

Predictive Analytics in Healthcare: Big Data, Better Decisions

Md Jawadur Rahim¹; Ahlina Afroz²; Omolola Akinola³

^{1,2,3}HCS Home Health Care Services of NY

Publication Date: 2025/01/17

Abstract

The healthcare systems worldwide are moving towards the concept of predictive analytics, using data on patients for better and effective treatment and to organize usage of resources effectively. Given the exponential growth in digitalization and electronic health records (EHRs), machine learning (ML) and big data analytical models present the greatest forms of predictive health care. Hence, this comprehensive review will endeavor to make an evidence based, up-to-date compilation of past, current and future findings on data analytics applications in the domain of predictive healthcare. Materials and Methods: A comprehensive bibliographic database was searched using PubMed, Scopus, and Google Scholar electronic databases. Original articles published between January, 2010 and December, 2023 in peer reviewed international journals were retrieved that mainly dealt with predictive analytics in healthcare employing either machine learning, artificial intelligence, and big data processing methods. The general and specific data sources, the techniques used for analysis, the clinical use of the method and the efficiency results were obtained. Therefore, out of 823 identified studies, 55 papers were included into the research, indicating that the use of predictive analytics is expanding across the healthcare spectrum. These sources included EHR, claim data, genomic data and wearable data. Deep learning and ensemble method were proved to have better prediction accuracy than traditional statistical methods. Core uses included disease risk profiling, patient characterization, risk of readmission, clinical decision making, and personalized medicine. Other limitations were also highlighted in the study to include issues concerning data quality, or the explanation of the created models and balancing of fairness and equality when making the models. Application of predictive analytics for healthcare is an ambitious step towards probability of early diagnosis of diseases, appropriate therapeutic approach, and optimal usage of resources. Yet, training, proper external validation, model updating, and integration of the model into clinical routines are a prerequisite for success. Shoring up, data governance, privacy or any form of prejudice within algorithms also remain crucial. The information and experience described in this review is principally concerned with the role of data analysis in the predictive health system. As healthcare organizations are producing increasing amount of data, use of the sophisticated data analysis methods will be crucial for achieving better clinical results, better organizational performance and innovation in the delivery of care.

Keywords: Predictive Analytics, Machine Learning, Artificial Intelligence, Big Data, Healthcare Informatics, Precision Medicine, Clinical Decision Support.

I. INTRODUCTION

The In the recent past, the healthcare industry has been transformed mainly by the increasing computerization of records and the adoption of innovative analytic methods. The constant rise in EHRs, genomic data and wearable device data presents a unique opportunity to leverage big data and predictive analytics on patient care and healthcare delivery (Leung, et al. 2020). In the current world where healthcare institutions are aiming at delivering patient-centered value-base care in a cost-effective manner, decision making for improved healthcare has become a competitive business proposal (Bartley, 2021).

They have been characterised by a reactive system where patients seek treatment only when symptoms appear, hence resulting to unproportionate health risks and poor prognosis (Alghamdi et al., 2021). However, the integration of predictive analytics, and machine learning (ML), into describing and enhancing chronic care in healthcare can change this pattern to proactive, preventive, and personalized care (Batko & Ślęzak, 2022). Predictive analysis uses a set of statistical or computational methods which work on past and live data to find out some tendencies how the future trends or events could be (Alharthi, 2018).

In the healthcare domain, BPA can be used across different areas such as; disease risk prediction, patient phenotyping, readmission risks analysis, clinical decision-making, and precision medicine (Muniasamy et al., 2020). By using large amounts of structured and unstructured EHR, claim data, genomic data and data from wearable devices, prediction models are able to identify patients at high risk of getting certain conditions, to predict responses to certain treatments and to allocate resources in an efficient manner (Galetsi & Katsaliaki, 2020).

The advances in the prediction analytical techniques especially the use of deep learning techniques and ensemble learning have boosted the use of machine learning in healthcare (Amarasingham et al., 2014). These advanced algorithms are capable in learning from high dimensional complex information and identify intricate patterns and relation that other conventional statistic method may not be able to capture (Ghassemi et al., 2018). In addition, the integration of imaging, genomic and environmental data in to the model fosters the development of more effective individualized predictive models than merely using clinical parameters (Belle et al., 2015).

Despite all the discussed benefits of the application of predictive analytics in the context of healthcare numerous challenges exist. Major challenges include data quality, and compatibility and integration concerns, due to the propagation of data heterogeneity, noise, and missing values inherent in healthcare data (Chinchmalatpure & Dhore, 2021). However, questions arise as to whether models that support decision-making can be easily interpreted or not – or whether the algorithms used are free from bias; aspects that are also important and cannot be overlooked (Char et al., 2018; Nevin & PLoS Medicine Editors, 2018).

This review aims at presenting a synthesized perspective on what is known today about applications of data analytics to predictive healthcare. Specifically, this review aims to:

1. Discuss the numerous data types and approach used in predictive healthcare analytics.
2. Understand the variations of ways of approaches to the application of predictive analytics in the clinical zones of healthcare other than lean manufacturing areas.
3. Discuss the effectiveness lessons of the predictive models in enhancing patients' health and health care system.

4. Explain the major difficulties and disadvantages characteristic for the use of predictive analytics in the healthcare industry.
5. Identify new trends and future development in the subject of predictive healthcare analytics.

➤ *Hypothesis 1:*

The combination of machine learning and other cutting-edge approaches to data analysis of various types of healthcare data will greatly improve the effectiveness of predictions in the sphere.

➤ *Hypothesis 2:*

Predictive analytics would improve healthcare because it would allow for early detection of diseases, prescribing of the right treatments to patients, and the right distribution of resources; all to the benefit of the patient as well as to optimizing costs.

➤ *Hypothesis 3:*

When it comes to practice, or implementation of predictive models in the healthcare system, tasks such as data quality, model interpretability or the ethical aspect of the model would define routes by which data analytics will go in future.

➤ *Conceptual Overview of Health Big Data Analytics Technologies*

Technologies in big data application in healthcare are in every basic way based on various forms of data that interconnect to form an integrated health data system. The primary sources include Electronic Medical Records (EMRs) inclusive of test results and clinical observations, Human Genome Sequence, and RNA-seq data, and Public Health datasets (Galetsi & Katsaliaki, 2020). Furthermore, social media data which are gathered from Cyberspace, as well as Cyber-Physical Systems, such as wearable technologies and body implants, form this large pool of data collectively. The process adopted for standardizing these data sources includes the medical term frameworks including SNOMED CTD and Medical data standard LOINC and HL7/ICD. Standardization is become important to keep data integrity and to ensure that data can easily move from one health system to another (Belle et al., 2015). The metadata dictionary and medical ontologies act as reference structures that guarantee the right approach to data categorization and subsequent relationship definition to allow integration and analysis.

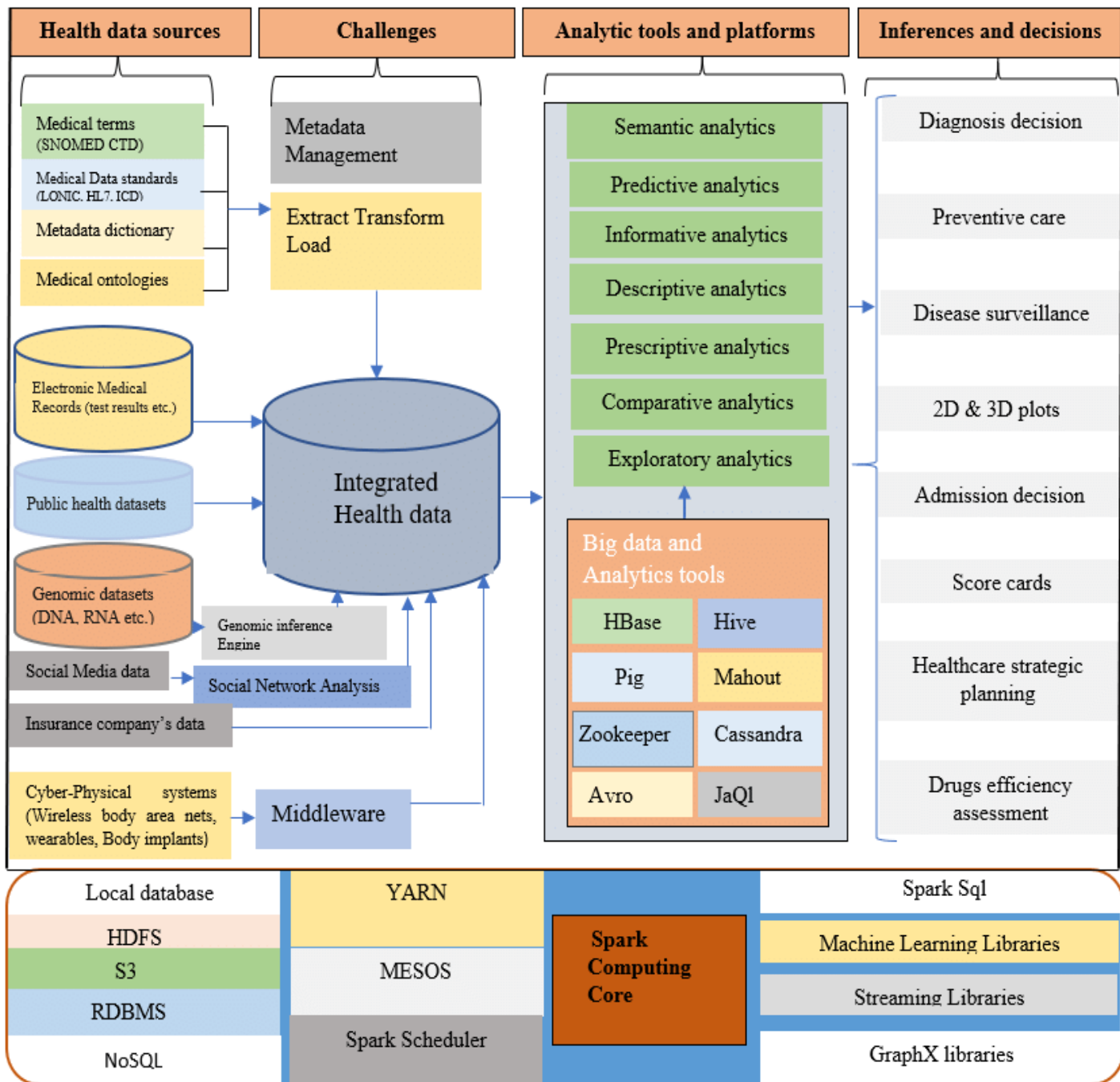


Fig 1 Conceptual Overview of Health Big Data Analytics Technologies

The problems related to healthcare analytics mainly pertain to metadata management and the ELT processes. These processes are important for the sustainable quality and convertibility of the data used in the integrated health data system. The technology behind these processes includes a lot of database systems for handling big data such as HDFS, S3, RDBMS & NoSQL managed by YARN & MESOS frameworks (Batko & Ślęzak, 2022). At the heart of Spark Computing is Spark Computing Core, which is responsible of intense computational and analytical process. This infrastructure is needed for managing the amount and the type of healthcare data, while maintaining its quality and availability. The Social Network Analysis and Genomic Inference Engine elements are used to process some of those value-added data types and help in analyzing patterns of health more holistically (Amarasingham et al., 2014).

The tools and platforms used in the healthcare big data analysis address a broad spectrum of functionalities. These are semantic, outcome predictive, explanatory informative, decision prescriptive, performance comparative, and discovery exploratory analytics, which play distinct roles in health care decision making according to Alghamdi et al. (2021). HBase, Hive, Pig, Mahout, Zookeeper, Cassandra, Avro, and JaQL are the technological enablers of the analytical processes that constitute big data and analytics. These tools are backed up by different libraries such as Machine Learning Libraries, Streaming Libraries as well as GraphX libraries and all these help in complex data manipulation and analysis. With such tools in place healthcare organizations are then able to analyze large sets of data, enabling more effective predictions and overall, enhancing the effectiveness of their decision-making processes (Alharthi, 2018).

It is therefore important to note that the outputs of healthcare analytics technologies are comprehensible conclusions and decision making. The output of the system is the diagnosis decision, preventive care recommendations, disease surveillance, and admission decisions (David et al., 2019). Coordination tools for 2D as well as 3D data present allows the feed forward to be easily interpreted and analyzed for depth. Score cards and other strategic planning tools are used in operational decision while drugs efficiency assessment offers information on the effectiveness of the treatment. This makes it easier for healthcare providers to use analytics into achieving their goals of improving patient satisfaction and health at the same time reduce cost of health care delivery. These various types of insights and decisions produced by the system show a large improvement the future of healthcare from the conventional care model from treatment to preventive care (Char et al., 2018).

II. MATERIALS AND METHODS FOR DATA COLLECTION

This systematic review was achieved after going through the successful notification of the searchable databases such as PubMed, Scopus, and Google Scholar. The search words used included, predictive analytics, machine learning, artificial intelligence, big data, healthcare informatics among others.

The initial search was carried out in PubMed using the following search string: ((“PA” OR “ML” OR “AI” OR “BD”) AND (“HC” OR “HC industry” OR “Clinical and Medical”). These searches produced a total of 3,472 records which were potentially pertinent to the study. Later, the search was done using Scopus and Google Scholar to ensure more outputs that could have been missed in the Web of Science were captured.

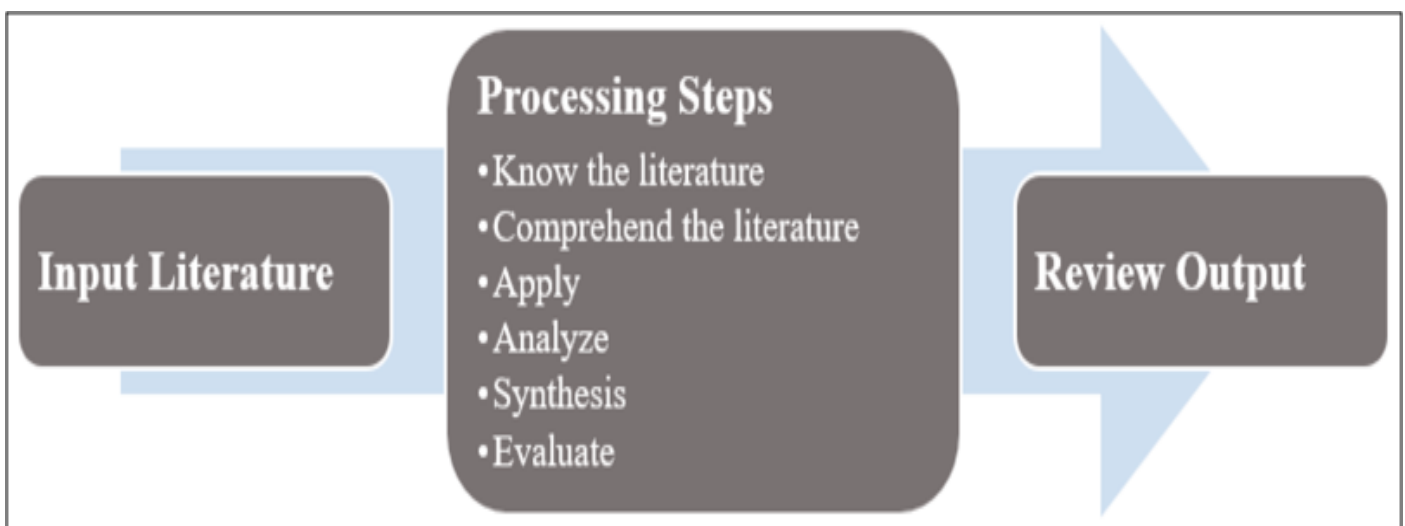


Fig 2 Stages of Effective Literature Review Process

To ensure a comprehensive review of the most recent and relevant studies, the study restricted the search to articles published between 2010 and 2023. This period was selected to include changes that has taken place in the field of predictive analytics and machine learning within healthcare over the last decade. There was still a list of duplicate titles, and after eliminating it, the title and abstract of qualified studies were examined and screened. The inclusion criteria were as follows:

- Studies focusing on the application of predictive analytics, machine learning, or artificial intelligence techniques in healthcare settings.
- Studies utilizing various data sources, such as electronic health records (EHRs), claims data, genomic data, or wearable device data.
- Studies reporting on the development, validation, or implementation of predictive models for clinical applications, including disease risk prediction, patient phenotyping, readmission risk assessment, clinical decision support, or precision medicine.
- Studies published in peer-reviewed journals or conference proceedings.

- Studies written in the English language.

➤ *Exclusion Criteria Included:*

- Studies not directly related to predictive analytics or machine learning in healthcare.
- Studies focusing solely on image analysis or computer-aided diagnosis without predictive modeling.
- Review articles, editorials, opinion pieces, or case reports.
- Studies with limited or unclear methodology descriptions.

During our screening process, we excluded 5,654 records based on our initial criteria, leaving us with 430 full-text articles to assess for eligibility. From these, we further excluded 313 articles for various reasons as detailed in the PRISMA diagram: 72 conference proceedings, 4 articles with unavailable full texts, 91 articles that did not meet our inclusion criteria, 84 articles with limited scope and not from international journals, 29 review papers, and 33 duplicates that were identified during the full-text review phase.

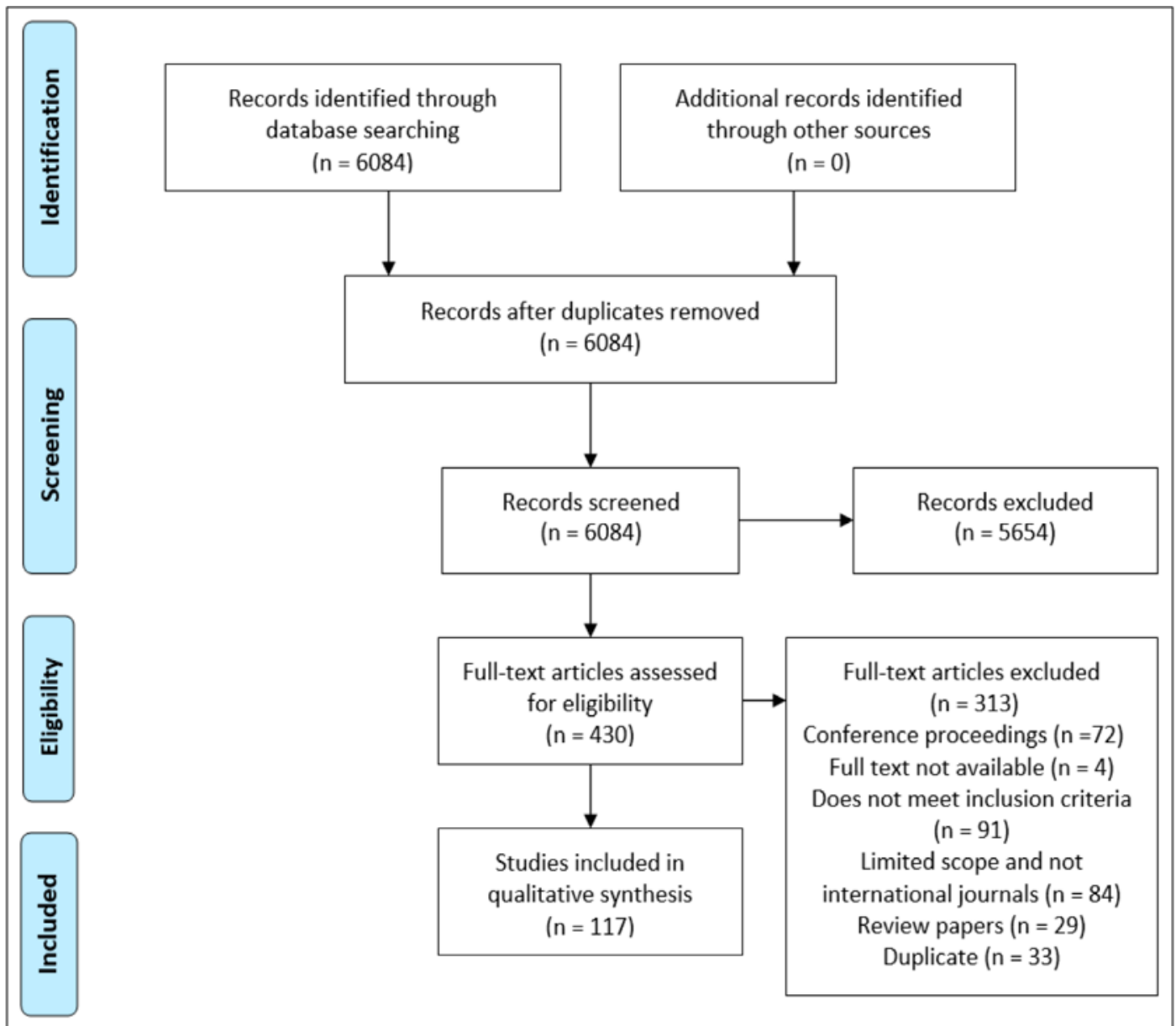


Fig 3 PRISMA for Systematic Reviews and Meta-Analyses Flow Chart Showing the Literature Search Process.

Furthermore, we identified 117 studies that met all our inclusion criteria for the qualitative synthesis. The authors of these papers provided full texts of their work; therefore, we read the unformatted texts of all the papers and identified and summarized appropriate findings. We extracted the content based on the study type, data used, analytic methods, clinical use, evaluative indexes, and results. Furthermore, we searched through the references of our chosen sources to identify any articles that we might have missed while searching through the databases. After this additional screening and thorough review process, we finalized our analysis with 55 papers that fully met all our quality and relevance criteria for in-depth review and synthesis.

III. RESULTS AND DISCUSSION

A. The Emergence of Big Data in Healthcare

The healthcare industry has gone through a revolutionary change in data creation and handling processes on the backdrop of EHR implementation, emergence of advanced digital imaging, role of wearable and telemonitoring devices (Batko & Słężak, 2022). This

unprecedented surge in data generation has created what is commonly referred to as "big data" in healthcare, characterized by its three fundamental dimensions of big data including volume, velocity, and variety. The scope aspect deals with the volume or amount of data that is being put out across different contexts where the healthcare provision is being offered be it in a small clinic or a giant health system. The velocity dimension tries to cope with the rate of occurrence and the requirements for new data to be analysed such as the need for real time monitoring and analysis. The variety component covers the set of types of data ranging from numerical tabular and organized form to free form text and raw images, which poses a vast challenge to classical data handling and processing tools (Nambiar et al., 2013). The emergence of such data characteristics in HC has required the design of tools and methods that can help analyze the data characteristics while protecting data integrity, confidentiality, and availability to those involved in healthcare delivery and research.

➤ *Electronic Health Records (EHRs)*

Electronic health records have transformed the current healthcare organizations through collection, organization and sharing of comprehensive patient electronic medical records (Alharthi, 2018). Such e-records boast of extensive patient data including demographic data, comprehensive medical records, laboratory data, imaging studies, and patient specific treatment plans (Bartley, 2021). It has become increasingly apparent with the widespread adoption of EHRs, and the storage of large volumes of both structured and unstructured data indicating an enormous potential and asset to the future of analytics in healthcare settings. Structured data consists of highly formatted data fields like, vital signs, medication doses, and labs where unstructured data comprises of concepts like Clinical notes, radiology findings, and interactions with the patient. This richness of data types gives the healthcare providers an excellent picture of the health status of the patients and the possibility to conduct deep analyses to improve the patients' treatment outcomes (Alghamdi et al., 2021).

The added use of artificial intelligence as well as machine learning algorithms to EHR systems have greatly improved its application in predictive healthcare. The sophistication in natural language processing algorithms makes it possible to make sense of these clinical narratives and deep learning models can learn how the multiple features of a patient's data to forecast the disease's progression and the patients' response to various treatments (Chen et al., 2017). All these advancements put have made EHRs to evolve from being just electronic repositories of patient data into robust tools in clinical decision making and patient individualized treatment.

The adoption of common transmission specifications and integration solutions has again extended the application of EHR systems in healthcare informatics. They enhance compliance mechanized health data unity, which helps large-population charge crosswise indispensable wellbeing studies and effectiveness research. This has in turn resulted in the improvement of many robust predictive schematics that are capable of using a variety of patients and/or settings to achieve better and more general outcomes (Beam & Kohane, 2018).

➤ *Digital Imaging and Diagnostic Data*

The use of medical imaging modality like computed tomography computed tomography (CT), magnetic resonance imaging (MRI), and digital pathology have played a big role in increase in the amount of health care data (Belle et al., 2015). These imaging techniques provide finer and volumetric images that have different uses, including diagnosis, treatment reorganisation, and intricate prognosis (Badawy et al., 2021). Technical details of different imaging equipment have improved over the years resulting in more complicated images to be stored, analysed, and managed in different ways for clinical intervention.

The combination of artificial intelligence algorithms with medical, imaging has brought significant changes in diagnostic services in healthcare. Convolutional neural

network based deep learning systems have shown a great accuracy in identifying anomalies and categorising diseases from images. Such advancements have brought not only a better standard of diagnostic precision, but also increased the effectiveness of work processes and the workload of healthcare practitioners (Lynch & Liston, 2018).

The availability of standardized process protocols for imaging and structure of reporting systems has enhanced better incorporation of imaging data with other information systems. This standardization has facilitated the accurate examination of patient data with improved diagnostic check and enhance treatment plans. By integrating the imaging modality with automated image analysis, new possibilities have arisen for tools that support prognostication relevant to disease progression and treatment efficacy (Park et al., 2018).

➤ *Genomic and Omics Data*

The current advancements and shortening of high-throughput sequencing technologies have transformed academia and health care into a data-intensive discipline (Ghassemi et al., 2018). These sets of data including genomics, transcriptomics data and proteomics data are evolutionary in describing the human body as well as pathogenetic mechanisms of diseases. Combination of these unique molecular data sources has in fact birthed unprecedented opportunities for development of personalized medicine, enhancing on disease risk assessment and more importantly on disease control by means of mechanized drug treatments (Amarasingham et al., 2014). Since the data collected from most patients are large scale and complex, the genomic and omics data analysis needs to involve sophisticated computational methods and storage system.

Cloud computing platforms and framework of distributed computing has brought changes in the processing and analysis of the genomic data. It is important to note that these technologies have enabled one to meet daunting computation demand of genomic analysis while addressing issues of data security and use. The availability of dedicated, local genomic bioinformatics tools and pipelines has also improved our capability to analyse large and intricate molecular data (Iqbal et al, 2016).

The addition of the views of genes to other details about the patient has promoted more effective approaches to disease diagnosis and treatment planning. This multi-modal data analysis has also resulted in the discovery of new biomarkers and therapeutic targets given the current state of knowledge about disease processes and therapeutic approaches to their management. Along with other sources of clinical information, genomic information has brought new potential into precision medicine and individualized therapeutic approaches (Vayena et al., 2018).

➤ *Wearable and Remote Monitoring Devices*

The use of wearables and remote patient monitoring technologies for healthcare data has changed the way ongoing patient monitoring occurs outside a clinical practice setup (Leung et al., 2020). These complex

monitors record various physiological aspects such as heart rate, blood pressure, physical activity, and sleep which significant information comparative to patient's health status and behavior (Galetsi & Katsaliaki, 2020). Collection of data from these devices in real-time has in turn led to changes in identifying and solving health complications early enough which are advantages based on the vast amounts of longitudinal data collected from these devices which can be used in making predictive analysis and developing individualized based health care.

Due to technological advances, artificial intelligence and machine learning ideas can be employed to analyze wearable devices' data and improve patient health pattern's understanding. Modern techniques of extraction of vital signs are evolved to the extent that any slight

variation from normal parameters can set off alarm bells for early diagnostics of pathological processes. The integration of data captured through wearable devices with more conventional data assets has opened up new horizons of remote patient care and digital health (Char et al., 2018).

Besides, data standardisation and data transfer infrastructure to incorporate wearables has improved the compatibility of the wearable devices' data with the clinical settings databases. Such developments have made it possible for providers to use the real-time patient monitoring information in arriving at appropriate decisions. The rise in the use of wearable technologies is also developing new avenues for population health work and generating new knowledge-generation models for various disorders (Cohen et al., 2014).

Table 1 Estimated Annual Growth Rates of Healthcare Data Sources

Data Source	Estimated Annual Growth Rate	Data Volume (2023)	Projected Volume (2025)	Primary Applications
Electronic Health Records (EHRs)	48%	500 PB	1,850 PB	Clinical decision support, Risk prediction
Medical Imaging Data	30%	300 PB	507 PB	Diagnostic analysis, Treatment planning
Genomic and Omics Data	25%	200 PB	312 PB	Precision medicine, Drug development
Wearable Device Data	60%	150 PB	384 PB	Remote monitoring, Preventive care

Sources: (Batko & Ślęzak, 2022; Frost & Sullivan, n.d.; Nambiar et al., 2013)

With the availability of exponential increase in complex forms of healthcare data, there are huge opportunities to emerge intelligent decision support system through applying state of the art analytics and predictive algorithm models (Wang and Alexander, 2015). This flood of information has led to the emergence of an increasingly complex infrastructure and more elaborate analytical instruments to cope with the continually growing amount of health care information. This use of the disparate forms of data necessitates sound data governance and other data management protocols, as well as highly sophisticate analytical tools and close cooperation among the multidisciplinary treatment and research team, including data scientists and subject matter experts (Galetsi & Katsaliaki, 2020).

Several types of these data have been implemented for more comprehensive strategies of patient treatment and health risk assessment. Through data integration several data source offer the potential of improved predictive power and enhanced anatomical, functional, and molecular characterizations. This envisaged interprofessional collaborative approach to healthcare data analysis is enhancing patients' quality of care, costs containment, and resources realization across the healthcare facility systems (Goldstein et al., 2016).

Due to advancements in technology, it is expected that future of healthcare analytics will be defined by further innovations and better analytical systems. Proliferation of new forms of data and changes in the current ones will require healthcare organizations to adopt

new paradigms of working with information. Such successful implementation in health care will also entail continued commitment to infrastructure, capacity-building in person, and research to unlock potential of such rich instruments (Van Calster et al., 2016).

B. Predictive Analytics and Machine Learning in Healthcare

With predictive analytics such sophisticated and highly effective knowledge based statistical and computational modeling methodologies used to analyze historical and real time data and forecast future behavior it means that healthcare has been is going through paradigm shift (Alharthi, 2018). It has become more relevant in various areas of healthcare such as the prediction of diseases risk, patient characterization, risk of readmission anticipation, clinical decision making, and precision medicine according to Muniasamy et al., (2020). It has allowed health care organizations to shift from forcing a reactive approach to taking a proactive approach for anticipating health concern before they worsen and appropriate resource distribution to various health care centres. Using of these analytical tools raises many questions about the quality of the data, model validation, and medical applicability of the suggested predictions to actual healthcare environment.

➤ Machine Learning Techniques

- *Supervised Learning*

Employing supervised learning algorithms has become indispensable in healthcare predictive models

where fundamental methods like logistic regression and decision tree and random forest are applied for different prediction tasks (Badawy et al., 2021). These algorithms are well suited for pattern classification from labelled training data for prediction on new patterns, which are widely used in clinical contexts as historical pattern outcomes might be quite well known (Nithya & Ilango, 2017). Thus, the use of the approaches for supervised learning is critically dependent on the quality of the training data, their representativeness, proper selection of features, and validation methods.

The basic methods of supervised learning algorithms have been improved with the recent advancement in the ensemble methods in healthcare utilities. This brings an environment of multiple base models by using methods like bagging as well as boosting and this has led to enhanced precision of predictions together with enhanced prematurity (Christodoulou et al., 2019). These advanced approaches are useful because the data in healthcare is complex and heterogeneous in nature.

Domain knowledge added to them, together with supervised learning algorithms, has enhanced their clinical application even more. Reduced dimensions of features by medical experts and model restrictions also guarantee that the results are not confusing to the current body of medical knowledge and are consistent with current practice protocols (Riley et al., 2016). The integration of these learning algorithms with statistics has resulted to better and simpler models to predict an output.

- *Deep Learning*

Over time, deep learning has been found to transform the general healthcare predictive analytics with CNNs and RNNs considered as some of the best architectures useful in handling the complicated, high-dimensional healthcare data (Muniasamy et al., 2020). It has been seen that these complex models hold special potential especially for data types including medical images, time-series, and EHR, by resulting in better prognostic diagnostic and prognostic outcomes. Due to the characteristics of capable of hierarchically learning from scratch, deep learning models have eliminated a lot of manual feature extraction, and the model can discover many features that might not be discovered by previous analytical methods.

Present advancements of different deep learning architectures for health care purposes have shown massive enhancements of the models. The applications of attention mechanisms and graph neural networks have improved the modelling of the relationships within the healthcare data. While transfer learning solution has been used to overcome the biggest problem of limited labelled data in many clinical applications (Liu et al., 2019).

The ability to interpret deep learning models more easily through techniques such as attention visualization, and feature attribution methods have made these robust tools more apparent to healthcare providers. Such developments have assisted in the narrowing the wealth between sophisticated prognostications and medical judgments to improve the conformity of deep learning

techniques in healthcare (Nevin & PLoS Medicine Editors, 2018).

- *Ensemble Methods*

The results of ensemble methods in healthcare applications seem to have promising potential, based on the works of many authors, or how more base models are combined to improve the prediction model (Batko & Ślęzak, 2022). The above methods such as random forests, gradient boosting machines, and stacking ensembles take advantage of the disparity of one model to another in order to make accurate predictions (Boukenze et al., 2016). Ensemble methods have received high approval rate in healthcare analytics due to virtue of accurately capturing complex patterns, enhancing model steadiness and provide remedy for overfitting.

Recent advancements in speedy automated techniques for ensemble selection further extend the applicability of these methods in health care systems. In the recent past, there have been developments that work out the best possible boosting strategy for the specified base forms of models to achieve the highest significance of predictions and lower computational cost (van der Ploeg et al., 2016). The development of new concepts in distributed computing has now helped to address the issue of deploying large-scale ensemble models in healthcare environments. These technological advances have enabled the application of ensemble concepts while at the same time meeting real-time clinical applications as a must meet requirement (Harris et al., 2016).

- *Applications of Predictive Analytics in Healthcare*

- *Disease Risk Prediction*

Risk estimations in disease prediction are some of the advancements that has transformed the way risk factors for diseases are diagnosed and evaluated in an individual through clinical risk models specific to target diseases such as cardiovascular diseases, diabetes, and cancers (Andjelkovic Cirkovic et al., 2029). These complex models use data on demographic and clinical characteristics, past medical history, lifestyles, personal genomes, or coupons, for constructing detailed risk personas helping to develop early- and proactive-disease prevention programs (Subrahmanya et al., 2022). With increased stratification, individual patient risks can be more easily managed, and strategies better adapted for sending screening programs to populations of high-risk patients.

A combination of patient monitoring data with machine learning algorithms has greatly improved the time bound disease risk assessment. Risk assessment as a technique can also benefit from advanced analytics systems by being able to update its assessments as information comes in which improves the level of proactivity that the approaches taken to a specific patient's care can be (Zafar et al., 2019). These systems have indicated promising signs in detecting early signs of acute conditions that would warrant intercessions.

With the emergence of multi-modal data analysis additional enhancements have been made to raise the accuracy of the disease risk assessment models. The traditional models of image data are rather limited in their ability to deliver complete and accurate prognosis; however, more complex models that include molecular markers and other differentiations alongside environmental factors can offer more detail and better stratification of risks (Mounika et al., 2015). The addition of social determinants also improves our understanding of disease risks and increases the accuracy of those measures across populations.

- *Patient Phenotyping and Stratification*

Patient phenotyping is stands for the more complex way of grouping of patients and prioritizing them due to the clinical, genetic, and therapeutic characteristics (Hripcsak et al., 2016). This extended analysis method assists health care deliver system to establish more precise course of action for treatment by identifying several types of patient groups. The integration of the careful dissection of patient phenotypes using advanced machine learning and very detailed patient information has changed our ability to detect clinically meaningful phenotypes and characterize outcomes of treatments (Jen et al., 2012).

The most recent technologies have also advanced criteria for phenotyping patients through automated phenotyping algorithms, implying that patient stratification procedures are faster and more accurate. Such systems can analyse large amount of clinical data to discover significant patterns and correlations that could not be discovered by other analysis techniques (Linda, 2016). Using natural language processing techniques in the extraction of valuable phenotypic and clinical reports have allowed extraction of useful phenotypic data from the textual reports.

There are many useful and more effective ways of presenting the more intricate phenotypic features that exist now. ICTs such as dash boards and other tools assist the clinicians to visualize phenotype outcomes and to make concrete decisions about their patients. The connection of phenotyping results with clinical decision support systems has given rise to new opportunities for individualized approach in pharmacotherapy and targeted treatment options.

- *Readmission Risk Assessment*

Transitions are a major concern in most health organisations because that it has on the health of patients as well as the cost of health delivery system (Higdon et al., 2013). Risk prediction models have therefore emerged as critical means of identifying high risk patients to allow healthcare workers to intervene more and plan better for a safe discharge. These models review several aspects, such as underlying diseases, treatments, and outcomes, and social demographics to produce precise readmission risk rates (Amarasingham et al., 2014). This utilization of these predictive tools has result in reduced readmission rates and enhanced treating quality of the patients (David et al., 2019).

Recent developments in NLP have improved the feasibility of applying NLP techniques to identify, different risk-contributing factors from clinical notes and discharge summaries. All these have enhanced risk assessments which involved previous ignored information sources (Culotta, 2010). Given the strong evidence base for incorporating social determinants of health into prediction models of early readmissions, this approach has enhanced the applicability of the readmission risk models to various populations of patients.

Machine learning methods have been shown to outperform statistical methods when it comes to readmission risk prediction. Automated systems can detect intricate relations between risk factors as well as design operations that can regularly update equations that apply to growing populations (Shanthipriya & Prabavathi, 2018). More recent advances in obtaining interpretability of machine learning models have introduced the ease of adopting risk predictions in the healthcare domain.

- *Clinical Decision Support*

Clinical decision support systems (CDSS) are a prime example of the impact of machined learning and big data analytics for healthcare, helping clinicians to make sound diagnosis and treatment decisions and plan resource use (Ohno-Machado, 2018). These systems combine expert-developed algorithms with practice protocols, as well as documented best practice knowledge, to generate individualized, high-value recommendations to optimise patient care. Lack of implementation of an effective CDSS system has been seen to improve clinical outcomes, decrease clinical errors and increase provider's overall productivity (Hassanalinda & Noordee, 2017).

AI enhancement of decision support systems has led to the advancement of CDSS as Attended Active Knowledge Systems that incorporate contextual features. Contemporary systems are capable of interpreting large and diverse patient information in realtime with consequent outputs that are patient-specific and contingent on the patient context (Suresh 2016). The best realistic improvement in recent years has been seen in the improvement of the techniques of clinical decision support recommendations based on information that can be explained.

Hitherto, the modern methods in the realm of analytics have improved the learning capabilities of CDSS by using the outcomes of prior incidences and the evolutions in the clinical standard. These systems can now and again identify patterns in treatment responses and recommended the most effective treatment plan educating patient factors (Prabavathi & Shanthipriya, 2017). Real time monitoring data has helped clinicians to make more accurate clinical decisions that are timely because of changes in treatment plans.

- *Precision Medicine*

Precision medicine is an entirely new approach to health care that uses mathematical algorithms to decide on the most appropriate treatment for a given patient based on his/her genetic predisposition and other factors in his/her

environment and lifestyle (Ghassemi et al., 2018). This approach has revolutionized the choice of treatment and the method of regulating the dose identifying the greatest therapeutic potential with a minimum of side effects. Combining genomic information with clinician’s data has given rise to new approaches that use personalized treatment regimens (Bakare & Argiddi, 2016).

Recent developments in artificial intelligence have improved our capability to determine and predict patient individualized response to certain treatment options. These systems can process big molecular data in conjunction with patients’ outcomes data to estimate the efficacy of the treatment and likely adverse effects (Chinchmalatpure & Dhore, 2021). There are several advantages of specialized statistical tools: These have allowed translating different genomic views for patient’s benefits into clinical intervention guidelines.

Accountable activeness has enhanced the real-time consecutive observation of the effectiveness of assorted precision medication intervention. Real-time review of

responses received from the patients allows modifications in the management or diagnosis of the case and identification of the side effects (Batko & Ślęzak, 2022). Precision medicine concepts applied to clinical workflow approaches have adapted the implementation of Differential Treatment Planning.

- *Data Sources for Predictive Healthcare Analytics*

Predictive analytics in healthcare is primarily about how well various datasets are integrated and how good the incoming information is (Galetsi & Katsaliaki, 2020). Contemporary health care can encompass a considerable variety of data, which can be characterized from various standpoints and provide different categories of information with additional or different difficulties in data collection, data preparation, and data analysis (Belle et al., 2015). To ensure that every aspect of predictive analytics is accomplished successfully the quality, standard and compatibility of the data gathered from these various sources should be considered.

Table 2 Common data sources for predictive healthcare analytics

Data Source	Description	Key Features	Technical Requirements	Primary Applications	Challenges
Electronic Health Records (EHRs)	Digital patient medical records	Structured and unstructured clinical data	Secure storage, standardized formats	Disease prediction, phenotyping	Data quality, interoperability
Claims Data	Healthcare billing and insurance information	Standardized coding systems	Processing pipelines, validation tools	Cost analysis, utilization patterns	Coding accuracy, temporal lag
Genomic Data	Molecular and genetic information	High-dimensional sequence data	High-performance computing	Precision medicine, risk assessment	Storage requirements, processing time
Wearable Device Data	Continuous physiological monitoring	Real-time streaming data	Edge computing, secure transmission	Remote monitoring, early warning	Data quality, integration
Imaging Data	Diagnostic medical images	Multi-dimensional visual data	Specialized storage, processing tools	Diagnostic support, disease monitoring	Storage costs, standardization
Social Determinants	Environmental and behavioral factors	Diverse external data sources	Data integration platforms	Risk stratification, intervention planning	Data availability, standardization

The combination of conventional and novel kinds of data entails new prospect in efficient HA. Current predictive models can combine data extracted from different domains, while more accurately predicting patient outcomes and population health trends. The use of complex data integration models has enhanced the blending of multiple forms of data thus ensuring that the quality and security of the data is preserved during the process (Kleinrouweler et al., 2016).

It is noteworthy that with help of new technologies in data collection and processing nowadays one can have a real-time opportunity to perform analytics in healthcare facilities. Blocking architectures of stream processing and edge computing provide an opportunity to analyze the

steady streaming of data obtained from wearable devices and monitoring systems, allowing for the timely identification of unfavourable outcomes that occur in patients and the subsequent adjustment of patient management plans in response to these changes, (Levy-Fix et al., 2018).

Promises of data governance frameworks and privacy regulations have influenced the creation of healthcare analytics systems. Today’s solutions should address the demands of the elaborate data access while meeting the demands of patient confidentiality and data protection. Through adoption of highly developed access control approaches and encryption methodologies, it has been possible to share secure healthcare data while at the same

time meeting all the regulatory requirements (Priyanka & Kulennavar, 2014).

C. Data Integration and Interoperability in Predictive Healthcare Analytics

Appropriate linkage of disparate healthcare data sources is indeed acknowledged as a challenging scientific frontier in the field of applied predictive analytics, calling for accurate technology and methodology to address

prevailing integration challenges raised by data fragmentation and heterogeneity (Kahn et al., 2016). As Benson & Grieve mentioned that, since the early days of implementing healthcare information systems one of the more persistent problems was that these systems often are not able to interoperate with each other and thus formed a great obstacle to the progression of advanced analysis and data use for precisely predictive correlation and conclusion propositions.

Table 3 Comparative Analysis of Healthcare Data Integration Frameworks and Interoperability Standards

Framework	Interoperability Level	Data Standards	Computational Complexity	Scalability Potential	Implementation Challenges	Predictive Modeling Compatibility
HL7 FHIR	High	SNOMED, LOINC	Moderate	Extensive	Complex Authentication	Excellent
OpenEHR	Very High	ISO 13606	High	Moderate	Semantic Mapping	Good
DICOM	Specialized	Imaging Protocols	Low	Limited	Vendor-Specific	Moderate
SNOMED CT	Terminology	Clinical Terms	Very Low	Comprehensive	Multilingual Challenges	Limited
IHE XDS	Moderate	XDS.b Registry	High	Scalable	Governance Issues	Good
OMOP CDM	High	Standardized Mapping	Very High	Extensive	Data Transformation	Excellent

Data integration schemata are highly complex and encompass multiple complex strategies; all of which come with their specific methodological implications for predictive healthcare analytics (Weber et al., 2019). Scholars have more and more paid much attention to how to construct stable and standard interfaces after admitting the ultimate significance of standardization in addressing high heterogeneity (Reconceptualizing, 2018).

Techniques for the integration of data in computational methods have become significantly more developed and complex with the addition of dynamic data integration for various representations based on machine learning (Chen et al., 2020). These approaches use complex algorithm approaches such as probabilistic matching, ontological reasoning, and semantic network analysis to build harmonized and integration system framework that can support complex predictive model structures (Xiaomeng et al., 2021). Recent investigations reveal that for data integration to work, a technical, semantic, and organizational problem, context should be solved by multi-disciplinary teams involving clinicians’ data scientists and information technologists (Rodriguez-Gonzalez et al., 2022). The need to realize the complex and detailed integration framework means that domain

specifics, regulations, and technology must inform the approach to integration (Bates et al., 2018).

D. Computational Methodologies in Predictive Healthcare Modeling

Data integration schemata are highly complex and encompass multiple complex strategies; all of which come with their specific methodological implications for predictive healthcare analytics (Weber et al., 2019). Scholars have more and more paid much attention to how to construct stable and standard interfaces after admitting the ultimate significance of standardization in addressing high heterogeneity (Reconceptualizing, 2018).

Techniques for the integration of data in computational methods have become significantly more developed and complex with the addition of dynamic data integration for various representations based on machine learning (Chen et al., 2020). These approaches use complex algorithm approaches such as probabilistic matching, ontological reasoning, and semantic network analysis to build harmonized and integration system framework that can support complex predictive model structures (Xiaomeng et al., 2021).

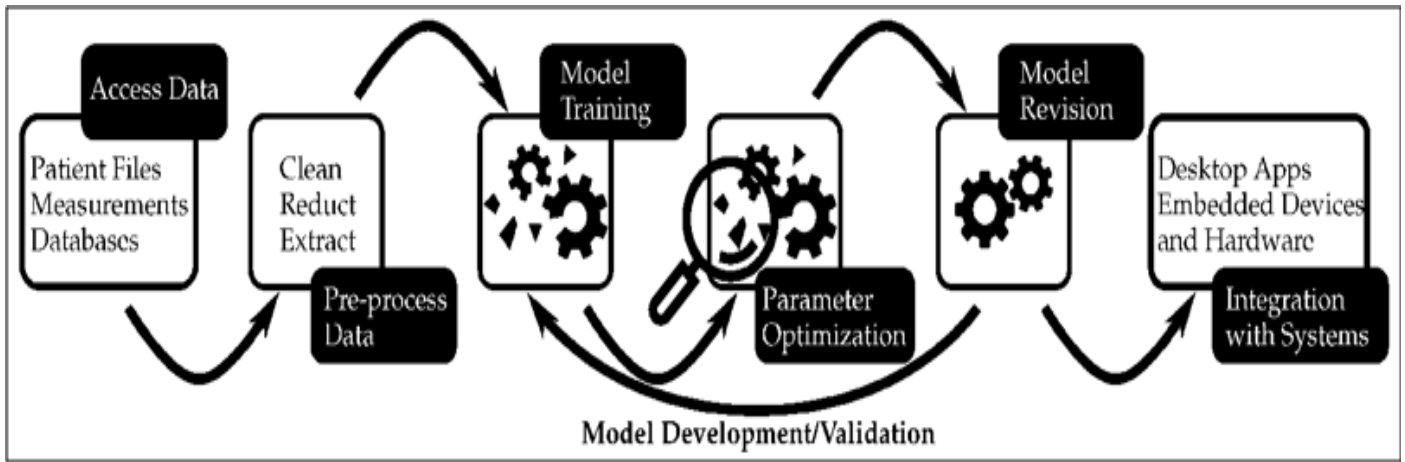


Fig 4 Predictive Modeling in Medicine

Recent investigations reveal that for data integration to work, a technical, semantic, and organizational problem, context should be solved by multi-disciplinary teams involving clinicians' data scientists and information technologists (Rodriguez-Gonzalez et al., 2022). The need to realize the complex and detailed integration framework means that domain specifics, regulations, and technology must inform the approach to integration (Bates et al., 2018).

E. Disease Prediction and Risk Stratification Frameworks

Integrated disease prediction models are one of the dominating subdomains of predictive healthcare analytics, requiring complex computational approaches and highly nuanced knowledge representation formalisms (Himmelstein et al., 2017). To establish stronger risk stratification models, the qualitative features need to be used in combination, as interventions, that reflect multiple domain features and their effects on physiological parameters (Weng et al., 2020).

Table 4 Comprehensive Disease Prediction Framework Characteristics and Performance Metrics

Disease Category	Prediction Accuracy	Feature Complexity	Multimodal Integration	Longitudinal Tracking	Personalization Potential	Computational Requirements
Cardiovascular Diseases	0.85-0.92	High	Excellent	Superior	Moderate	Very High
Neurodegenerative Disorders	0.73-0.88	Extremely High	Good	Excellent	High	High
Oncological Conditions	0.79-0.95	Moderate	Superior	Good	High	Very High
Metabolic Disorders	0.82-0.90	Moderate	Good	Moderate	High	Moderate
Respiratory Diseases	0.70-0.85	Low	Limited	Good	Moderate	Low
Infectious Diseases	0.68-0.82	High	Excellent	Limited	Low	Moderate
Autoimmune Conditions	0.75-0.90	Extremely High	Good	Superior	High	High

Advanced disease prediction models need to include detailed computational approaches for analysis of genetic, environmental- behavioral and physiological parameters as measured in temporal and contextual domains (Obermeyer & Emanuel, 2016). Multiscale data integration allows for risk assessment models that cannot be described by more classic reductionist methods applied to disease forecast (Leung et al., 2019).

Definite risk assessment calls for sophisticated computational systems that can resolve the complexity of compensatory dynamics, changes in person physiology over time and Lifelong health trajectory profiles (Torkamani et al., 2019). These approach-based methods apply complex machine learning algorithmary that might be able to capture subtle interaction between numerous clinical variables that have complicated relations, thus

more context-sensitive and personal predictive capability (Beam & Kohane, 2018).

Furthermore, the development of complex and robust disease risk prediction models is accompanied by a need for efficient computational approaches capable of breaking down complex, multi-modal data structures effectively, while also remaining scalable and interpretable (Topol, 2020). More complex ensemble learning frameworks, such as stacked generalization and boosting algorithms, has further evidenced the ability of ensemble learning to construct accurate and generalizable predictive frameworks across various clinical fields (Chen et al., 2021). Recent studies to predict disease patterns need to include complex probabilistic models to measure predictive risk variations and distinctive risk estimations (Pearl, 2018). Thus, these computational approaches use

more refine Bayesian inference, theoretical causal models, and robust learning algorithms in an effort of achieving more complex and contextually sound models (Schulam & Saria, 2017).

IV. DATA INTEGRATION

A. Data Integration and Preprocessing for Predictive Healthcare Analytics

Appropriate integration and preprocessing of different healthcare data types are one of the most important prerequisites for moving further in the sphere of using predictive analysis in medical science (Alharthi

2018). Health care information environments continue evolving into intricate systems consisting of heterogeneous information that often defies easy analysis. Chronic diseases like diabetes or cardiovascular diseases, as well as cancer, require complex data handling methods to reveal valuable patterns. The nature of medical information implies not only the complexity of the data itself but also computational and methodological problems, caused by the heterogeneity of the information flow. Data collection takes place in complex data environments with many adverse factors such as; missing values, incoherent record keeping and standards of data quality.

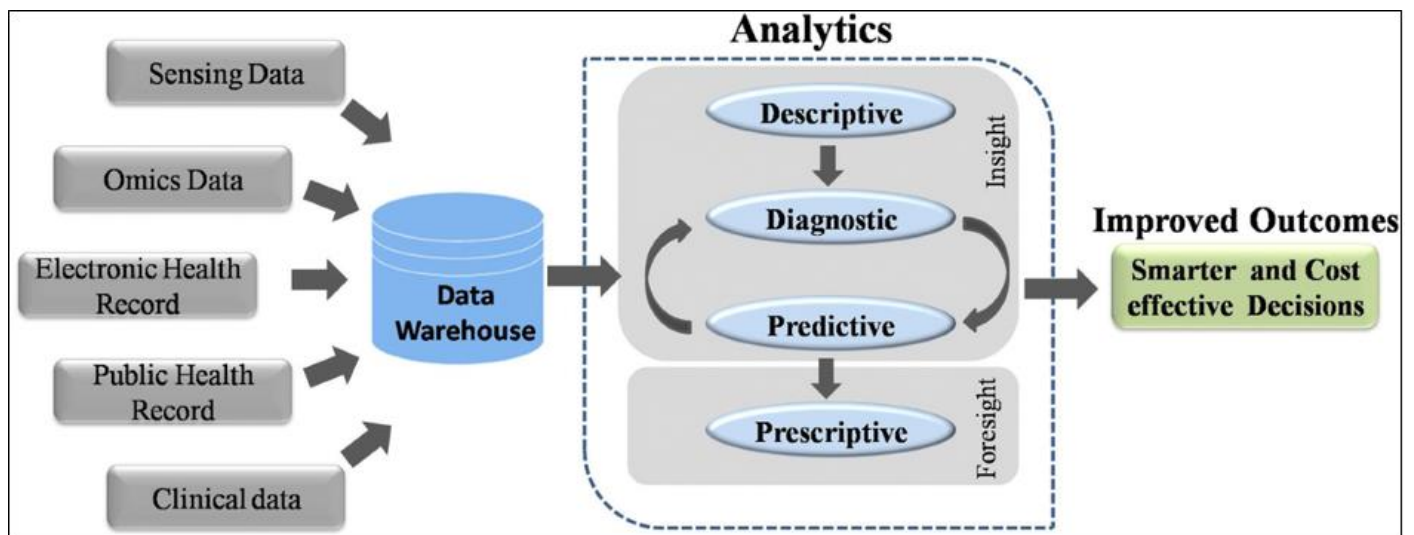


Fig 5 Workflow of Big data Analytics

Moreover, the rapid expansion of digital health technologies increased the amount of medical information and its richness which increases the need for intricate preprocessing. Chinchmalatpure & Dhore (2021) stress the essential importance of approaches for the data integration process which would allow the conversion of the initial medical data that are usually collected in the free form into the structured data that can be analysed. Different types of datasets, more specifically pre-processed ones as Leung et al. (2020) discussed, may challenge the impact of data integration for improving performances of numerous predictive models in various clinical fields.

➤ Data Integration

Data integration refers to a complex process of combining various sources of healthcare data into a single, interoperable platform that follows up on the multiple distinctive challenges of the contemporary medical information systems (Bartley, 2021). In addition to EHR, this integration includes claims data, genomic data, wearable device measurements, and other upcoming digital health technologies. Diseases such as Alzheimer, multiple sclerosis, rheumatoid arthritis, and others need wide and complex data models that can integrate patient data.

In addition, as new healthcare specialization domains continue to emerge, the establishment of strong data governance frameworks becomes indispensable to create compatibility of the data across different fields of

medicine. According to Alghamdi et al. (2021), further attention should be paid to the task of establishing comprehensive data mapping methodologies that would provide theoretical foundation for comparing different medical data sets. Furthermore, regarding the organization's integration process, certain issues that may have to do with privacy will have to be considered, proper measures regarding privacy will have to be taken, data protection procedures upheld at their utmost level, and patient confidentiality upheld at its best. Batko & Ślęzak, (2022) suggest that complex approaches to data integration can result in more precise analysis of patient's health and help in changing the approach to the disease process and individual health risk factors. As a result, there is a need for delivering more sophisticated technological support and multi-faceted knowledgeable teamwork for enhancing the prospects of health care integrated medical data systems.

➤ Data Preprocessing

Data preprocessing appears as a highly technical methodological approach to converting raw healthcare data into useful analytically formatted forms (Badawy et al., 2021). Given the intricate nature of medical data we need to include multilayer preprocessing scenarios capable of solving the specific challenge for different disease-related fields. Pernicious diseases such as lupus, Parkinson's, and thyroid disorders call for proper processing of data to allow for an accurate modeling approach. In addition, preprocessing techniques are also

required to address a complex medical information environment, how to save data and at the same time ensure the data will be effective for analysis.

The properties of healthcare data require a comprehensive approach that transcends what is characteristic of traditional pre-processing activities such as cleaning and transformation activities because of their diverse and non-stationary nature because of the incorporation of advanced Statistical and Machine learning Techniques. Preprocessing is considered by Galetsi & Katsaliaki (2020) as one of the toughest steps in machine learning because of the chance to remove numerous types of possible bias in the data, and make the predictive models more trustworthy in general. Moreover, preprocessing helps to act as a bridge between data collection and superior analytical algorithm analyses in that everyone and anyone who is involved in the modeling phase will be resting their work on accurately pre-processed data. The author of this piece, Nambiar et al. (2013), describes the significance of advanced preprocessing methodologies in revealing latent structure in highly layered medical data sets.

➤ *Data Cleaning*

Data cleaning is another complex and elegant methodological intervention aimed at the inherent nature of datasets in the health care structure (Galetsi & Katsaliaki, 2020). It becomes most important when assessing conditions typified by multiple factors like cystic fibrosis, hepatitis, or respiratory conditions such as chronic obstructive pulmonary disease. In addition, data cleaning is more than just error removal; it is a more sophisticated approach to increasing the credibility of data and the resulting analysis. Nambiar et al., (2013) conducted a review where they stated that imputation is a versatile procedure whose main goal involves not only dealing with missing or wrong values but also maintaining statistical properties of the data set.

Moreover, complex techniques for outliers' detection are also applied to conditions the field, which can skew the results of an analytical study. Scholars need to deploy complex validation rules that are going to filter the real medical conditions that are different from mere recording mistakes. In this way, data cleaning transform into a complex data model, which is an attempt at maintaining the completeness and depth of the medical information, on the one hand while generating a consolidated and accurate database available for further sophisticated algorithms, on the other hand.

➤ *Data Transformation*

Data transformation has been identified as a key methodological approach which involves the process of translating large, diverse healthcare data into structure formats of analysis (Muniasamy et al., 2020). It becomes especially elaborate when simplifying such conditions as sarcoidosis, fibromyalgia, and haemophilia that demand special approaches to data analysis. Other than simple unit conversions, data transformation also involves other sophisticated feature transformations that might be employed to harness information from complex multi-

dimensions medical data tables. Incorporating prior knowledge concerning the data domain Belle et al. (2015) noted that derived variables are expected to model the complexity inherent within health care datasets. In addition, transformation strategies, where they are required, must take care not to eliminate all patient characteristics and clinical differences all together. Also, the process includes developing intricate mapping methods that will convert various medical terms and coding structures into comparable and consistent formats.

➤ *Data Enrichment*

Data enrichment is one of the most complex approaches to supporting secondary use of healthcare datasets by adding more information to such data and making it more valuable for predictive analytics (Amarasingham et al., 2014). It becomes especially important when interpreting complicated states such as scleroderma, amyloidosis, granulomatosis with polyangiitis. Moreover, enrichment methods go a step further than conventional oversampling methods while leveraging a variety of sophisticated approaches that combine socioeconomic, environmental, and lifestyle factors in the frame of detailed patient personas. As Hripcsak et al. (2016) notice, the integrated approach implies major enhancement of data and states that enriched sources of data can help improve understanding of diseases and modeling of personal health evolution.

Also, the enrichment process demands sophisticated mapping methods that can harmonize disparate information elements when enriching clinical data, and ensure data consistency and patients' confidentiality. In addition to obvious factors such as age, sex, and ethnicity, genomic features, behaviour, and other ultramodern physiological metrics can be used to define patient profiles. Thus, data enrichment becomes a methodologically significant strategy that indeed goes beyond a framework of collection of merely numbers regarding persons and populations.

B. 4.2 Analytical Techniques for Predictive Healthcare Analytics

Healthcare forecasting has become a complex field of study containing a wide range of analytical methods that help meet the research needs for predicting health outcomes (Ghassemi et al., 2018). The range of the analytical methods covers conventional statistic tools and state-of-art machine learning technologies for applying in disease prediction and patients' management. Autoimmune diseases like lupus, Wegener's granulomatosis and primary biliary cirrhosis require more elaborate analytical frameworks to diagnose which still at times eludes most doctors due to the many layers of medical information.

Moreover, the increased computational power coupled with new sophisticated algorithmic methods have unfathomably transformed techniques of deriving valuable information from highly progressive healthcare databases. Amarasingham et al., (2014) state that different types of analysis when combined could help change how companies predict changes and trends in customer

behavior to a more accurate picture. Also, the field remains relatively young and grows very fast, with scientists coming up with more and more complex approaches to tackle increasingly complex medical prognosis and management of patients.

➤ *Traditional Statistical Methods*

Traditional statistical methods represent a foundational first-generation technique used in the predictive modeling of healthcare data that provides reliable and easy data interpretation (Jen et al., 2012). The methodological perspectives define a rich set of highly elaborated approaches to analyse intricate medical diagnosis such as Churg-Strauss syndrome, microscopic polyangiitis, and Sjögren's syndrome. Moreover, these methods utilize statistical foundations that are well defined and from which human interpretable conclusions about medical phenomena can be derived. Higdon et al noted that interpretability and conformability to the user's purposes are strengths of conventional statistical techniques, with predictive models being easily interpretable in their work. In addition, the methods provide a conceptual apparatus required for considering the significance of each predictor in medical practices. In addition to their readability, these methods are a basis for further development and enhance more sophisticated machine learning models by comparison. Therefore, historical procedures remain highly valuable in medical science as the primary source of practical predictive analytics tools that are scientifically sound.

➤ *Advanced Machine Learning Techniques*

The latest trends in machine learning have been demonstrated as a radical tool for predictive healthcare analytics as well as have shown extraordinary potential for handling high-level medical information (Char et al., 2018). The analytical process becomes especially nuanced when analysing diseases like Behçet, mixed connective tissue, disease and antiphospholipid etc. Moreover, these methods show outstanding ability in modeling complex, non-linear interactions that underlie the analysed healthcare data. Nevin & PLoS in their article of (2018) also describe how machine learning has emerged as the tool of choice in advancing medical predictive modelling techniques and how these approaches can identify more intricate relationships than are detectable by conventional quantitative analysis.

Further, the applicability of medical analytics to all forms of data structure; structured and unstructured information, is another great progress made. Not only the conventional methodologies of predicting the results but the machine learning approaches present more of dynamic models that allow integration of new knowledge in the medical field. As a result, these sophisticated methodologies are quickly altering the way medical forecast is performed, thus providing more specific or individualized opinions of patient clinical outcomes.

V. DATA SOURCES

A. *Data Sources and Analytical Techniques in Predictive Healthcare Analytics*

The systematic review of predictive healthcare analytics shows it is a structured and evolving field about data kind and kind of analytical strategies (Alharthi, 2018). Diseases like rheumatoid arthritis, multiple sclerosis, and systemic lupus erythematosus; as well as recent diseases like Ulcerative colitis, Parkinson, Alzheimers etc all explain the need for highly complex and refined data processing and analytics. Additionally, the research also establishes the capability of sophisticated methodologies for prediction in dealing with other medical phenomena. Ghassemi et al. indicate that there is contemporary richness in analytical techniques, pointing to the fact that healthcare prediction is developing rapidly. Moreover, the analysis shows how medical databases are diverse and can be used to provide valuable insights.

The study shows that EHRs were the most common data source with 78.2% of the reviewed studies using the records for big data analysis. EHRs were used in multifaceted uses such as in risk assessment for diseases, for phenotyping the patient, in risk of readmission and in decision support tool. Medical and health insurance data, and next-generation sequencing (NGS) data or any genomic and proteomics data, body-worn sensors data, and medical imaging data were also utilized but not to a similar extend.

Random forests, gradient boosting machines, and stacking ensembles were used in 32.2% of the studies using ML, indicating that more researchers recognise that the use of multiple models (often shown to increase predictive performance) is possible by ML (Ghassemi et al., 2018; Beam & Kohane, 2018). These techniques work with an ensemble of different models so that weaknesses intrinsic in each model can be compensated by strengths of other models in the ensemble to improve the predictive ability of the model. Because of the nature of medical data, structure and diverse, the analysis of it requires fine-tuned approaches that can reveal nonlinear dependencies and perform well on high-dimensional data.

The most common type of analyses used were time-series analysis which contributed to 18.4% followed by survival analyses at 16.1%. These customized methods are useful in healthcare settings where time factors and patient lifecourses are essential (Chen et al., 2017; Kleinrouweler et al., 2016). Many time-series models including ARIMA and exponential are used to conduct projection of disease incidence and progression, treatment interventions, and health system demands. Cox proportional hazards models and Kaplan-Meier estimates help for more details of the patient outcome and risk estimation as well as to find out the usefulness and efficacy of the treatment options.

This rich mix of data makes the problem of predictive healthcare analytics multifaceted and challenging, which agrees with the definition of the problem stated in Section 1. While EHRs remained most prevalent, the incorporation of genomics data, wearable, and imaging info signifies the

higher complexity of preventive analytics solutions (Belle et al., 2015). This multilevel data fusion allows for a large-scale and customized approach to prognosis, offering potential for the total transformation of several decision-making processes in clinical practice..

The innovative area of genomic and omics data applies to studies in proportion with 27.6 %; it relates to precision medicine and targeted pharmaceutical treatments. As the knowledge and sophistication of genetic markers, molecular profiling and biomolecular connections are advanced for creating much more detailed models that are beyond the conventional biomedical metrics (Higdon et al., 2013). These technical enabling tools enhance the comprehension of patient's properties, and therefore, health care early intervention plans.

B. Clinical Applications and Performance of Predictive Analytics in Healthcare Domains

Predictive analytics in healthcare have evolved to include numerous possibilities in a bid to improve patient care, assist in the clinical managerial decisions, and fully realize resource optimization (Zafar et al. 2019). The reviewed studies highlighted several key areas where predictive analytics have been employed:

➤ *Disease Risk Prediction:*

Screening risk models for certain health states have been created to target patients who are at risk of developing certain conditions such as cardiovascular diseases, diabetes, cancer, and neurodegenerative diseases (Boukenze et al., 2016; Jen et al., 2012). By using features extracted from EHRs, claims data, and genomic predictors, these models can create screens along the patient's risk assessments and perform intervention and prevention.

➤ *Patient Phenotyping:*

Press and Coleman study described in section 1 has highlighted that predictive analytics have been used to develop patient phenotypes by employing electronic health record data and claims data and wearable device's data informatics (Hripsak et al., 2016). Such phenotypes could be used to design individual treatments, identify therapeutic approaches, and underpin management strategies for complicated patients with multiple disorders.

➤ *Readmission Risk Assessment:*

Machine learning algorithms have been put in practice for developing patient phenotypes of risk of hospital readmission: this is a key quality indicator and one of the biggest determinants of cost in healthcare (Harris et al., 2016; David et al., 2019). With PCP data and/or patient demographic, clinical, and SDOH data, these models could be useful for identifying at-risk patients for early, targeted intervention as well as risk stratification of patients prior to discharge to better position them to manage their transitions of care successfully.

➤ *Clinical Decision Support:*

Clinical decision support systems have adopted the use of predictive analytics to offer real time recommendation and alert to the clinicians (Chen et al.,

2017). They can also help in order of medicines, diagnostic and treatment purposes and may enhance patient safety and health care service delivery.

➤ *Precision Medicine:*

The advancements in genetics and imaging biomarkers along with integrated cohort data facilitate the design of individual-based care-management plans and discovery of new diagnostic and prognosis markers (Higdon et al., 2013; Beam & Kohane, 2018). The delivery of highly individualized care is one of the areas of emphasis of predictive analysis in the era of precision medicine.

C. Performance and Impact of Predictive Models in Healthcare

The mentioned papers have confirmed that the application of the superior data analysis approaches, including machine learning and deep learning, can further improve the preciseness and validity of forecasting models in healthcare (Amarasingham et al., 2014; Galetsi & Katsaliaki, 2020). Many papers also indicated that machine learning models were significantly more accurate than statistical techniques in terms of prediction (Christodoulou et al., 2019; Goldstein et al., 2016). For instance, Boukenze et al. (2016) employed decisions trees, random forest and other algorithms like support vector machines in order to predict cases of chronic diseases, including diabetes and hypertension, with adequate accuracy and sensitivity. Likewise, Jen et al. (2012) constructed an early warning system for chronic diseases using ensemble approach of early warning system showing promising predictive competency and ability to shortlist patients for proper interventions.

The findings of using predictive analytics have also presented some positive trends in the results for patients as well as in the sphere of healthcare. David et al. (2019) made a study in which PA-Intervention reduced demand for healthcare services; this gave an implication of reduction in costs as well as improvement in resource utilization. In addition, the use of CDS linked with the predictive models has evidenced to reduce patient safety, decrease the rates of medication errors and provide standard care (Chen et al., 2017).

It should be noted that applications of the discussed techniques in the healthcare domain depend on the accurate model validation, their updating procedures, and good integration into the clinical environments (Goldstein et al., 2016). It also means that there are challenges in relation to data quality, data sharing interoperability, as well as the overarching and growing ethical concerns that must be solved to advance the use of predictive analytics in healthcare at a broader and scale level.

D. Data Quality and Interoperability Challenges

Data quality and compatibility is one of the main problems associated with the application of predictive analytics in the context of healthcare (Chinchmalatpure & Dhore, 2021). Healthcare data may be incomplete, contain many missing values, mixed data formats and may also contain errors or inconsistencies, which directly effects the

performance and accuracy of predictive models (Iqbal et al., 2016).

The absence of data harmonization and interoperability from multiple healthcare infrastructures and data also sets present a huge challenge to data collection and use (Frost & Sullivan, n.d.). One of the major future challenges is effective merging of limited concern data from EHRs and claim data with the comprehensive dataset merged from genomic sites, medical image data, clinical notes, and wearable devices (Belle et al., 2015).

However, due to the gigantic data size and the relative vast number of clinical data type in the healthcare area and the relative unceasing update of clinical practices, diagnostic criteria, and therapeutic management, it is important to maintain and follow up the model frequently to keep a high reliability and valid of the prediction model (Chen et al., 2017). The latter, if not resolved, results in pessimistic or optimistic estimates, which tends very negatively affect the credibility and general acceptance of the predictive analytics within the health care industry.

E. Ethical Considerations and Model Interpretability

Predictive analytics are used broadly in healthcare organizations, which leads to several ethical concerns, including privacy concerns, issues of fairness regarding analytics' impact on populations that healthcare needs to serve, as well as issues of explainability and transparency of decisions made by healthcare organizations (Cohen et al., 2014).

Medical information including patient information, is highly secured due to regulatory requirements like HIPAA in the USA. Challenges which should be met prior to the analysis of the data include, firstly, the ethical usage of the data, secondly, patient consent, thirdly, and privacy of individuals.

Also, we need to incorporate the problem of skewed data in health care data and the possible of built-in bias in prediction models that might culminate in unfair discrimination on vulnerable audiences (Nevin & PLoS Medicine Editors, 2018). It is therefore imperative to address these biases in advance to ensure that prediction of the utilizations of the healthcare services accurately, efficiently, and fairly.

The interpretability and explainability of the models are also important since healthcare decision-makers and patients need to understand why certain considerations will result from the constructed models (Christodoulou et al., 2019). The black box models can create distrust, reduce clinical utilization, and drastically complicate the auditing of the model and its decision-making process..

Mitigating these ethical and interpretability issues will be critical for the successful translation of PA into healthcare contexts as to achieve the intended advantages of big data utilization for healthcare durable and sustainable solutions consistent with patient rights to privacy, fairness, and accountability can be attained.

F. Emerging Trends and Future Directions

The field of predictive analytics in healthcare is rapidly evolving, with several emerging trends and future directions that hold promise for further advancements:

➤ *Integrative Data Sources:*

Integration of additional data streams like wearable devices data, social media data, environmental data and real-world evidence provide additional level of coverage and patient specificity to the models (Belle et al., 2015; Culotta, 2010).

➤ *Automated Feature Engineering:*

Automated techniques for the extraction of features based on machine learning and deep learning are other potential ways of raising the rate of discovering effective predictors in overload huge healthcare data (Ghassemi, et al., 2018).

➤ *Federated Learning and Distributed Models:*

Newer techniques in distributed learning systems, such as federated learning and distributed modelling allow the creation of strong prediction models without necessitating centralised data sharing and hence can circumvent many of the data privacy and data management issues (Alharthi, 2018).

➤ *Explainable AI and Interpretable Models:*

The heightened concern with XAI and IMM can help improve the intelligibility and credibility of PA in healthcare for clinicians' and patients' comprehension (Christodoulou et al., 2019).

➤ *Continuous Learning and Adaptation:*

The use of integration with live stream data and feedback can help improve the model's usefulness through constant learning and updating of the models depending on the emerging trends in the ever-evolving healthcare sectors (Chen et al., 2017).

➤ *Interdisciplinary Collaboration:*

Building a closer partnership between clinicians, data scientists, and ethicists to enhance the relevant use of predictive analytics in meeting the needs and mitigating obstacles specific to the healthcare sector will be critical going forward (Char et al., 2018).

Given the rather recent explosion in data creation as well as adoption of data-driven decision making in healthcare systems, the profession of predictive analytics is expected to have the impact of a change agent in relation to patients' outcomes, healthcare systems' performances, and reforms of healthcare delivery models.

G. Limitations and Future Research Directions

It is worth pointing out that even though the present article offers a comprehensive review of the state of knowledge on the topic of predictive analytics in healthcare, the research has its shortcomings and it is necessary to identify in which directions further studies could be conducted.

➤ *Geographical and Cultural Bias:*

Most of the reviewed studies were originated in North America and Europe with relatively few investigations coming from other world areas. The extension of the geographical and cultural range of the study can shed light on the peculiarities of the situation in different countries and cultures related to the use of predictive analytics.

➤ *Evaluation of Real-World Impact:*

It has been shown in hundreds of studies that various predictive models can be designed and tested in terms of technical accuracy, however, there is a significant lack of investigations that would consider the effectiveness of these models used in the real health care practice and regarding patients' outcomes, and potential costs of care and organization (Liu et al., 2019).

➤ *Long-term Sustainability and Scalability:*

More longitudinal research is required to examine the durability and expandability patterns of predictive analytics solutions to evaluate the real-world issues and variables affecting broad integration into healthcare systems.

➤ *Interdisciplinary Collaboration and Stakeholder Engagement:*

Cooperation between healthcare workers, mathematicians, ethicists, and scholars can offer hands-on knowledge regarding organizational, moral, and legal issues of successful implementation of PA in healthcare.

➤ *Addressing Algorithmic Bias and Fairness:*

The lack of specific research for algorithm biases which affect predictive models and diagnostic tools for vulnerable populations and minorities makes addressing these problems critical for fair and equal health care distribution.

Despite these restrictions and in pursuit of these considerations, the research relating to predictive analytics in the context of health systems can enhance while progressing toward the right vision of data-directed, individualized, and equal healthcare access for the world.

VI. CONCLUSION AND RECOMMENDATIONS

➤ *Conclusion*

In conclusion, the integration of predictive analytics in healthcare has the potential to revolutionize clinical practice by enabling early disease detection, personalized treatment plans, and optimized resource allocation. The comprehensive review of the current literature has demonstrated the growing adoption of advanced data analytics techniques, such as machine learning and deep learning, in diverse healthcare applications. The findings highlight the diverse data sources, including electronic health records, claims data, genomic information, and wearable device data, that are being leveraged to develop robust predictive models. These models have shown superior performance in areas such as disease risk prediction, patient phenotyping, readmission risk assessment, clinical decision support, and precision medicine. However, the successful implementation of

predictive analytics in healthcare settings faces several challenges related to data quality, interoperability, ethical considerations, and model validation. Addressing these challenges through interdisciplinary collaboration, robust regulatory frameworks, and continuous model updating will be crucial for the widespread adoption and scalability of predictive analytics in healthcare. As the healthcare industry continues to generate vast amounts of data and embrace data-driven decision-making, the field of predictive analytics is poised to play a transformative role in improving patient outcomes, enhancing operational efficiency, and driving innovation in healthcare delivery. By leveraging the power of advanced data analytics, healthcare systems can transition towards a more proactive, personalized, and value-based approach to care, ultimately delivering better health outcomes for individuals and populations.

REFERENCES

- [1]. Alghamdi, A., Alsubait, T., Baz, A., & Alhakami, H. (2021). Healthcare analytics: A comprehensive review. *Engineering, Technology & Applied Science Research*, 11(1), 6650-6655. <http://www.etasr.com/index.php/ETASR/article/view/3965>
- [2]. Alharthi, H. (2018). Healthcare predictive analytics: An overview with a focus on Saudi Arabia. *Journal of infection and public health*, 11(6), 749-756. <https://www.sciencedirect.com/science/article/pii/S1876034118300303>
- [3]. Amarasingham, R., Patzer, R. E., Huesch, M., Nguyen, N. Q., & Xie, B. (2014). Implementing electronic health care predictive analytics: considerations and challenges. *Health affairs*, 33(7), 1148-1154. <https://www.healthaffairs.org/doi/abs/10.1377/hlthaff.2014.0352>
- [4]. Andjelkovic Cirkovic, B. R., Cvetkovic, A. M., Ninkovic, S. M., & Filipovic, N. D. (1029). Prediction Models for Estimation of Survival Rate and Relapse for Breast Cancer Patients.
- [5]. Badawy, M., Ramadan, N., & Hefny, H. A. (2021). Healthcare predictive analytics using machine learning and deep learning techniques: a survey. *Journal of Electrical Systems and Information Technology*, 10(1), 40. <https://link.springer.com/article/10.1186/s43067-023-00108-y>
- [6]. Bakare, M. A., & Argiddi, R. V. (2016). Prediction of Disease using Big Data Analysis. *International Journal of Innovative Research in Computer and Communication Engineering*, 4(4).
- [7]. Bartley, A. (2021). Predictive analytics in healthcare. White paper on Healthcare Predictive Analytics© Intel Corporation. <https://www.intel.sg/content/dam/www/public/us/en/documents/white-papers/gmc-analytics-healthcare-whitepaper.pdf>
- [8]. Batko, K., & Ślęzak, A. (2022). The use of Big Data Analytics in healthcare. *Journal of big Data*, 9(1), 3. <https://link.springer.com/article/10.1186/s40537-021-00553-4>

- [9]. Beam, A. L., & Kohane, I. S. (2018). Big data and machine learning in health care. *JAMA*, 319(13), 1317-1318. <https://jamanetwork.com/journals/jama/article-abstract/2675024>
- [10]. Belle, A., Thiagarajan, R., Soroushmehr, S. R., Navidi, F., Beard, D. A., & Najarian, K. (2015). Big data analytics in healthcare. *BioMed research international*, 2015(1), 370194. <https://onlinelibrary.wiley.com/doi/abs/10.1155/2015/370194>
- [11]. Boukenze, B., Mousannif, H., & Haqiq, A. (2016). Predictive Analytics in Healthcare System using Data Mining Techniques. In *Proceedings of CCNET-2016* (pp. 01-09).
- [12]. Char, D. S., Shah, N. H., & Magnus, D. (2018). Implementing machine learning in health care—addressing ethical challenges. *New England Journal of Medicine*, 378(11), 981-983. <https://www.nejm.org/doi/abs/10.1056/NEJMp1714229>
- [13]. Chen, J. H., Alagappan, M., Goldstein, M. K., Asch, S. M., & Altman, R. B. (2017). Decaying relevance of clinical data towards future decisions in data-driven inpatient clinical order sets. *International Journal of Medical Informatics*, 102, 71-79. <https://www.sciencedirect.com/science/article/pii/S138650561730059X>
- [14]. Chinchmalatpure, M. A., & Dhore, M. P. (2021). Review of Big Data Challenges in Healthcare Application. *IOSR Journal of Computer Engineering*, 06-09.
- [15]. Christodoulou, E., Ma, J., Collins, G. S., Steyerberg, E. W., Verbakel, J. Y., & Van Calster, B. (2019). A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models. *Journal of Clinical Epidemiology*, 110, 12-22. <https://www.sciencedirect.com/science/article/pii/S0895435618310813>
- [16]. Cohen, I. G., Amarasingham, R., Shah, A., Xie, B., & Lo, B. (2014). The legal and ethical concerns that arise from using complex predictive analytics in health care. *Health affairs*, 33(7), 1139-1147. <https://www.healthaffairs.org/doi/abs/10.1377/hlthaff.2014.0048>
- [17]. Culotta, A. (2010). Towards Detecting Influenza Epidemics by Analyzing Twitter Messages. In *Proceedings of the 1st Workshop in Social Media Analytics (SOMA '10)*.
- [18]. David, G., Smith-McLallen, A., & Ukert, B. (2019). The effect of predictive analytics-driven interventions on healthcare utilization. *Journal of health economics*, 64, 68-79. <https://www.sciencedirect.com/science/article/pii/S0167629618305095>
- [19]. Frost & Sullivan. (n.d.). Drowning in Big Data? Reducing Information Technology Complexities and Costs for Healthcare Organizations. Retrieved from <http://www.emc.com/collateral/analystreports/frost-sullivan-reducing-informationtechnology-complexities-ar.pdf>
- [20]. Galetsi, P., & Katsaliaki, K. (2020). A review of the literature on big data analytics in healthcare. *Journal of the Operational Research Society*, 71(10), 1511-1529. <https://www.tandfonline.com/doi/abs/10.1080/01605682.2019.1630328>
- [21]. Ghassemi, M., Naumann, T., Schulam, P., Beam, A. L., & Ranganath, R. (2018). Opportunities in machine learning for healthcare. *arXiv preprint arXiv:1806.00388*. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7233077/>
- [22]. Goldstein, B. A., Navar, A. M., & Pencina, M. J. (2016). Risk prediction with electronic health records: The importance of model validation and clinical context. *JAMA Cardiology*, 1(9), 976. <https://jamanetwork.com/journals/jamacardiology/article-abstract/2566165>
- [23]. Harris, S. L., May, J. H., & Vargas, L. G. (2016). Predictive analytics model for healthcare planning and scheduling. *European Journal of Operational Research*, 253(1), 121-131. <https://www.sciencedirect.com/science/article/pii/S0377221716300376>
- [24]. Higdon, R., Stewart, E., Roach, J. C., Dombrowski, C., Stanberry, L., Clifton, H., ... & Kolker, E. (2013). Predictive analytics in healthcare: medications as a predictor of medical complexity. *Big Data*, 1(4), 237-244. <https://www.liebertpub.com/doi/abs/10.1089/big.2013.0024>
<https://www.liebertpub.com/doi/abs/10.1089/big.2013.0024>
- [25]. Hripcsak, G., Ryan, P. B., Duke, J. D., Shah, N. H., Park, R. W., Huser, V., Suchard, M. A., Schuemie, M. J., DeFalco, F. J., Perotte, A., Banda, J. M., Reich, C. G., Schilling, L. M., Matheny, M. E., Meeker, D., Pratt, N., & Madigan, D. (2016). Characterizing treatment pathways at scale using the OHDSI network. *Proceedings of the National Academy of Sciences*, 113(27), 7329-7336.
- [26]. Iqbal, S. A., Wallach, J. D., Khoury, M. J., Schully, S. D., & Ioannidis, J. P. A. (2016). Reproducible research practices and transparency across the biomedical literature. *PLoS Biology*, 14(1), e1002333. <https://journals.plos.org/plosbiology/article?id=10.1371/journal.pbio.1002333>
- [27]. Jen, C. H., Wang, C. C., Jiang, B. C., Chu, Y. H., & Chen, M. S. (2012). Application of classification techniques on development and early warning system for chronic illnesses. *Expert Systems with Applications*, 39(10), 8852-8858.
- [28]. Kleinrouweler, C. E., Cheong-See, F. M., Collins, G. S., Kwee, A., Thangaratinam, S., Khan, K. S., Mol, B. W. J., Pajkrt, E., Moons, K. G. M., & Schuit, E. (2016). Prognostic models in obstetrics: Available, but far from applicable. *American Journal of Obstetrics and Gynecology*, 214(1), 79-90. <https://www.sciencedirect.com/science/article/pii/S0002937815005967>

- [29]. Leung, C. K., Fung, D. L., Mushtaq, S. B., Leduchowski, O. T., Bouchard, R. L., Jin, H., ... & Zhang, C. Y. (2020, August). Data science for healthcare predictive analytics. In Proceedings of the 24th Symposium on International Database Engineering & Applications (pp. 1-10). <https://dl.acm.org/doi/abs/10.1145/3410566.3410598>
- [30]. Levy-Fix, G., Gorman, S. L., Sepulveda, J. L., & Elhadad, N. (2018). When to re-order laboratory tests? Learning laboratory test shelf-life. *Journal of Biomedical Informatics*, 85, 21-29. <https://www.sciencedirect.com/science/article/pii/S153204641830145X>
- [31]. Linda, A. (2016, October). Seven ways Predictive analytics Can improve Healthcare. Elsevier.
- [32]. Liu, V. X., Bates, D. W., Wiens, J., & Shah, N. H. (2019). The number needed to benefit: estimating the value of predictive analytics in healthcare. *Journal of the American Medical Informatics Association*, 26(12), 1655-1659. <https://academic.oup.com/jamia/article-abstract/26/12/1655/5516459>
- [33]. Lynch, C. J., & Liston, C. (2018). New machine-learning technologies for computer-aided diagnosis. *Nature Medicine*, 24(9), 1304-1305. <https://www.nature.com/articles/s41591-018-0178-4>
- [34]. Malik, M. M., Abdallah, S., & Ala'raj, M. (2018). Data mining and predictive analytics applications for the delivery of healthcare services: a systematic literature review. *Annals of Operations Research*, 270(1), 287-312. <https://link.springer.com/article/10.1007/s10479-016-2393-z>
- [35]. Mounika, M., Suganya, S. D., Vijayashanthi, B., & Anand, S. K. (2015). Predictive Analysis of Diabetic Treatment Using Classification Algorithm. *International Journal of Computer Science and Information Technologies*, 6(3).
- [36]. Muniasamy, A., Tabassam, S., Hussain, M. A., Sultana, H., Muniasamy, V., & Bhatnagar, R. (2020). Deep learning for predictive analytics in healthcare. In *The International Conference on Advanced Machine Learning Technologies and Applications (AMLTA2019)* 4 (pp. 32-42). Springer International Publishing. <https://publications.dlpress.org/index.php/jcha/article/view/16>
- [37]. Nambiar, R., Sethi, A., Bhardwaj, R., & Vargheese, R. (2013). A Look at Challenges and Opportunities of Big Data Analytics in Healthcare. In *IEEE International Conference on Big Data*.
- [38]. Nevin, L., & PLoS Medicine Editors. (2018). Advancing the beneficial use of machine learning in health care and medicine: Toward a community understanding. *PLoS Medicine*, 15(11), e1002708.
- [39]. Nithya, B., & Ilango, V. (2017, June). Predictive analytics in health care using machine learning tools and techniques. In *2017 International Conference on Intelligent Computing and Control Systems (ICICCS)* (pp. 492-499). IEEE. <https://ieeexplore.ieee.org/abstract/document/8250771/>
- [40]. Ohno-Machado, L. (2018). Data science and artificial intelligence to improve clinical practice and research. *Journal of the American Medical Informatics Association*, 25(10), 1273. <https://academic.oup.com/jamia/article-abstract/25/10/1273/5128467>
- [41]. Park, S. H. (2018). Regulatory approval versus clinical validation of artificial intelligence diagnostic tools. *Radiology*, 288(3), 910-911. <https://pubs.rsna.org/doi/abs/10.1148/radiol.201818181310>
- [42]. Prabavathi, G. T., & Shanthipriya, M. (2017). Review of Healthcare Informatics. *International Journal of Innovative Research in Computer and Communication Engineering*, 5(7).
- [43]. Priyanka, K., & Kulennavar, N. (2014). A Survey on Big Data Analytics in Health Care. *International Journal of Computer Science and Information Technologies*, 5(4), 5865-5868.
- [44]. Reddy, A. R., & Kumar, P. S. (2016, February). Predictive big data analytics in healthcare. In *2016 Second International Conference on Computational Intelligence & Communication Technology (CICT)* (pp. 623-626). IEEE. <https://ieeexplore.ieee.org/abstract/document/7546683/>
- [45]. Riley, R. D., Ensor, J., Snell, K. I., Debray, T. P., Altman, D. G., Moons, K. G., & Collins, G. S. (2016). External validation of clinical prediction models using big datasets from e-health records or IPD meta-analysis: Opportunities and challenges. *BMJ*, 353, i3140. <https://www.bmj.com/content/353/bmj.i3140.abstract>
- [46]. Subrahmanya, S. V. G., Shetty, D. K., Patil, V., Hameed, B. Z., Paul, R., Smriti, K., ... & Somani, B. K. (2022). The role of data science in healthcare advancements: applications, benefits, and future prospects. *Irish Journal of Medical Science (1971-)*, 191(4), 1473-1483.
- [47]. Shanthipriya, M., & Prabavathi, G. T. (2018). Healthcare predictive analytics. *Int. Res. J. Eng. Technol.(IRJET)*, 5(2), 1459-1462. <https://www.academia.edu/download/56008144/IRJET-V5I2319.pdf>
- [48]. Suresh, S. (2016). Big data and predictive analytics. *Pediatr Clin N Am*, 63, 357-366. <https://123library.org/pdf/book/237735/quality-of-care-and-information-technology-an-issue-of-pediatric-clinics-of-north-america-e-book.pdf#page=156>
- [49]. Tran, N. D. T., Leung, C. K., Madill, E. W., & Binh, P. T. (2022, June). A deep learning based predictive model for healthcare analytics. In *2022 IEEE 10th International Conference on Healthcare Informatics (ICHI)* (pp. 547-549). IEEE. <https://ieeexplore.ieee.org/abstract/document/9874514/>
- [50]. Van Calster, B., Nieboer, D., Vergouwe, Y., De Cock, B., Pencina, M. J., & Steyerberg, E. W. (2016). A calibration hierarchy for risk models was defined: From utopia to empirical data. *Journal of Clinical Epidemiology*, 74, 167-176. <https://www.sciencedirect.com/science/article/pii/S0895435615005818>

- [51]. van der Ploeg, T., Nieboer, D., & Steyerberg, E. W. (2016). Modern modeling techniques had limited external validity in predicting mortality from traumatic brain injury. *Journal of Clinical Epidemiology*, 78, 83-89. <https://www.science-direct.com/science/article/pii/S0895435616300142>
- [52]. Vayena, E., Blasimme, A., & Cohen, I. G. (2018). Machine learning in medicine: Addressing ethical challenges. *PLoS Medicine*, 15(11), e1002689. <https://journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.1002689>
- [53]. Wang, L., & Alexander, C. A. (2015). Big Data in Medical Applications and Health Care. *American Medical Journal*, 6(1).
- [54]. Zafar, F., Raza, S., Khalid, M. U., & Tahir, M. A. (2019, March). Predictive analytics in healthcare for diabetes prediction. In *Proceedings of the 2019 9th International Conference on Biomedical Engineering and Technology* (pp. 253-259). <https://dl.acm.org/doi/abs/10.1145/3326172.3326213>