# An Overview of Milestones of Big Data Analytics in Clinical and Medical Analysis

**Manu M R, B Balamurugan**

*Abstract: The technological advancements make changes during availability of knowledge in a huge way. As the volume of data is increasing exponentially, there is a need for better management of data to research and industry. This data, referred to as Big Data, is now employed by various organizations to extract valuable information which may reanalyzed computationally to reveal patterns, trends and associations revealing the human interaction and behavior for making various industrial decisions But the data must be optimized, integrated, secured and visualized to make any effective decision. Analyzing of the large volume of data is not beneficial always unless it is analyzed properly. The existing techniques are insufficient to analyze the large Data and identify the frequent services accessed by the cloud users. Various services can be integrated to provide a better environment to work in emergency cases pretty earlier. Using these services, people become widely vulnerable to exposure. The data is large and provides an insight in to future predictions, which could definitely prevent maximum medical cases from happening. But without big data analytics techniques and therefore the Hadoop cluster, this data remains useless. Through this paper, we'll explain how real time data may be useful to research and predict severe*

*Index Terms - Big Data, Medical Big Data mining, EHR, EMR, Hadoop, Hbase, Mapreduce, HDFS, Hive, Cassandra*

*Keywords: The Existing Techniques Are Insufficient To Analyze The Large Data And Identify The Frequent Services Accessed By The Cloud Users.*

## I. INTRODUCTION

Big Data has gained a huge attention past few years, ever since the data is digitized. The information available over the distributed architecture has evolved instant accessibility across the globe. This helps in providing evidence based medicine to the patients instantly, using live analysis to simulate spreading or viral diseases, providing proactive care to the patients and reducing readmission rates to the hospital. Big Data, a term that refers to dataset that is so huge that performing efficacious analysis using the traditional methods becomes very difficult. Analysis of this data uses newly researched technologies and distributed architecture that makes extraction of value from the dataset possible. Its properties are pointed out by 3V's, namely volume, variety and velocity. Recent studies point out that big data cannot be fully defined by 3V's and thus veracity, valence and value were added to complement its explanation. The data is crucial at it is complex to analysis without any bias The challenges of Big Data involves captures, storages ,transfer and analysis of data both in visualized as well as analytical behaviour. The Big Data analytics also provides business intelligence, statistical analysis, statistical forecasting, optimization, text analysis and predictive modelling. The traditional data analysis involves only centralized data base architecture which will be complex and not economical. The below table shows the properties of Big Data as an overview .

**Table 1: Properties of Big data**

| Volume | Variety | Velocity | Veracity | Valence | Value |
|---|---|---|---|---|---|
| Big data implies tremendous amount of data which is generated every single second of the day in our digitized world | Data can come in ever increasing different forms such as audio, video, image, text and what not. These non-ending forms contribute to this property of the data. | The speed of data generation and the speed by which data is transported from one port to another port. The data generation speed is much more than the data storage speed which again is very huge than speed at which insights are gained from that data. | It refers to quality factor of the data which varies greatly. It caters to the question that whether the data which is being stored and mined is useful to the real life problem that is being dealt with. | This points to how well can big data bond with each other which forms connections between otherwise disparate datasets. | The data incurred should possess real-time value and generate insights so as to sustain the changing trends of input, output, processing and analysis of data. |

**Manu M R\*,** Research Scholar, School of Computing Science and Engineering, Galgotias University, Greater Noida (U.P), India.
**Dr. B Balamurugan,** Associate Professor, School of Computing Science and Engineering, Galgotias University, Greater Noida (U.P), India.

## II. RELATED WORKS

### A) Software Present and Their Issues:

Different software [3] are used for various functions of Big Data analytics. Some related software are:

1) Virtualization: Xen, Oracle Virtual Box, KVM,VMware are some of the examples but they take time to start as they are heavyweight. Docker is there for Linux. It is lightweight and more efficient to interact with the host.

2) Big Data Storage: It is more convenient to store huge amount of data on the cloud rather than on disk. Hadoop Distributed File System (HDFS) provides high throughput data access to data. Some new NoSQL databases such as HBase, Cassandra, MongoDB are scalable in dealing with big semi structured/unstructured data.

3) Big Data Computation Model: Apache Hadoop written in Java is an open-source software framework. It

provides MapReduce parallel processing framework and job scheduling and resource management. Other software packages which are incorporated into Hadoop ecosystem are HBase, Hive, Pig etc. Apache Spark and Flink increased the performance by caching data in memory.

4) Big Data Analytics User Interface: Since most of the data analyst are not a computer scientist, User Interface plays a key role. Data Bricks, which is built on Spark, provides a web based cloud computing platform. Cloud Flows is an open source web-based platform

### B) Frame Works

Different frameworks consist of a set of systems which are deployed over multiple parallel nodes. They allow huge computations on reduced infrastructure. Apache Hadoop is the most known tool for processing of Big Data. Many open source components developed by Apache forms a Hadoop ecosystem.

Different techniques [4] for Big Data Analysis are:

1) HADOOP: Open-source software for Big Data. Normally, it deals with a very big amount of distributed

heterogeneous data. It offers storage and computation. It stores data using MapReduce technology.

2) MAPREDUCE: A Google designed technology for processing large data. It has two main parts: Map and Reduce. Map function distributes the data into several clusters for parallel processing and Reduce function clubs all the cluster of the same type to one final result.

3) HDFS: It is the core of Hadoop. It stores and manages data of large files. It splits data into blocks and then allocates it on servers in different locations.

4) Hive: A data warehouse tool that allows managing and requesting distributed data. It uses the SQL-like language called HiveQL.

5) HBase: It is a Hadoop Database. It is inspired by Google's Big Table. It manages and processes big tables in an efficient way.

6) CASSANDRA: A Facebook developed tool which is a column-oriented NoSQL database. It supports MapReduce and allows access of data of large volumes.

### C) Integration of Tools

Various tools can be integrated to get a platform which is more efficient to work on. This will redefine the workflow. We can also visualize the data and make our system light-weight.

They are described as follows:

1) New Data Analytics Flow: SPARK Notebook combines Scala code, SQL ,Mark-up or even JavaScript in a collaborative manner. It is good to verify algorithms but it is notgood for code reusing and sharing. Three main components are

Notebook, Widgets and Workflow. Various steps are:

i) Verify the new idea with Notebook which may include several iterations.

ii) Integrate the code into reusable widget through widget component in Web IDE in which input and output need to be defined by a template.

iii )Different widgets could be combined to form complicated applications.

By this, applications could be reused and shared with others, and the size of widget makes the maintenance easier.

## III. PROPOSED SYSTEM

### a) An overview of Medical Big Data mining

The data is most crucial part in the day to day life. The capturing of right data and its storage is an challenging task. The processing and storage capability with high accuracy of variability of data and satisfying the above challenges are accomplished by big datamining .Big datamining made a milestone in health care industry ,which in turns the origin of Medical big datamining. The process of uncover hidden pattern in medical records from a medical repository and transform it for discovering or diagnosis the disease is called as medical big data Ming.

The medical Big data mining widely used for recognizing and discovering of the hidden pattern in Medical records for easily diagnosis of the diseases. The medical data mining mainly processing in three stages data pre-processing, datamining and result evaluation. The data pre-processing involves the processing of raw data and information from medical records and its given for datamining process which involves different technique for synthesis of pre-processed data with algorithms and the final result can be evaluated with different parametric methodology. The data pre-processing is preliminary analysis which directly affects the accuracy and efficiency of datamining. For recognizing of various disease like cardiologic or malignant types can be deeply analysed using different technique such as classification and association and implemented with help of artificial neural networks and genetic algorithmic schemes.

The datamining techniques now effectively organized in industry as well as ecommerce. The medical data used in the datamining processes to evaluate and to be authentic .The medical data in the senses the data got form medical database and to be refined for manipulation purpose. The data authenticity and fraud abuse analysis can be effectively made by the datamining process and which intense make a effectiveness in the diagnosis and treatment of the diseases. The metabolic disease that is a common occurrence in the current generation can be rectified and analyses with the machine learning technique implemented with the help of trained data test by the datamining process. The medical datamining process can be diagrammatically represented below in which each stage is important in the processing ,from the initial data collection from the medical data repository to the final result.

The medical data collection stage represents the segmented data and analysed data from the repository and its given for pre-processing stage which includes different algorithmic steps and give for datamining processing with machine learning technique implementation. In this stage an accurate analysis with the help of trained data set can be formed and given for analysis phase .In this step it analysed with pre-trained data and an intermediate result is executed and at last with final processing driven to final results
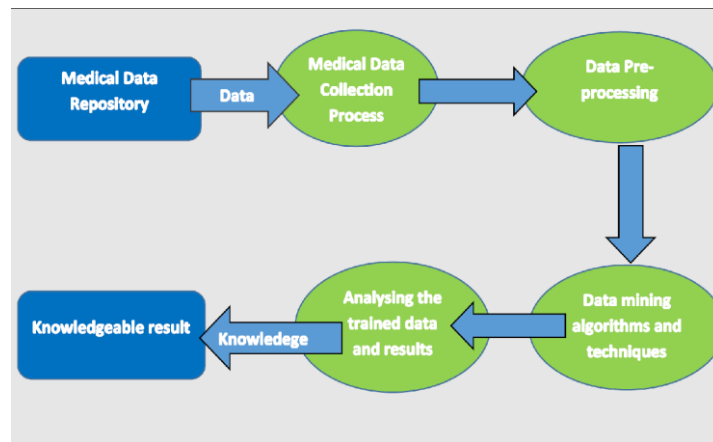


**Figure 1: Medical Big datamining process**

### b) Evolution of Electronic Medical record

The EHR-Electronic health record derived for patient centric treatment and outcome is a structured data collection automated by valuable mining tools .The EHR currently play a big role in medical filed which digitalized all the patient information and records suitable for variant diseases. The evolution of ERH derived for Federal government based on incentive programs.The medical records written in Late 1600 BC were Egyptian mnemonics which were based of 48 medical cases which includes injury, fracture and different tumor details which were handwritten and not in a properly formatted, Later on 1862 the papyrus text acquired by Edwin smith came for the medical record preservation. The legalization of medical records for insurance derived on 1880 and easily identified the malpractice of the records. At the end of 1898 patients record moving from retrospective document to case report. The end of 1960 the documents are structured in digital format which in turns the formation of EHR which called as Problem Oriented Medical Record (POMR).The record contain physical examination ,laboratory data and complex problem list which also includes the discharge summary record. In 1971 Locked Cooperation created a company 'Eclipsys Cooperation 'which computerized all the physical records which allowing the ordering of whole medical records includes the X-ray and all laser records .EHR now known to be the veteran health information and technology which is an electronic system that record and transfer the patient information and clinical operation .
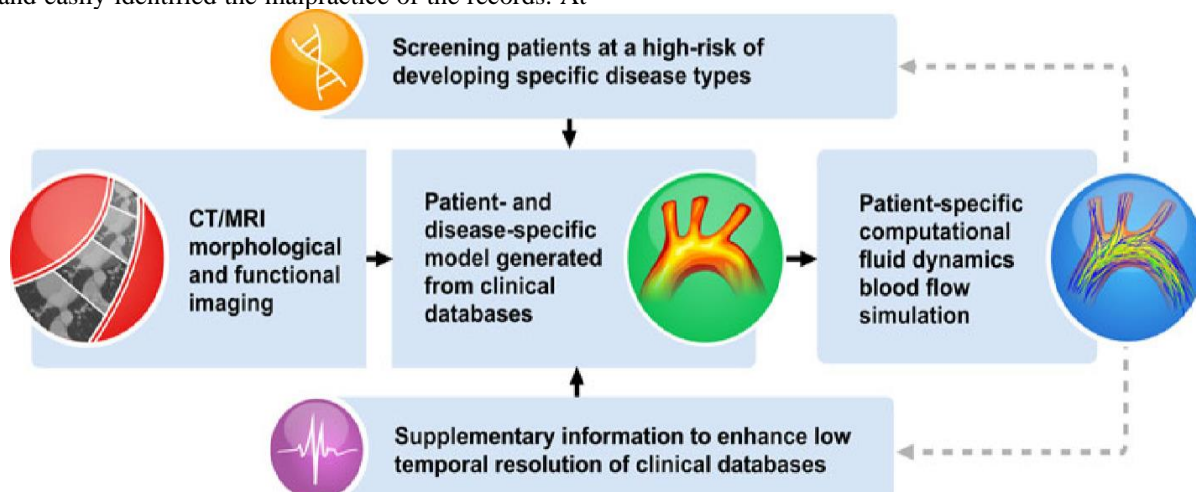


**Figure 2: Integration of imaging, modelling, and real-time sensing for the management of disease progression and planning of intervention procedures. This example of thoracic aortic dissection illustrates how risk stratification and subject-specific haemodynamic modelling substantiated with long-term continuous monitoring are used to guide the clinical decision process.**

### c) Differentiation of Electronic Health Care Record and Electronic Medical Record

The Electronic health care records and Electronic records differ in their word itself "health "and "medical" .One concerning with the life routine and healthcare information where as other deeply refer to the medical diagnosis and medical historic over view of the patient. Traditionally electronic medical record is a digital way of maintain g the patient medical record which is prominently used for tracking and historical background survey record. The same digital version is also used in Electronic healthcare records also but when concerning with EMR it have entire patient historical data which includes stream line sharing feature of real-time information.
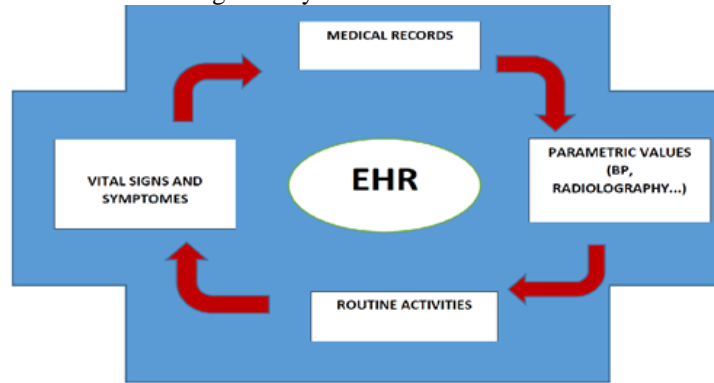


**Figure 3: Electronic Health Record (EHR) Process diagram**

The Electronic healthcare record also provides authorization and authenticity of healthcare feature while sharing. The Electronic medical record is not dealt with the sharing feature but it is mainly includes single patient centric digital data records. The Electronic healthcare record data mainly featured with life style analysis with primary set data. A better follows of pre-stored healthcare information is maintained in Electronic health care record. The Electronic medical record is a parametric oriented such as quantitative information.

The EMR and EHR both made the medical and healthcare industry a tremendous changes and effected with fastness and accuracy of test analysis and diagnosis of disease. The both will avoid the duplication of record and time saving for recognizing disease.
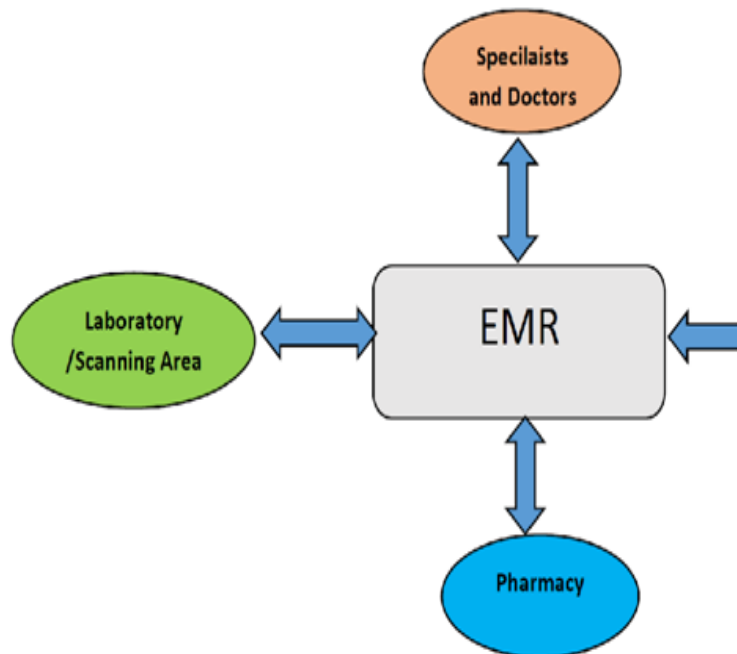


**Figure 4: Electronic Medical Record (EHR) System diagram**

### d) Data Mining Algorithms and Techniques in health care

Data Mining in medicine is an emerging field of great importance to provide a prognosis and deeper understanding of disease classification, specifically in Mental Health areas. The main objective of this paper is to present a review of the existing research works in the literature, referring to the techniques and algorithms of Data Mining in Mental Health, specifically in the most prevalent diseases such as: Dementia, Alzheimer, Schizophrenia and Depression.

419

Academic databases that were used to perform the searches are Google Scholar, IEEE Xplore, PubMed, Science Direct, Scopus and Web of Science, taking into account as date of publication the last 10 years, from 2008 to the present. Several search criteria were established such as 'techniques' AND 'Data Mining' AND'Mental Health','algorithms' AND 'Data Mining' AND 'dementia' AND 'schizophrenia' AND 'depression', etc. selecting the papers of greatest interest. A total of 211 articles were found related to techniques and algorithms of Data Mining applied to the main Mental Health diseases. 72 articles have been identified as relevant works of which 32% are Alzheimer's, 22% dementia, 24% depression, 14% schizophrenia and 8%bipolar disorders. Many of the papers show the prediction of risk factors in these diseases. From the review of the research articles analysed, it can be said that use of

Data Mining techniques applied to diseases such as dementia, schizophrenia, depression, etc. can be of great help to the clinical decision, diagnosis prediction and improve the patient's quality of life.

**e) Main techniques and algorithms of datamining used in the review**

The Data Mining techniques have recently become a predominant field of research with wide applications in medical

healthcare, financial services, telecommunications, natural sciences, etc. It is a process to discover useful models in data, with the aim of interpreting existing behaviours or predicting future results
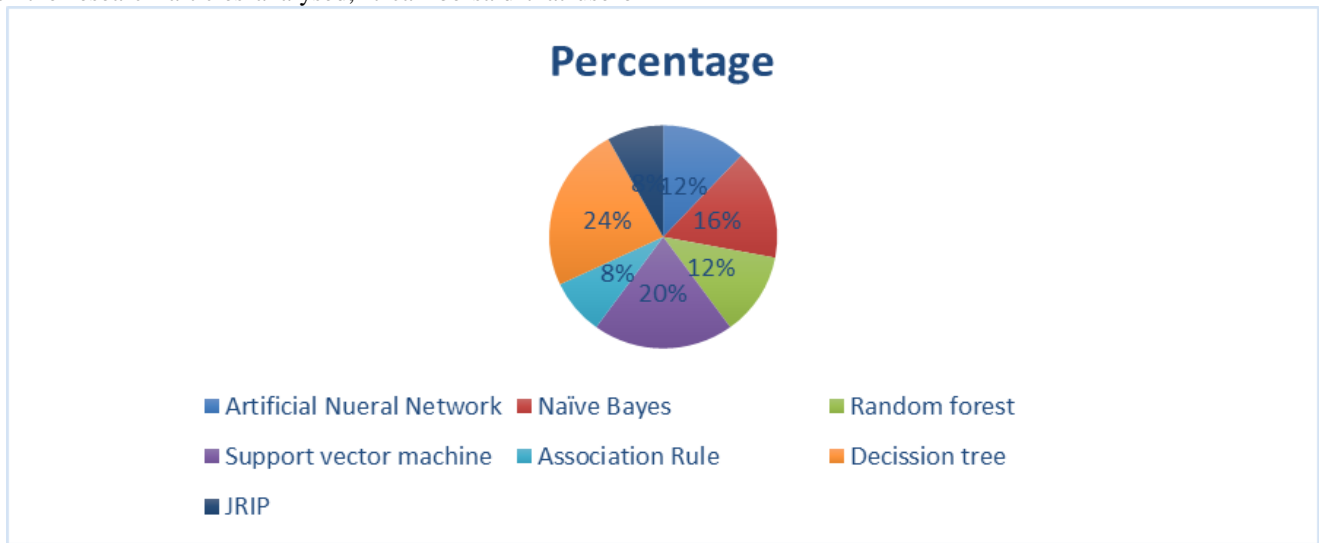


**Figure 5: Percentage of data mining technique and algorithm applied for Alzheimer's disease**

An Efficient Intelligent Traffic System is proposed in this paper to reduce the traffic congestion problem before it happens. This system will decrease the traffic queue size and provides an alternative route for the vehicles to avoid traffic and to achieve free flow of vehicles. An intelligent traffic light control system is deployed to prevent the traffic congestion before it occurs and based on an alert signal traffic route will be deviated. This can help the travelers to have free flow traffic. Traffic jam can be avoided and the condition of the traffic flows in many of the metropolitan cities can be improved. The Euler's approach used to convert map to graph was tested on a metropolitan city graph and the results are found to be satisfactory. Finally the overall framework is statistically proven to be better than the related traffic congestion models.

## IV. CONCLUSION

As the volume data is increasing exponentially, the need for analysing the data is also becoming important. This can helping taking future decisions. The data should be secured so that the industries can rely on them. Better tools will increase the number of users. Visualization plays an important part to study the data. Implementation all these features may lead to increase in costing factor which needs to be resolved. The medical innovation in Big data analytics and processing made a drastic change in clinical field. The

future enhancement involves in clinical oncology prediction with big data analytics techniques.

## REFERENCES

1. Rayan DasoriyaA Review of Big Data Analytics over Cloud in 2017 IEEE International Conference on Consumer Electronics-Asia (ICCE-Asia)Ali, S. S. M., George, B., Vanajakshi, L., & Venkatraman, J. (2012). A multiple inductive loop vehicle detection system for heterogeneous and lane-less traffic. *IEEE Transactions on Instrumentation and Measurement*, *61*(5), 1353-1360.
2. Ali, S. S. M., George, B., & Vanajakshi, L. (2013). An efficient multiple-loop sensor configuration applicable for undisciplined traffic. *IEEE Transactions on Intelligent Transportation Systems*, *14*(3), 1151-1161.
3. Botta, A., De Donato, W., Persico, V., & Pescapé, A. (2016). Integration of cloud computing and internet of things: a survey. *Future Generation Computer Systems*, *56*, 684-700.
4. Chen, C., Petty, K., Skabardonis, A., Varaiya, P., & Jia, Z. (2001). Freeway performance measurement system: mining loop detector data. *Transportation Research Record: Journal of the Transportation Research Board*, (1748), 96-102..
5. Cai, Y., Zhang, W., & Wang, H. (2010, March). Measurement of vehicle queue length based on video processing in intelligent traffic signal control system. In *Measuring Technology and Mechatronics Automation (ICMTMA), 2010 International Conference on* (Vol. 2, pp. 615-618). IEEE.

6. Cheung, S., Coleri, S., Dundar, B., Ganesh, S., Tan, C. W., & Varaiya, P. (2005). Traffic measurement and vehicle classification with single magnetic sensor. *Transportation research record: journal of the transportation research board*, (1917), 173-181.

7. Chiu, S., & Chand, S. (1993, December). Self-organizing traffic control via fuzzy logic. In *Decision and Control, 1993., Proceedings of the 32nd IEEE Conference on* (pp. 1897-1902). IEEE.Coifman, B. (2001).

8. Improved velocity estimation using single loop detectors. *Transportation Research Part A: Policy and Practice*, *35*(10), 863-880.Coifman, B. (2002). Estimating travel times and vehicle trajectories on freeways using dual loop detectors. *Transportation Research Part A: Policy and Practice*, *36*(4), 351-364.Coifman, B., Beymer, D., McLauchlan, P., & Malik, J. (1998). A real-time computer vision system for vehicle tracking and traffic surveillance. *Transportation Research Part C: Emerging Technologies*, *6*(4), 271-288.

9. De Lima, G. R. T., Silva, J. D. S., & Saotome, O. (2010). Vehicle inductive signatures recognition using a Madaline neural network. *Neural Computing and Applications*, *19*(3), 421-436. Guerrero-Ibáñez, J., Zeadally, S., & Contreras-Castillo, J. (2018). Sensor technologies for intelligent transportation systems. *Sensors*, *18*(4), 1212.

10. Guerrero-Ibanez, A., Contreras-Castillo, J., Buenrostro, R., Marti, A. B., & Muñoz, A. R. (2010, June). A policy-based multi-agent management approach for intelligent traffic-light control. In *Intelligent Vehicles Symposium (IV), 2010 IEEE*(pp. 694-699). IEEE.

11. Hartenstein, H., & Laberteaux, K. (Eds.). (2009). *VANET: vehicular applications and inter-networking technologies* (Vol. 1). John Wiley & Sons. [Online]. Available: http://www.nhtsa.gov/ He, W., Yan, G., & Da Xu, L. (2014). Developing vehicular data cloud services in the IoT environment. *IEEE Transactions on Industrial Informatics*, *10*(2), 1587-1595.

12. Houbraken, M., Logghe, S., Schreuder, M., Audenaert, P., Colle, D., & Pickavet, M. (2017). Automated Incident Detection Using Real-Time Floating Car Data. *Journal of Advanced Transportation*, *2017*. https://ccl.northwestern.edu/netlogo/download.shtml

13. Huang, D. Y., Chen, C. H., Hu, W. C., Yi, S. C., & Lin, Y. F. (2012). Feature-based vehicle flow analysis and measurement for a real-time traffic surveillance system. *Journal of Information Hiding and Multimedia Signal Processing*, *3*(3), 279-294. Hunter, T., Herring, R., Abbeel, P., & Bayen, A. (2009). Path and travel time inference from GPS probe vehicle data. *NIPS Analyzing Networks and Learning with Graphs*, *12*(1).

14. Iscaro, G., & Nakamiti, G. (2013, February). A supervisor agent for urban traffic monitoring. In *Cognitive Methods in Situation Awareness and Decision Support (CogSIMA), 2013 IEEE International Multi-Disciplinary Conference on* (pp. 167-170). IEEE.

15. Jia, Z., Chen, C., Coifman, B., & Varaiya, P. (2001). The PeMS algorithms for accurate, real-time estimates of g-factors and speeds from single-loop detectors. In *Intelligent Transportation Systems, 2001. Proceedings. 2001 IEEE* (pp. 536-541). IEEE.

16. Jost, J. (2013). *Compact Riemann surfaces: an introduction to contemporary mathematics*. Springer Science & Business Media.Kaewkamnerd, S., Chinrungrueng, J., Pongthornseri, R., & Dumnin, S. (2010, June). Vehicle classification based on magnetic sensor signal. In *Information and Automation (ICIA), 2010 IEEE International Conference on* (pp. 935-939). IEEE. Ki, Y. K., & Baik, D. K. (2006). Vehicle-classification algorithm for single-loop detectors using neural networks. *IEEE Transactions on Vehicular Technology*, *55*(6), 1704-1711.

17. Lamas-Seco, J. J., Castro, P. M., Dapena, A., & Vazquez-Araujo, F. J. (2015). Vehicle classification using the discrete fourier transform with traffic inductive sensors. *Sensors*, *15*(10), 27201-27214.Liu, H. X., He, X., & Recker, W. (2007). Estimation of the time-dependency of values of travel time and its reliability from loop detector data. *Transportation Research Part B: Methodological*, *41*(4), 448-461.Lu, X. Y., Varaiya, P., Horowitz, R., Guo, Z., & Palen, J. (2012). Estimating traffic speed with single inductive loop event data. *Transportation Research Record: Journal of the Transportation Research Board*, (2308), 157-166.

18. Meta, S., & Cinsdikici, M. G. (2010). Vehicle-classification algorithm based on component analysis for single-loop inductive detector. *IEEE Transactions on Vehicular Technology*, *59*(6), 2795-2805.

19. Oh, C., Park, S., & Ritchie, S. G. (2006). A method for identifying rear-end collision risks using inductive loop detectors. *Accident Analysis & Prevention*, *38*(2), 295-301.

20. Pan, X., Guo, Y., & Men, A. (2010, January). Traffic surveillance system for vehicle flow detection. In *Computer Modeling and Simulation, 2010. ICCMS'10. Second International Conference on* (Vol. 1, pp. 314-318). IEEE.

21. Salama, A. S., Saleh, B. K., & Eassa, M. M. (2010, November). Intelligent cross road traffic management system (ICRTMS). In *Computer Technology and Development (ICCTD), 2010 2nd International Conference on* (pp. 27-31). IEEE. Samadi, S., Rad, A. P., Kazemi, F. M., & Jafarian, H. (2012). Performance evaluation of intelligent adaptive traffic control systems: A case study. *Journal of transportation technologies*, *2*(03), 248.

22. Taghvaeeyan, S., & Rajamani, R. (2014). Portable roadside sensors for vehicle counting, classification, and speed measurement. *IEEE Transactions on Intelligent Transportation Systems*, *15*(1), 73-83.

23. Tao, F., Zuo, Y., Da Xu, L., & Zhang, L. (2014). IoT-based intelligent perception and access of manufacturing resource toward cloud manufacturing. *IEEE Transactions on Industrial Informatics*, *10*(2), 1547-1557.

24. Tong, M., & Tang, M. (2010, September). LEACH-B: An improved LEACH protocol for wireless sensor network. In *Wireless Communications Networking and Mobile Computing (WiCOM), 2010 6th International Conference on* (pp. 1-4). IEEE.

25. Wang, M., Shan, H., Lu, R., Zhang, R., Shen, X., & Bai, F. (2015). Real-time path planning based on hybrid-VANET-enhanced transportation system. *IEEE Transactions on Vehicular Technology*, *64*(5), 1664-1678.

26. Wang, M., Liang, H., Zhang, R., Deng, R., & Shen, X. (2014). Mobility-aware coordinated charging for electric vehicles in VANET-enhanced smart grid. *IEEE Journal on Selected Areas in Communications*, *32*(7), 1344-1360.

27. Wang, S., Li, R., & Guo, M. (2018). Application of nonparametric regression in predicting traffic incident duration. Transport, 33(1), 22-31. https://doi.org/10.3846/16484142.2015.1004104

28. Xiao, J., Gao, X., Kong, Q. J., & Liu, Y. (2014). More robust and better: a multiple kernel support vector machine ensemble approach for traffic incident detection. *Journal of Advanced Transportation*, *48*(7), 858-875.

29. Yao, B., Hu, P., Zhang, M., & Jin, M. (2014). A support vector machine with the tabu search algorithm for freeway incident detection. *International Journal of Applied Mathematics and Computer Science*, *24*(2), 397-404.

30. Zheng, B., Sayin, M. O., Lin, C. W., Shiraishi, S., & Zhu, Q. (2017, November). Timing and security analysis of VANET-based intelligent transportation systems. In *Computer-Aided Design (ICCAD), 2017 IEEE/ACM International Conference on*(pp. 984-991). IEEE.Zhu, L., & Jin, S. (2011, July). Speed estimation with single loop detector using typical effective vehicle length. In *Multimedia Technology (ICMT), 2011 International Conference on* (pp. 4096-4099). IEEE.

## AUTHORS PROFILE

**Mr. Manu M R,** is currently working as a Computer Science Teacher in Ministry of Education, Abudabi, UAE.He worked as an Assistant Professor, in School of Computing Science and Engineering, Galgotias University, NCR Delhi, India. He has completed ME in Computer Science and Engineering from Anna University Taramani Campus, Tamil Nadu, India and currently pursuing Ph.D in Computer Science and Engineering from Galgotias University NCR Delhi, India. His area of Interest is Big Data, Networks and Network security. He has undergone different research projects in networks specialization and published 20 papers in various international and national journals and submitted 4 Patent in the field of Artificial intelligence. He is currently writing Monograph and book chapters in CRC Press, Springer, Elsevier publishers.