



Big data analysis with parallel and distributed computing in cloud environment

Mahboob Alam *, Pradeep Kumar Harsha K.G., Vinooth P, Mukesh Raj
Department of Computer Science and Engineering
JSS Academy of Technical Education Noida
** Corresponding Author: mahboobru@gmail.com*

6599

Abstract

Analysis with large statistics or huge records analysis has come to be energetic research vicinity. It is very difficult the use of cutting-edge methodologies and statistics analysis software program equipment for one individual PC to handle extremely huge datasets with efficiency. The cloud based computing and parallel computing structures are taken to be a higher answer to perform massive facts analysis. The parallel computing concept is entirely reliant on dividing a huge trouble in smaller parts, every one of those is completed by using an individual standalone processor in my opinion. Further, those procedures are done simultaneously in an allotted but parallel way. These are some of the conventional strategies to deal with large records problems. The median one is the allotted procedure primarily reliant upon the paradigm of data parallelism, where a large dataset under consideration is broken down into n subsets by employing manual methods, and an equal number of algorithms are run for the respective n subsets. The eventual result may be received from an aggregation of outputs generated by the n algorithms. Another approach is a procedure based on map reduce, which runs beneath the platform of cloud computation. This process is basically the map & reduce methods, where the first one is responsible for filtering and sorting while the second one does an operation of summary to generate the final result. In this paper, we attempt to examine the overall difference of performances between the mapreduce and the dispensed techniques over huge datasets with respect to the analysis efficiency and accuracy. The experimentation is primarily dependent upon the 4 big size datasets that are used for the records classification troubles. The results exhibit that the performance of the mapreduce based technique in terms of classification is rather robust irrespective of the number of pc nodes used, and is preferable over the baseline unmarried machine and distributed approaches besides magnificent imbalanced dataset. Similarly, Mapreduce method calls for the minimal computing value to process those huge datasets.

Keywords: Big Data, Parallel Computing, Distributed Computing, Cloud, MapReduce.

DOI Number: 10.48047/NQ.2022.20.15.NQ88658

NeuroQuantology2022;20(15): 6599-6608

1. Introduction

Attributable to the recognition and advancement of associated net and statistics era, large quantities of information are produced in our everyday life. Massive flows of statistics, exabytes of records, are captured daily. Evidently, the technology of massive information is here already. Further related to the statistics length big information has different properties together with range & speed. The previous approach was that records can also be generated from a huge style of unstructured, partially-established, and records on the contrary

the last one points to the need of a real-time system and evaluation. Consequently, massive facts analytics with the aid of gadget getting to know and statistics analysis techniques has come to be a significant research trouble. Analysis involving huge information and massive data analysis could be extremely tough for assessing with the modern methods & statistics analysis SW program equipment because of the big length and convolutions. alternatively, the usage of a unmarried personal computer (pc) to ex-adorable the statistics analysis project over massive order data files demand extremely high computation expenses. Hence it becomes very



important to introduce efficiency in the environment of computing for handling big data chunks. Consistent with other research work, the answers in a response to the challenge of analysis huge order data files could be primarily built upon the concept of cloud and parallel computation set ups. Theoretically, parallelism lays emphasis on segmenting the selected (massive) hassle into simpler ones, all of these are then finished through a single unmarried processor individually, enabling a calculation made up of various computations to be executed simultaneously in a parallel & distributed guynner. However this introduces some research problems for dispensed records analysis & distributed machining setup getting to know. Being specific, statistically, the paradigm of data parallelism, known as the distributed technique in this document, are taken into consideration for handling huge order data. In information parallelism, the big order datasets are divided between a number of processor cores, each of them then does an equal amount of processing and calculation over a pre-determined delegated segment.

The concept of distributed technique is being put to use in the ensemble classier. This means, every single classifier is trained with a part of a particular set earmarked for training. Subsequently, to classify any newly introduced variable, we take a look at case, the check case is fed to each one of the pre-trained classifier, this in turn lets those classifiers operate with parallelism while incorporating distributed computation. In the end, the classification consequences generated by means of those classifiers are merged through a few aggregate strategies, together with voting & also casting of weighed vote to arrive at the ultimate result. Lately, cloud computing has enhanced the principle of parallelism to manage efficiently, the utility & intake of computational assets accross a pc cluster. In order to deal with the massive order datafiles issues, the mapreduce computation is commonly implemented the usage of hadoop1, a robust parallel programing enviornment. The mapreduce method combines the map & reduce procedures, where the previous in line with- bureaucracy sorting&

filtering, whereas the other one is a précis operation in an effort to generate the output.

In line with the previous dialogue, massive data analysis may be efficiently executed through the frequently used mapreduce & distributed techniques. Both techniques need a quite a few processors for performing a few analysis tasks parallel. However 1 difference among the 2 techniques is the resource management of the computational peripherals. In the traditional distributed method, you may partition a massive order dataset into m subsets for the m nodes to carry out the analysis assignment. However, mapreduce method mechanically does manage the intake of computing assets of range laptop nodes whilst dealing with huge order data file. This eliminates the need to segment it further into m subsets. Only the putting of the wide variety of pc nodes to manage the dataset is needed. However, the consumption of resources among nodes in not necessarily equitable anymore. It however raises a significant research query, that hasn't been raised earlier than this: Is there any difference in the way the methodologies of mapreduce and distributed operate on big scale datasets with respect to analysis of efficiency & accuracy?

Hence the objective of this investigation is to contribute in having a look at the analysis of the performance of these two methodologies on the two parameters concerned with the classification processes, namely the accuracy and the processing instances. Similarly, since using greater nodes no longer assures accuracy or efficiency, as a consequence of the associated overheads, one of a kind numbers of computer nodes might be used for comparing. This paper relaxation is prepared as described below. The second section outlines the associated literature for disbursed & map reduces- based totally facts analysis. The third & the fourth sections gift respectively the experimentation procedures & results. At last, a few conclusions are proposed in the fifth section.



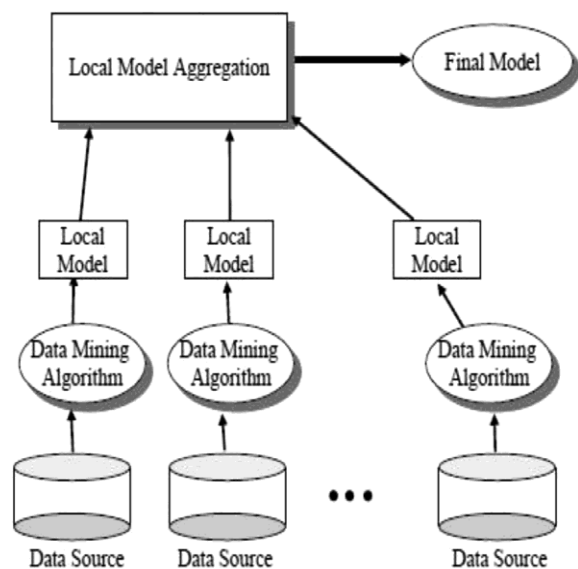


Fig.1. Distributed data analysis framework

2. Literature Review

Assigned registering can allude to utilizing dispersed frameworks to clear up computational difficulties. Particularly, a challenge is split in various duties, all of them are then resolved using 1 or greater computer systems (or processors) running in concurrency. Besides all the processors are capable of communicating with other ones by passing messages.

Within the conventional statistics analysis technique, the information is usually centralized, thereafter a particular algorithm is selected to systemizing & analyzing the records underneath a single computing platform. but, in case of massive facts problem or huge order statistics analysis, this isn't always this simplified, & the appearing of the records analysis obligations underneath the disbursed computing platform has come to be an critical region of research investigation. However, the aim of disbursed facts analysis is to carry out the statistics analysis responsibilities based totally at the distributed resources, inclusive of the records, systems, & data analysis algorithms. Figure 1 suggests a standard allotted records analysis framework wherein different statistics assets can be homogenous and/or heterogeneous. Each information analysis has a set of rules for

handling the respective data source underneath a unmarried computation setup main to a neighborhood model. afterwards, such nearby models are congregated on the way for producing the final version.

For solving the big statistics troubles, the paradigm of information parallelism is also an option. Taken a huge scale data file k , this can be subdivided into smaller m subsets, denoted as $d_1, d_2, d_3, \dots, d_n$, wherein every subset can also Includes one-of-a-kind count of facts samples and every subset might can possibly have mirrored(duplicated) sample of facts. After this, a chosen information analysis algorithm is applied in all of these neighborhood machines, which in my opinion is carried out over every subset. Finally with the help of combination component, these n analysis effects are combined to generate the final OP. This distributed method differs from the conventional one in the sense that, the latter plays the facts analysis algo over d directly on a single device. Hence obviously, whilst the size of the data file increases to become excessively large, time required for processing by the conventional technique significantly will increase, however the dispensed method is able to resolve this large scale dataset trouble efficiently.

Practically, majority of users, commonly, at the best have constrained computational assets for performing huge data analysis. Advantage that this approach offers is that just one unmarried computer is needed for addressing every subset at a time, and that device plays the identical undertaking for one of kind subsets n instances, ensuing in n exceptional fashions.



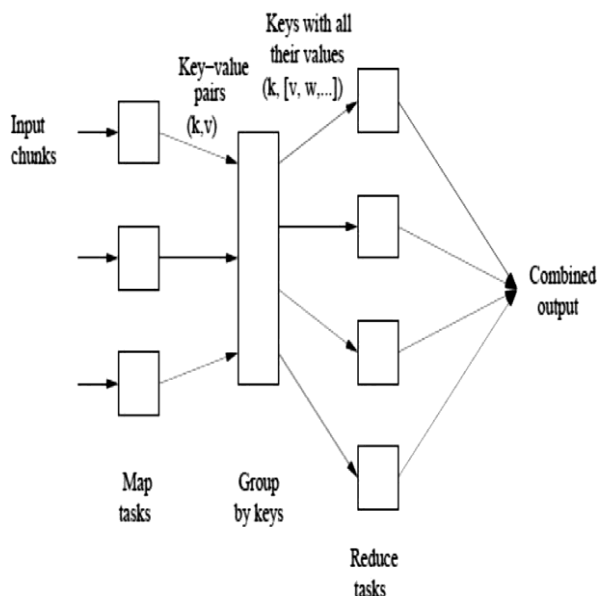


Fig.2. MapReduce computation procedure

The incredible effects obtained by using these distinct fashions are subsequently merged for arriving at the final value of result.

Although, one issue with the dispensed records analysis approach is that it requires manual segmentation or partitioning of the big records for the selected statistics. Thereafter the allocation of the computing nodes has to be done manually too for performing the analysis operation over those repartitioned subsets.

3. Big Data Analysis Procedure

The general execution got utilizing the disbursed & mapreduce techniques over huge scale data files with respect to the analysis efficiency & accuracy is tested by using evaluating 3 large records analysis strategies, particularly the centralized baseline allotted & the mapreduce processes.

3.1 The analysis procedure of the baseline (BL) big data

This method of baseline is done for large statistics analysis on one specific unmarried system & the records are centralized for analysis purposes. Figure 3 suggests an instance of the usage of the guide (Support

vector machine) classification approach for a dataset. To start with, the tenfold pass validation method is employed to break up the authentic dataset in 2 portions, nineteenth for education set & one-tenth for testing set. Subsequently, the training set is put to use for assembling the Support vector machine classifier. In the end, the trying out set is input into the Support vector machine for the end result of classification.

3.2 The distributed and parallel big data analysis procedure

This is slightly different from the procedure of Baseline; it works on the old divide & conquers technique to arrive at results. Here n subsets are created by dividing the parent dataset for n computing nodes & the support vector machine algorithm is used on each of these nodes. We have contrasted a varying number of nodes to determine the implications the node count has on the efficiency and the accuracy of the analysis has been chosen from a set of 5 values, each a multiple of 10, ranging from 10-50.

Figure 4 depicts an instance of analysis procedure over distributed big data involving five nodes. Firstly, we break down the given dataset into two segments that are earmarked for training and testing, each with a share of 90 and 10 percent respectively. This is achieved using the tenfold cross validation strategy. Thereafter we break it further down to 5 unduplicated subsets, each of which is put to use for training the Support Vector Machine classifier. This leads to the construction of 5 which results in five individual Support Vector Machine. The next step is to feed the set earmarked for testing into the 5 Support Vector Machine classifiers at once. We the proceed to use the majority voting combination method for combining the separate outputs generated by our 5 Support Vector Machines, to generate the ultimate classification result. We measure & assess the training & classification times along



with the accuracy of classification. One point worth noting here is that the exact same platform for computation is used for the distributed procedure as the one for the baseline.

3.3 The MapReduce(MR) based big data analysis procedure

Figure five depicts an instance of the usage of 5 laptop nodes in keeping with- shape massive data analysis based totally at the framework of mapreduce. Observe here, the surroundings of cloud computing, aiding this procedure are simulated by a laptop

Thereafter the individual generation of results and their eventual merging of them is done as already discussed.

server. The selected huge scale dataset is placed within the Hadoop distributed file system based totally on a master & m virtual machines allocated for processing of data and evaluation venture, where m is a multiple of 10, ranging from 10 to 50 for evaluation. Listed below are the SW & the HW specifications.

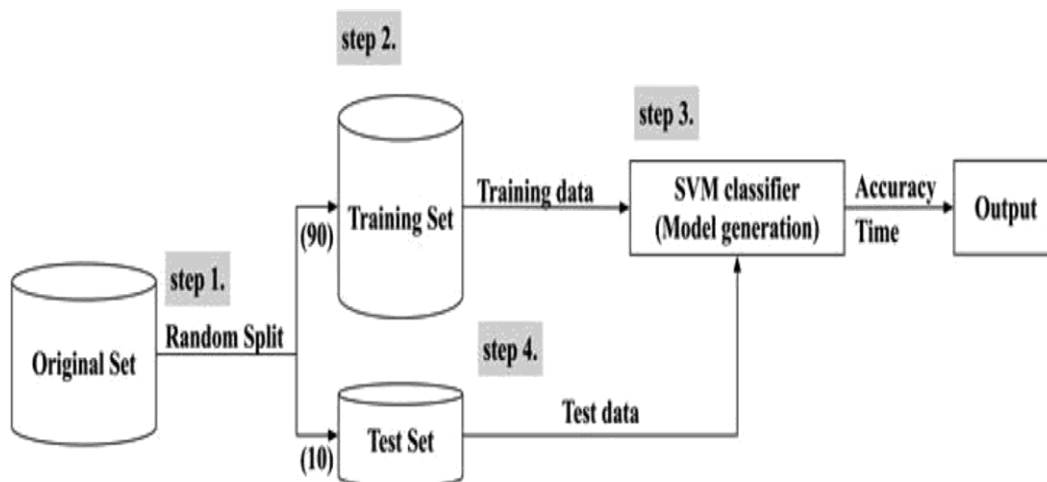


Fig. 3. BL Big Data analysis procedure

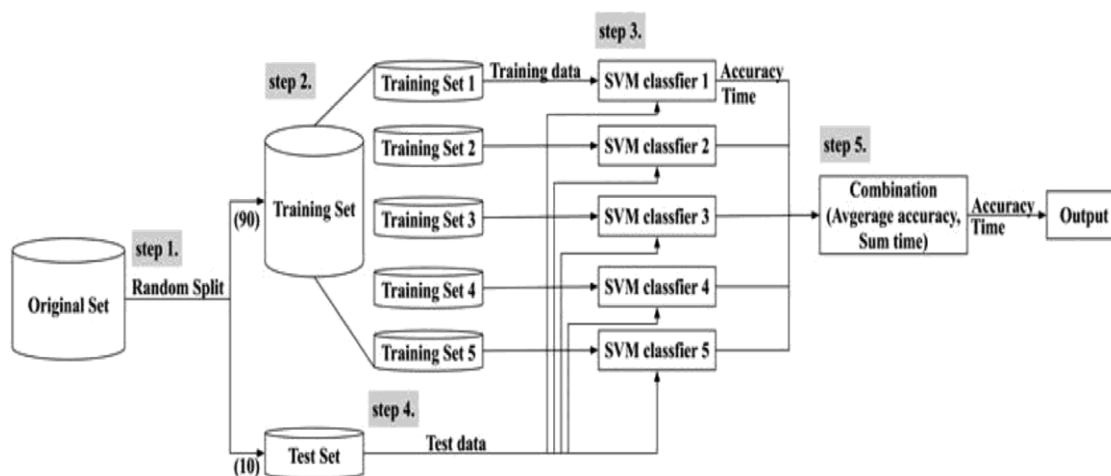


Fig.4. Distributed and parallel big data analysis procedure.



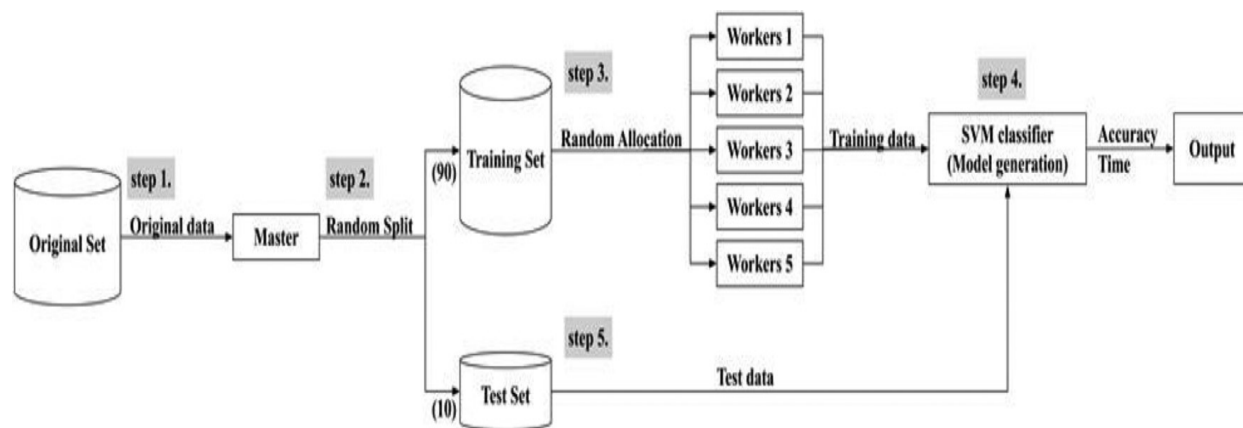


Fig.5. MR based big data analysis procedure.

Manufacturer	Dell Inc.
Model	PowerEdge T610
Processor	Intel(R) Xeon(R) CPU E5620 @ 2.40 GHz
Host CPU	8 CPUs
Memory	196.0 GB
Host operating system	VMkernel
Virtualisation software	VMware vSphere Hypervisor (ESXi)
Number of virtual machines	1 to 51 (MASTER-1, WORKERS-all of the rest)
Guest operating system	CentOS 6.5
Guest CPU	Single-Core
Guest memory	3.0–4.0 GB
MapReduce environment	Hadoop 2.2.0
Data analysis environment	Spark 0.8.1

Table 1. Software & Hardware specification

Unlike the dispersed large facts analysis method, each of the five workers might also address one-of-a-kind segment of the nine-tenth schooling set, that's managed through the grasp mechanically. Consequently, the 5 worker possess one-of-a-kind computational complexities at some point of the analysis project. In different phrases, the allotted based technique focuses on segmenting the dataset as such, while the mapreduce primarily based technique is handy for handling the number of persons.

4. Result and Analysis

To compare the overall performance of the three special huge facts analysis processes, 4 massive scale datasets that cowl exceptional area troubles are used. They are the kdd cup2 2004 i.e prediction of protein homology & 2008 (breast most cancers), cover- type3 and character activity4 datasets. The second table features the fundamental infor- mation related with these 4



datasets. The previous 2 datasets belongs to the two-class classification issue.

Besides, every single dataset is split into 90% schooling and 10% trying out sets based totally on the validation strategy of 10 folds for schooling & trying out the single vector machine classifier, respectively. Then various performance across various parameters like accuracy, performance time for activities like training etc are examined for assessment. The same are listed below in a section.

Datasets	Feature count	Sample count	Class count
Breast cancer	117	102,294	2
Protein homolog y	74	145,751	2
Cover type	54	581,012	7
Person activity	8	164,860	11

Table 2. composition of the four datasets.

No.of nodes	Physical surroundings	Virtual surroundings
1	BL procedure	MR
10	Distributed and parallel Procedure	based
20		procedure
30		
40		
50		

Table 3. Surrounding conditions for the three big data analysis procedure

Environmental settings of the 3 procedures are listed in the third table for assessment. For the fundamental process, merely a single laptop is put to use for the massive data analysis challenge. The allotted manner relies only upon segmenting the schooling set into a number of subsets which is a multiple of 10, ranging from 10 to 50. Where the association of computing nodes & subsets in one to one for the classifier education undertaking. On the contrary, the mapre-duce based system educates the classifier

by using a computer sever with a configuration of 1, 10, 20, 30, 40, and 50 digital machines.

4.1 Results from the datasets related to breast cancer

The accuracy of classification of our 3 huge statistics analysis approaches over the breast most cancers dataset are shown in figure 3. Apparently the Support Vector Machine classifier based upon the baseline & allotted procedures can offer enhanced overall performances vis-a-vis mapreduced one. Mainly, the Support Vector Machine obtained the usage stats of the distributed & baseline tactics(with around ten to fifty nodes) each one of them delivered an accuracy figure north of 99.38 percent; on the other hand Support Vector Machine received by way of mapreduced based totally technique most effective produces around fifty eight(58%) precision.

One of the factors that can be attributed to substandard performance from the mapreduced primarily based method might be that there are various class imbalances in the dataset with high dimensionalities for breast cancer, wherein the dataset consists of ninety 94% and 06% information for the benign and malignant instructions, respectively. Considering the framework of mapreduced, the Support Vector Machine classifier doesn't learn efficaciously to differentiate among the 2 instructions.

In comparison, fig. 7 indicates the training costs & the costs associated with running Support Vector Machine in allotted & mapreduced based totally procedures. We haven't compared the baseline procedures in this as it takes 10,225 s to perform the task.

Talking about the figure 7, we will study the fact that the computation prices drops with the increase in the laptop node, primarily based on disbursed procedures, however no noticeable decrease is observed in the time it takes for processing when we use 30-50 nodes. In particular, just around 3 minutes are needed to check & train the Support Vector Machines while using 50 nodes & hence the very best fee of classification accuracy is generated (99.38%).



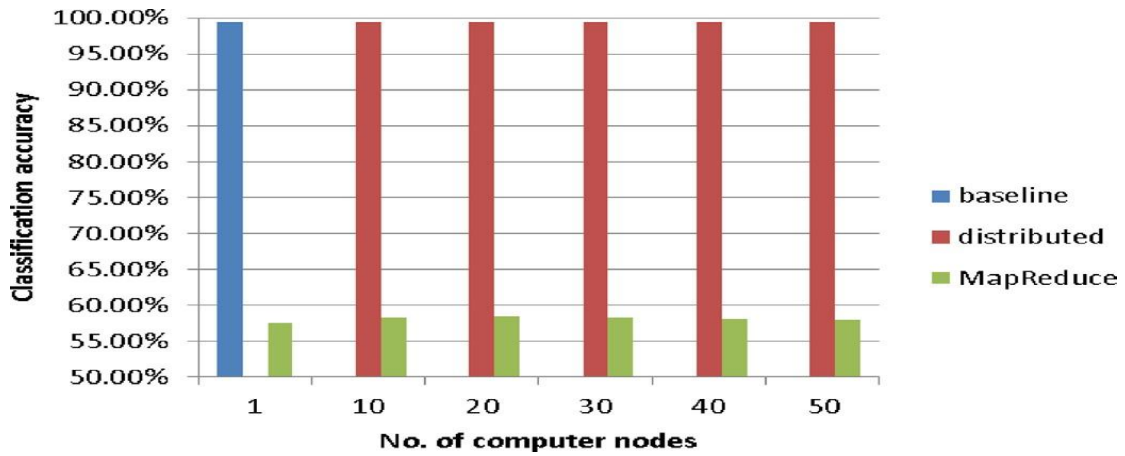


Fig.6. Big data analysis procedures accuracy.

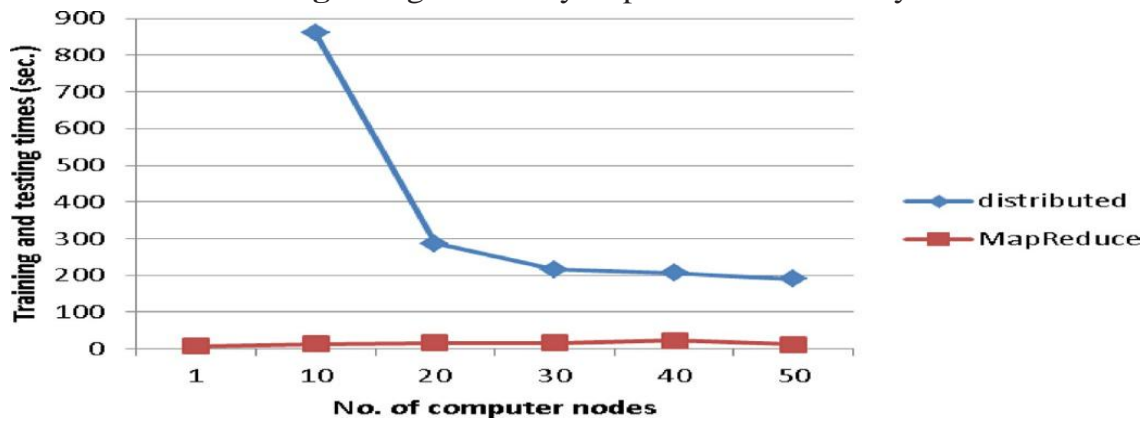


Fig.7. Distributed & MapReduce based procedures training & testing.

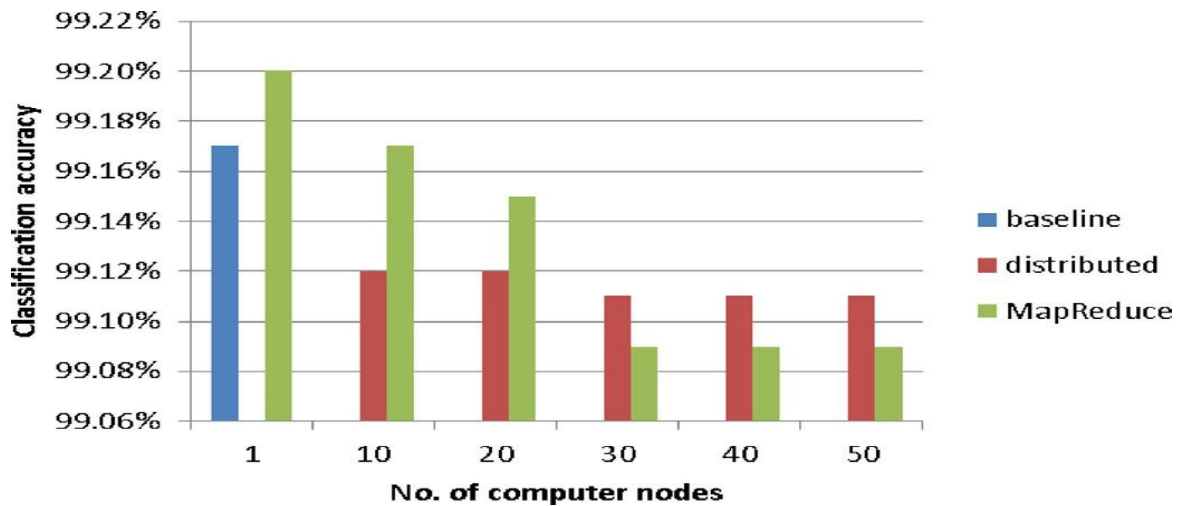


Fig.8. Procedures accuracy over the protein homology dataset.



5. Conclusion

Big information analysis are possibly resolved with efficiency underneath a computational environment which runs on parallelism. In standard, unusual techniques can be used.

Madain one of which works entirely on the disbursed manner, that specializes in the parallel processing of the information prospect to break down a given huge scale items into some of sub items, every item is dealt with by a particular searning method carried out by a common machine.

Final result is received with aid of combining generated outputs by means of gaining knowledge of fashions. The second one is the procedure based upon Map Reduce, in which the range of maps might be consumer-defined, however actually are con- trolled via a master to mechanically control the application & consumption of assets for a pc group. Further, reduce feature combines the map outputs to generate result.

Our aim here is to have a look at the analysis of performing efficiency of the allotted & map reduce reliant strategies over bulk information issues. Our experimentation effects are primarily dependent on 4 massive scale datafiles, exhibit that map reduce primarily based methodology is quite stable in terms of analysis accuracy regardless of the number of laptop nodes used and it is able to permit the svm clas- sifier to attain the maximum fee of classification accuracy aside from elegance imbalance dataset. Similarly, minimum processing time is needed for educating & testing the Single vector machine, however adding wide variety of laptop nodes might increase the processing times by a little. Hence it is established that using 1-10 laptop nodes is the optimum choice.

Future Work:

Several problems can be taken into consideration in potential research. primarily, larger scale datasets with diverse quantity of records, distinct count of capabilities (that is dimensionality), & vary- ent feature kinds which includes categoric, numeric & combined

information sorts may be used for similarly comparisons. 2nd, further for the construction of the Support Vector Machine classifiers, other techniques of classification's performance with respect to the 3 selected tactics may be tested. However, to sum up it'd be beneficial to analyze impact of utilizing specific HW environments for computing on the 3 varied strategies.

References

1. Wang, C. J., Ng, C. Y., & Brook, R. H. (2020). Response to COVID-19 in Taiwan: big data analytics, new technology, and proactive testing. *Jama*, 323(14), 1341-1342.
2. Iqbal, R., Doctor, F., More, B., Mahmud, S., & Yousuf, U. (2020). Big data analytics: Computational intelligence techniques and application areas. *Technological Forecasting and Social Change*, 153, 119253.
3. Tsai, C. F., Lin, W. C., & Ke, S. W. (2016). Big data mining with parallel computing: A comparison of distributed and MapReduce methodologies. *Journal of Systems and Software*, 122, 83-92.
4. Abaker, I., Hashem, T., Yaqoob, I., Anuar, N.B., Mokhtar, S., Gani, A., Khan, S.U., (2015). The rise of "big data" on cloud computing: review and open research issues. *Inf. Syst.* 47, 98–115.
5. Bi, X., Zhao, X., Wang, G., Zhang, P., Wang, C., (2015). Distributed extreme learning machine with kernels based on MapReduce. *Neurocomputing* 149, 456–463.
6. Cano, J.R., Herrera, F., Lozano, M., (2003). Using evolutionary algorithms as instance selection for data reduction: an experimental study. *IEEE Trans. Evol. Comput.* 7 (6), 561–575.
7. Coulouris, G., Dollimore, J., Kindberg, T., Blair, G., (2011). *Distributed Systems: Concepts And Design*, 5th ed. Addison-Wesley.



8. Dean, J., Ghemawat, S., (2010). Map reduce: a flexible data processing tool. *Commun.ACM* 53 (1), 72–77.
9. Fan, W., Bifet, A., (2012). Analysis big data: current status, and forecast to the future. *ACM SIGKDD Explor. Newslett.* 14 (2), 1–5.
10. Fernandez, A., del Rio, S., Lopez, V., Bawakid, A., del Jesus, M.J., Benitez, J.M., Herrera, F., (2014). Big data with cloud computing: an insight on the computing environment, MapReduce, and programming frameworks. *Wiley Interdiscip. Rev.* 4 (5), 380–409.
11. Galar, M., Fernandez, A., Barrenechea, E., Bustince, H., Herrera, F., (2012). A review on ensembles for the class imbalance problem: bagging-, boosting-, and hybrid-based approaches. *IEEE Trans. Syst. Man Cybern. Part C* 42 (4), 463–484.
12. García, S., Derrac, J., Cano, J.R., Herrera, F., (2012). Prototype selection for nearest neighbor classification: taxonomy and empirical study. *IEEE Trans. Pattern Anal. Mach. Intell.* 34 (3), 417–435.
13. Gottlieb, A., Almasi, G., (1989). *Highly Parallel Computing*. Benjamin-Cummings Publishing.
14. Kittler, J., Hatef, M., Duin, R.P.W., Matas, J., (1998). On combining classifiers. *IEEE Trans. Pattern Anal. Mach. Intell.* 20 (3), 226–239.
15. Kohavi, R., (1995). A study of cross-validation and bootstrap for accuracy estimation and model selection. In: *International Joint Conference on Artificial Intelligence*, pp. 1137–1143.
16. Lopez, V., del Rio, S., Benitez, J.M., Herrera, F., (2015). Cost-sensitive linguistic fuzzy rule based classification systems under the MapReduce framework for imbalanced big data. *Fuzzy Sets Syst.* 258, 5–38.
17. Mayer-Schonberger, V., Cukier, K., (2014). *Big Data: A Revolution That Will Transform How We Live, Work, And Think*. Eamon Dolan/Mariner Books.
18. Park, B., Kargupta, H., (2002). Distributed data analysis: algorithms, systems, and applications. In: Ye, N. (Ed.), *Data Analysis Handbook*. Oxford University Press, pp. 341–358.
19. Peteiro-Barral, D., Guijarro-Berdinas, B., (2013). A survey of methods for distributed machine learning. *Prog. Artif. Intell.* 2, 1–11.
20. Pyle, D., (1999). *Data Preparation For Data Analysis*. Morgan Kaufmann.
21. Qian, J., Lv, P., Yue, X., Liu, C., Jing, Z., (2015). Hierarchical attribute reduction algorithms for big data using MapReduce. *Knowl. Based Syst.* 73, 18–31.
22. Rajaraman, A., Ullman, J.D., (2011). *Analysis Of Massive Datasets*. Cambridge University Press.
23. Triguero, I., Peralta, D., Bacardit, J., Garcia, S., Herrera, F., (2015). MRPR: a MapReduce solution for prototype reduction in big data classification. *Neurocomputing* 150, 331–345.
24. Wilson, D.R., Martinez, T.R., (2000). Reduction techniques for instance-based learning algorithms. *Mach. Learn.* 38, 257–286.
25. Wu, X., Zhu, X., Wu, G.-Q., Ding, W., (2014). Data analysis with big data. *IEEE Trans. Knowl. Data Eng.* 26 (1), 97–107.



26. Zaki, M.J., (2000). Parallel and distributed data analysis: an introduction. Lect. Notes Comput. Sci. 1759, 1–23.
27. Zheng, L., Li, L., Duan, L., Lu, K., Shi, Z., Wang, M., Wu, W., Luo, P., (2012). Distributed data analysis: a survey. Inf. Technol. Manage. 13, 403–409.
28. Zhou, Z.-H., Chawla, N.W., Jin, Y., Williams, G.J., (2014). Big data opportunities and challenges: discussions from data analytics perspectives. IEEE Comput. Intell. Mag. 9 (4), 62–74.

