

Open Research Issues, Tools and Challenges in Big Data Analytics

Prof. Prachi N. Deshmukh Prof. Rasika S. Badre Prof. Shruti G. Taley

Abstract- A huge storage area of terabytes of data is generated each day from modern information systems and digital technologies such as Internet of Things and cloud computing. Study of these massive data requires a lot of efforts at multiple levels to extract knowledge for decision making. Therefore, big data analysis is a current area of study and improvement. The basic goal of this paper is to explore the possible impact of big data challenges, open research issues, and various tools related with it. As a result, this article provides a platform to explore big data at various stages. Additionally, it opens a new possibility for researchers to extend the result, based on the challenges and open research issues.

Keywords- Big data analytics; Hadoop; Massive data; Structured data; Unstructured Data.

I. INTRODUCTION

In digital world, information are generated from different sources and the rapid transition from digital technologies has led to growth of big data. It provides evolutionary breakthroughs in many fields with gathering of large datasets. In general, it refers to the collection of large and difficult datasets which are difficult to process using established database management tools or data processing applications. These are obtainable in structured, semi-structured, and unstructured format in petabytes and outside. Formally, it is distinct from 3Vs to 4Vs. 3Vs refers to volume, velocity, and variety. Volume refers to the large amount of data that are being generated everyday whereas velocity is the rate of growth and how fast the data are collected for being analysis. Variety provide information about the types of data such as structured, unstructured, semi-structured etc. The fourth V refers to veracity that includes accessibility and responsibility. The prime goal of big data analysis is to method data of high volume, velocity, variety, and veracity using various usual and computational intelligent techniques [11]. Generally, Data warehouses have been used to control the big dataset. In this

case extracting the exact knowledge from the accessible big data is foremost issue. The key problem in the analysis of big data is the need of organization between database systems as well as with analysis tools such as data mining and statistical analysis. These challenges generally occur when we wish to execute knowledge discovery and demonstration for its practical applications. Additionally, the study on difficulty theory of big data will help understand required characteristics and formation of difficult patterns in big data, simplify its representation, gets better knowledge abstraction, and guide the design of computing models and are interested in disseminating the findings of big data. This paper focuses on challenges in big data and its available techniques.

II. OPEN RESEARCH ISSUES IN BIG DATA ANALYTICS

Big data analytics and data science are becoming the research main point in industries and academia. Data science aims at researching big data and knowledge extraction from data. Applications of big data and data science contain information science, uncertainty modeling, uncertain data analysis, machine learning, statistical learning, pattern recognition, data warehousing, and signal processing. Main focus of this section is to discuss open research issues in big data analytics. The research issues pertaining to big data analysis are classified into three big categories namely internet of things (IoT), cloud computing, bio inspired computing, and quantum computing.

A. IoT for Big Data Analytics

Internet has simplified global interrelations, the ability of businesses, cultural revolutions and an incredible number of personal characteristics. Currently, machines are getting in on the act to control incalculable autonomous gadgets via internet and create Internet of Things (IoT). Thus, appliances are appropriate the user of the internet, just like humans with the web browsers. Internet of Things is attracting the notice of recent researchers for its most capable opportunities and

challenges. It has an essential economic and societal impact for the future construction of information, network and communication technology. The new regulation of the future will be ultimately, everything will be connected and intelligently restricted. The idea of IoT is becoming more important to the practical world due to the improvement of mobile devices, fixed and ubiquitous communication technologies, cloud computing, and data analytics. In a broader sense, just like the internet, Internet of Things enables the devices to exist in a myriad of places and facilitates applications ranging from minor to the crucial. A number of diversified technologies such as computational intelligence, and big data can be incorporated simultaneously to improve the data management and knowledge discovery of large scale automation applications. Much research in this direction has been carried out by Mishra, Lin and Chang [12]. Knowledge achievement from IoT data is the biggest challenge that big data professionals are facing. Therefore, it is essential to develop communications to analyze the IoT data. An IoT device generates constant streams of data and there researchers can develop tools to extract meaningful information from these data using machine learning techniques. Accepting these streams of data generated from IoT devices and analyzing them to acquire important information is a challenging issue and it leads to big data analytics. Machine learning algorithms and computational intelligence techniques is the only clarification to handle big data from IoT perspective. Key technologies that are related with IoT are also discussed in many research papers [17]. Figure 1 depicts an overview of IoT big data and knowledge discovery process.

B. Cloud Computing for Big Data Analytics

The improvement of virtualization technologies have complete supercomputing more accessible and affordable. Computing infrastructures that are unknown in virtualization software make systems to act like a true computer, but with the flexibility of specification data such as a number of processors, disk space, memory, and operating system. The use of these virtual computers is identified as cloud computing which has been one of the most robust big data techniques. Big Data and cloud computing technologies are residential with the importance of developing a scalable and on demand accessibility of resources and data. Cloud computing harmonize massive data by on demand contact to configurable computing resources through virtualization techniques. The profit of utilizing the Cloud computing include offering resources when there is a demand and pay only for the resources which is needed to develop the product. Concurrently, it improves availability and cost reduction. So Cloud computing helps in developing a business model for all varieties of applications with infrastructure and tools. Big data application using cloud computing should maintain data analytic and development. The cloud environment should propose tools that allow data scientists and business analysts to interactively and collaboratively explore knowledge acquisition data for further processing and extracting fruitful outcome. This can assist to solve large applications that may arise in various domains. In addition to this, cloud computing should also allow scaling of tools from virtual technologies into new technologies similar to spark, R, and other types of big data processing techniques.

C. Quantum Computing for Big Data Analysis

A quantum computer has memory that is exponentially bigger than its physical size and can manage an exponential set of inputs simultaneously [10]. This exponential improvement in computer systems might be achievable. If an actual quantum computer is accessible now, it could have solved problems that are exceptionally difficult on recent computers, of course today's large data problems. The main technical difficulty in building a quantum computer could soon be possible. Quantum computing provides a way to join the quantum mechanics to method the information. In traditional computer, information is obtainable by long strings of bits which instruct either a zero or a one. On the other hand a quantum computer uses quantum bits or qubits. The difference between qubit and bit is that, a qubit is a quantum system that encodes the zero and the one into two discernible quantum states. Therefore, it can be capitalized on the phenomena of superposition and entanglement. It is because qubits perform quantumly. For example, 100 qubits in

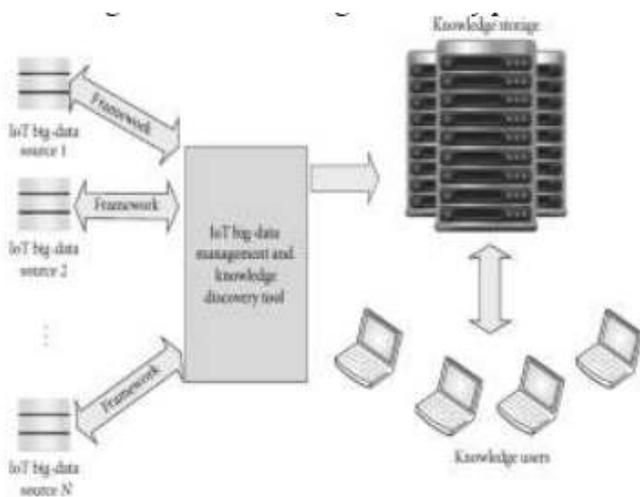


Fig. 1: IoT Big Data Knowledge Discovery

quantum systems need to 2100 complex values to be stored in a common computersystem. It means that many big data problems canbe solve much more rapidly by superior scalequantum computers compare with typicalcomputers. Hence it is a challenge for thisinvention to built a quantum computer and assistquantum computing to explain big data problems.

III. TOOLS FOR BIG DATA PROCESSING

Large numbers of tools are easy to get toprocess big data. In this section, we discuss somecurrent techniques for analyzing big data withconsequence on three important promising toolsnamely MapReduce, Apache Mahout, and Dryad.

A. Apache Hadoop and MapReduce

The most well-known software stage for big dataanalysis is Apache Hadoop and Mapreduce. Itconsists of hadoop kernel, mapreduce, hadoopdistributed file system (HDFS) and apache hive etc.Map reduce is a programming model for dealingout huge datasets is based on divide and conquermethod. The divide and conquer method isimplemented in two steps such as Map step andReduce Step. Hadoop works on two kinds of nodessuch as master node and worker node. The masternode divides the input into smaller sub problemsand then distributes them to worker nodes in mapstep. Thereafter the master node join the outputs forall the subproblems in reduce step. Moreover,Hadoop and MapReduce ability as a well-builtsoftware support for solving big data problems. It isalso cooperative in fault-tolerant storage and highthroughput data giving out.

B. Apache Mahout

Apache mahout goal to give scalable andcommercial machine learning techniques for largescale and intelligent data analysis applications.Core algorithms of mahout including clustering,classification, pattern mining, regression,dimensionality reduction, evolutionary algorithms,and batch based mutual filtering execute on top ofHadoop step through map decrease framework. Thegoal of mahout is to build a vibrant, reactive,diverse community to facilitate discussions on theproject and potential use cases. The basic objectiveof Apache mahout is to provide a tool for enrichingbig challenges. The different companies those whohave implemented scalable machine learningalgorithms are Google, IBM, Amazon, Yahoo, Twitter, and facebook [4].

C. Dryad

It is any more popular programmingmodel for implementing parallel and distributedprograms for handling large situation

bases ondataflow graph. It consists of a cluster ofcomputing nodes, and an user use the devices of acomputer cluster to run their program in adistributed way. Indeed, a dryad user use thousandsof machines, each of them with several processorsor cores. The main benefit is that users do not needto know anything about correspondprogramming. A dryad application runs acomputational directed graph that is composed ofcomputational vertices and communicationchannels. Therefore, dryad provides a large numberof functionality including generating of workgraph, scheduling of the machines for the availableprocesses, evolution failure handling in the cluster,set of performance metrics, visualizing the work,invoking user defined policies and dynamicallyupdating the work graph in response to thesestrategy decisions without knowing the semanticsof the vertices [5].

IV. CHALLENGES IN BIG DATA ANALYTICS

Here the challenges of big data analyticsare classify into four large categories namely datastorage space and analysis; knowledge discoveryand computational complexities; scalability andvisualization of data; and information security. Wediscuss these issues for a short time in thefollowing subsections.

A. Data Storage and Analysis

In current years the size of informationhas grown exponentially by various means such asmobile devices, aerial sensory technologies, remotesensing, radio frequency identification readers etc.Thus, the first test for big data analysis is storagespace mediums and high input/output speed. Insuch cases, the data accessibility must be on the topmain concern for the knowledge discovery anddemonstration. The prime reason is being that, itmust be accessed easily and rapidly for furtherstudy. In past decades, analyst use hard disk drivesto accumulate data but, it slower casualinput/output performance than sequentialinput/output. To overcome this limitation, the ideaof solid state drive (SSD) and phase changememory (PCM) was introduced. However theavailable storage technologies cannot grasp therequired performance for processing big data.Data reduction, data selection, featureselection is a crucial job particularly when dealingwith large datasets. This presents an exceptionalchallenge for researchers. It is because, existingalgorithms may not always respond in an sufficienttime when dealing with these high dimensionaldata. Automation of this technique and developingnew machine learning algorithms to make sureconsistency is a key challenge in recent years.

B. Knowledge Discovery and Computational Complexities

Knowledge discovery and demonstration is a major concern in big data. It includes a number of sub fields such as authentication, archiving, management, preservation, in order retrieval, and representation. There are numerous tools for knowledge discovery and demonstration such as fuzzy set [9], rough set [19], soft set [3], near set [8], formal model analysis [15], main component analysis [7] etc to name a few. Furthermore a lot of hybridized techniques are also developed to process real life problems. All these techniques are trouble dependent. Further some of these techniques may not be suitable for big datasets in a sequential computer. At the same time some of these techniques has good quality characteristics of scalability over similar computer. As the size of big data keeps growing exponentially, the existing tools may not be well-organized to method these data for obtaining significant information. The most admired approach in case of big dataset management is data warehouses and data marts. Data warehouse is mostly dependable to accumulate data that are sourced from operational systems while data mart is based on a data warehouse and facilitates analysis. Analysis of huge dataset requires more computational complexities. The main apprehension is to handle inconsistencies and uncertainty current in the datasets. In general, systematic modeling of the computational complexity is used. It may be difficult to set up a comprehensive mathematical system that is broadly related to Big Data. But a domain specific data analytics can be done simply by understanding the particular complexities. A series of such development could suggest big data analytics for different areas. Much research and survey has been carried out in this direction using machine learning techniques with the least memory requirements. The basic purpose in these research is to minimize computational cost giving out and complexity [13], [2], [14].

C. Scalability and Visualization of Data

The mainly central challenge for big data analysis techniques is its scalability and security. It is necessary to grow sampling, on-line, and multi-resolution analysis techniques. Incremental techniques have excellent scalability property in the side of big data analysis. As the information size is scaling much quicker than CPU speeds, there is a natural dramatic shift in processor technology being fixed with growing number of cores [1]. This move in processors leads to the growth of parallel computing. The reason of visualizing data is to current them additional effectively using some techniques of graph theory. Graphical visualization provides the link between data with proper explanation. However, online marketplace like flipkart, amazon, e-bay have millions of users and billions of

generate to solve each month. This generates a lot of information.

D. Information Security

In big data analysis massive sum of data are connected, analyzed, and mined for meaningful patterns. Preserving responsive information is most important issue in big data analysis [6]. Therefore, information protection is becoming a big data analytics problem. Security of big data can be improved by using the techniques of authentication, authorization, and encryption. Different protection measures that big data applications appearance are scale of network, variety of different devices, real time security monitoring, and need of interruption system [18], [26]. The security challenge caused by big data has attracted the attention of information security. Therefore, interest has to be particular to develop a multi level protection rule model and avoidance system.

V. CONCLUSION

In current years information are generated at a amazing pace. Analyzing these information is challenging for a general man. To this conclusion in this paper, we study the various research issues, challenges, and tools worn to analyze these big data. From this study, it is appreciate that every big data stage has its individual focus. Some of them are intended for batch processing whereas some are high-quality at real-time analytic. Each big data stage also has exact functionality. Different techniques available for the analysis include statistical analysis, machine learning, data mining, intelligent analysis, cloud computing, quantum computing, and data stream processing. We regard as in prospect researchers will give more focus to these techniques to explain troubles of big data successfully and resourcefully.

REFERENCES

- 1 A. Jacobs, The pathologies of big data, Communications of the ACM, 2(8) (2009), pp.36-44.
- 2 Changwon. Y, Luis. Ramirez and Juan. Liuzzi, Big data analysis using modern statistical and machine learning methods in medicine, International Neurology Journal, 18 (2014), pp.50-57.
- 3 D. Molodtsov, Soft set theory first results, Computers and Mathematics with Applications, 37(4/5) (1999), pp.19-31.
- 4 G. Ingersoll, Introducing apache mahout: Scalable, commercial friendly machine learning for building intelligent applications, White Paper, IBM Developer Works, (2009), pp.1-18.

- 5 H. S. G. Fox and J. Qiu, Performance model for parallel matrix multiplication with dryad: Dataflow graph runtime, Second International Conference on Cloud and Green Computing, 2012, pp.675-683.
- 6 H. Zhu, Z. Xu and Y. Huang, Research on the security technology of big data information, International Conference on Information Technology and Management Innovation, 2015, pp.1041-1044.
- 7 I. T. Jolliffe, Principal Component Analysis, Springer, New York, 2002.
- 8 J. F. Peters, Near sets. General theory about nearness of objects, Applied Mathematical Sciences, 1(53) (2007), pp.2609-2629.
- 9 L. A. Zadeh, Fuzzy sets, Information and Control, 8 (1965), pp.338- 353.
- 10 M. A. Nielsen and I. L. Chuang, Quantum Computation and Quantum Information, Cambridge University Press, New York, USA 2000.
- 11 M. K. Kakhani, S. Kakhani and S. R. Biradar, Research issues in big data analytics, International Journal of Application or Innovation in Engineering & Management, 2(8) (2015), pp.228-232.
- 12 N. Mishra, C. Lin and H. Chang, A cognitive adopted framework for IoT big data management and knowledge discovery prospective, International Journal of Distributed Sensor Networks, 2015, (2015), pp. 1-13
- 13 O. Y. Al-Jarrah, P. D. Yoo, S. Muhaidat, G. K. Karagiannidis and K. Taha, Efficient machine learning for big data: A review, Big Data Research, 2(3) (2015), pp.87-93.
- 14 P. Singh and B. Suri, Quality assessment of data using statistical and machine learning methods. L. C. Jain, H. S. Behera, J. K. Mandal and D. P. Mohapatra (eds.), Computational Intelligence in Data Mining, 2 (2014), pp. 89-97.
- 15 R. Wille, Formal concept analysis as mathematical theory of concept and concept hierarchies, Lecture Notes in Artificial Intelligence, 3626 (2005), pp.1-33.
- 16 X. Jin, B. W. Wah, X. Cheng and Y. Wang, Significance and challenges of big data research, Big Data Research, 2(2)(2015), pp.59-64.
- 17 X. Y. Chen and Z. G. Jin, Research on key technology and applications for internet of things, Physics Procedia, 33, (2012), pp. 561-566.
- 18 Z. Hongjun, H. Wenning, H. Dengchao and M. Yuxing, Survey of research on information security in big data, Congresso da sociedade Brasileira de Computacao, 2014, pp.1-6.
- 19 Z. Pawlak, Rough sets, International Journal of Computer Information Science, 11 (1982), pp.341-356.