by the recognizable proof of the individual things that are visit in the database and afterward extending them to bigger item sets as long as adequately those thing sets show up regularly enough in the database. These continuous item sets that are dictated by Apriori can be utilized for the

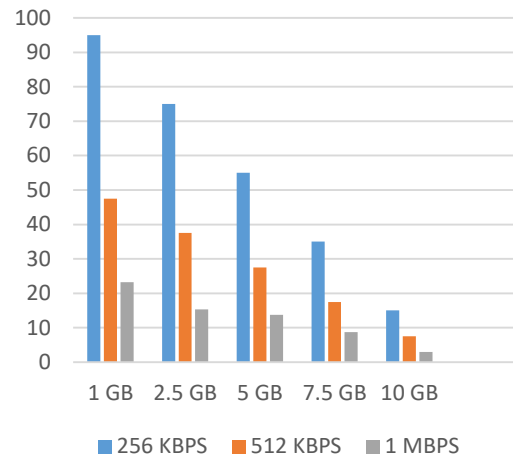assurance of affiliation rules which at that point feature general patterns.

Customary affiliation administer mining calculations just create countless incessant principles, however these standards don't give valuable responses to what the high utility guidelines are. We build up an original thought of best K objective-coordinated information mining, which centers around mining the best K high utility shut examples that specifically bolster a given business objective. To affiliation mining, we add the idea of utility to catch very alluring factual examples and present a level-wise thing set mining calculation. With both positive and negative utilities, the ant monotone pruning technique in Apriori calculation never again holds. Accordingly, we build up another pruning methodology in view of utilities that permit pruning of low utility itemsets to be finished by methods for a weaker yet ant monotonic condition. Our exploratory outcomes demonstrate that our calculation does not require a client determined least utility and thus is powerful by and by.

| Size (GB) | Time % based mining (256kbps) | Time % based mining (512kbps) | Time % based mining (1mbps) |
|---|---|---|---|
| 1 | 95 | 47.5 | 23.25 |
| 2.5 | 75 | 37.5 | 15.35 |
| 5 | 55 | 27.5 | 13.75 |
| 7.5 | 35 | 17.5 | 8.75 |
| 10 | 15 | 7.5 | 3.75 |

Table 2: The Apriori Time Sequence Representation of the mining operation with respect to the size in Repositories

The Underlying Graph representation which describes the graphical stands of the mining operation with respect to its input.

## Apriori Clustering Algorithm



Graph 2: Graphical Representation of the Apriori clustering algorithm with statistical data

## SUPPORT VECTOR MACHINE FOR BIG DATA

When it comes to machine learning, support vector machines that are also known as support vector networks are basically supervised learning models that come with associated learning algorithms which then analyze data that are used for the analysis of regression and classification.

An SVM model is created that is a representation of the examples as points in space, that are further mapped so that the examples of the separate categories are then divided by a clear gap that is ought to be as wide as possible.
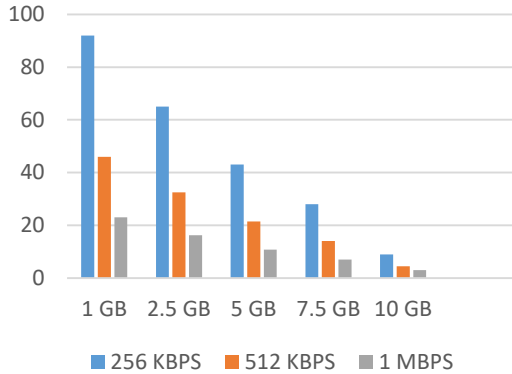
This is a supervised learning, since a dataset is used to first teach the SVM about the classes. Only then is the SVM capable of classifying new data

| Size (GB) | Time % based mining (256kbps) | Time % based mining (512kbps) | Time % based mining (1mbps) |
|---|---|---|---|
| 1 | 92 | 46 | 23 |
| 2.5 | 65 | 32.5 | 16.25 |
| 5 | 43 | 21.5 | 10.75 |
| 7.5 | 28 | 14 | 7 |
| 10 | 9 | 4.5 | 2.25 |

Table 3: The Support Vector Machine Time Sequence Representation of the mining operation with respect to the size in Repositories

The Underlying Graph representation which describes the graphical stands of the mining operation with respect to its input.



Graph 3: Graphical Representation of the SVM clustering algorithm with statistical data

Support vector machine (SVM) takes in a hyperplane to group information into 2 classes. At an abnormal state, SVM plays out a comparative errand like C4.5 aside from SVM doesn't utilize choice trees by any means.

## K-NEAREST NEIGHBOR FOR BIG DATA

The k-nearest neighbors calculation (k-NN) is a sort of apathetic learning or case based learning and is considered as a non-parametric strategy that is utilized for characterization and regression.\In both the specified cases, the info comprises of the k nearest preparing cases in the element space and the yield relies upon whether the calculation is being utilized for arrangement or relapse. This kNN Algorithm is considered and is additionally among the most straightforward of all machine learning calculations.

1. kNN can get very computationally expensive when trying to determine the nearest neighbors on a large dataset.
2. Noisy data can throw off kNN classifications.
3. Features with a larger range of values can dominate the distance metric relative to features that have a smaller range, so feature scaling is important.
4. Since data processing is deferred, kNN generally requires greater storage requirements than eager classifiers.
5. Selecting a good distance metric is crucial to kNN's accuracy.

### K-NN Identifier and Classifier Algorithm

1. Calculate "d(x, $x_i$)" i =1, 2, ….., **n**; where **d** denotes the Euclidean distance between the points.
2. Arrange the calculated **n** Euclidean distances in non-decreasing order.
3. Let **k** be a +ve integer, take the first **k** distances from this sorted list.
4. Find those **k**-points corresponding to these **k**-distances.
5. Let $k_i$ denotes the number of points belonging to the $i^{th}$ class among **k** points i.e. k ≥ 0
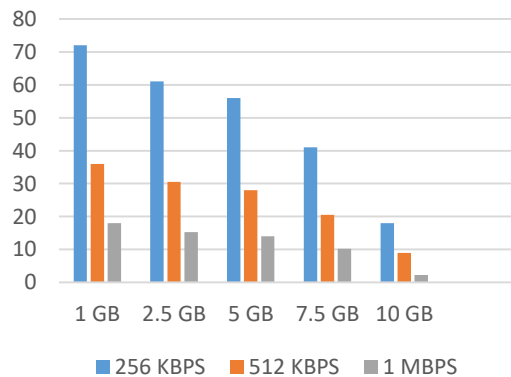6. If $k_i > k_j$ ∀ i ≠ j then put x in class i.

Closest neighbor is an extraordinary instance of k-closest neighbor class. Where k esteem is 1 (k = 1). For this situation, new information point target class will be allocated to the first nearest neighbor.

| Size (GB) | Time % based mining (256kbps) | Time % based mining (512kbps) | Time % based mining (1mbps) |
|---|---|---|---|
| 1 | 72 | 36 | 18 |
| 2.5 | 61 | 30.5 | 15.25 |
| 5 | 56 | 28 | 14 |
| 7.5 | 41 | 20.5 | 10.25 |
| 10 | 18 | 9 | 4.5 |

Table 4: The K-NN Time Sequence Representation of the mining operation with respect to the size in Repositories

The Underlying Graph representation which describes the graphical stands of the mining operation with respect to its input.



Graph 4: Graphical Representation of the K-NN clustering algorithm with statistical data

## PAGE RANK ALGORITHM FOR BIG DATA

PageRank (PR) that was named after Larry Page who is one of the originators of Google is a calculation that is utilized by Google Search to rank the sites in their web index comes about. PageRank, that is the principal calculation that was utilized by the organization isn't the main calculation that is being utilized by Google to arrange internet searcher comes about, however it is the best-known method for estimating the significance of site pages.

PageRank is a link analysis algorithm designed to determine the relative importance of some object linked within a network of objects. In spite of the fact that the exact importance of a PageRank number isn't uncovered by Google, we can get a feeling of its relative significance.

| Website | PageRank |
|---|---|
| twitter.com | 10 |
| facebook.com | 9 |
| reddit.com | 8 |
| stackoverflow.com | 7 |
| tumblr.com | 6 |
| crucial.com | 5 |
| programmingzen.com | 4 |
| dearblogger.org | 3 |

Fig 1: Statistical Page Ranking based on Click Points

 PageRank is its robustness due to the difficulty of getting a relevant incoming link and have a graph or network and want to understand relative importance, priority, ranking or relevance

### Page Ranking Strategies

1. Ranks Page based on the number of other pages that link to it
2. Gives an indication of the relative importance of a page
3. An Appropriate SERP Listing
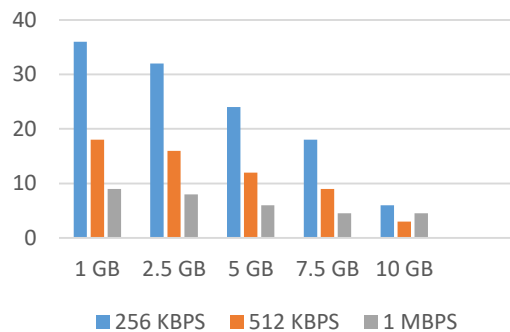4. Calculated by nature and number of backlinks

| Size (GB) | Time % based mining (256kbps) | Time % based mining (512kbps) | Time % based mining (1mbps) |
|---|---|---|---|
| 1 | 36 | 18 | 9 |
| 2.5 | 32 | 16 | 8 |
| 5 | 24 | 12 | 6 |
| 7.5 | 18 | 9 | 4.5 |
| 10 | 6 | 3 | 1.5 |

Table 4: Page Rank Time Sequence Representation of the mining operation with respect to the size in Repositories

The Underlying Graph representation which describes the graphical stands of the mining operation with respect to its input.
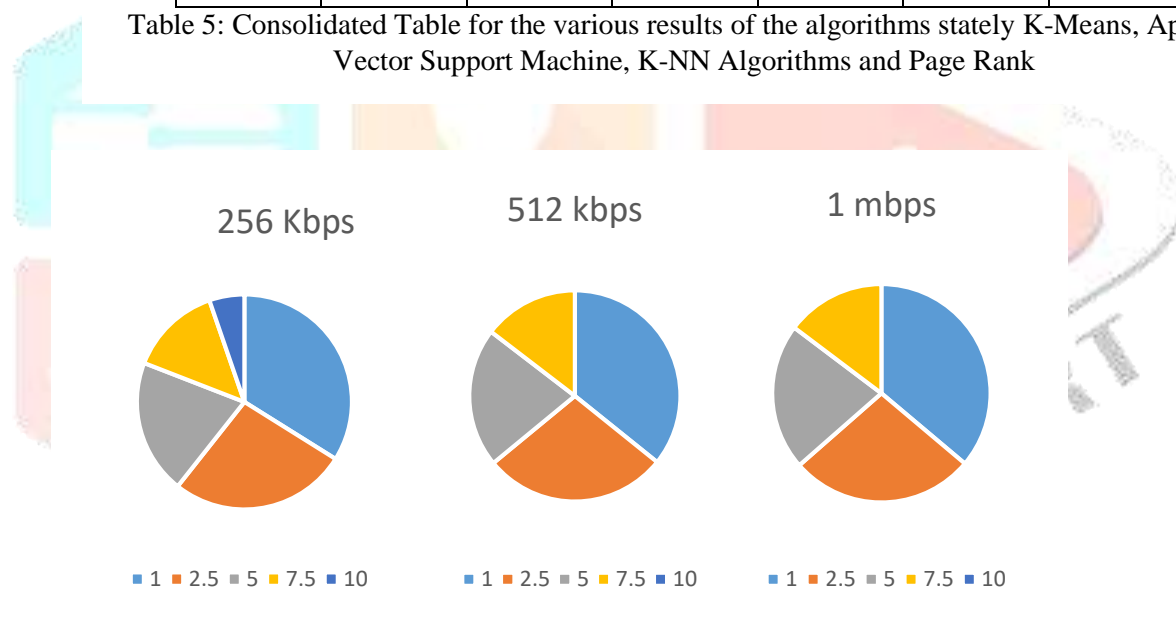


Pagerank Clustering Algorithm

These algorithms are taken into consideration with the

| Size (GB) | K-MEANS | | | APRIORI | | | SUPPORT VECTOR MACHINE | | |
|---|---|---|---|---|---|---|---|---|---|
| | 256 kbps | 512 kbps | 1 mbps | 256 kbps | 512 kbps | 1 mbps | 256 kbps | 512 kbps | 1 mbps |
| 1 | 90 | 45 | 22.5 | 95 | 47.5 | 23.25 | 92 | 46 | 23 |
| 2.5 | 71 | 35.5 | 17.75 | 75 | 37.5 | 15.35 | 65 | 32.5 | 16.25 |
| 5 | 52 | 26 | 13 | 55 | 27.5 | 13.75 | 43 | 21.5 | 10.75 |
| 7.5 | 35 | 17.5 | 8.5 | 35 | 17.5 | 8.75 | 28 | 14 | 7 |
| 10 | 12 | 6 | 3 | 15 | 7.5 | 3.75 | 9 | 4.5 | 2.25 |

| Size (GB) | K-NN ALGORITHM | | | PAGE RANK | | |
|---|---|---|---|---|---|---|
| | 256 kbps | 512 kbps | 1 mbps | 256 kbps | 512 kbps | 1 mbps |
| 1 | 72 | 36 | 18 | 36 | 18 | 9 |
| 2.5 | 61 | 30.5 | 15.25 | 32 | 16 | 8 |
| 5 | 56 | 28 | 14 | 24 | 12 | 6 |
| 7.5 | 41 | 20.5 | 10.25 | 18 | 9 | 4.5 |
| 10 | 18 | 9 | 4.5 | 6 | 3 | 1.5 |

Table 5: Consolidated Table for the various results of the algorithms stately K-Means, Apriori, Vector Support Machine, K-NN Algorithms and Page Rank



256 Kbps   512 kbps   1 mbps

■ 1 ■ 2.5 ■ 5 ■ 7.5 ■ 10    ■ 1 ■ 2.5 ■ 5 ■ 7.5 ■ 10    ■ 1 ■ 2.5 ■ 5 ■ 7.5 ■ 10

Web pages on the World Wide Web link to each other. If rayli.net links to a web page on CNN, a vote is added for the CNN page indicating rayli.net finds the CNN web page relevant.

The table is created on the basis of the various results obtained from the different data mining algorithms

1. K-means Algorithm
2. Apriori Algorithm
3. Support Vector Machine Algorithm
4. K-NN Algorithm
5. Page Rank Algorithm

## IV. CONCLUSION

This paper concludes with various strategic mining algorithms with their statistical results in the point of the performance of the mining operations. Enormous Data Mining is one of the developing innovation which manages the capacity of tremendous information's as the stores. The Big information is utilized as a part of the different fields for the scattered use. The acceptance of different information mining calculations in the field has made the mining things more proficient in the stores. Measurably reviewing all the conceivable calculations which creates the mining activity more doable is

investigated in this paper and finished up with the best discretionary calculation which makes the mining task more effective. The calculations like K-implies Clustering, Apriori Algorithm, Support Vector Machine Algorithm, k-NN Algorithm and Page Rank calculations are mulled over which makes the huge information mining activity more proficient. The diverse outcomes are acquired from the calculations concerning the generous and basic mining. The Statistical Result review is delivered and deduced in this paper

## V.REFERENCES

[1] Wei Fan, Albert Bifet. Mining big data: current status, and forecast to the future, ACM SIGKDD Explorations Newsletter, Volume 14 Issue 2, December 2012

[2] Sneha Gupta, Manoj S. Chaudhari, Big Data Issues and Challenges, nternational Journal on Recent and Innovation Trends in Computing and Communication ISSN: 2321-8169, Volume: 3, Issue: 2

[3] Stephen Kaisler ; Frank Armour ; J. Alberto Espinosa ;William Money. Big Data: Issues and Challenges Moving Forward, 46th Hawaii International Conference on System Sciences (HICSS), 2013, ISSN :1530-1605

[4] Big Data computing and clouds: Trends and future direction by Rajkumar Buyya

[5] Marcos D. Assunção et.al. Big Data computing and clouds: Trends and and Distributed Compu

[6] Xindong Wu, , Ding, Data Mining Knowledge And D January 2014

[7] Deepak S. Tan Analysis Using Had Advanced Research Technology (IJARC 2015

[8] Kale Suvarna Journal of Compute eETECME October 20

[9] Privacy Preservin Partitioned Data by Jai SIGKDD 02, Edmonto

[10] Privacy-Preservin Vertically Partitioned I Kwang Raymond Choc by IEEE Transaction Volume 11;Issue 8, Au

[11] Distributed Decision Tree Algorithm and Its Implementation on Big Data Platforms by JingXiang Chen, Taowang, Ralph Abbey, Joseph Pimgenot published in IEEE International Conference on Data Science and Advanced Analysis 2016

[12] Optimal Online Deterministic Algorithms and Adaptive Heuristics for Energy and Performance Efficient Dynamic Consolidation of Virtual Machines in Cloud Data Centers by Anton Beloglazov and Raj Kumar Buyya published in Concurrency and Computation: Practice and Experience , 2012

[13] Content-Based Scheduling of Virtual Machines (VMs) in the Cloud by Sobir Bazarbayev,William H Sanders, Matti Hiltunen and Kaustubh Joshi Published in Distributed Computing Systems (ICDCS), 2013 IEEE 33rd International Conference .

[14] Dynamic Resource Management in Cloud Data Centers for Server Consolidation by Alexander Ngenzi , Dr.Selvarani R and Dr.Suchithra R Nair published in www.arXiv.org/ArXiv:1505.00577

[15] Resourc-Aware Scheduling for Data Centers with Heterogenous Servers by Tony T Tran, Peter Yun Zhang, Heyse Li, Douglas G Down, J. Christopher Beck Published in Google and Natural Sciences and Engineering Research Council of Canada (NSERC), http://hdl.handle.net/1807/69000

[16] Greedy Scheduling of Tasks with Time Constraints for Energy-efficient cloud computing data centers published at Journal of Cloud Computing:Advanced, Systems and Applications 2015 4:5 DOI 10.1189/s13677-015-0031-y

**Mrs.R.Karunia krishnapriya** received the Engineering in 2000 and the Master of comp from Vellore Institute of Technology, Tan Ph.D. degree in Computer science and Eng Tamilnadu,India.She has published more th conference and many workshops from many and also conducted many workshops and co for their research works for her teaching fiel mining, text mining, web and cloud minin society of India and Indian society for techn

**Dr. G.Vinodhini** is an Assistant Professor a and Engineering in Annamalai University, She received the B.Tech Degree in Inf University in 2003 and M.E Degree in Comp from Annamalai University, Tamil Nadu, a in Computer Science and Engineering fron research interests include Text mining, Num Business Intelligence. She published mor indexed reputed journals like Elsevier, Spr citation index of more than 370. She is als in area of data mining. She is a Life memb Indian Society for Technical Education.

**Dr.K.Sakthivel** received the B.Sc. degree i of Computer Applications in 1997 Tiruchitrappalli, Tamil Nadu, India,and M Engineering in 2005 from Anna University, degree in Information and Communica Chennai in 2012. Currently working as pro