# APPLICATIONS OF HEALTHCARE SYSTEMS USING MACHINE LEARNING AND BIG DATA ANALYTICS

## Abstract

Big data Massive informational volumes have amazing potential. Due to the huge potential that is incorporated over there, it was given extra focus during the past two decades. Big data is a label used when the size of the data itself becomes a part of the problem. A common strategy for big data analysis on computing clusters is divide and conquer.Utilizing Big Data analytics has the potential to enhance patient results, progress, and individualized healthcare, strengthen the bond between healthcare providers and patients, as well as curtail medical expenses. This article introduces the concept of healthcare data, the integration of big data within healthcare systems, and the various applications and benefits derived from employing Big Data analytics within the healthcare domain. Computer clusters with the shared-nothing architecture are the major computing platforms for big data processing and analysis. Additionally, we highlight the technological advancements related to big data in healthcare, including aspects like cloud computing and real-time data processing. Furthermore, we address the hurdles associated with implementing Big Data analytics in healthcare systems.

**Keywords:** Big data; Big Data analytics; Healthcare; Spark; Machine Learning.

## Authors

**Sruthi Yenninti**
Assistant Professor
Department of CSE (AI&ML, DS)
Anil Neerukonda Institute of Technology & Sciences
Visakhapatnam, Andhra Pradesh, India
sruthi9814@gmail.com

**N. Sivaganga Kumari**
Assistant Professor
Department CSE (AI&ML, DS)
Anil Neerukonda Institute of Technology & Sciences
Visakhapatnam, Andhra Pradesh, India
pandusivaganga9@gmail.com

**S. Joshua Johnson**
Assistant Professor
Department of CSE (AI&ML,DS)
Anil Neerukonda Institute of Technology & Sciences
Visakhapatnam, Andhra Pradesh, India
joshua.sirasapalli@gmail.com

# I. INTRODUCTION

An overwhelming volume of data is now being generated from business transactions, computer simulations, mobile devices, healthcare systems, sensors, satellites, social media, and so on. This massive quantity of data can be used to produce high-value information for decision support, forecasting, business intelligence, research on data-intensive science, and other fields of application. Big data analytics is the process of extracting valuable insights from large datasets. In healthcare, big data analytics can be used to improve patient care, reduce costs, and improve outcomes. There are many different ways that big data analytics can be used in healthcare. Here are a few examples: Identifying patients at risk for disease. Big data analytics can be used to identify patients who are at risk for developing certain diseases. This information can then be used to develop preventive measures or to provide early treatment. Tracking the effectiveness of treatments. Big data analytics can be used to track the effectiveness of different treatments. This information can help doctors make better decisions about which treatments to prescribe. Identifying new patterns and trends. Big data analytics can be used to identify new patterns and trends in healthcare. This information can help researchers develop new treatments and therapies. Improving the efficiency of healthcare delivery. Big data analytics can be used to improve the efficiency of healthcare delivery. For example, big data analytics can be used to schedule appointments, track patient records, and manage inventory. Big data analytics has the potential to revolutionize healthcare. By providing doctors and researchers with more information, big data analytics can help improve patient care, reduce costs, and improve outcomes.

'Big data' represents vast quantities of information with the capacity for remarkable outcomes. Over the last two decades, it has gathered significant attention due to the immense latent potential it holds. Diverse sectors, both public and private, are generating, storing, and analyzing big data to enhance the quality of their services. Within the healthcare realm, substantial sources of big data encompass hospital records, patient medical histories, medical test results, and Internet of Things-connected devices. Biomedical research is also a notable contributor to the extensive pool of big data relevant to public health. Properly managing and dissecting this data is essential for deriving meaningful insights; otherwise, attempting to extract solutions through big data analysis to search for a needle in a haystack. Each stage of big data handling presents its own set of challenges, which can be effectively surmounted by employing high-performance computing solutions tailored for big data analysis.

The healthcare system is a multifaceted structure designed with the primary objective of preventing, diagnosing, and treating health-related concerns or deficiencies in individuals. Hence, in order to furnish pertinent solutions for advancing public health, healthcare providers must possess the requisite infrastructure to methodically generate and analyze big data. A streamlined approach to managing, analyzing, and interpreting big data has the potential to be transformative opportunities for modern healthcare. This is precisely why various sectors, including healthcare, are taking decisive measures to translate this potential into enhanced services and financial gains. Through the seamless integration of biomedical and healthcare data, contemporary healthcare organizations could potentially reshape medical therapies and personalized medicine, leading to a revolutionary shift in the field.

Healthcare advancements have led to the creation of mobile and web applications that enable patients to submit queries about their symptoms to medical professionals via a central server. These applications might include providing patients with rapid assistance, immediate first aid instruction, or suitable referrals for additional medical care. Furthermore, a mobile cloud computing (MCC)-based healthcare framework has been developed. This platform collects and analyses real-time biomedical data from people in various settings, including blood pressure and ECG readings. The ability to conduct in-depth information analysis to improve the implementation of effective strategies has been rendered possible by advanced information technology, which is crucial in the collection of large amounts of healthcare data.

1. **Usage of Big Data Analytics in Healthcare Systems:** As outlined in Table 1, big data exhibits notable attributes in terms of volume, velocity, variety, variability, value, complexity, and sparseness[1], [2]. In the realm of healthcare, big data holds considerable promise with applications encompassing disease surveillance, epidemic management, clinical decision support, and population health oversight [3]. These potential yields significant advantages, including the early detection of ailments. The integration of Big Data analytics into intelligent healthcare systems introduces innovative electronic and mobile health solutions that amplify efficiency and curtail medical expenditures[4].

**Table 1**: **Aspects of Bigdata**

| Aspects | Description | Examples in Healthcare |
|---|---|---|
| Volume | Data size | Treatment plans, multiple conditions, and cohorts of patients |
| Velocity | Data generation rate (batches, streams, infrequent intervals) | Sensing and diagnostics transmitting patients' status and behaviors |
| Variety | Various formats and data types (numbers, text,images) | Clinical, medical, and omics data and images from various patients under diverse conditions |
| Variability | Data changes with time | Health data from wearable sensors |
| Veracity | Imprecise or untruthful data | Clinician notes about patients' states, patients' feedback |
| Value | Inherent value (often achieved through data mining) | Analyzing numerous patients' feedback and identifying the side effects of a drug |
| Complexity | Hierarchies, linkages between items and recurrent structure of data | Multi-pharmacy, multi-morbidity |
| Sparseness | The low density of useful information (due to null values, missing data, etc.) | Many missing data of patient feedback on progress and symptoms |
| Visualization | Visual content is the process of Amalgamating the medical data gathered from various resources and transforming it. | EHRs/EMRs, remote monitoring, Hospital management, diagnostic centers, laboratories, and pharmaceutical measures. |

Big Data analytics accuracy is of the utmost significance. Abbreviations, spelling errors, and illegible notes may be included in personal health records (PHRs). During ambulatory monitoring, measurements made in unregulated and less reliable surroundings could differ from the information acquired by qualified practitioners in clinical settings.

## II. CONCEPTUAL OVERVIEW OF HEALTH BIG DATA ANALYTICS

Healthcare big data analytics involves the systematic analysis of large and diverse healthcare-related datasets to extract valuable insights, patterns, and knowledge. It combines advanced technologies, analytical techniques, and domain expertise to drive evidence-based decisions, improve patient care, and enhance healthcare processes[5], [6]. Here's a conceptual overview of the key components of health big data analytics:

1. **Health Big Data Sources:** Healthcare big data sources encompass a wide range of structured and unstructured data generated within the healthcare ecosystem[7]–[9]. These data sources contribute to a comprehensive understanding of patient health, medical processes, and outcomes. Here are some key healthcare big data sources:

   - **Electronic Health Records (EHRs):** EHRs contain patient medical histories, diagnoses, medications, treatment plans, lab results, and other clinical information. These digitized records provide a longitudinal view of a patient's health and are a primary source for healthcare analytics.

   - **Medical Imaging Data:** Medical imaging, including X-rays, MRIs, CT scans, and ultrasounds, generates large volumes of image data. These images are crucial for diagnosing and monitoring conditions and can be analyzed for pattern recognition and anomaly detection.

   - **Genomic Data:** Genomic data includes information about a person's genetic makeup. Advances in DNA sequencing technology have led to large-scale genetic databases, enabling personalized medicine, disease risk prediction, and drug development.

   - **Health Sensor Data:** Wearable devices, biosensors, and IoT-enabled devices capture real-time patient data like heart rate, blood pressure, glucose levels, and physical activity. This continuous monitoring facilitates remote patient management and health trend analysis.

   - **Healthcare Claims Data:** Claims data from insurance providers contain details about medical procedures, diagnoses, treatments, and associated costs. Analyzing claims data can reveal patterns in healthcare utilization and costs.
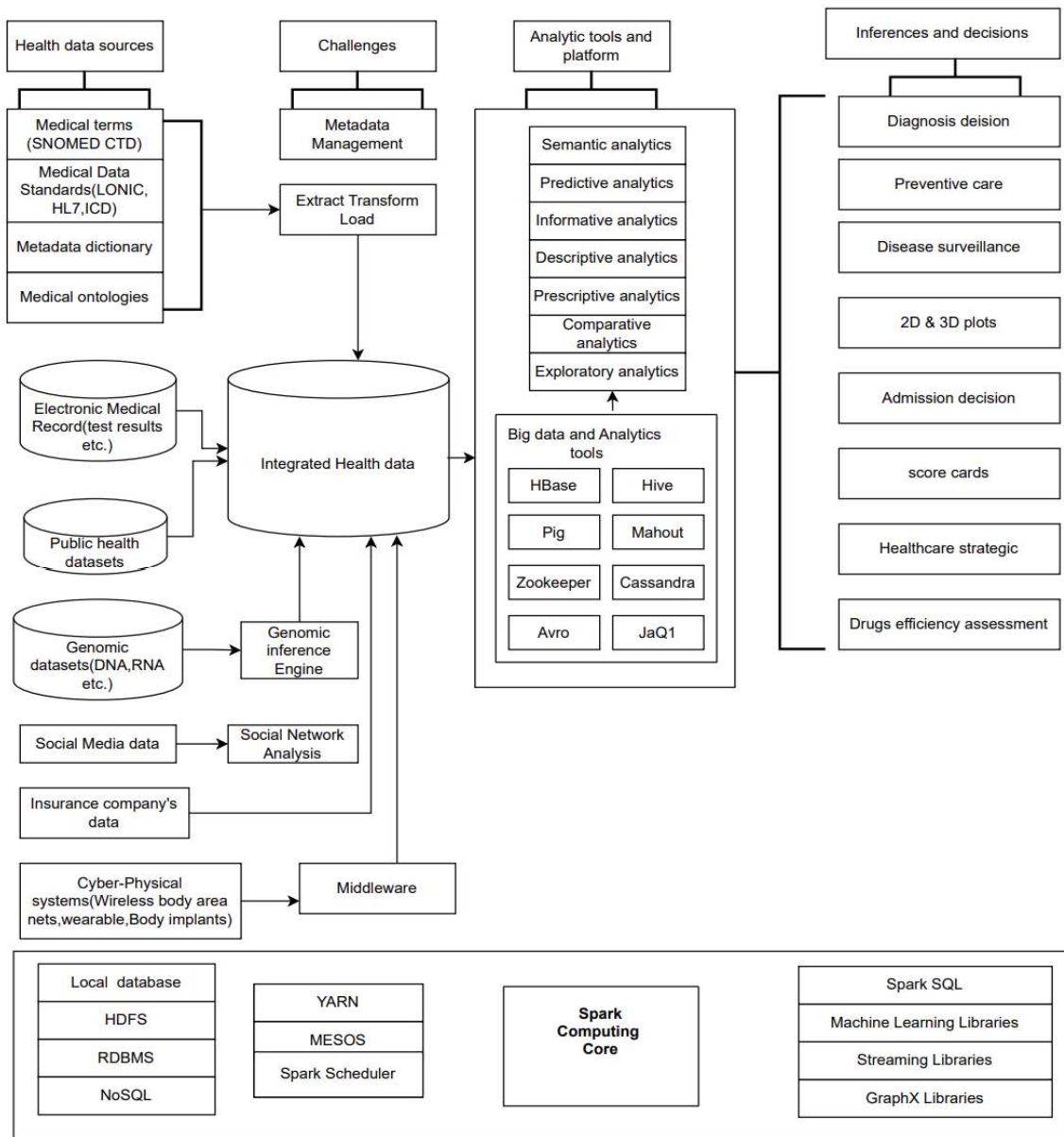
**Figure 1:** Overview of Health Big Data in Spark

- **Social Determinants of Health Data:** Factors like socioeconomic status, education, living conditions, and access to healthcare services impact health outcomes. Incorporating these social determinants of health into analysis provides a holistic perspective.

- **Patient Reported Data:** Patients provide information about their symptoms, quality of life, and experiences through surveys and questionnaires. This data offers insights into patient perspectives and can guide treatment plans.

- **Clinical Trial Data:** Data from clinical trials, including experimental treatments and outcomes, contribute to medical research and drug development. Analyzing clinical trial data helps assess treatment efficacy and safety.

- **Public Health Data:** Government agencies and health organizations collect data on disease prevalence, outbreaks, vaccinations, and environmental factors. This data aids in tracking and controlling public health concerns.

- **Medical Literature and Research Data:** Research articles, medical journals, and conference proceedings contain valuable insights from medical professionals and researchers. Text mining techniques can extract knowledge from these unstructured sources.

- **Pharmaceutical Research Data:** Pharmaceutical companies generate data from drug development, including preclinical and clinical trials. This data informs drug safety, efficacy, and regulatory compliance.

- **Healthcare Surveys and Feedback:** Patient and provider surveys, feedback, and reviews offer insights into healthcare experiences, patient satisfaction, and areas for improvement.

- **Hospital Operational Data:** Data on hospital operations, including bed occupancy, resource utilization, and patient flow, helps optimize healthcare delivery and resource allocation.

- **Electronic Prescription Data:** Prescription records provide information about prescribed medications, dosage, and adherence. Analyzing this data aids in understanding treatment patterns and outcomes.

- **Biometric** Data: Biometric identifiers such as fingerprints, retinal scans, and facial *recognition* are used for patient identification and secure access to medical records.

The combination of these diverse data sources allows healthcare professionals, researchers, and data scientists to gain a comprehensive understanding of patient health, medical trends, treatment effectiveness, and population health. However, managing and analyzing healthcare big data comes with challenges related to data privacy, security, interoperability, and the need for advanced analytical tools.

2. **Integration of Big Data Challenges:** Integrating health data from various sources to create a comprehensive and interoperable system presents several challenges. These challenges span technical, ethical, regulatory, and operational domains. Here are some key challenges associated with integrated health data. data security, managing variability in data quality and standards, achieving semantic interoperability, navigating complex regulatory compliance (e.g., Health Insurance Portability and Accountability Act (HIPAA), General Data Protection Regulation (GDPR)), and determining data ownership. Scalability, real-time integration, and addressing cultural resistance further compound these issues[10]. Addressing these challenges requires coordinated efforts, clear data

governance, and a focus on ethical data use to harness the potential benefits of integrated health data effectively.

3. **Analytical Tools:** Analytical tools are software or platforms designed to process, analyze, visualize, and interpret data. They help transform raw data into actionable insights, enabling better decision-making[11]. Some commonly used analytical tools include:

- Structured Query Language (SQL) databases manage structured data and allow users to query, manipulate, and retrieve information using standardized SQL commands.
- Business Intelligence (BI) tools like Tableau, Power BI, and QlikView offer interactive dashboards and visualizations for data exploration and reporting.
- Statistical Analysis System (SAS) software provides advanced analytics, data management, and predictive modeling capabilities for business and research.
- R is an open-source programming language and environment for statistical computing and graphics, widely used for data analysis and visualization.
- Python is a versatile programming language with libraries like NumPy, pandas, and scikit-learn, making it popular for data analysis, machine learning, and data visualization.
- Hadoop is an open-source framework that stores and processes large datasets across distributed computing clusters using MapReduce programming.
- Apache Spark is a fast and general-purpose cluster-computing framework for big data processing, capable of handling batch and real-time analytics.
- Machine Learning Platforms like TensorFlow, PyTorch, and scikit-learn offer tools for creating and deploying machine learning models.

**In Healthcare, Big Data and Analytical Tools Play A Transformative Role:**

- **Patient Care:** Analyzing patient data helps in personalized treatment plans and disease management.
- **Drug Discovery:** Analyzing genetic and molecular data accelerates drug development.
- **Healthcare Operations:** Optimizing resource allocation and improving patient flow in hospitals.
- **Predictive Analytics:** Identifying patient readmission risks or disease outbreaks.
- **Remote Monitoring:** Analyzing data from wearables and sensors for proactive healthcare.
- **Epidemiological Studies:** Tracking disease patterns and trends at a population level.

4. **Inferences and Decisions:** Inferences and decisions are the key outcomes of analyzing data using various techniques and tools. They form the basis for informed actions, strategies, and improvements across diverse fields[6]. Here's a closer look at inferences, decisions, and their significance:

- Inferences are conclusions drawn from data analysis. They involve interpreting patterns, trends, relationships, and correlations within the data. Inferences are critical

for understanding the underlying meaning of data and for making predictions about future events or outcomes. Inferences can be both descriptive and predictive:

- Descriptive Inferences describe the current state of affairs, providing insights into what has happened. For example, understanding the distribution of patients across different age groups in a healthcare dataset.

- Predictive Inferences involve making predictions based on historical data, statistical models, and patterns observed. For instance, predicting patient readmission rates based on past medical history and treatment.

- Decisions are actions taken based on the insights gained from data analysis. They are influenced by the inferences drawn from the data and are aimed at achieving specific goals. Decisions can be categorized into two main types:

  ➢ Operational Decisions are day-to-day decisions that impact routine activities. For example, adjusting the number of nurses on duty based on patient admission patterns in a hospital.
  ➢ Strategic Decisions are broader decisions that shape long-term directions and goals. For instance, allocating resources to invest in new medical technologies based on trends in patient needs and healthcare advancements.

**The Significance of Inferences And Decisions Play A Pivotal Role in Various Domains:**

- **Healthcare:** Inferences drawn from patient data enable personalized treatment plans, early disease detection, and informed medical interventions. Decisions based on these insights improve patient outcomes and optimize healthcare processes.

- **Business and Marketing:** Inferences about customer behavior and preferences drive targeted marketing campaigns. Decisions about product development, pricing, and distribution are guided by market trends and consumer insights.

- **Finance:** Financial data analysis leads to inferences about investment trends and risk assessment. Financial decisions regarding portfolio allocation and investment strategies are driven by these insights.

- **Manufacturing**: Analyzing production data helps infer efficiency and quality issues. Decisions related to optimizing production schedules and resource allocation are made based on these insights.

- **Education:** Inferences about student performance guide educational strategies and interventions. Decisions regarding curriculum adjustments and teaching methods are based on these insights.
- **Public Policy:** Inferences about social trends, economic indicators, and public health data inform policy decisions related to infrastructure, healthcare, education, and more.

## III. MACHINE LEARNING TECHNIQUES FOR HEALTH CARE ANALYTICS

Machine learning techniques have made significant contributions to healthcare analytics by enabling the extraction of insights, predictions, and patterns from vast amounts of medical data. Here are some machine learning techniques commonly used in healthcare analytics:

1. **Predictive Modeling:** Predictive modeling involves training a model to predict outcomes based on input features. In healthcare, this can be used for disease diagnosis, patient risk assessment, and treatment effectiveness prediction. For example, logistic regression is commonly used to predict the likelihood of a disease based on patient characteristics like age, gender, and medical history.

2. **Diagnosis and Image Analysis**

   - **Convolutional Neural Networks (CNN):** CNNs are deep learning models designed for image analysis. In healthcare, they're used to identify patterns in medical images like X-rays, MRIs, and CT scans. Layers of convolutions learn to recognize features like edges, textures, and shapes, making them invaluable for detecting abnormalities such as tumors or fractures.

   - **Transfer Learning:** This technique leverages pre-trained CNN models (often trained on massive image datasets like ImageNet) and fine-tunes them on medical image datasets. This approach benefits from the general features learned by the model and adapts them to specific medical image analysis tasks.

   - **Recurrent Neural Networks (RNN):** RNNs are suitable for sequential data, such as time series from patient monitors or electronic health records. Long Short-Term Memory (LSTM), a type of RNN, can capture temporal dependencies in patient data for tasks like predicting disease progression.

3. **Natural Language Processing (NLP)**

   - **Named Entity Recognition (NER):** NER identifies entities like drugs, diseases, and medical procedures in text data. In electronic health records, NER can help extract valuable information for clinical research or decision support.

   - **Sentiment Analysis:** NLP models can determine the sentiment or emotional tone in patient reviews, social media posts, or doctor-patient interactions. This insight can aid in understanding patient satisfaction and feedback.

   - Text Classification: Text classification involves categorizing medical notes, research articles, or other textual data into relevant categories. This assists in indexing and retrieving information from vast amounts of unstructured text.

4. **Clustering and Patient Stratification**

- **K-Means:** K-Means clustering groups patients with similar features. This can aid in identifying patient subgroups for personalized treatment plans or disease stratification.

- **Hierarchical Clustering:** This technique arranges patients in a tree-like structure, revealing relationships between subgroups. It can help uncover finer distinctions within larger patient populations.

- **Latent Dirichlet Allocation (LDA):** LDA is applied to text data, like medical articles, to identify latent topics. This aids researchers in understanding prevalent themes and trends in medical literature.

5. **Time Series Analysis**

- **ARIMA:** Autoregressive Integrated Moving Average models are used for time series forecasting. In healthcare, they can predict patient admissions, disease outbreaks, or resource demand.

- **Long Short-Term Memory (LSTM):** LSTMs are specialized for sequential data, making them suitable for predicting patient outcomes based on historical electronic health record data.

6. **Anomaly Detection**

- **Isolation Forest:** This algorithm isolates rare instances (anomalies) by constructing decision trees. In healthcare, it can detect unusual medical events or outlier patient records.

- **One-Class SVM:** This method identifies anomalies in a dataset by finding a boundary that best encompasses the majority of the data points. It's useful for detecting outliers in cases where the majority class is well-defined.

7. **Reinforcement Learning**

- **Treatment Planning:** Reinforcement learning can optimize treatment plans by learning from patient responses over time. It considers the sequential nature of treatment decisions and adjusts strategies accordingly.

8. **Clinical Trial Design:** Reinforcement learning can aid in designing adaptive clinical trials that dynamically adjust based on patient responses. This optimizes trial efficiency and maximizes patient benefit.

9. **Causal Inference**

- **Propensity Score Matching:** This technique helps estimate causal effects in observational data by matching treated and control subjects with similar propensity scores. It's used to assess the impact of treatments or interventions.

- **Instrumental Variables:** In healthcare, where conducting controlled experiments is often not feasible, instrumental variables can be used to address confounding variables and estimate causal effects of interventions.

10. **Survival Analysis:**

- **Kaplan-Meier Estimator:** This non-parametric statistic estimates the survival function of a population. It's used to analyze time-to-event data, such as patient survival times after a diagnosis.

- **Cox Proportional Hazards Model:** This model assesses how different factors (covariates) affect the hazard rate over time. It's often employed to analyze factors influencing patient survival times.

11. **Healthcare IoT and Wearables:**

- **Sensor Data Analysis:** Machine learning can process data from wearable devices and sensors to monitor patient health remotely, detect anomalies, and trigger interventions when necessary.

    These techniques collectively empower healthcare professionals to make more informed decisions, improve patient outcomes, and optimize healthcare processes by leveraging the wealth of data available in the field. It's important to note that the successful deployment of these techniques requires careful consideration of ethical concerns, data privacy regulations, and validation against clinical standards.

## IV. CONCLUSION

    Big data analytics are essential to the handling, governance, and application of healthcare systems' big data, analytics, and artificial intelligence. In this article, a machine learning strategy using little data that was developed from big data was discussed.Medical facilities work with both structured and unstructured data, which originates from databases, transactions, unstructured email and document content, devices, and sensors. In addition to the therapeutic setting, analytics are used in the administrative, business, and commercial sectors. It amply demonstrated how data-driven these decisions are. What has been analyzed in the literature is supported by the study's findings. The advantage of data-based healthcare is attracting medical facilities.

Big Data analytics has the ability to improve healthcare globally and have an impact. The defining of tactics used by medical facilities to promote and implement such solutions, as well as the advantages they obtain from using Big Data analysis, will be the focus of future studies on the use of Big Data in medical facilities.

**Future Work:** Big Data analytics in healthcare systems has its limitations. The issues of big data in practically every field are data collection, storage, sharing, searching, and analysis. Big Data analytics in healthcare systems also face problems related to standards for healthcare data, data security, and privacy, data quality, real-time processing, integration of heterogeneous or divergent data, and data security.

There are practically unlimited opportunities for future healthcare research. These suggest a strategy that may be applied in different healthcare applications and provide mechanisms to identify "patients" with relation to the usage of Big Data Analytics to diagnose certain illnesses. Big Data Analytics may also be utilized for research on the spread of diseases, the effectiveness of cancer treatments, or studies in psychology and psychiatry, such as emotion recognition utilizing deep learning and neural network algorithms.

## REFERENCES

[1] D. De Silva, F. Burstein, H. F. Jelinek, and A. Stranieri, "Addressing the Complexities of Big Data Analytics in Healthcare: The Diabetes Screening Case," Australasian Journal of Information Systems, vol. 19, Sep. 2015, doi: 10.3127/ajis.v19i0.1183.

[2] M. A. Srinuvasu, A. Koushik, and E. B. Santhosh, "Big Data Challenges and Solutions," International Journal of Computer Sciences and Engineering, vol. 5, no. 10, pp. 250–255, Oct. 2017, doi: 10.26438/ijcse/v5i10.250255.

[3] S. Sabharwal, S. Gupta, and K. Thirunavukkarasu, "Insight of big data analytics in healthcare industry," in 2016 International Conference on Computing, Communication and Automation (ICCCA), IEEE, Apr. 2016, pp. 95–100. doi: 10.1109/CCAA.2016.7813696.

[4] M. I. Pramanik, R. Y. K. Lau, H. Demirkan, and Md. A. K. Azad, "Smart health: Big data enabled health paradigm within smart cities," Expert Syst Appl, vol. 87, pp. 370–383, Nov. 2017, doi: 10.1016/j.eswa.2017.06.027.

[5] S. Dash, S. K. Shakyawar, M. Sharma, and S. Kaushik, "Big data in healthcare: management, analysis and future prospects," J Big Data, vol. 6, no. 1, p. 54, Dec. 2019, doi: 10.1186/s40537-019-0217-0.

[6] G. Harerimana, B. Jang, J. W. Kim, and H. K. Park, "Health Big Data Analytics: A Technology Survey," IEEE Access, vol. 6, pp. 65661–65678, 2018, doi: 10.1109/ACCESS.2018.2878254.

[7] Z. Sun and Y. Huo, "The Spectrum of Big Data Analytics," Journal of Computer Information Systems, vol. 61, no. 2, pp. 154–162, Mar. 2021, doi: 10.1080/08874417.2019.1571456.

[8] S. Dash, S. K. Shakyawar, M. Sharma, and S. Kaushik, "Big data in healthcare: management, analysis and future prospects," J Big Data, vol. 6, no. 1, p. 54, Dec. 2019, doi: 10.1186/s40537-019-0217-0.

[9] M. Mohammad Yousef, "Big Data Analytics in Health Care: A Review Paper," International Journal of Computer Science and Information Technology, vol. 13, no. 2, pp. 17–28, Apr. 2021, doi: 10.5121/ijcsit.2021.13202.

[10] K. Batko and A. Ślęzak, "The use of Big Data Analytics in healthcare," J Big Data, vol. 9, no. 1, p. 3, Dec. 2022, doi: 10.1186/s40537-021-00553-4.

[11] T. Ramesh and V. Santhi, "Exploring big data analytics in health care," International Journal of Intelligent Networks, vol. 1, pp. 135–140, 2020, doi: 10.1016/j.ijin.2020.11.003.