

A Survey on optimising difficulty measures and maximising reliability in Analytics and Big Data Computations on Cloud <https://doi.org/10.56343/STET.116.011.003.002> <http://stetjournals.com>

K.Selvamani* and S.Ramamoorthy

*Dy.Manager-HR-CorpSuper Auto Forge Private Limited, Kolapakkam,Vandalooore, Chennai-127, Tamilnadu, India.

Department of Computer Applications, Dr. MGR University, Maduravoyal, Chennai-600 092, Tamilnadu, India.

Abstract

Currently society has been surrounded with various developing activities in utilizing generated data. Most Popular and demanding concepts of big data analytics, have become more computational than reliable. Study on big data analytics ensures to provide, data management and related frameworks, various business models related to user interactions by identifying available gaps in the technology, provides challenges and open issues using cloud related to big data and analytics. Parallely, cloud computing is a solution to reduce the problems by enabling the services by modifying the system to the exact demand. This paper provides deep understanding of approaches and environments in big data environment related to cloud and workflow of the basic phases of big data analytics and timely state of every stage of cloud related analytics by providing cloud computational trust models.

Key words: Data Mining, Analytics, Cloud computing , Difficulty measures in Big Data, Computational trust models.

Received : December 2017

Revised and Accepted : February 2018

INTRODUCTION

Big Data

Managing of data is very important for computing huge amounts of data. Algorithm and designs depend on the efficient use of map reduce and architecture of hadoop-based structure for analytics of highly time-synchronized datasets. "Big data analytics" the center of the latest frontier of data analysis is described by the two circumstances "1. complex data" and "2. Big Data". The computerized machinery of data has led to the tremendous growth of information, that makes applicable to excessive new sources of data originating from organizations, enterprises, government and all other fields. This leads us to the point of issue, how to examine and process this data to identify valuable information and acknowledge improved decision making (Marcos *et al.*, 2015). The individual various innovative methods form a broad range of conditions of big data analytics classified as adaptive data analysis for big data, computing techniques and computational techniques for big data. Computational intelligence and complex adaptive system.

Data mining

Data mining tools express outcome of future rage and behavior, acknowledging businesses to make

proactive, knowledge-based decisions. Data mining specifications rely on:

Associated connections - Observing for patterns on events relating to each other.

Path analysis- Observing for patterns where one event shows path to another later event.

Classification- Observing for new patterns

Clustering- Discovering and plainly documenting groups

Forecasting- Determining patterns that show high priority to be reasonable predictions relates the future is known as predictive analytics.

Cloud computing

Cloud computing is the effectiveness of transferring a variety of IT services, which are altogether different from each other. This range creates a particular awareness of what cloud computing is between users.

Cloud computing services

Cloud computing services there are three major divisions in cloud computing services namely, :

1. Infrastructure-as-a-service (IaaS), solutions distribute framework in the form of virtual hardware, data storage, and networking and maintains very large potential for flexibility.

*Corresponding Author :
email: selvamani62@gmail.com

2. Platform-as-a-service (PaaS), cloud services offers platforms, various tools, and other services to users, the services offered to manage operating system and middleware with resources to transfer.

3. Software-as-a-service (SaaS) solutions provide services to software applications and are shared among many users.

LITERATURE SURVEY

Data obtained from different sources, like databases, streams, and repositories, are used to develop models. High volumes and different structure of the data needs pre-processing tasks for integrating data, removal of noise and purifying it. The processed data is used to train typical methods and to find its specifications. Despite decision makers would like to groundwork their foundational decisions and actions on awareness obtained from this data (Davenport and Harris, 2007) making the perception of data, deriving unclear patterns, and making use of those patterns to predict future performance.

Knowledge discovery in data (KDD) Fayyad *et al.*, 1996 targets to derive unclear, not accessible information using protected and very broad analysis and explanation. Data mining (Written *et al.*, 2010) more particularly, targets to explore earlier unidentified communications between irrelevant attributes on data sets by performing designed approaches from various areas containing machine learning concepts, database systems for storage, and statistics for implementing and analysis.

Analytics includes approaches of KDD, data mining techniques, text mining methods, statistical approaches and measurable analysis, descriptive and predicting designs, and progressive and collective determination initiative decisions and processes (Davenport *et al.*, 2010; King, 2015). Analytics results can be categorized as descriptive and predictive. Descriptive analytics deals with the historical data to recognize patterns and develop management information in the description, it mix up with past behavior of designing. Predictive analytics try to foresee the expected data by analyzing existing and historical data (Davenport and Harris, 2010). Operating analytics on large quantity of data needs prominent approaches to store, filter, reconstruct and extract the data. Most of the demands of expanding data management results on cloud environments have been acknowledged approximately (Wikipedia, 2017; Grossman, 2009; Attention, 2013) and solutions to accomplish data analysis on the cloud. High volumes and distinct types of the data requires pre-processing schedules for integration of data, avoiding noise, and purifying it. The processed data is used to train a

method and to evaluate its parameters. Once the method is evaluated, it will be authenticated before its utilization. Generally, this period demands the use of the standard input data and definite approaches to confirm the designed model. Finally, the design is expanded and applicable for data as it reaches. This stage is known as model scoring, it develops predictions, instructions, and suggestions. The obtained solutions are made clear and understood to evaluate, that results in achieving new models. Analytics solutions are required to recognize the various collective cloud deployment approaches maintained by resourceful enterprises, where clouds are particular. Private extend on a private network and are trained by the organization alone or by observers.

Regardless of demanding analytics and big data, practical implementations of them are complicated and time exhausting enterprise. As (Balmin *et al.*, 2013) describes abundant data computing infrastructure; and yield for consultative timings of analysts who deal with the organization to get proper recognition of its business, storage, and with analytics (Barga *et al.*, 2012). This effort focuses on technical problems and analyse current existing effort on providing solutions to generate analytics effectiveness for big data on the cloud. Assuming the traditional analytics progress, are main issues in the development of an analytics solution. Working with big data deals with many of the challenges of cloud analytics claims to managing data, data integration, and processing of data. Recent work concentrates on the concerns like formats of data, representation of data, a data repository for storage, accessibility of data, privacy in providing data protection, and quality of data.

This paper concentrates on the important and most challenging growth of data in institutions and organizations that come across to enhance the quality of administering resolutions and reliability in performance. Current work deals with the challenges on cloud environments by providing models and specifying data models on the cloud. Security is absolutely a key challenge for maintaining analytics solutions on public clouds. Security and interpretation of data accuracy and exactness (Yu *et al.*, 2013) are significant outlook of this survey. This also describes the functioning of the vm player in the cloud making use of big data and the different circumstances in computational measures.

This paper also deals with services offered by the cloud between users and cloud service providers by establishing trust factor, various cloud computational trust models providing key management technologies based on encryption which seems to be more reliable (Pearson and Benameur, 2010).

Models between the Cloud and Big Data

Frequently and most available models for better enhancement of big data analytics to generate solutions on cloud services is PaaS and SaaS. IaaS is mostly not used for data that resides in high-level analytics applications but essentially for handling data for storage and needs for data computations, Cloud computing approaches helps in increasing the possible scalable analytics solutions (Sun *et al.*, 2011). Cloud computing is a part of distributed computing family that assures assets in the form of end user services for instance (SaaS), framework like (IaaS) and a platform services like (PaaS). The cloud computing model continuously provoke big database service along with (AaaS, BDaaS) named as (DaaS) "Database as a Service" means applications easily deployed to have database services in any applicational environment (Wang, 2012). BDaaS is a version to provide service

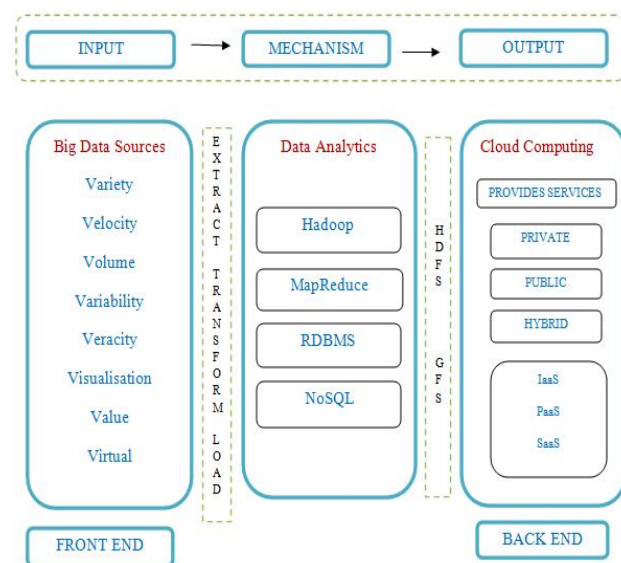


Fig.1. Connectivity among big data and cloud computing.

identical to SaaS or IaaS. Huge data as a service mostly depend on cloud repository to manage continues connected data access to the enterprise the information and considers that to host in the cloud (Reda *et al.*, 2013).

Figure 1 shows the input of the big data sources as the front end by extracting, transforming and loading of data and output services provided by the cloud, it provides services like IaaS, SaaS, PaaS on the back end related to the public, private and Hybrid basis. The entire process of mechanism deals with the Data Analytics like Hadoop, MapReduce, RDBMS, NoSQL.

Virtual Machine between the cloud and big data

Virtual Machine is a software application that resembles an environment for virtual computing to run the operating system and correlates applications with different distinct or similar virtual machines installed in a single machine. Distributed systems, network computing and parallel programming are old versions for prominent enabling components of the cloud in the virtual technology. Use of virtual technology, single virtual machine can also host many number of virtual machines (Ahson *et al.*, 2010). This technology supports the capacity to reduce the workload on devices and also consolidate them in single physical server.

Virtual technology is one of the most user friendly platform for big data plus traditional applications. Considering the applications in big data, it reduces and manages big data framework, generating quick results that is cost-effective (Techtarget, 2012). Current virtual data establishes a wide dimension of sources comprising multi dimensional storages, like web and various data services, XML specific documents, analytical devices, and applications from outdoor as well as indoor. NoSQL plays a role in data storage for present day sources that holds the virtual data (Vmware, 2012).

Data inside the technology world are designed in terms of letters, words, numbers, symbols or images, but with the improvement of tools in multitasking technology, a possibility of data variation in content and source (Zomaya and Sakr, 2017). There exist lots of variations on traditional data and big data in terms of volumes, data generation rate, type of data, sources provides the data and storage of data. Traditional data generated MB and GB volumes of data taking long periods of time in centralized data resources and are stored only in RDBMS because the data is purely structured. Whereas big data provides PBs and ZBs volumes of data in the more rapid rate of generation on multiple resources in distributed environments and are stored in HDFS, NoSQL for semi-structured and unstructured data.

Big Data Computational difficulty measures

There are many interesting choices in mining-related to big data. Several challenges are kept in front while analyzing big data sets. Difficulties at different stages involve: Attaining of data, storage, searching of data, sharing of data, analysis of data, data management, and data visualization. Irrespective of that there are security measures and privacy issues while distribution of data over applications (Rahman *et al.*, 2011). Truth is that the size of data gradually keeps on growing and the capacity to handle the data and explore beyond zettabytes is highly complex. Following are the some technological issues

Big Data management

Working in big data, data scientists are facing real-time challenges. Firstly comes when a collection of data, integrating and cache. Data sets generated from various sources deal with the support of hardware and software requirements.

It is very necessary and important to manage bigdata effectively for extraction of reliable understanding and to optimize overheads. Big data management is to provide clean data for reliability and to accumulate data coming from disparate sources by ensuring security and privacy in encoding data.

Big Data cleaning

Traditional data management works with these steps but when comes to big data it concentrates on managing complexity with respect to cleaning related to nature of data (velocity, volume, and variety) and process in combination of different applications. On the other part, incomplete data has noises, errors (Lopez and Xavier, 2012). The challenge is to remove those noises, errors and to mention how much data is reliable and useful.

Big Data aggregation

Applications, networks, warehouses distributed on different platforms is another challenge to make them synchronize. Aggregating internal data with external resources includes third parties. This leads to maximization of predictive models in analytics.

Imbalanced systems capacities

The Primary considerable issue is the capacity related to CPU performance. With an increase of performance for every 18months by Moore's law, the performance of disk also increasing at the same rate and also I/O operations with different performance. Imbalance system capacities bring to slow access and that in turn leads to effects in performance and scalability. This will slow down performance of the system.

Imbalanced Big Data

Another technological issue comes when we classify imbalanced dataset. Real-world applications produce classes with individual distinct distributions. In particular, many individual domains have more uneven distributions. Dealing with multi-class tasks with various misclassification.

Finding a way to solve this problem scientists have introduced a binary classification for two class and multi-class classification for more than two instances, another way is decomposition and ensemble methods (DEM). It follows a procedure called decomposition, it decomposes multi-class classification problem into

binary classification and predictions based on aggregative strategies.

Cloud Computing Trust Mechanisms

Trust mechanisms are the best way to provide security to the system, it provides reliable policies to security and access control (Subhulakshmi *et al.*, 2016). This develops highly reliable computing resources to end users as well as servers dynamically. Following are some of the trust mechanisms for cloud computing:

Reputation-Based trust

These are based on the direct connections and it is very useful to the cloud users in opting cloud services without any specific requirement.

SLA Based trust

Service level agreement provides services by maintaining legal contacts between cloud user and cloud services providers. although it provides most of the services to CSP, it lacks for stakeholders. Visible elements are provided services and invisible elements are kept secure. they do not provide enough efficiency to perform SLA.

Domain Based trust

It mainly depends on the transactions among entities. It deals with the validation mechanism that handles every domain. When a single domain is authenticated by a unity, then it will be accepted by all other domains because of authentication.

Platform-Based trust models

It contains policies for establishing platforms for applications. It provides a particular trust assurance step and estimates using cloud services.

Authentication based trust models

These models are based on key encryption techniques maintaining data with confidentiality, integrity, and availability. Cloud provides services trusted virtual environment module, mutual trust-based access control, grid and cloud trust model, hierarchical attribute set based encryption, trusted platform software stack and improved trusted cloud computing platform.

CONCLUSION

Big data analytics is very much challenging and time demanding job that desires for valuable software, high computational infrastructure and tasks being used by industries, organizations, gaining an advantage over their opponents. It is a platform that is supported by various available processing analytical tools. It is a part that acts as a part of cloud computing providing great solutions to problems by specifying resources on

demand with costs corresponding to the data usage. This is not achieved by the constraints of inconsistent processing time among the run time of the analytics. This paper provides the various circumstances of analytics, big data analytics, and cloud computing trust mechanisms. Various computational difficulty measures in big data to achieve reliable data and computational trust mechanisms for ensuring security using key management technologies and the connection between cloud users and cloud service providers. This ensures that the data can be made more reliable using computing trust models by overcoming the difficulty measures.

REFERENCES

- Ahson., Syed, A. and ILYAS. 2010. Mohammad (ed.) *Cloud computing and software services: theory and techniques*. CRC Press. PMID:20134077
<https://doi.org/10.1201/EBK1439803158>
- Attention, 2013. shoppers:Store is tracking your cell, New York Times. URL <http://www.nytimes.com/2013/07/15/business/attention-shopper-stores-are-tracking-your-cell.html>.
- Balmin, A., Beyer, K., Ercegovic, V., Ozen, J.M.F., Pirahesh, H. and Shekita, E. 2013. A platform for extreme Analytics, *IBM J. Res. Dev.* 57P.3-4.
<https://doi.org/10.1147/JRD.2013.2242693>
- Barga, R.S., Ekanayake, J. and Lu, W. 2012. Project Daytona: Data Analytics as a cloud service, In: A.Kementsietsidis, M.A.V. Salles, *Proceedings of the International Conference of Data Engineering (ICDE 2012)*, IEEE Computer Society, P.1317-1320.
<https://doi.org/10.1109/ICDE.2012.136>
- Davenport, T.H. and Harris, J.G.2007. *Competing on Analytics: The New Science of Winning*, Harvard Business review Press.
- Davenport, T.H., Harris, J.G. and Morison, R. 2010. *Analytics at work: Smarter Decisions, Better Results*, Harvard Business Reviews Press.
- Grossman, R.L. 2009. what is analytic infrastructure and why should you care? *ACMSIG KDD Explorations Newsletter* 11 (1) : P. 5-9.
<https://doi.org/10.1145/1656274.1656277>
- King, E.A. 2015. How to buy data mining: A framework for avoiding costly project pitfalls in predictive analytics, *DMReview* 15(10).
- Lopez and Xavier, 2012. Big Data and advanced spatial analytics" In : *Proceedings of the 3rd international Conference on Computing for Geospatial research and Applications, P.5*, ACM.
<https://doi.org/10.1145/2345316.2345322>
- Marcos, D., Assuncao, Rodrigo N. Calheiros, Silvia Bianchi, Macro Netto, A.S.and RajKumarBuyya. 2015. Big Data computing and clouds: Trends and future directions *J. Parallel Distrib. Comput.* P.79-80.
<https://doi.org/10.1016/j.jpdc.2014.08.003>
- Pearson, S. and Benameur, A. 2010. Privacy, security and trust issues arising from cloud computing. In : *Cloud computing technology and Science,IEEE Second International Conference* 693-702.
<https://doi.org/10.1109/CloudCom.2010.66>
- Rahman, M., Li, X. and Palit, H. 2011. Hybrid Heuristic for Scheduling Data Analytics workflow Applications in hybrid Cloud Environment, In: *Proceedings of the IEEE International Symposium on Parallel and Distributed Processing Workshops and Phd Forum(IPDPSW)*, P.966-974.
<https://doi.org/10.1109/IPDPS.2011.243>
- Reda, K., Febretti, A., Knoli, A., Aurisano, J., Leigh, J.,Johnson, A., Papka,M. and Hereid, M.E. 2013.Visualizing Large, Heterogeneous Data inHybrid-reality Environments, *IEEE computer graphics and Applications* 33(4):38-48, PMID:24808058
<https://doi.org/10.1109/MCG.2013.37>
- Subhulakshmi, R., Suryagandhi, S., Mathubala, R. and Sumathi, P. 2016. An evaluation on cloud Computing Research Challenges and its Novel Tools, *International Journal of Advanced research in Basic Engineering Science and technology(IJARBEST)* volume 2, Special Issue 19.
- Sun, X., Gan, B. and Zhang, Y. 2011. An, H.Cao,C.Guo, towards delivering analytical solutions in cloud: Businee models and technical challenges, In : *Proceedings of the IEEE 8th international Conference on e-Business Engineering(ICEBE 2011)*, IEEE Computer Society, Washington, USA,P.347-351.
<https://doi.org/10.1109/ICEBE.2011.81>
- TechTarget, 2012. <http://searchio.techtarget.com/definition/big-data-as-a-service-bdas>.
- Vmware, 2012. <https://www.vmware.com/asean/solutions/big-data.html>.
- Wang, H. 2012. Integrity Verification of Cloud_Hosted Data Analytics computations In : *Proceedings of the 1st International workshop on cloud Intelligence(cloud-I 2012)*, ACM, New York, USA.
<https://doi.org/10.1145/2347673.2347678>
- Wikipedia contributors., 4 Aug 2017. "Virtual Machine". Wikipedia, *The Free Encyclopedia*. wikipedia,web. 15.
- Written,I.H., Frank,E. and Hall, M.A. 2010. *DataMining: Practica Machine Learning Tools and Techniques*, third ed., Morgan Kaufmann.
- Yu,P.S. 2013. On mining big data in J.Wang, H.Xiong, Y. Ishikawa, J.Xu, J. Zhou, web. Age Information management, In:*Lecture notes in Computer Science*, vol. 7923, Springer-Verlag, Berlin, Heidelberg.
- Zomaya, A.Y., and Sakr, S.2017. *Handbook of Big Data Technologies*, Springer.
<https://doi.org/10.1007/978-3-319-49340-4>