

TERM PAPER PROJECT

BASIC BIOMEDICAL ENGINEERING

BIG DATA ANALYTICS IN HEALTHCARE

NATIONAL INSTITUTE OF TECHNOLOGY, RAIPUR

492010, CHHATTISGARH, INDIA



SUBMITTED TO = DR. SAURABH GUPTA SIR

SUBMITTED BY = PIYUSH KUMAR SAHU

ROLL NO. = 21111037

BRANCH = BIOMEDICAL ENGINEERING

SECTION = 'A'

E-mail ID = varundvnpihu@gmail.com

PHONE NO. = 9179984034

SEMESTER = 1ST



Figure 1:

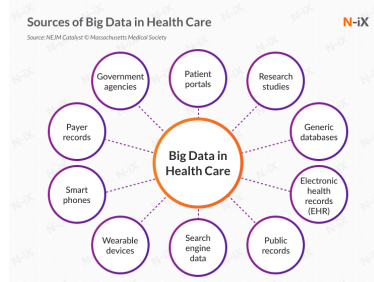


Figure 2:

Abstract

‘Big data’ is massive amounts of information that can work wonders. It has become a topic of special interest for the past two decades because of a great potential that is hidden in it. Various public and private sector industries generate, store, and analyze big data with an aim to improve the services they provide. In the healthcare industry, various sources for big data include hospital records, medical records of patients, results of medical examinations, and devices that are a part of internet of things.

Biomedical research also generates a significant portion of big data relevant to public healthcare. This data requires proper management and analysis in order to derive meaningful information. Otherwise, seeking solution by analyzing big data quickly becomes comparable to finding a needle in the haystack. There are various challenges associated with each step of handling big data which can only be surpassed by using high-end computing solutions for big data analysis. That is why, to provide relevant solutions for improving public health, healthcare providers are required to be fully equipped with appropriate infrastructure to systematically generate and analyze big data. An efficient management, analysis, and interpretation of big data can change the game by opening new avenues for modern healthcare.

1 Introduction

Today, we are facing a situation wherein we are flooded with tons of data from every aspect of our life such as social activities, science, work, health, etc. In a way, we can compare the present situation to a data deluge. The technological advances have helped us in generating more and more data, even to a level where it has become unmanageable with currently available technologies. This has led to the creation of the term ‘big data’ to describe data that is large and unmanageable.

In order to meet our present and future social needs, we need to develop new strategies to organize this data and derive meaningful information. One such special social need is healthcare. Like every other industry, healthcare organizations are producing data at a tremendous rate that presents many advantages and challenges at the same time. In this review, we discuss about the basics of big data including its management, analysis and future prospects especially in healthcare sector.

2 The Data Overload

Every day, people working with various organizations around the world are generating a massive amount of data. The term “digital universe” quantitatively defines such massive amounts of data

created, replicated, and consumed in a single year. International Data Corporation (IDC) estimated the approximate size of the digital universe in 2005 to be 130 exabytes (EB). The digital universe in 2017 expanded to about 16,000 EB or 16 zettabytes (ZB). IDC predicted that the digital universe would expand to 40,000 EB by the year 2020. To imagine this size, we would have to assign about 5200 gigabytes (GB) of data to all individuals. This exemplifies the phenomenal speed at which the digital universe is expanding. The internet giants, like Google and Facebook, have been collecting and storing massive amounts of data.

These observations have become so conspicuous that has eventually led to the birth of a new field of science termed ‘Data Science’. Data science deals with various aspects including data management and analysis, to extract deeper insights for improving the functionality or services of a system (for example, healthcare and transport system). Additionally, with the availability of some of the most creative and meaningful ways to visualize big data post-analysis, it has become easier to understand the functioning of any complex system. As a large section of society is becoming aware of, and involved in generating big data, it has become necessary to define what big data is. Therefore, in this review, we attempt to provide details on the impact of big data in the transformation of global healthcare sector and its impact on our daily lives.

3 Defining Big Data

As the name suggests, ‘big data’ represents large amounts of data that is unmanageable using traditional software or internet-based platforms. It surpasses the traditionally used amount of storage, processing and analytical power. The ‘big’ part of big data is indicative of its large volume. In addition to volume, the big data description also includes velocity and variety. Velocity indicates the speed or rate of data collection and making it accessible for further analysis; while, variety remarks on the different types of organized and unorganized data that any firm or system can collect, such as transaction-level data, video, audio, text or log files. These three Vs have become the standard definition of big data. Although, other people have added several other Vs to this definition, the most accepted 4th V remains ‘veracity’.

The term “big data” has become extremely popular across the globe in recent years. Almost every sector of research, whether it relates to industry or academics, is generating and analyzing big data for various purposes. The most challenging task regarding this huge heap of data that can be organized and unorganized, is its management. Implementation of artificial intelligence (AI) algorithms and novel fusion algorithms would be necessary to make sense from this large amount of data. Indeed, it would be a great feat to achieve automated decision-making by the implementation of machine learning (ML) methods like neural networks and other AI techniques. However, in absence of appropriate software and hardware support, big data can be quite hazy. We need to develop better techniques to handle this ‘endless sea’ of data and smart web applications for efficient analysis to gain workable insights.

4 Healthcare as a Big Data Repository

Healthcare is a multi-dimensional system established with the sole aim for the prevention, diagnosis, and treatment of health-related issues or impairments in human beings. The major components of a healthcare system are the health professionals (physicians or nurses), health facilities (clinics, hospitals for delivering medicines and other diagnosis or treatment technologies), and a financing institution supporting the former two. The health professionals belong to various health sectors like dentistry, medicine, midwifery, nursing, psychology, physiotherapy, and many others. Healthcare is required at several levels depending on the urgency of situation.

Professionals serve it as the first point of consultation (for primary care), acute care requiring skilled professionals (secondary care), advanced medical investigation and treatment (tertiary care) and highly uncommon diagnostic or surgical procedures (quaternary care). At all these levels, the health professionals are responsible for different kinds of information such as patient’s medical history (diagnosis and prescriptions related data), medical and clinical data (like data from imaging and laboratory examinations), and other private or personal medical data. Previously, the common practice to store such medical records for a patient was in the form of either handwritten notes or typed reports.

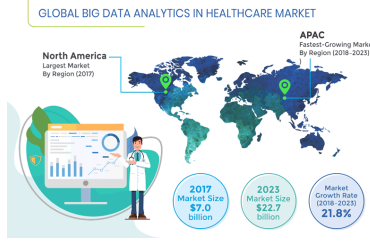


Figure 3:



Figure 4:

5 Electronic Health Records

EHRs have introduced many advantages for handling modern healthcare related data. Below, we describe some of the characteristic advantages of using EHRs. The first advantage of EHRs is that healthcare professionals have an improved access to the entire medical history of a patient. The information includes medical diagnoses, prescriptions, data related to known allergies, demographics, clinical narratives, and the results obtained from various laboratory tests. The recognition and treatment of medical conditions thus is time efficient due to a reduction in the lag time of previous test results. With time we have observed a significant decrease in the redundant and additional examinations, lost orders and ambiguities caused by illegible handwriting, and an improved care coordination between multiple healthcare providers. Overcoming such logistical errors has led to reduction in the number of drug allergies by reducing errors in medication dose and frequency. Healthcare professionals have also found access over web based and electronic platforms to improve their medical practices significantly using automatic reminders and prompts regarding vaccinations, abnormal laboratory results, cancer screening, and other periodic checkups.

EHRs enable faster data retrieval and facilitate reporting of key healthcare quality indicators to the organizations, and also improve public health surveillance by immediate reporting of disease outbreaks. EHRs also provide relevant data regarding the quality of care for the beneficiaries of employee health insurance programs and can help control the increasing costs of health insurance benefits. Finally, EHRs can reduce or absolutely eliminate delays and confusion in the billing and claims management area. The EHRs and internet together help provide access to millions of health related medical information critical for patient life.

6 Digitization of healthcare and big data

Similar to EHR, an electronic medical record (EMR) stores the standard medical and clinical data gathered from the patients. EHRs, EMRs, personal health record (PHR), medical practice management software (MPM), and many other healthcare data components collectively have the potential to improve the quality, service efficiency, and costs of healthcare along with the reduction of medical errors. The big data in healthcare includes the healthcare payer-provider data (such as EMRs, pharmacy prescription, and insurance records) along with the genomics-driven experiments (such as genotyping, gene expression data) and other data acquired from the smart web of internet of things (IoT)

The management and usage of such healthcare data has been increasingly dependent on information technology. The development and usage of wellness monitoring devices and related software that can generate alerts and share the health related data of a patient with the respective health care providers

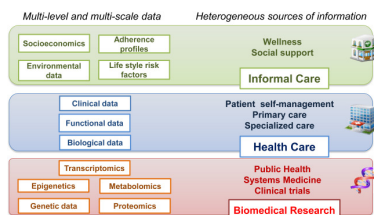


Figure 5:

has gained momentum, especially in establishing a real-time biomedical and health monitoring system. These devices are generating a huge amount of data that can be analyzed to provide real-time clinical or medical care. The use of big data from healthcare shows promise for improving health outcomes and controlling costs.

7 Big Data In Biomedical Research

A biological system, such as a human cell, exhibits molecular and physical events of complex interplay. In order to understand interdependencies of various components and events of such a complex system, a biomedical or biological experiment usually gathers data on a smaller and/or simpler component. Consequently, it requires multiple simplified experiments to generate a wide map of a given biological phenomenon of interest. This indicates that more the data we have, the better we understand the biological processes. With this idea, modern techniques have evolved at a great pace. For instance, one can imagine the amount of data generated since the integration of efficient technologies like next-generation sequencing (NGS) and Genome wide association studies (GWAS) to decode human genetics. NGS-based data provides information at depths that were previously inaccessible and takes the experimental scenario to a completely new dimension. It has increased the resolution at which we observe or record biological events associated with specific diseases in a real time manner.

8 Extracting information from EHR datasets

Emerging ML or AI based strategies are helping to refine healthcare industry’s information processing capabilities. For example, natural language processing (NLP) is a rapidly developing area of machine learning that can identify key syntactic structures in free text, help in speech recognition and extract the meaning behind a narrative. NLP tools can help generate new documents, like a clinical visit summary, or to dictate clinical notes. The unique content and complexity of clinical documentation can be challenging for many NLP developers. Nonetheless, we should be able to extract relevant information from healthcare data using such approaches as NLP.

AI has also been used to provide predictive capabilities to healthcare big data. For example, ML algorithms can convert the diagnostic system of medical images into automated decision-making. Though it is apparent that healthcare professionals may not be replaced by machines in the near future, yet AI can definitely assist physicians to make better clinical decisions or even replace human judgment in certain functional areas of healthcare.

9 Internet Of Things [IOT]

Healthcare industry has not been quick enough to adapt to the big data movement compared to other industries. Therefore, big data usage in the healthcare sector is still in its infancy. For example, healthcare and biomedical big data have not yet converged to enhance healthcare data with molecular pathology. Such convergence can help unravel various mechanisms of action or other aspects of predictive biology. Therefore, to assess an individual’s health status, biomolecular and clinical datasets need to be married. One such source of clinical data in healthcare is ‘internet of things’ (IoT).

In fact, IoT is another big player implemented in a number of other industries including healthcare. Until recently, the objects of common use such as cars, watches, refrigerators and health-monitoring devices, did not usually produce or handle data and lacked internet connectivity. However, furnishing

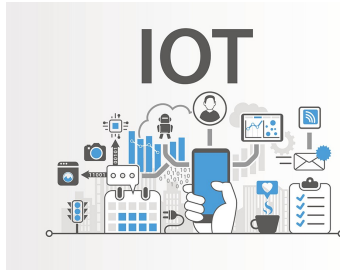


Figure 6:

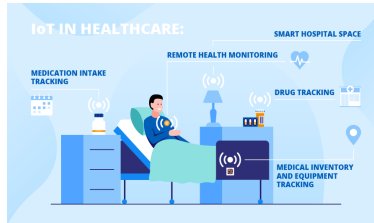


Figure 7:

such objects with computer chips and sensors that enable data collection and transmission over internet has opened new avenues. In fact, IoT has become a rising movement in the field of healthcare. IoT devices create a continuous stream of data while monitoring the health of people (or patients) which makes these devices a major contributor to big data in healthcare. Such resources can interconnect various devices to provide a reliable, effective and smart healthcare service to the elderly and patients with a chronic illness.

10 Advantages of IOT

Using the web of IoT devices, a doctor can measure and monitor various parameters from his/her clients in their respective locations for example, home or office. Therefore, through early intervention and treatment, a patient might not need hospitalization or even visit the doctor resulting in significant cost reduction in healthcare expenses. Some examples of IoT devices used in healthcare include fitness or health-tracking wearable devices, biosensors, clinical devices for monitoring vital signs, and others types of devices or clinical instruments. Such IoT devices generate a large amount of health related data.

In fact, big data generated from IoT has been quiet advantageous in several areas in offering better investigation and predictions. On a larger scale, the data from such devices can help in personnel health monitoring, modelling the spread of a disease and finding ways to contain a particular disease outbreak. The analysis of data from IoT would require an updated operating software because of its specific nature along with advanced hardware and software applications.

11 Mobile Computing and Mobile Health [mhealth]

In today's digital world, every individual seems to be obsessed to track their fitness and health statistics using the in-built pedometer of their portable and wearable devices such as, smartphones, smartwatches, fitness dashboards or tablets. With an increasingly mobile society in almost all aspects of life, the healthcare infrastructure needs remodeling to accommodate mobile devices . The practice of medicine and public health using mobile devices, known as mHealth or mobile health, pervades different degrees of health care especially for chronic diseases, such as diabetes and cancer Healthcare organizations are increasingly using mobile health and wellness services for implementing novel and innovative ways to provide care and coordinate health as well as wellness. Mobile platforms can improve healthcare by accelerating interactive communication between patients and healthcare providers.



Figure 8:



Figure 9:

12 Nature Of The Big Data In Healthcare

EHRs can enable advanced analytics and help clinical decision-making by providing enormous data. However, a large proportion of this data is currently unstructured in nature. An unstructured data is the information that does not adhere to a pre-defined model or organizational framework. The reason for this choice may simply be that we can record it in a myriad of formats.

Nonetheless, the healthcare industry is required to utilize the full potential of these rich streams of information to enhance the patient experience. In the healthcare sector, it could materialize in terms of better management, care and low-cost treatments. We are miles away from realizing the benefits of big data in a meaningful way and harnessing the insights that come from it. In order to achieve these goals, we need to manage and analyze the big data in a systematic manner.

13 Management And Analysis Of Big Data

Big data is the huge amounts of a variety of data generated at a rapid rate. The data gathered from various sources is mostly required for optimizing consumer services rather than consumer consumption. This is also true for big data from the biomedical research and healthcare. The major challenge with big data is how to handle this large volume of information. To make it available for scientific community, the data is required to be stored in a file format that is easily accessible and readable for an efficient analysis. In the context of healthcare data, another major challenge is the implementation of high-end computing tools, protocols and high-end hardware in the clinical setting.

Experts from diverse backgrounds including biology, information technology, statistics, and mathematics are required to work together to achieve this goal. The data collected using the sensors can be made available on a storage cloud with pre installed software tools developed by analytic tool developers. These tools would have data mining and ML functions developed by AI experts to convert the information stored as data into knowledge. Upon implementation, it would enhance the efficiency of acquiring, storing, analyzing, and visualization of big data from healthcare. The main task is to annotate, integrate, and present this complex data in an appropriate manner for a better understanding.

Heterogeneity of data is another challenge in big data analysis. The huge size and highly heterogeneous nature of big data in healthcare renders it relatively less informative using the conventional technologies. The most common platforms for operating the software framework that assists big data analysis are high power computing clusters accessed via grid computing infrastructures. Cloud computing is such a system that has virtualized storage technologies and provides reliable services. It offers high reliability, scalability and autonomy along with ubiquitous access, dynamic resource discovery and composability. Such platforms can act as a receiver of data from the ubiquitous sensors, as a computer to analyze and interpret the data, as well as providing the user with easy to understand web-based visualization.



Figure 10:

In IoT, the big data processing and analytics can be performed closer to data source using the services of mobile edge computing cloudlets and fog computing. Advanced algorithms are required to implement ML and AI approaches for big data analysis on computing clusters. A programming language suitable for working on big data (e.g. Python, R or other languages) could be used to write such algorithms or software. Therefore, a good knowledge of biology and IT is required to handle the big data from biomedical research. Such a combination of both the trades usually fits for bioinformaticians.

14 Commercial platforms for healthcare data analytics

In order to tackle big data challenges and perform smoother analytics, various companies have implemented AI to analyze published results, textual data, and image data to obtain meaningful outcomes. IBM Corporation is one of the biggest and experienced players in this sector to provide healthcare analytics services commercially. IBM's Watson Health is an AI platform to share and analyze health data among hospitals, providers and researchers.

Similarly, Flatiron Health provides technology-oriented services in healthcare analytics specially focused in cancer research. Other big companies such as Oracle Corporation and Google Inc. are also focusing to develop cloud-based storage and distributed computing power platforms. Interestingly, in the recent few years, several companies and start-ups have also emerged to provide health care-based analytics and solutions.

14.1 AYASDI

Ayasdi is one such big vendor which focuses on ML based methodologies to primarily provide machine intelligence platform along with an application framework with tried and tested enterprise scalability. It provides various applications for healthcare analytics, It is also capable of analyzing and managing how hospitals are organized, conversation between doctors, risk-oriented decisions by doctors for treatment, and the care they deliver to patients. It also provides an application for the assessment and management of population health, a proactive strategy that goes beyond traditional risk analysis methodologies.

14.2 IBM Watson

This is one of the unique ideas of the tech-giant IBM that targets big data analytics in almost every professional sector. This platform utilizes ML and AI based algorithms extensively to extract the maximum information from minimal input. IBM Watson enforces the regimen of integrating a wide array of healthcare domains to provide meaningful and structured data.

IBM Watson has been used to predict specific types of cancer based on the gene expression profiles obtained from various large data sets providing signs of multiple druggable targets. IBM Watson is also used in drug discovery programs by integrating curated literature and forming network maps to provide a detailed overview of the molecular landscape in a specific disease model.

15 Challenges associated with healthcare big data

Methods for big data management and analysis are being continuously developed especially for real-time data streaming, capture, aggregation, analytics (using ML and predictive), and visualization



Figure 11:

solutions that can help integrate a better utilization of EMRs with the healthcare.

However, the availability of hundreds of EHR products certified by the government, each with different clinical terminologies, technical specifications, and functional capabilities has led to difficulties in the interoperability and sharing of data. Nonetheless, we can safely say that the healthcare industry has entered into a ‘post-EMR’ deployment phase. Now, the main objective is to gain actionable insights from these vast amounts of data collected as EMRs. Here, we discuss some of these challenges in brief.

15.1 Storage

Storing large volume of data is one of the primary challenges, but many organizations are comfortable with data storage on their own premises. It has several advantages like control over security, access, and up-time. However, an on-site server network can be expensive to scale and difficult to maintain. It appears that with decreasing costs and increasing reliability, the cloud-based storage using IT infrastructure is a better option which most of the healthcare organizations have opted for. Organizations must choose cloud-partners that understand the importance of healthcare-specific compliance and security issues. Additionally, cloud storage offers lower up-front costs, nimble disaster recovery, and easier expansion.

15.2 Cleaning

The data needs to be cleansed or scrubbed to ensure the accuracy, correctness, consistency, relevancy, and purity after acquisition. This cleaning process can be manual or automatized using logic rules to ensure high levels of accuracy and integrity. More sophisticated and precise tools use machine-learning techniques to reduce time and expenses and to stop foul data from derailing big data projects.

15.3 Unified Formats

Patients produce a huge volume of data that is not easy to capture with traditional EHR format, as it is knotty and not easily manageable. It is too difficult to handle big data especially when it comes without a perfect data organization to the healthcare providers. A need to codify all the clinically relevant information surfaced for the purpose of claims, billing purposes, and clinical analytics. Therefore, medical coding systems like Current Procedural Terminology (CPT) and International Classification of Diseases (ICD) code sets were developed to represent the core clinical concepts.

15.4 Accuracy

Some studies have observed that the reporting of patient data into EMRs or EHRs is not entirely accurate yet, probably because of poor EHR utility, complex workflows, and a broken understanding of why big data is all important to capture well. All these factors can contribute to the quality issues for big data all along its lifecycle. The EHRs intend to improve the quality and communication of data in clinical workflows though reports indicate discrepancies in these contexts. The documentation quality might improve by using self-report questionnaires from patients for their symptoms.

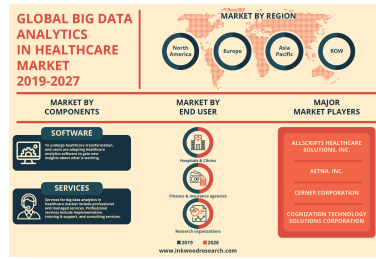


Figure 12:

15.5 Image Pre-Processing

Studies have observed various physical factors that can lead to altered data quality and misinterpretations from existing medical records. Medical images often suffer technical barriers that involve multiple types of noise and artifacts. Improper handling of medical images can also cause tampering of images for instance might lead to delineation of anatomical structures such as veins which is non-correlative with real case scenario. Reduction of noise, clearing artifacts, adjusting contrast of acquired images and image quality adjustment post mishandling are some of the measures that can be implemented to benefit the purpose.

15.6 Security

There have been many security breaches, hackings, phishing attacks, and ransomware episodes that data security is a priority for healthcare organizations. After noticing an array of vulnerabilities, a list of technical safeguards was developed for the protected health information (PHI). These rules, termed as HIPAA Security Rules, help guide organizations with storing, transmission, authentication protocols, and controls over access, integrity, and auditing. Common security measures like using up-to-date anti-virus software, firewalls, encrypting sensitive data, and multi-factor authentication can save a lot of trouble.

15.7 Visualisation

A clean and engaging visualization of data with charts, heat maps, and histograms to illustrate contrasting figures and correct labeling of information to reduce potential confusion, can make it much easier for us to absorb information and use it appropriately. Other examples include bar charts, pie charts, and scatterplots with their own specific ways to convey the data.

15.8 Data Sharing

Patients may or may not receive their care at multiple locations. In the former case, sharing data with other healthcare organizations would be essential. During such sharing, if the data is not interoperable then data movement between disparate organizations could be severely curtailed. This could be due to technical and organizational barriers. This may leave clinicians without key information for making decisions regarding follow-ups and treatment strategies for patients.

The biggest roadblock for data sharing is the treatment of data as a commodity that can provide a competitive advantage. Therefore, sometimes both providers and vendors intentionally interfere with the flow of information to block the information flow between different EHR systems. The healthcare providers will need to overcome every challenge on this list and more to develop a big data exchange ecosystem that provides trustworthy, timely, and meaningful information by connecting all members of the care continuum. Time, commitment, funding, and communication would be required before these challenges are overcome.

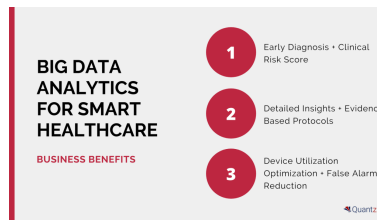


Figure 14:

5. However, in a short span we have witnessed a spectrum of analytics currently in use that have shown significant impacts on the decision making and performance of healthcare industry. The exponential growth of medical data from various domains has forced computational experts to design innovative strategies to analyze and interpret such enormous amount of data within a given timeframe. The integration of computational systems for signal processing from both research and practicing medical professionals has witnessed growth.
6. High volume of medical data collected across heterogeneous platforms has put a challenge to data scientists for careful integration and implementation. It is therefore suggested that revolution in healthcare is further needed to group together bioinformatics, health informatics and analytics to promote personalized and more effective treatments. Furthermore, new strategies and technologies should be developed to understand the nature (structured, semi structured, unstructured), complexity (dimensions and attributes) and volume of the data to derive meaningful information. The greatest asset of big data lies in its limitless possibilities.
7. The birth and integration of big data within the past few years has brought substantial advancements in the health care sector ranging from medical data management to drug discovery programs for complex human diseases including cancer and neurodegenerative disorders. We believe that big data will add-on and bolster the existing pipeline of healthcare advances instead of replacing skilled manpower, subject knowledge experts and intellectuals, a notion argued by many. One can clearly see the transitions of health care market from a wider volume base to personalized or individual specific domain.
8. In the coming year it can be projected that big data analytics will march towards a predictive system. This would mean prediction of futuristic outcomes in an individual's health state based on current or existing data (such as EHR-based and Omics-based). Similarly, it can also be presumed that structured information obtained from a certain geography might lead to generation of population health information. Taken together, big data will facilitate healthcare by introducing prediction of epidemics (in relation to population health), providing early warnings of disease conditions, and helping in the discovery of novel biomarkers and intelligent therapeutic intervention strategies for an improved quality of life.