

Challenges and OpenResearch Issues and Tools on Big Data Analytics

GameilSaadHamzh Ali¹, Dr.A.Nithya²

ABSTRACT

Modern information systems and digital technologies such as Internet of Things and cloud computing is generated a huge repository of terabytes of data each day. It requires a lot of efforts to analysis of these massive data at multiple levels to extract knowledge for decision making. Therefore, the current area of research and development isbig data analysis. The aim of this paper is to explore the potential impact of big data challenges, open research issues, and various tools associated with it. As a result, this article opens a new horizon for researchers to develop the solution, based on the challenges and open research issues.

Keywords:Big data analytics; Structureddata; Semi Structured data ,Unstructured Data

1- INTRODUCTION

Data are generated in digital world from several sourcesand the fast transition from digital technologies has led togrowth of big data. It provides evolutionary breakthroughs inmany fields with collection of large datasets. In generally,there are a collection ofcomplexand

largedatasets whichare difficult to process usingdata processing applications or traditional database managementtools . These are availablein structured, semi-structured, and unstructured format inetabytes and beyond. Formally, it is defined from four challenges refers tovelocity, volume,variety, and veracity. Velocity is the rate of growth and how fast the data are gathered for being analysiswhereasVolume refers tothe huge amount of data that are being generated every day. Variety provides informationabout the types of data such as structured, unstructured, semi structuredetc. The fourth V refers to veracity that includesavailability and accountability. The prime objective of big dataanalysis is to process data of high volume, velocity, variety, andveracity using various traditional and computational intelligenttechniques [1].Gandomi and Haider [2]discussedSome of these extraction methods for obtaininghelpful information . The following Figure 1 refers to the definition of bigdata. However exact definition for big data is not defined andthere is a believe that it is problem specific. This will support usin obtaining better decision making, insight discovery andoptimization while being innovative and cost-effective.It is expected that the growth of big data is estimated toreach 25 billion by 2015 [3]. From the perspective of

the information and communication technology, big data is a robust impetus to the next generation of information technology industries [4], which are broadly built on the third platform, mainly referring to big data, cloud computing, internet of things, and social business. Generally, Data warehouses have been used to manage the large dataset. In this case extracting the exact knowledge from the available big data is a main issue. Most of the presented approaches in data mining are not usually able to handle the large datasets successfully. In big data analysis the key problem is the lack of coordination between database systems as well as with analysis tools such as statistical analysis and data mining. These challenges generally arise when we wish to perform knowledge discovery and representation for its practical applications. A fundamental problem is how to quantitatively describe the essential characteristics of big data. There is a need for epistemological implications in describing data revolution [5]. Additionally, the study on complexity theory of big data will help understand essential characteristics and formation of complex patterns in big data, simplify its representation, get better knowledge abstraction, and guide the design of computing models and algorithms on big data [4]. Much research was carried out by various researchers on big data and its trends [6], [7], [8]. However, it is to be noted that all data available in the form of big data are not useful for analysis or decision making process. Industry and academia are interested in disseminating the findings of big data. This paper focuses on challenges in big data and its available techniques. Additionally, we state open research issues in big data. So, to elaborate this, the paper is divided

into following sections. Section 2 deals with challenges that arise during fine tuning of big data. Section 3 provides the open research issues that will help us to process big data and extract useful knowledge from it. Section 4 furnishes an insight to big data tools and techniques. Conclusion remarks are provided in section 5 to summarize outcomes.

2- CHALLENGES IN BIG DATA ANALYTICS

Big data in recent years has been accumulated in several domains like health care, public administration, retail, biochemistry, and other interdisciplinary scientific researches. Web-based applications encounter big data frequently, such as social computing, internet text and documents, and internet search indexing. Social computing includes social network analysis, online communities, recommender systems, reputation systems, and prediction markets where as internet search indexing includes ISI, IEEE Xplore, Scopus, Thomson Reuters etc. Considering this advantages of big data it provides a new opportunities in the knowledge processing tasks for the upcoming researchers. However opportunities always follow some challenges. To handle the challenges we need to know various computational complexities, information security, and computational method, to analyze big data. For example, many statistical methods that perform well for small data size do not scale to voluminous data. Similarly, many computational techniques that perform well for small data face significant challenges in analyzing big data. Various challenges that the health sector face was

being researched by much researchers [9], [10]. Here the challenges of big data analytics are classified into four broad categories namely storage and analysis of data; knowledge

discovery computational complexities; visualization and scalability of data; and information security. We discuss these issues briefly in the following subsections.

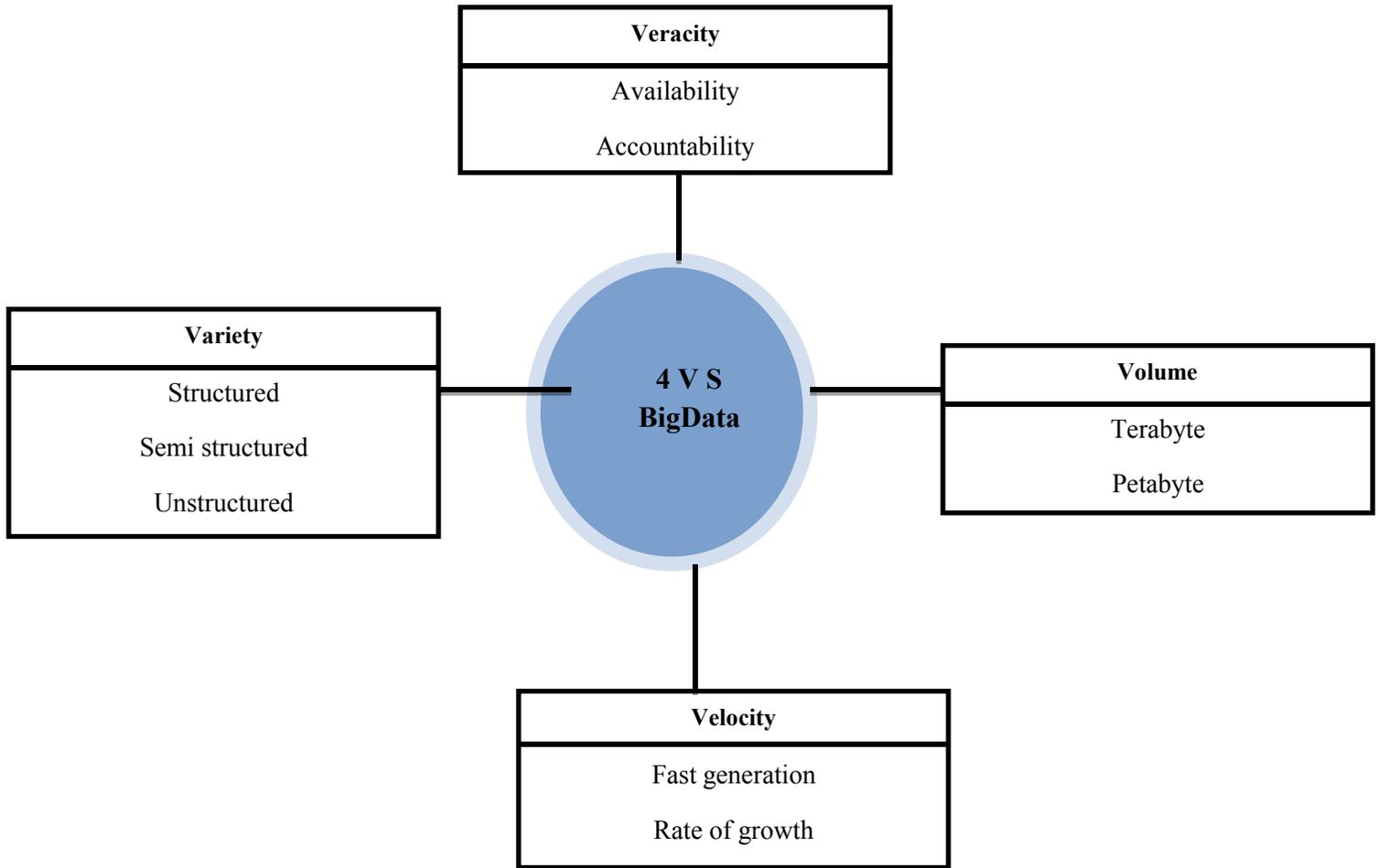


Fig. 1: Characteristics of Big Data[38]

A. Storage and Analysis of Data

In latest years the size of data has grown exponentially by several resources such as aerial sensory technologies, mobile devices, radio frequency identification readers, remote sensing etc. These data are stored on spending much

cost whereas they ignored or deleted finally because there is no enough space to store them. Therefore, the storage mediums and higher input/output speed is the first challenge for big data analysis. In such cases, the data accessibility must be on the top priority for the knowledge discovery and representation. The prime reason is

being that, it must be accessed easily and promptly for further analysis. In past decades, analyst use harddisk drives to store data but, it slower random input/output performance than sequential input/output. To overcome this limitation, the concept of solid state drive (SSD) and phase change memory (PCM) was introduced. However the available storage technologies cannot possess the required performance for processing big data. Another challenge with Big Data analysis is attributed to diversity of data. With the ever growing of datasets, data mining tasks has significantly increased. Additionally data selection, feature selection, data reduction is an essential task especially when dealing with large datasets. This presents an unprecedented challenge for researchers. It is because, existing algorithms may not always respond in an adequate time when dealing with these high dimensional data. Automation of this process and developing new machine learning algorithms to ensure consistency is a major challenge in recent years. In addition to all these Clustering of large datasets that help in analyzing the big data is of prime concern [11]. Recent technologies such as hadoop and mapReduce make it possible to collect large amount of semi structured and unstructured data in a reasonable amount of time. The key engineering challenge is how to effectively analyze these data for obtaining better knowledge. A standard process to this end is to transform the semi structured or unstructured data into structured data, and then apply data mining algorithms to extract knowledge. A framework to analyze data was discussed by Das and Kumar [12]. Das et al in their paper [13] was

discussed the similarly detail explanation of data analysis for public. The major challenge in this case is to pay more attention for designing storage systems and to elevate efficient data analysis tool that provide guarantees on the output when the data comes from different sources. Furthermore, design of machine learning algorithms to analyze data is essential for improving efficiency and scalability.

B. Knowledge Discovery and Computational Complexities

The main issue in big data is Knowledge discovery and representation. It includes a number of sub fields such as information retrieval, archiving, authentication, management, preservation, and representation. There are several tools for knowledge discovery and representation such as fuzzy set [14], rough set [15], soft set [16], near set [17], formal concept analysis [18], principal component analysis [19] etc. to name a few. Additionally many hybridized techniques are also developed to process real life problems. All these techniques are problem dependent. Further some of these techniques may not be suitable for large datasets in a sequential computer. At the same time some of the techniques has good characteristics of scalability over parallel computer. Since the size of big data keeps increasing exponentially, the available tools may not be efficient to process these data for obtaining meaningful information. The most popular approach in case of large dataset management is data warehouses and data marts. Data warehouse is mainly responsible to store data that are sourced from operational systems whereas data mart is based on a

datawarehouse and facilitates analysis. It requires more computational complexities' to Analysis of large. The major issue is to handle uncertainty and inconsistencies present in the datasets. In general, systematic modeling of the computational complexity is used. It may be difficult to establish a comprehensive mathematical system that is broadly applicable to Big Data. But a domain specific data analytics can be done easily by understanding the particular complexities. A series of such development could simulate 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 objective in this research is to minimize computational cost processing and complexities [20], [21], [22]. However, current big data analysis tools have poor performance in handling computational complexities, uncertainty,

C. Scalability and Visualization of Data

Scalability and security are the most important challenge for big data analysis techniques. In the recent decades researchers have paid attention to accelerate data analysis and its speed up processors followed by Moore's Law. For the previous, it is necessary to develop on-line, sampling, and multi-resolution analysis techniques. Incremental techniques have good scalability property in the aspect of big data analysis. As the data size is scaling much faster than CPU speeds, there is a natural dramatic shift in processor technology being embedded with increasing number of cores [23]. This shift in processors leads to the development of parallel computing. Real time applications like navigation, social networks, finance, internet search,

timeliness etc. requires parallel computing. The objective of visualizing data is to present them more adequately using some techniques of graph theory. Graphical visualization provides the link between data with proper interpretation. However, online marketplace like flipkart, amazon, e-bay have millions of users and billions of goods to sold each month. This generates a lot of data. To this end, some company uses a tool Tableau for big data visualization. It has capability to transform large and complex data into intuitive pictures. These help employees of a company to monitor latest customer feedback, visualize search relevance, and their sentiment analysis. However, current big data visualization tools mostly have poor performances in functionalities, scalability, and response time. We can observe that big data have produced many challenges for the developments of the hardware and software which leads to cloud computing, scalability, parallel computing, distributed computing, visualization process. To overcome this issue, we need to correlate more mathematical models to computer science.

D. Information Security

The huge amount of data in big data analysis are correlated, analyzed, and mined for meaningful patterns. All organizations have different policies to safe guard their sensitive information. The main issue in big data analysis is preserving sensitive information. There is a huge security risk associated with big data [24]. Therefore, information security is becoming a big data analytics problem. Big data security can be improved by using the techniques of authorization, authentication, and encryption. Various security

measures that big data applications face are scale of network, variety of different devices, real time security monitoring, and lack of intrusion system [25], [26]. The challenge of security caused by big data has attracted the attention of information security. Therefore, attention has to be given to develop a multi-level security policy model and prevention system. Although much research has been carried out to secure big data [25] but it requires lot of improvement. The major challenge is to develop a multi-level security, privacy preserved data model for big data.

3. RESEARCH ISSUES IN BIG DATA ANALYTICS

Data science and big data analytics are becoming the focal point of research in academia and industries. Data science aims at researching big data and knowledge extraction from data. Big data applications and data science include information science, uncertainty modeling, uncertain data analysis, machine learning, statistical learning, pattern recognition, data warehousing, and signal processing. Effective integration of technologies and analysis will result in predicting the future drift of events. 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 broad categories namely internet of things (IoT), cloud computing, bio inspired computing, and quantum computing. However it is not limited to these issues. More research issues related to health care big data can be found in Husing Kuo et al. paper [9].

A. IoT for Big Data Analytics

Internet has rearranged the art of businesses, global interrelations, cultural revolutions and an unbelievable number of personal characteristics. Currently, machines are getting in on the act to control innumerable autonomous gadgets via internet and create Internet of Things (IoT). Therefore, the user of the internet is becoming appliances, just like humans with the web browsers. Internet of Things is attracting the attention of recent researchers for its most promising opportunities and challenges. It has societal impact and imperative economic for the future construction of information, communication technology and network. The new regulation of future will be eventually; everything will be connected and intelligently controlled. The concept of IoT is becoming more pertinent to the realistic world due to the development of mobile devices, embedded and ubiquitous communication technologies, cloud computing, and data analytics. Moreover, IoT presents challenges in combinations of volume, velocity and variety. 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 trivial to the crucial. Conversely, it is still mystifying to understand IoT well, including definitions, content and differences from other similar concepts. Several diversified technologies such as computational intelligence, and big-data can be incorporated together 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 [27]. Knowledge acquisition from IoT data is the biggest challenge that big data professionals are facing. Thus, it is important to develop infrastructure to analyze the IoT data. An IoT device generates continuous streams of data and the researchers can develop tools to extract meaningful information from these data using machine learning techniques. Understanding these streams of data generated

from IoT devices and analyzing them to get meaningful information is a challenging issue and it leads to big data analytics. Machine learning algorithms and computational intelligence techniques is the only solution to handle big data from IoT perspective. Key technologies that are associated with IoT are also discussed in many research papers [28]. Figure 2 depicts an overview of IoT big data and knowledge discovery process.

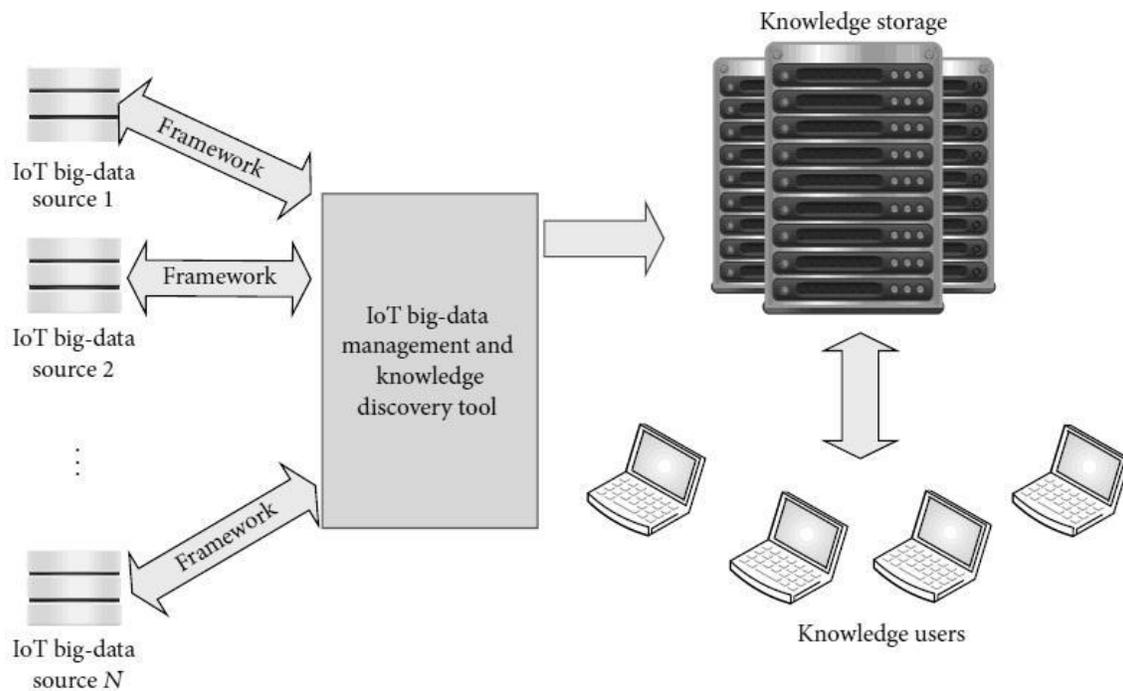


Fig. 2: IoT Big Data Knowledge Discovery[38]

Knowledge exploration systems have originated from theories of human information processing such as frames, rules, tagging, and semantic networks. In general, it consists of four segments such as knowledge acquisition, knowledge base, knowledge dissemination, and knowledge application. In phase of knowledge acquisition, various traditional and computational

intelligence techniques to discover the knowledge is discovered by using various traditional and computational intelligence techniques. The discovered knowledge is stored in knowledge bases and expert systems are generally designed based on the discovered knowledge. Knowledge dissemination is important for obtaining meaningful information

from the knowledgebase. Knowledge extraction is a process that searches documents, knowledge within documents as well as knowledgebases. The final phase is to apply discovered knowledge in various applications. It is the ultimate goal of knowledge discovery. The knowledge exploration system is necessarily iterative with the judgment

of knowledge application. There are many issues, discussions, and researches in this area of knowledge exploration. It is beyond the scope of this survey paper. For better visualization, knowledge exploration systems are depicted in Figure 3.

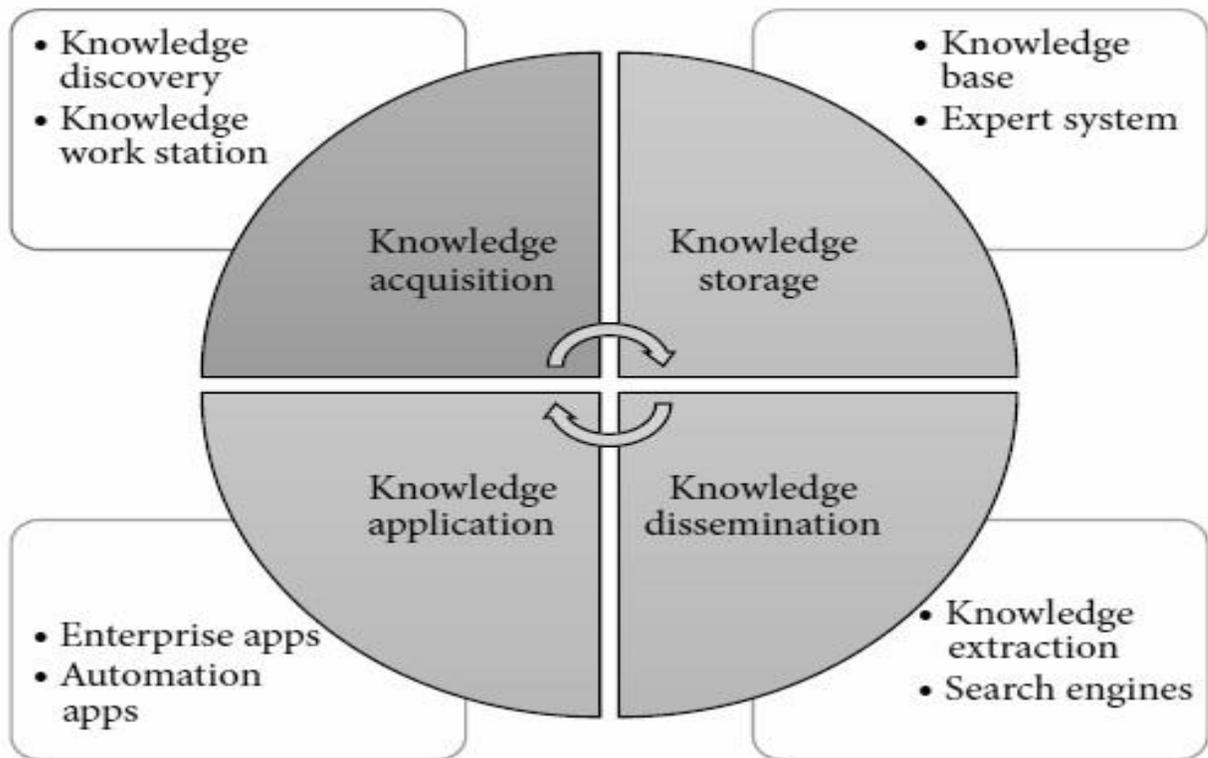


Fig. 3: IoT Knowledge Exploration System^[38]

B. Cloud Computing for Big Data Analytics

Virtualization technologies development have made supercomputing more accessible and affordable. Computing infrastructures that are hidden in virtualization software make systems to behave like a true computer, but with the flexibility of specification details such as number

of processors, disk space, memory, and operating system. The use of these virtual computers is known as cloud computing which has been one of the most robust big data techniques. Big Data and cloud computing technologies are developed with the importance of developing a scalable and on-demand availability of resource and data.

Cloud computing harmonize massive data by on demand access to configurable computing resources through virtualization techniques. The benefits 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. Simultaneously, it improves availability and cost reduction. Open challenges and research issues of big data and cloud computing are discussed in detail by many researchers which highlights the challenges in data management, data variety and velocity, data storage, data processing, and resource management [29], [30]. So Cloud computing helps in developing a business model for all varieties of applications with infrastructure and tools. Application of big data using cloud computing should support data development and analytic. The cloud environment should provide tools that allow data scientists and business analysts to interactively and collaboratively explore knowledge acquisition data for further processing and extracting fruitful results. This can help to solve large applications that may arise in various domains. In addition to this, cloud computing should also enable scaling of tools from virtual technologies into new technologies like spark, R, and other types of big data processing techniques. Big data forms a framework for discussing cloud computing options. Depending on special need, user can go to the marketplace and buy infrastructure services from cloud service providers such as Google, Amazon, IBM, software as a service (SaaS) from a whole crew of companies such as NetSuite, Cloud9, Jobscience etc. Another advantage of cloud computing is

cloud storage which provides a possible way for storing big data. The obvious one is the time and cost that are needed to upload and download big data in the cloud environment. Else, it becomes difficult to control the distribution of computation and the underlying hardware. But, the major issues are privacy concerns relating to the hosting of data on public servers, and the storage of data from human studies. All these issues will take cloud computing and big data to a high level of development.

4. BIG DATA PROCESSING TOOLS

Big numbers of tools are available to process big data. In this section, we discuss some current techniques for analyzing big data with emphasis on three important emerging tools namely MapReduce, Apache Spark, and Storm. Most of the available tools concentrate on batch processing, stream processing, and interactive analysis. Most batch processing tools are based on the Apache Hadoop infrastructure such as Mahout and Dryad. Stream data applications are mostly used for real time analytic. Some examples of large scale streaming platform are Storm and Splunk. The interactive analysis process allow users to directly interact in real time for their own analysis. For example Dremel and Apache Drill are the big data platforms that support interactive analysis. These tools help us in developing the big data projects. An extraordinary list of big data tools and techniques is also discussed by much researchers [6], [34]. The typical work flow of big data project discussed by Huang et al is highlighted in this section [35] and is depicted in Figure 4.

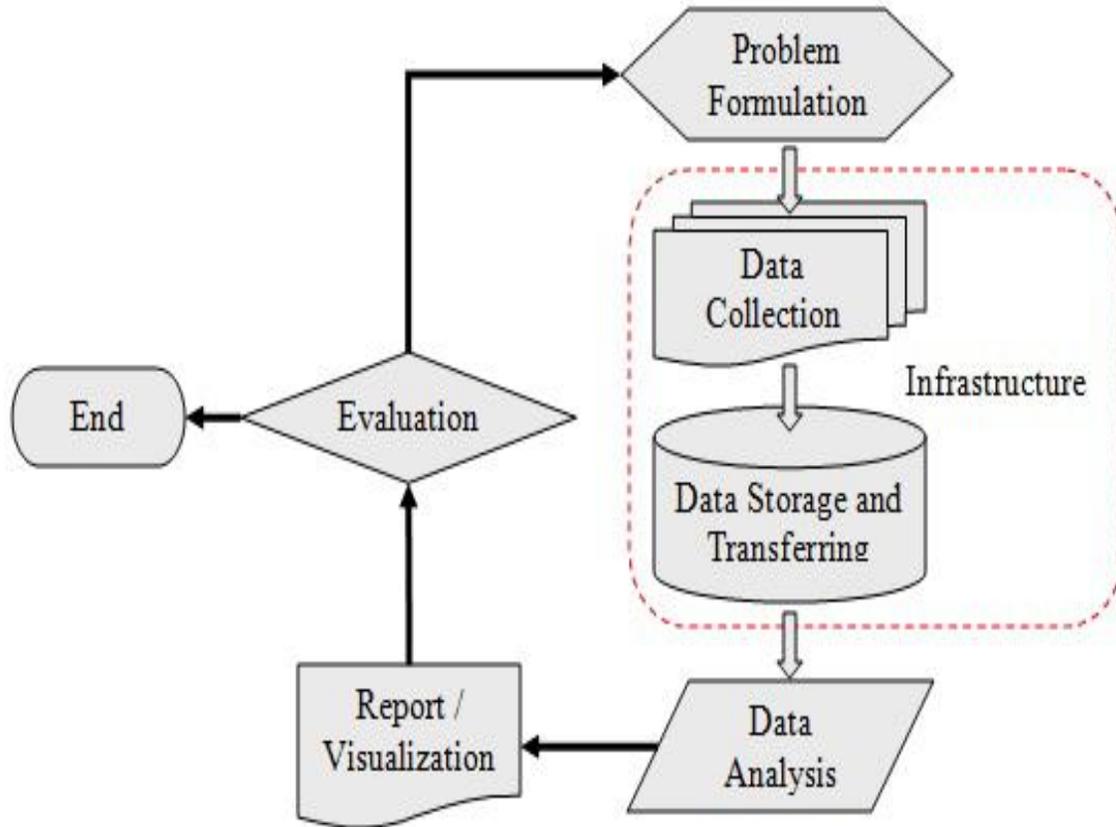


Fig. 4: Workflow of Big Data Project[38]

A. Apache Hadoop and MapReduce

The most recognized software platform for big data analysis is Apache MapReduce and Hadoop. It involves Hadoop kernel, MapReduce, Hadoop distributed file system (HDFS) and Apache Hive etc. MapReduce is a programming model for processing large datasets based on divide and

conquer method. The method of divide and conquer is implemented in two steps such as Map step and Reduce Step. Hadoop works on two kinds of nodes such as master node and worker node. The master node divides the input into smaller sub-problems and then distributes them to worker nodes in the map step. Thereafter the master

node combines the outputs for all the sub problems in reduce step. Moreover, Hadoop and MapReduce work as a powerful software framework for solving big data problems. It is also helpful in fault-tolerant storage and high throughput data processing.

B. Apache Mahout

The goal of Apache Mahout is to provide scalable and commercial machine learning techniques for large scale and intelligent data analysis applications. Core algorithms of Mahout including clustering, classification, pattern mining, regression, dimensionality reduction, evolutionary algorithms, and batch based collaborative filtering run on top of Hadoop platform through map reduce framework. The Mahout aims to build a responsive, vibrant, diverse community to facilitate discussion on the project and potential use cases. The basic objective of Apache Mahout is to provide a tool for alleviating big challenges. The different companies those who have implemented scalable machine learning algorithms are Google, IBM, Amazon, Yahoo, Twitter, and Facebook [36].

C. Apache Spark

Apache Spark is an open source big data processing framework built for speed processing, and sophisticated analytics. It is easy to use and was originally developed in 2009 in UC Berkeley's AMPLab. It was open sourced in 2010 as an Apache project. By using Spark you can quickly write applications in Java, Scala, or Python. In addition to map reduce operations, it supports SQL queries, streaming data, machine learning, and graph data processing. Spark runs on top of existing Hadoop distributed file system (HDFS) infrastructure to provide enhanced and additional functionality. Spark consists of components namely driver program, cluster manager and worker nodes. The driver program serves as the starting point of execution of an application on the Spark cluster. The cluster manager allocates the resources and the worker nodes to do the data processing in the form of tasks. Each application will have a set of processes called executors that are responsible for executing the tasks. The major advantage is that it provides support for deploying Spark applications in an existing Hadoop clusters. Figure 5 depicts the architecture diagram of Apache Spark. The various features of Apache Spark are listed below:

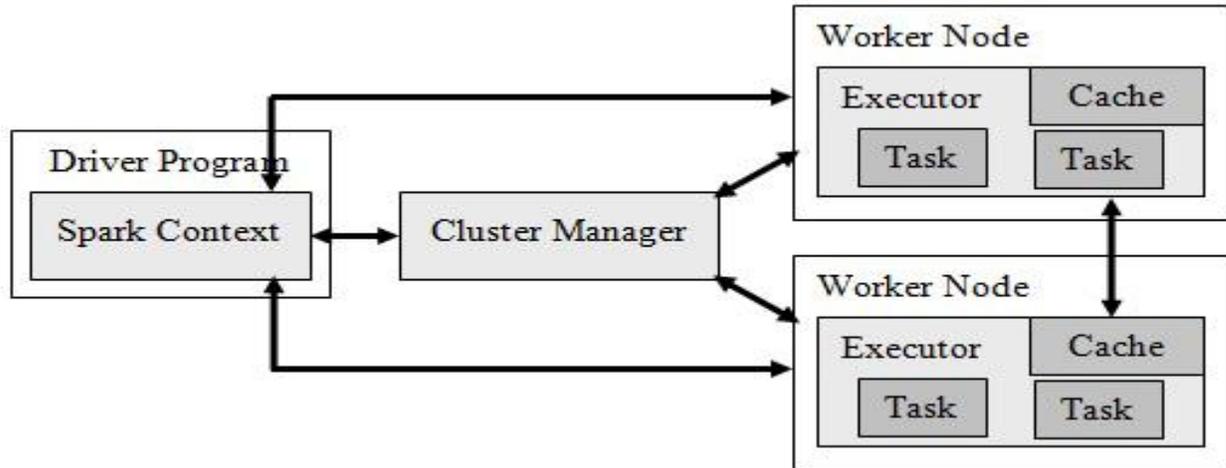


Fig. 5: Architecture of Apache Spark[38]

- The prime focus of spark includes resilient distributed datasets (RDD), which store data in-memory and provide fault tolerance without replication. It supports iterative computation, improves speed and resource utilization.
- The foremost advantage is that in addition to MapReduce, it also supports streaming data, machine learning, and graph algorithms.
- Another advantage is that, a user can run the application program in different languages such as Java, R, Python, or Scala. This is possible as it comes with higher-level libraries for advanced analytics. These standard libraries increase developer productivity and can be seamlessly combined to create complex workflows.
- Spark helps to run an application in Hadoop cluster, up to 100 times faster in memory, and 10 times faster when running on disk. It is possible because

of the reduction in number of read or write operations to disk.

- It is written in scala programming language and runs on java virtual machine (JVM) environment. Additionally, it supports java, python and R for developing applications using Spark.

5. SUGGESTIONS FOR FUTURE WORK

The amount of data collected from various applications all over the world across a wide variety of fields today is expected to double every two years. It has no utility unless these are analyzed to get useful information. This necessitates the development of techniques which can be used to facilitate big data analysis. The development of powerful computers is a boon to implement these techniques leading to automated systems. The transformation of data into knowledge is by no means an easy task for high performance large-scale data processing,

including exploiting parallelism of current and upcoming computer architectures for data mining. Moreover, these data may involve uncertainty in many different forms. Many different models like fuzzy sets, rough sets, soft sets, neural networks, their generalizations and hybrid models obtained by combining two or more of these models have been found to be fruitful in representing data. These models are also very much fruitful for analysis. More often than not, big data are reduced to include only the important characteristics necessary from a particular study point of view or depending upon the application area. So, reduction techniques have been developed. Often the data collected have missing values. These values need to be generated or the tuples having these missing values are eliminated from the data set before analysis. More importantly, these new challenges may comprise, sometimes even deteriorate, the performance, efficiency and scalability of the dedicated data intensive computing systems. The later approach sometimes leads to loss of information and hence not preferred. This brings up many research issues in the industry and research community in forms of capturing and accessing data effectively. In addition, fast processing while achieving high performance and high throughput, and storing it efficiently for future use is another issue. Further, programming for big data analysis is an important challenging issue. Expressing data access requirements of applications and designing programming language abstractions to exploit parallelism are an immediate need [38]. Additionally, machine learning concepts and tools are gaining popularity

REFERENCES

among researchers to facilitate meaningful results from these concepts. Research in the area of machine learning for big data has focused on data processing, algorithm implementation, and optimization. Many of the machine learning tools for big data are started recently and need drastic change to adopt it. We argue that while each of the tools has their advantages and limitations, more efficient tools can be developed for dealing with problems inherent to big data. The efficient tools to be developed must have provision to handle noisy and imbalance data, uncertainty and inconsistency, and missing values.

6. CONCLUSION

In recent years data are generated at a dramatic pace. Analyzing these data is challenging for a general man. To this end in this paper, we survey the various research issues, challenges, and tools used to analyze these big data. From this survey, it is understood that every big data platform has its individual focus. Some of them are designed for batch processing whereas some are good at real-time analytic. Each big data platform also has specific functionality. Different techniques used for the analysis include statistical analysis, machine learning, data mining, intelligent analysis, cloud computing, quantum computing, and data stream processing. We believe that in future researchers will pay more attention to these techniques to solve problems of big data effectively and efficiently.

- [1] M. K.Kakhani, S. Kakhani and S. R.Biradar, Research issues in bigdata analytics, *International Journal of Application or Innovation inEngineering & Management*, 2(8) (2015), pp.228-232.
- [2] A. Gandomi and M. Haider, Beyond the hype: Big data concepts, methods,and analytics, *International Journal of Information Management*,35(2) (2015), pp.137-144.
- [3] C. Lynch, Big data: How do your data grow?, *Nature*, 455 (2008),pp.28-29.
- [4] X. Jin, B. W.Wah, X. Cheng and Y. Wang, Significance and challengesof big data research, *Big Data Research*, 2(2) (2015), pp.59-64.
- [5] R. Kitchin, Big Data, new epistemologies and paradigm shifts, *BigData Society*, 1(1) (2014), pp.1-12.
- [6] C. L. Philip, Q. Chen and C. Y. Zhang, Data-intensive applications,challenges, techniques and technologies: A survey on big data, *InformationSciences*, 275 (2014), pp.314-347.
- [7] K. Kambatla, G. Kollias, V. Kumar and A. Gram, Trends in big dataanalytics, *Journal of Parallel and Distributed Computing*, 74(7) (2014),pp.2561-2573.
- [8] S. Del. Rio, V. Lopez, J. M. Bentez and F. Herrera, On the use ofmapreduce for imbalanced big data using random forest, *InformationSciences*, 285 (2014), pp.112-137.
- [9] MH. Kuo, T. Sahama, A. W. Kushniruk, E. M. Borycki and D. K.Grunwell, Health big data analytics: current perspectives, challengesand potential solutions, *International Journal of Big Data Intelligence*,1 (2014), pp.114-126.
- [10] R. Nambiar, A. Sethi, R. Bhardwaj and R. Vargheese, A look atchallenges and opportunities of big data analytics in healthcare, *IEEEInternational Conference on Big Data*, 2013, pp.17-22.
- [11] Z. Huang, A fast clustering algorithm to cluster very large categoricaldata sets in data mining, *SIGMOD Workshop on Research Issues onData Mining and Knowledge Discovery*, 1997.[12] T. K. Das and P. M. Kumar, Big data analytics: A framework forunstructured data analysis, *International Journal of Engineering andTechnology*, 5(1) (2013), pp.153-156.
- [13] T. K. Das, D. P. Acharjya and M. R. Patra, Opinion mining about aproduct by analyzing public tweets in twitter, *International Conferenceon Computer Communication and Informatics*, 2014.
- [14] L. A. Zadeh, Fuzzy sets, *Information and Control*, 8 (1965), pp.338-353.
- [15] Z. Pawlak, Rough sets, *International Journal of Computer InformationScience*, 11 (1982), pp.341-356.
- [16] D. Molodtsov, Soft set theory first results, *Computers and Mathematicswith Applications*, 37(4/5) (1999), pp.19-31.
- [17] J. F.Peters, Near sets. General theory about nearness of objects,*Applied Mathematical Sciences*, 1(53) (2007), pp.2609-2629.
- [18] R. Wille, Formal concept analysis as mathematical theory of conceptand concept hierarchies, *Lecture Notes in Artificial Intelligence*, 3626(2005), pp.1-33.
- [19] I. T.Jolliffe, *Principal Component Analysis*, Springer, New York, 2002.[20] O. Y. Al-Jarrah, P. D. Yoo, S. Muhaidat, G. K. Karagiannidis andK. Taha, Efficient machine learning for big data: A review, *Big DataResearch*, 2(3) (2015), pp.87-93.
- [21] Changwon. Y, Luis. Ramirez and Juan.Liuzzi, Big data analysisusing modern statistical and machine learning methods in medicine,*International Neurourology Journal*, 18 (2014), pp.50-57.
- [22] P. Singh and B. Suri, Quality assessment of data using statistical andmachine learning methods. L. C.Jain, H. S.Behera, J. K.Mandal andD. P.Mohapatra (eds.), *Computational Intelligence in Data Mining*, 2(2014), pp. 89-97.
- [23] A. Jacobs, The pathologies of big data, *Communications of the ACM*,52(8) (2009), pp.36-44.
- [24] H. Zhu, Z. Xu and Y. Huang, Research on the security technology of bigdata information, *International Conference on Information Technologyand Management Innovation*, 2015, pp.1041-1044.
- [25] Z. Hongjun, H. Wenning, H. Dengchao and M. Yuxing, Survey ofresearch on information security in big data, *Congresso da sociedadeBrasileira de Computacao*, 2014, pp.1-6.
- [26] I. Merelli, H. Perez-sanchez, S. Gesing and D. D.Agostino, Managing,analysing, and integrating big data in medical bioinformatics: openproblems and future perspectives, *BioMed Research International*, 2014,(2014), pp.1-13.
- [27] N. Mishra, C. Lin and H. Chang, A cognitive adopted frameworkfor iot big data management and knowledge discovery prospective,*International Journal of Distributed Sensor Networks*, 2015, (2015), pp.1-13
- [28] X. Y.Chen and Z. G.Jin, Research on key technology and applicationsfor internet of things, *Physics Procedia*, 33, (2012), pp. 561-566.
- [29] M. D. Assuno, R. N. Calheiros, S. Bianchi, M. a. S. Netto and R. Buyya, Big data computing and clouds: Trends and future directions, *Journalof Parallel and Distributed Computing*, 79 (2015), pp.3-15.
- [30] I. A. T. Hashem, I. Yaqoob, N. BadrulAnuar, S. Mokhtar, A. Gani andS. Ullah Khan, The rise of big data on cloud computing: Review andopen research issues, *Information Systems*, 47 (2014), pp. 98-115.
- [31] L. Wang and J. Shen, Bioinspired cost-effective access to big data,*International Symposium for Next Generation Infrastructure*, 2013, pp.1-7.
- [32] C. Shi, Y. Shi, Q. Qin and R. Bai Swarm intelligence in big dataanalytics, H. Yin, K. Tang, Y. Gao, F. Klawonn, M. Lee, T. Weise,B. Li and X. Yao (eds.), *Intelligent Data Engineering and AutomatedLearning*, 2013, pp.417-426.
- [33] M. A. Nielsen and I. L.Chuang, *Quantum Computation and QuantumInformation*, Cambridge University Press, New York, USA 2000.
- [34] M. Herland, T. M. Khoshgoftaar and R. Wald, A review of data miningusing big data in health informatics, *Journal of Big Data*, 1(2) (2014),pp. 1-35.
- [35] T. Huang, L. Lan, X. Fang, P. An, J. Min and F. WangPromises andchallenges of big data computing in health sciences, *Big Data Research*,2(1) (2015), pp. 2-11.
- [36] G. Ingersoll, Introducing apache mahout: Scalable, commercial friendlymachine learning for building intelligent applications, *White Paper,IBM Developer Works*, (2009), pp. 1-18.
- [37] H. Li, G. Fox and J. Qiu, Performance model for parallel matrix multiplicationwith dryad: Dataflow graph runtime, *Second InternationalConference on Cloud and Green Computing*, 2012, pp.675-683.
- [38] D. P. Acharjya ,Kausar Ahmed P , Survey on Big Data Analytics: Challenges, Open Research Issues and Tools, *International Journal of Advanced Computer Science and Applications*, Vol. 7(2), 2016

GameilSaadHamzh Ali: MSc Computer Science , Ph.D.
Research Scholar , Department of Computer Science,
RathnavelSubramaniam college of Arts & Science, India

Dr.A. Nithya : MSc ,MPhil ,Ph.D. , Research Supervisor,
Department of Computer Science, RathnavelSubramaniam
collegeof Arts &Science,India