



2019

## Big Data Analytics for Healthcare Industry: Impact, Applications, and Tools

Sunil Kumar

*the Directorate of Livestock Farms, Guru Angad Dev Veterinary and Animal Sciences University, Ludhiana, India.*

Maninder Singh

*the Department of Computer Science, Punjabi University, Patiala, India.*

Follow this and additional works at: <https://tsinghuauniversitypress.researchcommons.org/big-data-mining-and-analytics>



Part of the [Computer Engineering Commons](#), [Computer Sciences Commons](#), and the [Data Science Commons](#)

### Recommended Citation

Sunil Kumar, Maninder Singh. Big Data Analytics for Healthcare Industry: Impact, Applications, and Tools. *Big Data Mining and Analytics* 2019, 2(1): 48-57.

This Research Article is brought to you for free and open access by Tsinghua University Press: Journals Publishing. It has been accepted for inclusion in *Big Data Mining and Analytics* by an authorized editor of Tsinghua University Press: Journals Publishing.

# Big Data Analytics for Healthcare Industry: Impact, Applications, and Tools

Sunil Kumar\* and Maninder Singh

**Abstract:** In recent years, huge amounts of structured, unstructured, and semi-structured data have been generated by various institutions around the world and, collectively, this heterogeneous data is referred to as big data. The health industry sector has been confronted by the need to manage the big data being produced by various sources, which are well known for producing high volumes of heterogeneous data. Various big-data analytics tools and techniques have been developed for handling these massive amounts of data, in the healthcare sector. In this paper, we discuss the impact of big data in healthcare, and various tools available in the Hadoop ecosystem for handling it. We also explore the conceptual architecture of big data analytics for healthcare which involves the data gathering history of different branches, the genome database, electronic health records, text/imagery, and clinical decisions support system.

**Key words:** big data; healthcare; Hadoop; MapReduce

## 1 Introduction

Every day, data is generated by a range of different applications, devices, and geographical research activities for the purposes of weather forecasting, weather prediction, disaster evaluation, crime detection, and the health industry, to name a few. In current scenarios, big data is associated with core technologies and various enterprises including Google, Facebook, and IBM, which extract valuable information from the huge volumes of data collected<sup>[1-3]</sup>. An era of open information in healthcare is now under way. Big data is being generated rapidly in every field including healthcare, with respect to patient care, compliance, and various regulatory requirements. As

the global population continues to increase along with the human lifespan, treatment delivery models are evolving quickly, and some of the decisions underlying these fast changes must be based on data<sup>[4]</sup>. Healthcare shareholders are promised new knowledge from big data, so called both for its volume as well as its complexity and range. Pharmaceutical-industry experts and shareholders have begun to routinely analyze big data to obtain insight, but these activities are still in the early stages and must be coordinated to address healthcare delivery problems and improve healthcare quality. Early systems for big-data analytics of healthcare informatics have been established across many scenarios, e.g., the investigation of patient characteristics and determination of treatment cost and results to pinpoint the best and most cost-effective treatments<sup>[4]</sup>. Health informatics is described as the assimilation of healthcare sciences, computing sciences and information sciences in the study of healthcare information. Health informatics involves data acquisition, storage, and retrieval to provide better results by healthcare providers. In the healthcare system, data is characterized by its heterogeneity and

---

• Sunil Kumar is with the Directorate of Livestock Farms, Guru Angad Dev Veterinary and Animal Sciences University, Ludhiana, India. E-mail: sunilkapoorldh@gmail.com.

• Maninder Singh is with the Department of Computer Science, Punjabi University, Patiala, India. E-mail: singhmaninder25@yahoo.com.

\* To whom correspondence should be addressed.

Manuscript received: 2018-05-16; accepted: 2018-08-02

variety as a result of the linking of a diverse range of biomedical data sources including, for example, sensor data, imagery, gene arrays, laboratory tests, free text, and demographics<sup>[5]</sup>. Most data in healthcare system (e.g., doctor's notes, lab test results, and clinical data) is unstructured and is not stored electronically, i.e., it exists only in hard copies and its volume is increasing very rapidly. Currently, there is a major focus on the digitization of these vast stores of hard copy data. The revolutions of data size are actually creating a problem in order to achieve this goal<sup>[6]</sup>. The various terminologies and models that have been developed to resolve the problems associated with big data focus on solving four issues known as the four Vs, namely: volume, variety, velocity, and veracity. The various classes of data in healthcare applications include Electronic Health Records (EHR), machine generated/sensor data, health information exchanges, patient registries, portals, genetic databases, and public records. Public records are major sources of big-data in the healthcare industry and require efficient data analytics to resolve their associated healthcare problems. According to a survey conducted in 2012, healthcare data totaled nearly 550 petabytes and will reach nearly 26 000 petabytes in 2020<sup>[5]</sup>. In light of the heterogeneous data formats, huge volume, and related uncertainties in the big-data sources, the task of realizing the transformation of raw data into actionable information is daunting. Being so complex, the identification of health features in medical data and the selection of class attributes for health analytics demands highly sophisticated and architectural specific techniques and tools.

## 2 Big Data Analytics in Health Informatics

The main difference between traditional health analysis and big-data health analytics is the execution of computer programming. In the traditional system, the healthcare industry depended on other industries for big data analysis. Many healthcare shareholders trust information technology because of its meaningful outcomes—their operating systems are functional and they can process the data into standardized forms. Today, the healthcare industry is faced with the challenge of handling rapidly developing big healthcare data. The field of big data analytics is growing and has the potential to provide useful insights for the healthcare system. As noted above, most of the massive amounts of data generated by this system is saved in hard copies, which must then be digitized<sup>[7]</sup>. Big

data can improve healthcare delivery and reduce its cost, while supporting advanced patient care, improving patient outcomes, and avoiding unnecessary costs<sup>[8]</sup>. Big data analytics is currently used to predict the outcomes of decisions made by physicians, the outcome of a heart operation for a condition based on patient's age, current condition, and health status. Essentially, we can say that the role of big data in the health sector is to manage data sets related to healthcare, which are complex and difficult to manage using current hardware, software, and management tools. In addition to the burgeoning volume of healthcare data, reimbursement methods are also changing<sup>[9]</sup>. Therefore, purposeful use and pay based on performance have emerged as important factors in the healthcare sector. In 2011, organizations working in the field of healthcare had produced more than 150 exabytes of data<sup>[10]</sup>, all of which must be efficiently analyzed to be at all useful to the healthcare system<sup>[11]</sup>. The storage of healthcare related data in EHRs occurs in a variety of forms. A sudden increase in data related to healthcare informatics has also been observed in the field of bioinformatics, where many terabytes of data are generated by genomic sequencing<sup>[11]</sup>. There are a variety of analytical techniques available for interpreting medical, which can then be used for patient care<sup>[12]</sup>. The diverse origins and forms of big data are challenging the healthcare informatics community to develop methods for data processing. There is a big demand for technique that combines dissimilar data sources<sup>[13]</sup>.

A number of conceptual approaches can be employed to recognize irregularities in vast amounts of data from different datasets. The frameworks available for the analysis of healthcare data are as follows:

- **Predictive Analytics in Healthcare:** For the past two years, predictive analysis has been recognized as one of the major business intelligence approaches, but its real world applications extend far beyond the business context. Big data analytics includes various methods, including text analytics and multimedia analytics<sup>[14]</sup>. However, one of the most crucial categories is predictive analytics which includes statistical methods like data mining and machine learning that examine current and historical facts to predict the future. Predictive methods which are being used today in the hospital context to determine if patient may be at risk for readmission<sup>[15]</sup>. This data can help doctors to make important patient care decisions. Predictive analysis requires an understanding and use of machine learning, which is widely applied in this

approach.

- **Machine Learning in Healthcare:** The concept of machine learning is very similar to that of data mining<sup>[4]</sup>, both of which scan data to identify patterns. Rather than extracting data based on human understanding, as in data mining applications, machine learning uses that data to improve the program's understanding. Machine learning identifies data patterns and then alters the program function accordingly<sup>[16]</sup>.

- **Electronic Health Records:** EHR represents the most widespread health application of big data in healthcare. Each patient has his/her own medical records, with details that include their medical history, allergies diagnosis, symptoms, and lab test results. Patient records are shared in both public and private sectors with healthcare providers via a secure information system. These files are modifiable, in that doctors can make changes over time and add new medical test results, without the need for paper work or duplication of data.

### 3 Four Vs of Big Data in Healthcare

Four primary attributes (shown in Fig. 1) that are

associated with big data: volume, velocity, variety, and veracity.

- **Volume:** Big data is a term to referring to huge volumes of collected data. There is no fixed threshold for the volume of this data. Typically, the term is used with respect to massive-scale data which must be managed, stored, and analyzed using traditional databases and data processing architecture<sup>[14]</sup>. The volume of data generated by modern IT and the healthcare system has been growing and is driven by the reduced costs of data storage and processing architectures and the need to extract valuable insights from data to improve business processes, efficiencies, and services to consumers<sup>[4]</sup>.

- **Velocity:** Velocity, which represents primary reason for the exponential growth of data, refers to how fast data is collected<sup>[14]</sup>. Healthcare systems are generating data at increasingly higher speeds. In the volume and variety of the structured or unstructured data collected, the velocity of the generation of this data after processing requires a decision based on its output.

- **Variety:** Variety refers to the form of the data, i.e., unstructured or structured, text, medical imagery, audio, video, and sensor data. Structured data information includes clinical data (patient record data), which

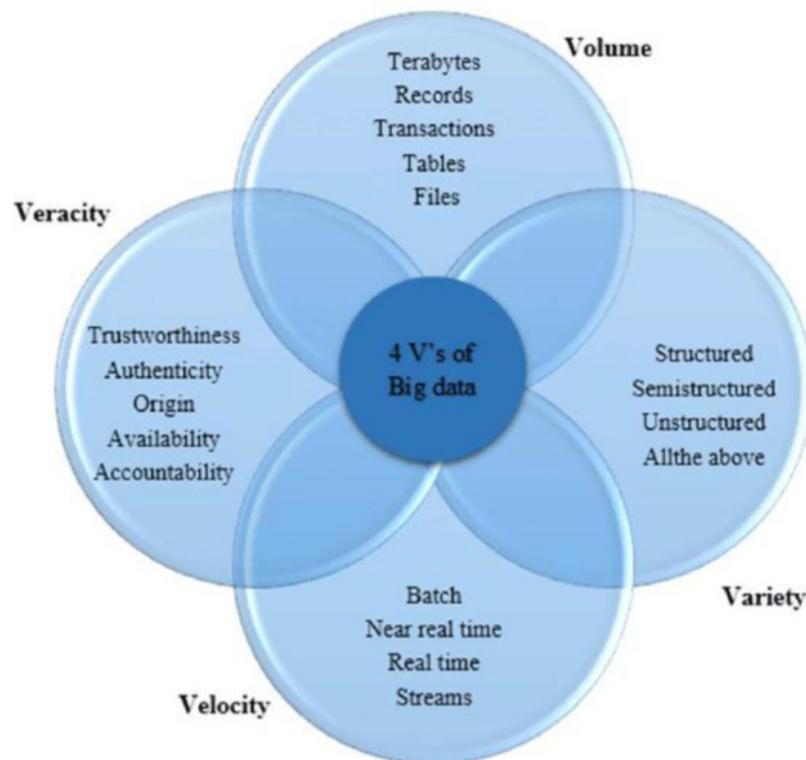


Fig. 1 Big data's four Vs in healthcare.

must simply be collected, stored, and processed by a particular device. Structured data comprises just 5% to 10% of healthcare data. Unstructured or semi-structured data includes e-mails, photos, videos, audios, and other health related data such as hospital medical reports, physician’s notes, paper prescriptions, and radiograph films<sup>[13]</sup>.

- **Veracity:** The veracity of data is the degree of assurance that the meaning of data is consistent. Different data sources vary in their levels of data credibility and reliability<sup>[9]</sup>. The outcomes of big-data analytics must be credible and error-free, but in healthcare, unsupervised machine learning algorithms make decisions that are used by automated machines based on data that may be worthless or misleading<sup>[4]</sup>. Healthcare analytics are tasked with extracting useful insights from this data to treat patients and make the best possible decisions.

#### 4 Impact of Big Data on the Healthcare System

The potential of big data is that it could revolutionize outcomes regarding the most suitable or accurate patient diagnosis and the accuracy information used in the health informatics system<sup>[15]</sup>. As such, the investigation of huge amounts of information will have a powerful effect on medicinal services framework in five respects, or “pathways” (shown in Fig. 2). Improving outcomes for patients with respect to these pathways, as described below, will be the focus of the

healthcare system and will directly impact the patient.

- **Right Living:** Right living refers to the patient living a better and healthier life<sup>[15]</sup>. By right living, patients could manage themselves by making the best decisions for themselves, based on the utilization of information mining better choices and enhancing their wellbeing. By choosing the right path for their daily health, regarding their diet, preventive care, exercise, and other activities of daily living, patients can play an active role in realizing a healthy life<sup>[16]</sup>.

- **Right Care:** This pathway ensures that patients receive the most appropriate treatment available and that all providers obtain the same data and has the same objectives to avoid redundancy of planning and effort<sup>[17]</sup>. This aspect has become more viable in the era of big data.

- **Right Provider:** Healthcare providers in this pathway can obtain an overall view of their patients by combining data from various sources such as medical equipment, public health statistics, and socioeconomic data<sup>[15]</sup>. The accessibility of this information enables human service providers to conduct targeted investigations and develop the skills and abilities to identify and provide better treatment options to patients<sup>[18]</sup>.

- **Right Innovation:** This pathway recognizes that new disease conditions, new treatments, and new medical will continue to evolve<sup>[15]</sup>. Likewise, advancements in the provision of patient services, for

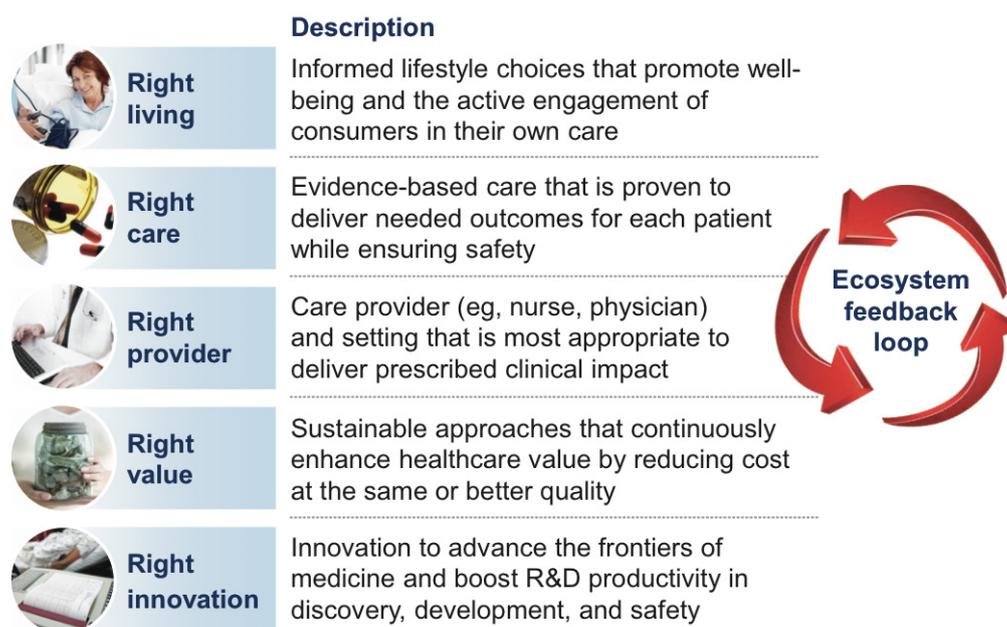


Fig. 2 New pathways that have impact on big data.

example, upgrading medications and the efficiency of research and development efforts, will enable new ways to promote wellbeing and patient health via national social insurance system<sup>[17]</sup>. The availability of early trial data is important for stakeholders. This data can be used to explore high-potential targets and identify techniques for improving traditional clinical treatment methods.

- **Right Value:** To improve the quality and value of health-related services, providers must pay careful and ongoing attention to their patients. Patients must obtain the most beneficial results identified by their social insurance system<sup>[18]</sup>. Measures that could be taken to ensure the intelligent use of data includes, for example, identifying and destroying data misrepresentation, manipulations, and waste, and improving resources<sup>[19]</sup>.

## 5 Hadoop-Based Applications for Health Industry

In light of the fact that healthcare data exists primarily in printed form, there is a need for the active digitization of print form data. The majority of this data is also unstructured, so it is a major challenge for this industry to extract meaningful information regarding patient care, clinical operations, and research. The collection of software utilities known as the Hadoop ecosystem can help the healthcare sector to manage this vast amount of data. The various applications of the Hadoop ecosystem in the healthcare sector are as follows:

- **Treatment of Cancer and Genomics:** We know that human DNA contains three billion base pairs. To fight cancer, it is vital that large amounts of data are efficiently organized. The patterns of cancer mutations and their reactions vary based on individual genetics, which explains the non-curability of some cancer. Oncologists have determined that in recognizing the patterns of cancer, it is important to provide specific treatment for specific cancers, based on the patient's genetic makeup. The Hadoop technology MapReduce facilitates the mapping of three billion DNA base pairs to determine the appropriate cancer treatment for each particular patient. Arizona State University is working on project to develop a healthcare model that takes individual genomic data and selects a treatment based on identification of the patient's cancer gene. This model provides basis for treatment through big data analysis to improve the chances of saving patients lives.

- **Monitoring of Patient Vitals:** Hospital staff throughout the world connect their work output using

big-data technology. Various hospitals around the globe use Hadoop-based components in the Hadoop Distributed File System (HDFS), including the Impala, HBase, Hive, Spark, and Flume frameworks, to convert the huge amount of unstructured data generated by sensors that take patient vital signs, heartbeats per minute, blood pressure, blood sugar level, and respiratory rate. Without Hadoop, these healthcare staff could not analyze this unstructured data being generated by patient healthcare systems. In Atlanta, Georgia, there are 6200 Intensive Care Units (ICUs) for pediatric healthcare, where children can stay for more than one month depending on their problem. These ICUs are equipped with a sensor technology that tracks the child's health status with respect to heartbeat, blood pressure, and other vital signs. If any problem occurs, an alert is automatically generated to medical staff to ensure the child's safety.

- **Hospital Network:** Several hospitals use the Hadoop ecosystem's NoSQL database to collect and manage their huge amounts of real-time data from diverse sources related to patient care, finances, and a payroll, which helps them identify high-risk patients while also reducing day-to-day expenditures.

- **Healthcare Intelligence:** Hadoop technology also supports the healthcare intelligence applications used by hospitals and insurance companies. Hadoop ecosystem's Pig, Hive, and MapReduce technologies process large datasets related to medicines, diseases, symptoms, opinions, geographic regions, and other factors to extract meaningful information (e.g., desired age) for insurance companies.

- **Prevention and Detection of Frauds:** In the early faces of big data analytics, health-based insurance groups utilize multiple paths to identify fraud activity and establish methods to prevent medical fraud. With Hadoop, companies use applications based on a prediction model to identify those committing fraud via data regarding their previous health claims, voice recordings, wages, and demographics. Hadoop's NoSQL database is also helpful in preventing fraud related to medical claims at an early stage by the use of real-time Hadoop based health applications, authentic medical claim bills, weather forecasting data, voice data recordings, and other data sources.

## 6 Big Data Analytics Architecture for Health Informatics

Currently, the main focus in big-data analytics is

to gain an in-depth insight and understanding of big data rather than to collect it<sup>[20]</sup>. Data analytics involves the development and application of algorithms for analyzing various complex data sets to extract meaningful knowledge, patterns, and information. In recent years, researchers have begun to consider the appropriate architectural framework for healthcare systems that utilize big-data analytics, one of which uses a four-layer architecture that comprises a transformation layer, data-source layer, big data platform layer, and analytical layer<sup>[14]</sup>. In this layered system, data originates from different sources and has various formats and storage systems. Each layer has a specific data-processing functionality for performing specific tasks on the HDFS, using the MapReduce processing model. The other layers perform other tasks, i.e., report generation, query passing, data mining processing, and online analytical processing.

The main requirement in big-data analytical processing is to bundle the data at high speed to minimize the bundling time. The next priority in big-data analytical processing is to efficiently update and transform queries at a constant time<sup>[21]</sup>. The third requirement in the big-data analytical processing is to utilize and efficiently manage the storage area space. The last specification of big-data analytics is to efficiently become familiar with the rapidly progressing workload notations. Big-data analytics frameworks differ from traditional healthcare processing systems with respect to how they process big data<sup>[22]</sup>. In the current health care system, data is processed using traditional tools installed in a single stand-alone system like a desktop computer. In contrast, big data is processed by clustering and scans multiple nodes of clusters in the network<sup>[23]</sup>. This processing is based on the concept of parallelism to handle large medical data sets<sup>[24]</sup>. Freely available frameworks, such as Hadoop, MapReduce, Pig, Sqoop, Hive, and HBase Avro, all have ability to process the health related data sets for healthcare systems.

Big-data technologies broadly refer to scientific innovations that mimic those used for large datasets<sup>[25]</sup>. In the first component is the requirement for big data sources for processing. In the second component clusters with a centralized big-data processing infrastructure are at the peak of high performance<sup>[24]</sup>. It has been observed that the tools mainly available for big-data analytics processing provide data security, scalability, and manageability with the help of the MapReduce paradigm. In the

third component, big data analytics applications have a storage domain to integrate accessed databases that use different applications<sup>[26]</sup>. In the fourth component, are the most popular big-data analytics applications in healthcare systems, which include reports, Online Analytical Processing (OLAP), queries, and data mining.

As shown in Fig. 3, healthcare data come from a range of sources including EHRs, genome databases, genome data files, text and imagery (unstructured data sources), clinical decision support systems, government related sources, medical test labs and pharmacies, and health insurance companies. These data are frequently available in different scheme tables, and are in ASCII/text and stored at various locations.

In the next section, we describe the various big-data Hadoop-based processing tools that support the development of health-based applications for the health industry.

## 7 Hadoop's Tools and Techniques for Big Data

To manage unstructured big data that does not fit into any database, special tools are needed. To examine this type of big dataset, the IT sector uses the Hadoop platform for a wide variety of methods that have been developed to record, organize, and analyze this type of data<sup>[27,28]</sup>. More efficient tools are needed to extract meaningful output from big data. Most of the tools are implemented in the Apache Hadoop architecture including MapReduce, Mahout, Hive, and others<sup>[29]</sup>. Below, we discuss the various tools used in processing healthcare big datasets.

- **Apache Hadoop:** The name Hadoop has evolved to mean many different things<sup>[23]</sup>. In 2002, it was established as a single software project to support a web search engine. Since that time, it has grown into an ecosystem of tools and applications that are used to analyze large amounts and types of data<sup>[30]</sup>. Hadoop can no longer be considered to be a monolithic single project, but rather an approach to data processing that radically differs from the traditional relational database model<sup>[23]</sup>. A more practical definition of the Hadoop ecosystem and framework is the following: open source tools, libraries, and methodologies for "big data" analysis in which a number of data sets are collected from different sources, i.e., Internet images, audios, videos, and sensor records as both structured and unstructured data to be processed<sup>[22]</sup>. Figure 4

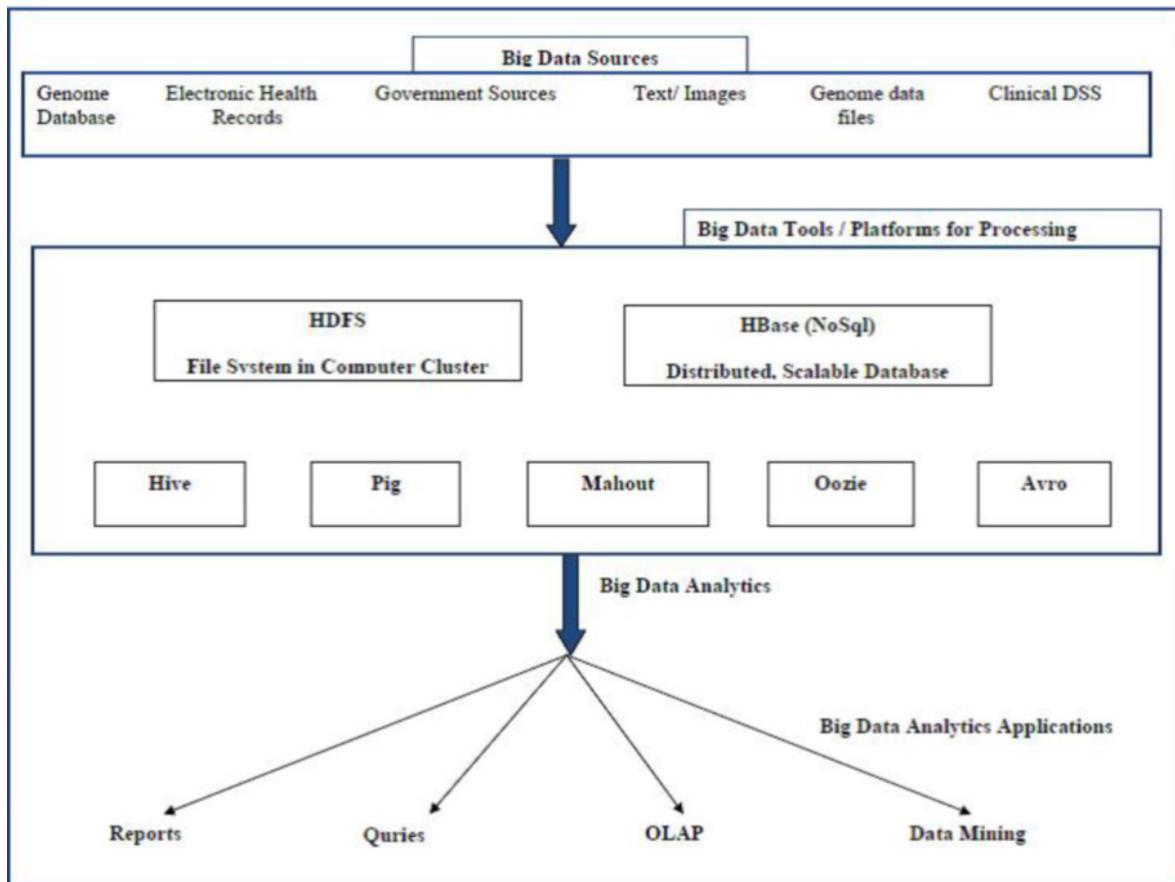


Fig. 3 Conceptual architecture of big data analytics for health informatics.

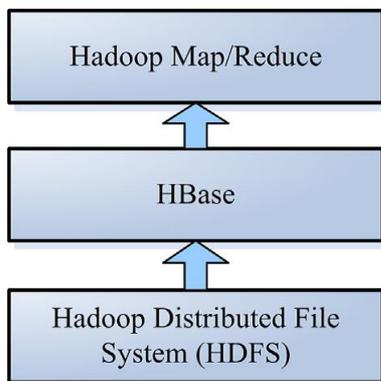


Fig. 4 Hadoop system architecture.

shows the physical layout architecture of Hadoop which consists of MapReduce, HBase, and HDFS.

- **HDFS:** The HDFS was designed for processing big data<sup>[21]</sup>. Although it can support many users simultaneously, HDFS is not designed as a true parallel file system. Rather, the design assumes a large file write-once/ read-many model that enables other optimizations and relaxes many of the concurrency and coherency overhead requirements of a true parallel file

system<sup>[31]</sup>. HDFS is designed for data streaming by which large amounts of data are read from disk in bulk<sup>[25]</sup>. The HDFS block size is 64 MB or 128 MB. There are two types of nodes: a name node and multiple data node(s). A single name node manages all the metadata needed to store and retrieve the actual data from the data nodes<sup>[13]</sup>. No data is actually stored on the name node. Files are stored as blocks in proper sequence and these blocks are equal in size<sup>[24]</sup>. The features of HDFS are its distributed nature and reliability. Storage of metadata and file data is separated. Metadata is stored in name node and application data is stored in data node.

- **MapReduce:** Apache Hadoop is often associated with MapReduce computing. The MapReduce computation model is a very powerful tool used in many health applications and is more common than most users realize. Its underlying concept is very simple<sup>[25]</sup>. In MapReduce, there are two stages: a mapping stage and a reducing stage. In the mapping stage, a mapping procedure is applied to input data. The reducing phase is implemented when counting

is complete<sup>[26]</sup>. The MapReduce programming phase also has two stages: a mapping stage that accepts input in key value pairs and generates output in key value pairs and a second reducing stage, in which each phase consists of key-value pairs as input and output<sup>[12]</sup>. There is a fixed size data segment division step in Hadoop which is called input splits<sup>[20]</sup>. The Map function generates the value pairs and the key, which are stored in the mapper. Any keys that are the same are merged. A simplified view of MapReduce is shown in Fig. 5.

- **Apache Hive:** Hive is a data warehousing layer at the top of Hadoop, in which analyses and queries can be performed using SQL-like procedural language<sup>[32]</sup>. Apache Hive can be used to perform ad-hoc queries, summarization, and data analysis. Hive is considered to be a de facto standard for SQL based queries over petabytes of data using Hadoop and offers the features easy data extraction, transformation, and access to the HDFS comprising data files or other HBase storage system<sup>[33]</sup>.

- **Apache Pig:** Apache Pig is one of the available open-source platforms being used to better analyze big data. Pig is an alternative to the MapReduce programming tool<sup>[34]</sup>. First developed by the Yahoo web service provider as a research project, Pig allows users to develop their own user-define functions and supports many traditional data operations such as join, sort, filter, etc.

- **Apache HBase:** HBase is a column-oriented NoSQL database used in Hadoop<sup>[35]</sup>, in which user can store large numbers of rows and columns. HBase has the functionality of random read/write operations. It also supports record level updates, which is not possible using HDFS<sup>[36]</sup>. HBase provides parallel data storage via the underlying distributed file systems across commodity servers. The file system of choice is

typically HDFS, due to the tight integration of HBase and HDFS<sup>[33]</sup>. If there is need for a structured low-latency view of the high-scale data stored via Hadoop, then HBase is the correct choice. Its open-source code scales linearly to handle petabytes of data on thousands of nodes.

- **Apache Oozie:** To run a complex system or tight system design or if there are a number of interconnected stations with data dependencies between them, there is a need for sophisticated technique called Apache Oozie. Apache Oozie can handle and run multiple jobs related to Hadoop. Oozie has two portions: workflow engines that store and execute workflow collections of Hadoop-based jobs and a coordinator engine that processes workflow jobs based on how they are designed in the process schedule. Oozie is designed to construct and manage Hadoop jobs as workflow in which the output of one job serves as the input for a subsequent job<sup>[37]</sup>. Oozie is not a substitute for the Yarn scheduler. Oozie workflow jobs are represented as Directed Acyclic Graphs (DAGs) of actions<sup>[28]</sup>. Oozie plays the role of a service in the cluster and clients submit their jobs for proactive or reactive execution.

- **Apache Avro:** Avro is a serialization format that makes it possible for data to be exchanged between programs written in any language<sup>[38]</sup>. It is often used to connect Flume data flows. The Avro system is schema-based, where the role of a scheme is to perform the read and write operations with the language being independent. Avro serializes the data that have a built-in schema<sup>[33]</sup>. It is a framework for the serialization of persistent data and remote procedure calls between Hadoop nodes and between client programs and Hadoop services.

- **Apache Zookeeper:** Zookeeper is a centralized system used by applications to maintain a healthcare system and provide organizing and other elements

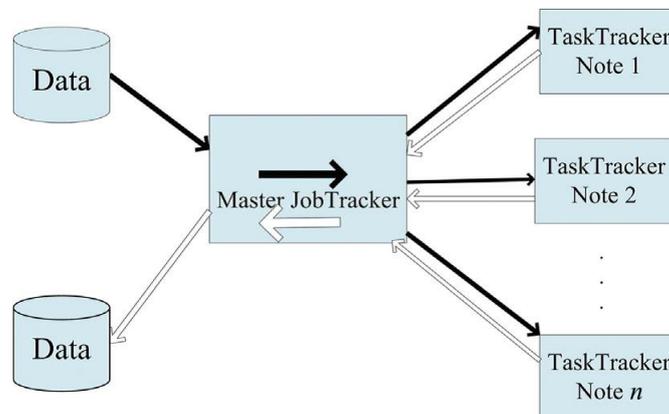


Fig. 5 MapReduce procedure.

on and between nodes<sup>[39]</sup>. It maintains the common objects needed in large cluster environments, including configuration information and the hierarchical naming space. These services can be used by different applications to coordinate the distributed processing of Hadoop clusters. Zookeeper also ensures application reliability<sup>[40]</sup>. If an application master dies, zookeeper generates a new application master to resume the tasks.

- **Apache Yarn:** Hadoop Yarn is a distributed shell application and is an example of a Hadoop non-MapReduce application built on top of Yarn<sup>[41]</sup>. Yarn has two components, a Resource Manager (RM) that handles all the resources within a cluster that are required for the tasks and Node Manager (NM), located on every host in a cluster and handles the available resources on the independent host. Both components handle the scheduling of jobs and manage the containers, memory management, CPU throughput, and I/O system which run the dedicated application code.

- **Apache Sqoop:** Apache Sqoop is a powerful tool that performs the functionality of extracting the data from Relational Database Management System (RDMS) and inputting it into Hadoop architecture for query processing. To do so, this process uses the MapReduce paradigm or other standard level tools, e.g., Hive<sup>[42]</sup>. Once placed in HDFS, the data can be used by Hadoop applications.

- **Apache Flume:** Apache Flume is a highly reliable service for accurately collecting data and moving large volumes of data from independent machines to HDFS<sup>[43]</sup>. Often data transport involves a number of flume agents that may traverse a series of machines and locations. Flume is often used for log files, data generated by social media, and email messages.

## 8 Conclusion

In this paper, we have provided an in-depth description and a brief overview of big data in general and in healthcare system, which plays a significant role in healthcare informatics and greatly influences the healthcare system and the big data four Vs in healthcare. We also proposed the use of a conceptual architecture for solving healthcare problems in big data using Hadoop-based terminologies, which involves the utilization of the big data, generated by different levels of medical data and the development of methods for analyzing this data and to obtain answers to medical questions. The combination of big data and healthcare

analytics can lead to treatments that are effective for specific patients by providing the ability to prescribe appropriate medications for each individual, rather than those that work for most people. As we know, big data analytics is in the early stage of development and current tools and methods cannot solve the problems associated with big data. Big data may be viewed as big systems, which present huge challenges. Therefore, a great deal of research in this field will be required to solve the issues faced by the healthcare system.

## References

- [1] A. Gandomi and M. Haider, Beyond the hype: Big data concepts, methods and analytics, *International Journal of Information Management*, vol. 35, no. 2, pp. 137–144, 2015.
- [2] A. O’Driscoll, J. Daugelaite, and R. D. Sleator, “Big Data”, Hadoop and cloud computing in genomics, *Journal of Biomedical Informatics*, vol. 46, no. 5, pp. 774–781, 2013.
- [3] C. L. P. Chen and C. Y. Zhang, Data-intensive applications, challenges, techniques and technologies: A survey on big data, *Information Sciences*, vol. 275, pp. 314–347, 2014.
- [4] M. Herland, T. M. Khoshgoftaar, and R. Wald, A review of data mining using big data in health informatics, *Journal of Big Data*, vol. 1, no. 1, p. 2, 2014.
- [5] D. H. Shin and M. J. Choi, Ecological views of big data: Perspective and issues, *Telematics and Informatics*, vol. 32, no. 2, pp. 311–320, 2015.
- [6] B. Saraladevi, N. Pazhaniraja, P. V. Paul, M. S. Basha, and P. Dhavachelvan, Big data and Hadoop-A study in security perspective, *Procedia Computer Science*, vol. 50, pp. 596–601, 2015.
- [7] X. Wu, X. Zhu, G. Q. Wu, and W. Ding, Data mining with big data, *IEEE transactions on Knowledge and Data Engineering*, vol. 26, no. 1, pp. 97–107, 2014.
- [8] S. Sharma and V. Mangat, Technology and trends to handle big data: Survey, in *Proc. 5th International Conference on Advanced Computing & Communication Technologies*, 2015, pp. 266–271.
- [9] R. Mehmood and G. Graham, Big data logistics: A health-care transport capacity sharing model, *Procedia Computer Science*, vol. 64, pp. 1107–1114, 2015.
- [10] D. P. Augustine, Leveraging big data analytics and Hadoop in developing India healthcare services, *International Journal of Computer Applications*, vol. 89, no. 16, pp. 44–50, 2014.
- [11] J. A. Patel and P. Sharma, Big data for better health planning, in *Proc. International Conference on Advances in Engineering and Technology Research*, 2014, pp. 1–5.
- [12] A. E. Youssef, A framework for secure healthcare systems based on big data analytics in mobile cloud computing environments, *International Journal of Ambient Systems and Applications*, vol. 2, no. 2, pp. 1–11, 2014.

- [13] MAPR, Healthcare and life science use cases, <https://mapr.com/solutions/industry/healthcare-and-lifescience-use-cases/>, 2018.
- [14] W. Raghupathi and V. Raghupathi, Big data analytics in healthcare: Promise and potential, *Health Information Science and Systems*, vol. 2, no. 1, p. 3, 2014.
- [15] J. Sun and C. K. Reddy, Big data analytics for healthcare, in *Proc. 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2013, pp. 1525–1525.
- [16] C. Mike, W. Hoover, T. Strome, and S. Kanwal. Transforming health care through big data strategies for leveraging big data in the health care industry, <http://ihealthtran.com/iHT2 BigData 2013.pdf>, 2013.
- [17] J. Anuradha, A brief introduction on big data 5Vs characteristics and Hadoop technology, *Procedia Computer Science*, vol. 48, pp. 319–324, 2015.
- [18] M. Viceconti, P. J. Hunter, and R. D. Hose, Big data, big knowledge: Big data for personalized healthcare, *IEEE Journal of Biomedical and Health Informatics*, vol. 19, no. 4, pp. 1209–1215, 2015.
- [19] Y. Sun, H. Song, A. J. Jara, and R. Bie, Internet of things and big data analytics for smart and connected communities, *IEEE Access*, vol. 4, pp. 766–773, 2016.
- [20] A. Jain and V. Bhatnagar, Crime data analysis using Pig with Hadoop, *Procedia Computer Science*, vol. 78, pp. 571–578, 2016.
- [21] T. Jach, E. Magiera, and W. Froelich, Application of Hadoop to store and process big data gathered from an urban water distribution system, *Procedia Engineering*, vol. 119, pp. 1375–1380, 2015.
- [22] C. Uzunkaya, T. Ensari, and Y. Kavurucu, Hadoop ecosystem and its analysis on tweets, *Procedia-Social and Behavioral Sciences*, vol. 195, pp. 1890–1897, 2015.
- [23] S. G. Manikandan and S. Ravi, Big data analysis using Apache Hadoop, in *Proc. International Conference on IT Convergence and Security*, 2014, pp. 1–4.
- [24] V. Ubarhande, A. M. Popescu, and H. Gonzalez-Velez, Novel data-distribution technique for Hadoop in heterogeneous cloud environment, in *Proc. 9<sup>th</sup> International Conference on Complex, Intelligent, and Software Intensive Systems*, 2015, pp. 217–224.
- [25] S. Maitrey and C. K. Jha, Handling big data efficiently by using map reduce technique, in *Proc. International Conference on Computational Intelligence & Communication Technology*, 2015, pp. 703–708.
- [26] J. Dean and S. Ghemawat, MapReduce: Simplified data processing on large clusters, *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, 2008.
- [27] Cloudera, Whole genome research drives healthcare to Hadoop, <https://www.cloudera.com/content/dam/www/marketing/resources/solution-briefs/whole-genome-research-inhealthcare.pdf.landing.html>, 2018.
- [28] R. Misra, B. Panda, and M. Tiwary, Big data and ICT applications: A study, in *Proc. 2nd International Conference on Information and Communication Technology for Competitive Strategies*, 2016, p. 41.
- [29] A. G. Picciano, The evolution of big data and learning analytics in american higher education, *Journal of Asynchronous Learning Networks*, vol. 16, no. 3, pp. 9–20, 2012.
- [30] Apache Hadoop, <http://hadoop.apache.org/>, 2018.
- [31] A. Katal, M. Wazid, R. H. Goudar, and T. Noel, Big data: Issues, challenges, tools and good practices, in *Proc. 6<sup>th</sup> International Conference on Contemporary Computing*, 2013, pp. 404–409.
- [32] Apache Hive, <https://hive.apache.org/>, 2018.
- [33] K. K. Y. Lee, W. C. Tang, and K. S. Choi, Alternatives to relational database: Comparison of NoSQL and XML approaches for clinical data storage, *Computer Methods and Programs in Biomedicine*, vol. 110, no. 1, pp. 99–109, 2013.
- [34] Apache Pig, <https://pig.apache.org/>, 2018.
- [35] E. Dede, B. Sendir, P. Kuzlu, J. Weachock, M. Govindaraju, and L. Ramakrishnan, Processing Cassandra datasets with Hadoop-streaming based approaches, *IEEE Transactions on Services Computing*, vol. 9, no. 1, pp. 46–58, 2016.
- [36] Apache HBase, <http://hbase.apache.org/>, 2018.
- [37] Apache Oozie, <https://oozie.apache.org/>, 2018.
- [38] Apache Avro, <https://avro.apache.org/>, 2018.
- [39] Apache Zookeeper, <https://zookeeper.apache.org/>, 2018.
- [40] Apache Zookeeper, <https://www.ibm.com/analytics/hadoop/zookeeper>, 2018.
- [41] Apache Yarn, <https://yarn.apache.org/>, 2018.
- [42] Apache Sqoop, <https://sqoop.apache.org/>, 2018.
- [43] Apache Flume, <https://flume.apache.org/>, 2018.



**Sunil Kumar** received the M.Tech degree in Internet and Communication Technology (ICT) from Punjabi University, Patiala, Punjab in 2010. He is working in the Directorate of Livestock Farms, Guru Angad Dev Veterinary and Animal Sciences University, Ludhiana, India. He is also pursuing his PhD in computer science at Punjabi University, Patiala, India. His research includes wireless networks, ad-hoc networks, and big data analytics.



**Maninder Singh** is working as an assistant professor in the Department of Computer Science, Punjabi University, Patiala. He received the PhD degree in computer science from Punjabi University, Patiala in 2009. He is the member of various professional societies such as IEEE, ACM, CSTA, etc. His research includes wireless networks, ad-hoc networks, big data analytics, and machine learning.