

# Optimizing Healthcare Outcomes through Data-Driven Predictive Modeling

# Md Nagib Mahfuz Sunny<sup>1</sup>, Mohammad Balayet Hossain Sakil<sup>1</sup>, Abdullah Al Nahian<sup>1</sup>, Syed Walid Ahmed<sup>2</sup>, Md Newaz Shorif<sup>1</sup>, Jennet Atayeva<sup>3</sup>

<sup>1</sup>Department of Engineering & Technology, Trine University, Detroit, USA
 <sup>2</sup>MBBS, Children's Clinic of Michigan, Hamtramck, USA
 <sup>3</sup>Department of Graduate & Professional Studies, Trine University, Detroit, USA
 Email: Nagibmahfuz1996@gmail.com, balayet.me@gmail.com, nahian77@gmail.com, syed\_walid@hotmail.com, Newshorif2016@gmail.com, jennetatayeva27@gmail.com

How to cite this paper: Sunny, M.N.M., Sakil, M.B.H., Nahian, A.A., Ahmed, S.W., Shorif, M.N. and Atayeva, J. (2024) Optimizing Healthcare Outcomes through Data-Driven Predictive Modeling. *Journal of Intelligent Learning Systems and Applications*, **16**, 384-402.

https://doi.org/10.4236/jilsa.2024.164019

Received: September 13, 2024 Accepted: October 19, 2024 Published: October 22, 2024

Copyright © 2024 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

http://creativecommons.org/licenses/by/4.0/

# Abstract

This study investigates the transformative potential of big data analytics in healthcare, focusing on its application for forecasting patient outcomes and enhancing clinical decision-making. The primary challenges addressed include data integration, quality, privacy issues, and the interpretability of complex machine-learning models. An extensive literature review evaluates the current state of big data analytics in healthcare, particularly predictive analytics. The research employs machine learning algorithms to develop predictive models aimed at specific patient outcomes, such as disease progression and treatment responses. The models are assessed based on three key metrics: accuracy, interpretability, and clinical relevance. The findings demonstrate that big data analytics can significantly revolutionize healthcare by providing datadriven insights that inform treatment decisions, anticipate complications, and identify high-risk patients. The predictive models developed show promise for enhancing clinical judgment and facilitating personalized treatment approaches. Moreover, the study underscores the importance of addressing data quality, integration, and privacy to ensure the ethical application of predictive analytics in clinical settings. The results contribute to the growing body of research on practical big data applications in healthcare, offering valuable recommendations for balancing patient privacy with the benefits of data-driven insights. Ultimately, this research has implications for policy-making, guiding the implementation of predictive models and fostering innovation aimed at improving healthcare outcomes.

#### **Keywords**

Big Data Analytics, Predictive Analytics, Healthcare, Clinical Decision-Making, Data Quality, Privacy

# **1. Introduction**

The healthcare industry is undergoing a transformative shift driven by the prospects of big data analytics. The integration of vast, diverse, and rapidly growing health data from sources such as Electronic Health Records (EHRs), wearable devices, and genomic data has presented unprecedented opportunities to enhance patient care and clinical decision-making through predictive analytics [1] [2]. By leveraging historical data to forecast patient outcomes and optimize treatment strategies, healthcare systems aim to improve quality and reduce costs [3].

The development of big data analysis in healthcare began with the digitization of patient records through early Electronic Health Records (EHRs) in the 1960s and 1970s. This shift from paper to digital allowed for the aggregation of vast amounts of health data, laying the groundwork for advanced data analysis techniques. By the 2000s, big data technologies emerged, enabling the integration of EHRs with genomic data, medical imaging, and real-time monitoring from wearable devices. This evolution has since fueled the use of predictive analytics, transforming clinical decision-making by allowing for more personalized, data-driven healthcare interventions [4] [5].

The wealth of data available in the healthcare sector has the potential to revolutionize the way clinicians approach patient care. Big data analytics, which encompasses the collection, processing, and analysis of large, complex datasets, can uncover hidden patterns, correlations, and insights that were previously unattainable [6]. In the context of healthcare, big data analytics can be used to predict the likelihood of disease onset, forecast the progression of chronic conditions, and identify patients at high risk of adverse outcomes [7]. This data-driven decisionmaking can lead to more personalized and proactive treatment approaches, ultimately improving patient outcomes and reducing the strain on healthcare resources.

However, the effective implementation of big data analytics in healthcare faces significant challenges. Key issues include the fragmentation of data systems that complicate integration, variability in data quality that affects predictive accuracy, privacy concerns that limits data accessibility, and the opaque nature of complex machine learning models that can obscure the rationale behind predictions [8] [9]. These barriers must be addressed to ensure the reliable and trustworthy application of predictive analytics in clinical settings.

The primary aim of this study is to explore the application of big data analytics in predicting patient outcomes and to develop robust predictive models that enhance clinical decision-making. The specific objectives include: (1) reviewing the current state of big data analytics in healthcare, with a focus on predictive analytics; (2) developing predictive models for specific patient outcomes using machine learning algorithms; (3) evaluating the performance of these models in terms of accuracy, interpretability, and clinical relevance; (4) identifying and addressing challenges related to data quality, integration, and privacy in the context of predictive analytics; and (5) providing recommendations for the implementation of predictive models in clinical practice [10].

To guide this investigation, the following research questions will be addressed: (1) How can big data analytics be effectively utilized to predict patient outcomes in healthcare? (2) What machine learning models are most appropriate for various types of predictions within the healthcare context? (3) What are the primary challenges in implementing predictive analytics in healthcare, and how can these challenges be mitigated to ensure the accuracy and ethical use of predictive models? [11].

This study is significant as it contributes to the expanding field of big data analytics in healthcare by addressing critical challenges that impede the implementtation of predictive models. By focusing on data quality, integration, and privacy, this research aims to provide a framework that enhances the reliability and applicability of predictive analytics in clinical settings [12]. Furthermore, the findings have the potential to influence clinical practice and policy-making, as healthcare systems increasingly adopt data-driven approaches to personalized medicine and improved patient outcomes [13]. Importantly, the study also emphasizes the ethical considerations surrounding the use of patient data, fostering trust among patients and healthcare providers alike [14].

In conclusion, the integration of big data analytics into healthcare presents a transformative opportunity to enhance patient care and clinical decision-making. However, the effective implementation of predictive analytics requires addressing the challenges of data fragmentation, quality, privacy, and model interpretability. This study aims to contribute to the growing body of research in this field by developing robust predictive models, identifying and mitigating the associated challenges, and providing recommendations for the successful integration of big data analytics into clinical practice. By addressing these critical issues, this research can pave the way for a more data-driven, personalized, and efficient healthcare system that ultimately improves patient outcomes.

# 2. Review of Previous Study

The integration of big data analytics into the healthcare sector has garnered significant attention in recent years, with a growing body of research exploring its potential applications and limitations [1] [2]. This literature review aims to summarize the existing research on big data analytics and predictive analytics in healthcare, as well as identify the gaps that the current study seeks to address.

### 2.1. Big Data Analytics in Healthcare

The healthcare industry has been at the forefront of the big data revolution, as it

generates and collects vast amounts of data from various sources, including electronic health records (EHRs), medical imaging, genomic sequencing, and wearable devices [5] [6]. The exponential growth in healthcare data has created new opportunities for data-driven decision-making, with big data analytics emerging as a powerful tool to extract valuable insights and improve patient outcomes [7].

Existing research has demonstrated the potential of big data analytics in various healthcare applications. For instance, studies have shown that the analysis of EHR data can help identify risk factors for chronic diseases, predict the likelihood of hospital readmissions, and optimize resource allocation within healthcare systems [8] [9]. Similarly, the integration of genomic data and clinical data has enabled personalized medicine approaches, where treatment strategies are tailored to an individual's genetic profile [10] [11].

Furthermore, the advent of IoT-enabled (Internet of Things) medical devices and wearables has provided continuous streams of real-time patient data, enabling the development of predictive models for early disease detection and monitoring [12] [13]. Researchers have explored the use of machine learning algorithms to analyze these data sources and identify patterns that can inform clinical decisionmaking [14] [15].

Despite these advancements, the healthcare industry continues to face significant challenges in the effective implementation of big data analytics. One of the primary challenges is the fragmentation of data systems, which often hinders data integration and limits the ability to derive comprehensive insights [16] [17]. Additionally, concerns regarding data quality, privacy, and security have emerged as critical barriers to the widespread adoption of big data analytics in healthcare [18] [19].

#### 2.2. Predictive Analytics in Healthcare

Predictive analytics, a subset of big data analytics, has garnered particular attention in the healthcare domain. Predictive models leverage historical data to forecast future outcomes, such as the likelihood of disease onset, the progression of chronic conditions, and the response to medical interventions [20] [21]. These models have the potential to improve clinical decision-making, enhance patient outcomes, and optimize the allocation of healthcare resources.

Numerous studies have explored the application of predictive analytics in various healthcare scenarios. For instance, researchers have developed models to predict the risk of hospital-acquired infections, identify patients at high risk of readmission, and forecast the trajectory of chronic diseases like diabetes and heart failure [22] [23]. These studies have demonstrated the potential of predictive analytics to enhance clinical practice and improve patient care.

However, the implementation of predictive analytics in healthcare faces several challenges. One of the primary concerns is the interpretability of complex machine learning models, which can often be perceived as "black boxes" by clinicians [24] [25]. Healthcare professionals require transparent and explainable models to

understand the rationale behind the predictions and have confidence in the decision-making process [26] [27].

Additionally, the integration of predictive models into clinical workflows has been a significant challenge, as healthcare organizations often struggle to seamlessly incorporate these tools into their existing systems and processes [28] [29]. Addressing these integration barriers is crucial for ensuring the effective and widespread adoption of predictive analytics in healthcare.

#### 2.3. Current Technologies and Application Cases in Healthcare

In recent years, healthcare has embraced transformative technologies like machine learning (ML), artificial intelligence (AI), and Internet of Things (IoT)-enabled devices [30]. Machine learning algorithms are widely applied in predictive modeling to analyze vast datasets, enabling early detection of diseases like diabetes and cardiovascular conditions. For example, AI-driven predictive models have been used to forecast heart disease risks by analyzing patient EHRs, vital signs, and lifestyle factors. Similarly, IoT-enabled devices such as wearable sensors continuously collect real-time health metrics—like heart rate and glucose levels—allowing for ongoing monitoring and immediate intervention [31].

A case study in diabetes management demonstrated how ML models could predict blood sugar fluctuations and suggest personalized treatments, significantly reducing complications. In cardiovascular care, predictive models have been applied to identify high-risk patients for heart failure, enabling early preventative measures. These advancements highlight the practical applications of these technologies in enhancing patient outcomes and optimizing clinical decision-making [32] [33].

#### 2.4. Gaps in the Literature

Despite the growing body of research on big data analytics and predictive analytics in healthcare, several gaps remain that the current study seeks to address. First, while existing studies have developed and tested predictive models for specific healthcare scenarios, there is a need for a more comprehensive evaluation of model performance, considering not only accuracy but also interpretability and clinical relevance [34] [35]. Second, the literature has extensively documented the challenges associated with data quality, integration, and privacy in the context of big data analytics [36] [37], but there is a need for research that directly addresses these challenges and proposes practical solutions to enable the effective implementation of predictive analytics in healthcare. Third, much of the existing research has focused on the technical aspects of predictive analytics, without providing sufficient guidance on how these models can be effectively integrated into clinical workflows and decision-making processes [38] [39]. There is a need for research that bridges the gap between academic findings and practical clinical applications. Finally, as the use of predictive analytics in healthcare becomes more prevalent, there is a growing need to address the ethical implications, such as patient privacy, data ownership, and the potential for bias in algorithmic decisionmaking [40] [41]. Existing research has largely overlooked these ethical considerations, which must be addressed to ensure the responsible and trustworthy application of predictive analytics in healthcare.

The current study aims to address these gaps by developing and evaluating predictive models that are not only accurate but also interpretable and clinically relevant, while simultaneously recognizing and tackling the challenges associated with data quality, integration, and privacy. Furthermore, this research will provide recommendations for the effective implementation of predictive analytics in clinical practice, with a strong emphasis on the ethical considerations surrounding the use of patient data.

#### 3. Method

This study employs a quantitative research design to develop and evaluate predictive models for healthcare applications. The research approach consists of four key components: data sources, model development, model evaluation, and addressing challenges related to data quality, integration, and privacy.

#### **3.1. The Processing of the Analysis**

Data processing is a crucial step in ensuring the quality, integrity, and usability of the dataset for analysis. This study utilizes data from various sources, including electronic health records (EHRs), wearable devices, and genomic data, which must be meticulously prepared before applying Big Data Analytics illustration in **Figure 1**.



Figure 1. Process of big data analytics.

Data Capture: The first step involves capturing data from multiple sources. In this study, data is obtained from EHRs, wearable devices that monitor vital signs, and genomic databases that provide genetic information about patients. Each data source is characterized by different formats, structures, and levels of granularity, making it essential to standardize and synchronize the data. Real-time data collection from wearable devices is also incorporated to integrate up-to-date information on patient health metrics.

Data Preparation: Data preparation is a critical phase that involves cleaning, normalizing, and transforming the data to ensure it is accurate and consistent. Missing data points are addressed through imputation techniques. Numerical data, such as lab results and vital signs, are normalized to a common scale using min-max scaling, while categorical data, including medication and diagnosis codes, are encoded into numerical formats using one-hot encoding. This preparation is necessary to make the data compatible with advanced analytical models and algorithms.

Data Processing: The prepared data undergoes further processing, where it is transformed into a structured format that is ready for analysis. This includes integrating different data types from the EHRs, wearable devices, and genomic data into a cohesive dataset. Data processing also involves handling outliers through statistical techniques, ensuring that extreme values do not skew the analysis results. Standardizing the data allows for seamless integration across platforms and prepares it for deeper analysis through machine learning algorithms.

Big Data Analytics: Once the data is processed, it is fed into Big Data Analytics platforms, which apply sophisticated algorithms to extract meaningful insights. This includes big data mining techniques that identify patterns and correlations within the data, and machine learning models that predict treatment responses and patient outcomes. Algorithms are specifically tailored to handle healthcare data, addressing issues such as patient variability and data complexity.

Result of Big Data Analysis: The final output is a set of analytical results that provide actionable insights into clinical decision-making. These results help in optimizing patient care, enhancing treatment protocols, and improving overall healthcare management at the Children's Clinic of Michigan. Insights derived from the analysis can inform personalized treatment plans, predict patient risks, and support data-driven decision-making processes.

#### 3.2. Data Collection

The data sources and collection methods employed in this study focus on integrating diverse data types to evaluate Big Data analytics in clinical settings. The primary data sources include Electronic Health Records (EHRs), wearable device data, genomic data, and patient socio-economic and lifestyle factors.

The study utilized data extracted from EHRs at the Children's Clinic of Michigan, which contained detailed patient information, including demographics, diagnoses, medications, lab results, comorbidities, treatment responses, and visit frequency. Wearable device data provided insights into patients' real-time health metrics, such as heart rate, blood oxygen levels, and physical activity levels. Genomic data were included to explore potential genetic influences on treatment responses and disease progression.

Socio-economic and lifestyle factors, such as dietary habits and socio-economic status (SES), were collected through patient surveys and linked with EHR data. All data from these diverse sources were merged into a cohesive dataset using a centralized data management platform, ensuring data normalization and integration.

Patients included in the study were selected based on criteria such as being under active treatment at the clinic within the past five years and having complete data available from all sources. The timeframe for data collection ranged from January 2019 to December 2023. Ethical considerations were prioritized throughout the data collection process, including obtaining necessary approvals and securing patient consent where applicable.

The structured data collection process ensures a comprehensive and integrated approach, providing a robust foundation for evaluating the impact of Big Data analytics on patient outcomes at the Children's Clinic of Michigan.

#### 3.3. Model Analysis

The analysis involved the use of various machine learning models to predict treatment responses, frequent visits, and other clinical outcomes based on integrated patient data. Key features such as age, gender, diagnosis codes, treatment response, vital signs, and wearable device data were selected to build predictive models.

Once the features were selected, the data was preprocessed. This included handling missing values, normalizing continuous variables, and encoding categorical variables using techniques such as one-hot encoding. The data was then split into training and test sets, with 80% of the data used for training the models and the remaining 20% reserved for testing their performance. This split ensures that the models are trained on a robust dataset and can be evaluated on unseen data to assess their generalizability. The training process involved feeding the training dataset into each machine learning model and adjusting the model parameters to minimize prediction errors. Hyperparameter tuning was performed to optimize the performance of each model. Hyperparameters are settings that are not learned from the data but are set before the training process begins, such as the regularization strength in logistic regression, the number of trees in a random forest, and the learning rate in a neural network. The hyperparameters were tuned using techniques such as grid search and cross-validation to find the combination that yielded the best performance on the validation set. Cross-validation was particularly important to prevent overfitting, which occurs when a model performs well on the training data but poorly on new, unseen data. By training the model on different subsets of the data and validating it on the remaining subset, cross-validation provides a more reliable estimate of the model's performance.

#### 3.3.1. Machine Learning Models

The study employed three different machine learning models: Logistic Regression, Random Forest, and Neural Networks. Each model was chosen for its unique strengths and applicability to different aspects of the dataset.

#### • Logistic Regression:

Logistic Regression is a linear model commonly used for binary classification tasks. In this study, it was applied to predict binary outcomes such as treatment response (good or poor) and the likelihood of frequent hospital visits (yes or no). The logistic regression model can be expressed as:

$$\hat{y} = \sigma \left( X \beta + \epsilon \right) \tag{1}$$

where  $\hat{y}$  is the predicted probability of the binary outcome, *X* is the vector of input features,  $\beta$  is the vector of coefficients,  $\epsilon$  is the error term, and  $\sigma$  is the sigmoid function. The model was tuned using a grid search over a range of regularization parameters (*C*) and penalty types (L1 or L2).

#### • Random Forest:

Random Forest is an ensemble learning method that combines multiple decision trees to improve predictive performance and reduce overfitting. It is particularly effective for handling complex, high-dimensional data, and provides valuable insights into feature importance. The Random Forest algorithm builds multiple decision trees during training, each on a random subset of the data and features.

$$\hat{y} = \frac{1}{N} \sum_{i=1}^{N} h_i(X)$$
 (2)

where  $\hat{y}$  is the predicted outcome, *N* is the number of trees in the forest, and  $h_i(X)$  is the prediction from the *i*-th tree for input features *X*.

The final prediction is made by averaging the predictions of all the trees (in the case of regression) or by majority voting (in the case of classification). The model was tuned by adjusting the number of trees, the maximum depth of the trees, and the minimum number of samples required to split an internal node.

#### • Neural Networks:

Neural Networks are powerful models capable of capturing complex, non-linear relationships in data. In this study, they were applied to analyze patterns in wearable device data and genomic data, with a focus on predicting personalized treatment plans and identifying high-risk patients. The neural network model used in this study consisted of multiple layers: an input layer corresponding to the selected features, hidden layers with neurons and activation functions to capture non-linearity, and an output layer producing the final predictions.

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g\left(z^{[l]}\right)$$
(3)

where  $z^{[l]}$  is the linear transformation in layer l,  $W^{[l]}$  and  $\backslash b^{[l]}$  are the weights and biases for layer l,  $a^{[l-1]}$  is the activation from the previous layer, and g is the activation function (e.g., ReLU, Sigmoid).

The neural network was fine-tuned using a combination of grid search and adaptive learning techniques, focusing on achieving the best balance between accuracy and computational efficiency.

Through the application of these three machine learning models, the study successfully built predictive models that leveraged diverse patient data to enhance clinical decision-making. The models not only demonstrated strong predictive capabilities but also provided valuable insights into the factors influencing treatment outcomes and patient risk profiles, highlighting the potential of Big Data analytics in transforming clinical practice.

#### 3.3.2. Model Evaluation

Several measures were used to assess the model's performance after training, including accuracy, precision, recall, confusion matrix, and F1-score. These measurements are essential for determining how well the model classifies feelings.

1) Accuracy: This represents the overall correctness of the model and is calculated as the ratio of correctly predicted instances (both True Positives and True Negatives) to the total number of instances.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$
(4)

where:

- TP = True Positives;
- TN = True Negatives;
- FP = False Positives;
- FN = False Negatives.

2) Precision: Indicates the proportion of correct identifications. It is calculated as:

$$Precision = \frac{TP}{TP + FP}$$
(5)

3) Recall: Measures the ability of the model to identify all relevant instances and is calculated as:

$$\operatorname{Recall} = \frac{TP}{TP + FN} \tag{6}$$

4) F1-Score: The harmonic mean of Precision and Recall, providing a single metric that balances both concerns:

F1-Score = 
$$2 \times \frac{\text{Percison} \times \text{Recall}}{\text{Precidion} + \text{Recall}}$$
 (7)

The model's performance was further examined in the confusion matrix and classification report, which displayed the distribution of true versus predicted classes and offered insights into the model's strong and weak points.

# 4. Result and Finding

This study aimed to leverage Big Data analytics to enhance predictive models for patient outcomes in the healthcare sector. The findings provide valuable insights into the factors influencing treatment responses, visit frequencies, and other critical clinical outcomes from **Figure 2** and **Figure 3**. The patient cohort analysis revealed a diverse demographic profile, with a mean age of 15.51 years and a gender distribution of 59% females and 41% males. The study population exhibited a wide range of diagnosis codes, medications, lab results, and comorbidities, highlighting the complexity of the healthcare challenges faced. Notably, Vitamin D Deficiency and Asthma emerged as prevalent conditions. The majority of patients (81%) reported a good treatment response, indicating the effectiveness of the interventions.



Figure 2. Patient data distributions and relationships.





The data visualizations uncovered several key trends. The age distribution was skewed towards younger patients, with a substantial proportion under 10 years old. The Body Mass Index (BMI) followed a near-normal distribution, with most patients falling within a healthy weight range. While treatment responses were predominantly positive, the lack of correlation between age and treatment duration suggested the need to explore additional predictive factors.

These findings underscore the potential of Big Data analytics to inform clinical decision-making and optimize patient outcomes. The insights gained can guide the development of tailored interventions, particularly for pediatric populations and patients with specific health conditions. Future research should further investigate the complex interplay between patient demographics, comorbidities, and healthcare engagement patterns to refine predictive models and enhance the de-livery of personalized, data-driven healthcare.

#### **Model Performance Analysis**

The logistic regression model (**Table 1**) used to predict treatment responses achieved an overall accuracy of 85%. The model performed exceptionally well in classifying the Good response category, with a precision of 0.85, recall of 1.00, and an F1-score of 0.92 for 17 samples. However, the model failed to correctly classify the Not Given & Moderate category, with both precision and recall at 0.00, indicating that it could not distinguish this response type.

The confusion matrix further highlights the model's limitations, showing that it classified all Good responses correctly but misclassified all other responses as Good. This suggests that while logistic regression performs well for dominant classes, it struggles with minority or less frequent categories, which impacts the overall macro average scores (precision = 0.43, recall = 0.50, F1-score = 0.46).

Table 1. Logistic regression analysis.

395

	Precision	Recall	F1-score	Support			
Good	0.85	1.00	0.92	17			
Not Given & Moderate	0.00	0.00	0.00	3			
Logistic Regression Accuracy: 0.85							
Accuracy			0.85	20			
Macro avg	0.43	0.50	0.46	20			
Weighted avg	0.72	0.85	0.78	20			
<b>Confusion Matrix:</b> $\begin{bmatrix} 17 & 0 \\ 3 & 0 \end{bmatrix}$							

The Random Forest model (**Table 2**) outperformed the logistic regression, achieving an overall accuracy of 92%. The model demonstrated a strong performance in classifying the Good response, with a precision of 0.75 and a recall of 0.89, indicating effective identification of this category. For the Not Given category,

the model achieved a recall of 0.67 and a higher F1-score of 0.78, suggesting a more balanced performance in capturing less frequent responses.

The Moderate category had no representation in the support data, which led to some inconsistency in precision and recall values. However, the confusion matrix emphasizes that the Random Forest model managed to correctly classify more instances across all response categories compared to the logistic regression model, highlighting its robustness in handling diverse patient outcomes.

	Precision	Recall	F1-score	Support			
Good	0.75	0.89	0.67	17			
Not Given	0.56	0.67	0.78	3			
Moderate	0.78	0.78	0.76	0			
Random Forest Accuracy: 0.92							
Accuracy			0.92	20			
Macro avg	0.60	0.52	0.54	20			
Weighted avg	0.90	0.81	0.90	20			
	Confusion Matrix:	$\begin{bmatrix} 15 & 2 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$	0 0 3				

Table 2. Random forest model.

The feature importance (Figure 4) analysis from the Random Forest model identified Frequent visits as the most significant predictor of patient outcomes, underscoring the impact of healthcare engagement on treatment success. Length of treatment ranked second, reflecting its influence on outcomes, likely due to extended care for complex conditions. Body Mass Index (BMI) was the third most important feature, indicating its relevance in clinical assessments.





396

Age was moderately important, suggesting it contributes to outcomes but is not as critical as the more direct health engagement metrics. Finally, Gender was deemed the least important, implying minimal impact on predicting outcomes, which aligns with contemporary findings that often de-emphasize gender in favor of more specific health metrics.

The Neural Network model (Table 3) achieved an accuracy of 90%, with the highest performance for the Moderate response category (precision = 1.00, recall = 0.67, F1-score = 0.80). The Good category maintained high performance, similar to the other models, but struggled with the Not Given response, where precision and recall were zero, indicating no correct classifications.

The neural network's overall weighted averages (precision = 0.95, recall = 0.90, F1-score = 0.92) reflect its effectiveness, especially in recognizing complex patterns within the dominant classes. This suggests that the neural network model provides nuanced recognition of complex patterns, particularly in underrepresented response types, though it is slightly less accurate than the Random Forest model.

	Precision	Recall	F1-score	Support			
Good	0.75	0.89	0.94	17			
Not Given	0.00	0.00	0.00	0			
Moderate	1.00	0.67	0.80	3			
Neural Network Accuracy: 0.90							
Accuracy			0.90	20			
Macro avg	0.65	0.54	0.58	20			
Weighted avg	0.95	0.90	0.92	20			
	Confusion I	<b>Matrix:</b> $\begin{bmatrix} 16 & 1 \\ 0 & 0 \\ 1 & 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$				

Table 3. Neural network.

The comparative analysis of logistic regression, Random Forest, and neural network models reveals that while all models perform well for frequent response categories, logistic regression struggles with minority classes. Random Forest offers a balanced performance across all categories, making it highly suitable for clinical decision-making. Neural networks, though slightly less accurate than Random Forests, provide nuanced recognition of complex patterns, especially in underrepresented response types.

These insights are critical for refining predictive models and ensuring robust clinical application. The findings highlight the importance of considering various modeling techniques and their strengths in accurately predicting patient outcomes, which can ultimately inform and optimize healthcare decision-making.

# **5.** Conclusions

397

This study has provided valuable insights into the application of data-driven

predictive modeling in the healthcare domain. By analyzing a comprehensive dataset encompassing patient demographics, clinical characteristics, and healthcare utilization patterns, we have identified key factors that influence treatment outcomes and patient engagement.

The comparative analysis of logistic regression, Random Forest, and neural network models revealed that while all performed well in predicting dominant response categories, the Random Forest model offered the most balanced and robust performance across all classes, including minority responses. This finding highlights the importance of considering diverse modeling approaches to ensure accurate and reliable predictions that can inform clinical decision-making.

The feature importance analysis further underscored the critical role of healthcare engagement metrics, such as visit frequency and treatment duration, in driving positive patient outcomes. Additionally, body mass index (BMI) emerged as a significant predictor, emphasizing the need to incorporate comprehensive health assessments in predictive models.

These insights hold profound implications for the healthcare sector. By leveraging data-driven predictive analytics, clinicians can develop more personalized and targeted interventions, optimizing resource allocation and improving overall patient care. Moreover, the ability to accurately identify at-risk individuals or subgroups can facilitate proactive risk mitigation strategies, reducing the burden on healthcare systems and enhancing population health outcomes.

As the healthcare industry continues to embrace the transformative potential of big data and advanced analytics, this study serves as a valuable framework for integrating data-guided methods into clinical practice. By harnessing the power of predictive modeling, healthcare providers can navigate the complexities of modern medicine more effectively, ultimately enhancing the quality of life for patients and communities.

## 6. Future Research

Blockchain offers promising solutions for enhancing health data privacy and security, particularly in managing sensitive patient data. Additionally, advancements in neural networks, including deep learning models, are expected to further improve predictive accuracy by handling more complex datasets like genomic information and real-time monitoring from IoT-enabled devices. These technologies could revolutionize big data analytics in healthcare by fostering more secure, accurate, and personalized care for patients.

# **Authors' Contributions**

Md Nagib Mahfuz Sunny led the research, conducted the literature review, and implemented the machine learning models. Mohammad Balayet Hossain Sakil contributed to model development, evaluation, and optimization. Abdullah Al Nahian focused on data integration, preprocessing, and addressing privacy concerns. Syed Walid Ahmed provided clinical insights, assessing the medical relevance and impact of the models on patient care. Md. Newaz Shorif assisted in interpreting the results and making ethical considerations. Jennet Atayeva contributed to policy-related aspects and manuscript writing. All authors reviewed and approved the final manuscript

# **Conflicts of Interest**

The authors declare no conflicts of interest regarding the publication of this paper.

#### References

- Abedjan, Z., Boujemaa, N., Campbell, S., Casla, P., Chatterjea, S., Consoli, S., *et al.* (2019) Data Science in Healthcare: Benefits, Challenges and Opportunities. In: Consoli, S., Reforgiato Recupero, D. and Petković, M., Eds., *Data Science for Healthcare*, Springer, 3-38. <u>https://doi.org/10.1007/978-3-030-05249-2\_1</u>
- [2] Agrawal, R. and Prabakaran, S. (2020) Big Data in Digital Healthcare: Lessons Learnt and Recommendations for General Practice. *Heredity*, **124**, 525-534. <u>https://doi.org/10.1038/s41437-020-0303-2</u>
- [3] Ahmed, Z., Mohamed, K., Zeeshan, S. and Dong, X. (2020) Artificial Intelligence with Multi-Functional Machine Learning Platform Development for Better Healthcare and Precision Medicine. *Database*, **2020**, baaa010. https://doi.org/10.1093/database/baaa010
- Srividhya, G. (2022) Electronic Health Records: A Transitional View. In: Hemalatha, R.J., Akila, D., Balaganesh, D. and Paul, A., Eds., *The Internet of Medical Things* (*IoMT*) *Healthcare Transformation*, Wiley, 289-300. https://doi.org/10.1002/9781119769200.ch15
- [5] Panayides, A.S., Amini, A., Filipovic, N.D., Sharma, A., Tsaftaris, S.A., Young, A., et al. (2020) AI in Medical Imaging Informatics: Current Challenges and Future Directions. *IEEE Journal of Biomedical and Health Informatics*, 24, 1837-1857. https://doi.org/10.1109/jbhi.2020.2991043
- [6] Bruintjies, A. (2022) Factors Affecting Big Data Adoption in a Government Organisation in the Western Cape. *South African Journal of Information Management*, 26, a1690. <u>https://doi.org/10.4102/sajim.v26i1.1690</u>
- He, J., Baxter, S.L., Xu, J., Xu, J., Zhou, X. and Zhang, K. (2019) The Practical Implementation of Artificial Intelligence Technologies in Medicine. *Nature Medicine*, 25, 30-36. <u>https://doi.org/10.1038/s41591-018-0307-0</u>
- [8] Awrahman, B.J., Aziz Fatah, C. and Hamaamin, M.Y. (2022) A Review of the Role and Challenges of Big Data in Healthcare Informatics and Analytics. *Computational Intelligence and Neuroscience*, 2022, 5317760. <u>https://doi.org/10.1155/2022/5317760</u>
- [9] Price, W.N. and Cohen, I.G. (2019) Privacy in the Age of Medical Big Data. *Nature Medicine*, **25**, 37-43. <u>https://doi.org/10.1038/s41591-018-0272-7</u>
- [10] Shams, I., Ajorlou, S. and Yang, K. (2021) A Predictive Analytics Approach to Reducing 30-Day Avoidable Readmissions among Patients with Heart Failure, Acute Myocardial Infarction, Pneumonia, or COPD. *Health Care Management Science*, 24, 82-101. <u>https://doi.org/10.1007/s10729-020-09501-9.</u>
- [11] Wang, L. and Alexander, C.A. (2020) Big Data Analytics in Medical Engineering and Healthcare: Methods, Advances and Challenges. *Journal of Medical Engineering & Technology*, 44, 267-283. <u>https://doi.org/10.1080/03091902.2020.1769758</u>
- [12] Himanen, L., Geurts, A., Foster, A.S. and Rinke, P. (2019) Data-Driven Materials

Science: Status, Challenges, and Perspectives. *Advanced Science*, **6**, Article ID: 1900808. <u>https://doi.org/10.1002/advs.201900808</u>

- [13] Negro, A. (2021) Graph-Powered Machine Learning. Simon and Schuster.
- [14] Mansour, M., Saeed Darweesh, M. and Soltan, A. (2024) Wearable Devices for Glucose Monitoring: A Review of State-Of-The-Art Technologies and Emerging Trends. *Alexandria Engineering Journal*, 89, 224-243. https://doi.org/10.1016/j.aej.2024.01.021
- [15] Rayan, R.A., Tsagkaris, C., Zafar, I., Moysidis, D.V. and Papazoglou, A.S. (2022) Big Data Analytics for Health: A Comprehensive Review of Techniques and Applications. In: Kulkarni, A.J., *et al.*, Eds., *Big Data Analytics for Healthcare*, Elsevier, 83-92. <u>https://doi.org/10.1016/b978-0-323-91907-4.00002-9</u>
- [16] Lalmi, F. and Adala, L. (2021) Big Data for Healthcare: Opportunities and Challenges. *Studies in Computational Intelligence*, 935, 217-229. <u>https://doi.org/10.1007/978-3-030-62796-6\_12</u>
- [17] Wang, Y., Kung, L., Gupta, S. and Ozdemir, S. (2019) Leveraging Big Data Analytics to Improve Quality of Care in Healthcare Organizations: A Configurational Perspective. *British Journal of Management*, **30**, 362-388. <u>https://doi.org/10.1111/1467-8551.12332</u>
- [18] Macias, C.G. and Carberry, K.E. (2020) Data Analytics for the Improvement of Healthcare Quality. In: Giardino, A., Riesenberg, L. and Varkey, P., Eds., *Medical Quality Management*, Springer, 121-138. <u>https://doi.org/10.1007/978-3-030-48080-6\_6</u>
- [19] Kumar, A., Le, D.N., Dubey, A.K., Kumar, S.A. and Bhatia, S. (2022) Evolving Predictive Analytics in Healthcare: New AI Techniques for real-Time Interventions. Institution of Engineering Technology. <u>https://doi.org/10.1049/pbhe043e</u>
- [20] Eschenbrenner, B. (2019) Identifying Essential Factors for Deriving Value from Big Data Analytics in Healthcare. In: Nah, F.H. and Siau, K., Eds., *HCI in Business, Government and Organizations. Information Systems and Analytics. HCII* 2019, Springer, 189-198. <u>https://doi.org/10.1007/978-3-030-22338-0\_15</u>
- [21] Iqbal, R., Doctor, F., More, B., Mahmud, S. and Yousuf, U. (2020) Big Data Analytics: Computational Intelligence Techniques and Application Areas. *Technological Forecasting and Social Change*, **153**, Article ID: 119253. <u>https://doi.org/10.1016/j.techfore.2018.03.024</u>
- [22] Ravi, V. and Kumar Cherukuri, A. (2021) Handbook of Big Data Analytics Volume 2: Applications in ICT, Security and Business Analytics. Institution of Engineering Technology. <u>https://doi.org/10.1049/pbpc037g</u>
- [23] Pramanik, M.I., Lau, R.Y.K., Azad, M.A.K., Hossain, M.S., Chowdhury, M.K.H. and Karmaker, B.K. (2020) Healthcare Informatics and Analytics in Big Data. *Expert Systems with Applications*, **152**, Article ID: 113388. <u>https://doi.org/10.1016/j.eswa.2020.113388</u>
- [24] Mehta, N., Pandit, A. and Kulkarni, M. (2019) Elements of Healthcare Big Data Analytics. *Studies in Big Data*, 66, 23-43. <u>https://doi.org/10.1007/978-3-030-31672-3\_2</u>
- [25] Zafar, F., Raza, S., Khalid, M.U. and Tahir, M.A. (2019) Predictive Analytics in Healthcare for Diabetes Prediction. *Proceedings of the* 2019 9th International Conference on Biomedical Engineering and Technology, Tokyo, 28-30 March 2019, 253-259. <u>https://doi.org/10.1145/3326172.3326213</u>
- [26] Thangarasu, G. and Subramanian, K. (2019) Big Data Analytics for Improved Care Delivery in the Healthcare Industry. *International Journal of Online and Biomedical*

Engineering (iJOE), 15, 40-51. https://doi.org/10.3991/ijoe.v15i10.10875

- [27] Stiglic, G., Kocbek, P., Fijacko, N., Zitnik, M., Verbert, K. and Cilar, L. (2020) Interpretability of Machine Learning-Based Prediction Models in Healthcare. *WIREs Data Mining and Knowledge Discovery*, **10**, e1379. <u>https://doi.org/10.1002/widm.1379</u>
- [28] Weerasinghe, K., Scahill, S.L., Pauleen, D.J. and Taskin, N. (2022) Big Data Analytics for Clinical Decision-Making: Understanding Health Sector Perceptions of Policy and Practice. *Technological Forecasting and Social Change*, **174**, Article ID: 121222. https://doi.org/10.1016/j.techfore.2021.121222
- [29] Subramanian, M., Shanmuga Vadivel, K., Hatamleh, W.A., Alnuaim, A.A., Abdelhady, M. and Ve, S. (2021) The Role of Contemporary Digital Tools and Technologies in COVID-19 Crisis: An Exploratory Analysis. *Expert Systems*, **39**, e12834. <u>https://doi.org/10.1111/exsy.12834</u>
- [30] Ranjan, R. and Sahana, B.C. (2024) A Comprehensive Roadmap for Transforming Healthcare from Hospital-Centric to Patient-Centric through Healthcare Internet of Things (IoT). *Engineered Science*, **30**, Article 1175. <u>https://doi.org/10.30919/es1175</u>
- [31] Virginia Anikwe, C., Friday Nweke, H., Chukwu Ikegwu, A., Adolphus Egwuonwu, C., Uchenna Onu, F., Rita Alo, U., *et al.* (2022) Mobile and Wearable Sensors for Data-Driven Health Monitoring System: State-Of-The-Art and Future Prospect. *Expert Systems with Applications*, **202**, Article ID: 117362. https://doi.org/10.1016/j.eswa.2022.117362
- [32] Hosseini, M.M., Zargoush, M., Alemi, F. and Kheirbek, R.E. (2020) Leveraging Machine Learning and Big Data for Optimizing Medication Prescriptions in Complex Diseases: A Case Study in Diabetes Management. *Journal of Big Data*, 7, Article No. 26. <u>https://doi.org/10.1186/s40537-020-00302-z</u>
- [33] Fernandes, D. (2024) Prescriptive Analytics in Healthcare: Advanced Decision Making for Optimal Treatment. <u>https://www.theseus.fi/bitstream/handle/10024/863050/Fernandes\_Daniel.pdf?sequence=2</u>
- [34] Rehman, A., Naz, S. and Razzak, I. (2021) Leveraging Big Data Analytics in Healthcare Enhancement: Trends, Challenges and Opportunities. *Multimedia Systems*, 28, 1339-1371. <u>https://doi.org/10.1007/s00530-020-00736-8</u>
- [35] Rieke, N., Hancox, J., Li, W., Milletari, F., Roth, H.R., Albarqouni, S., *et al.* (2020) The Future of Digital Health with Federated Learning. *NPJ Digital Medicine*, **3**, Article No. 119. <u>https://doi.org/10.1038/s41746-020-00323-1</u>
- [36] Singh, S., Cha, J., Kim, T. and Park, J. (2021) Machine Learning Based Distributed Big Data Analysis Framework for Next Generation Web in IoT. *Computer Science* and Information Systems, 18, 597-618. https://doi.org/10.2298/csis200330012s
- [37] Li, B., Du, K., Qu, G. and Tang, N. (2023) Big Data Research in Nursing: A Bibliometric Exploration of Themes and Publications. *Journal of Nursing Scholarship*, 56, 466-477. <u>https://doi.org/10.1111/jnu.12954</u>
- [38] Et al., N.K. (2023) Harnessing the Power of Big Data: Challenges and Opportunities in Analytics. Tuijin Jishu/Journal of Propulsion Technology, 44, 363-371. https://doi.org/10.52783/tjjpt.v44.i2.193
- [39] Kaur, P. (2023) Internet of Things (IoT) and Big Data Analytics (BDA) in Healthcare. In: Lytras, M.D., Housawi, A.A. and Alsaywid, B.S., Eds., *Digital Transformation in Healthcare in Post-Covid-19 Times*, Elsevier, 45-57. https://doi.org/10.1016/b978-0-323-98353-2.00015-0
- [40] Saber, H., Somai, M., Rajah, G.B., Scalzo, F. and Liebeskind, D.S. (2019) Predictive

Analytics and Machine Learning in Stroke and Neurovascular Medicine. *Neurological Research*, **41**, 681-690. <u>https://doi.org/10.1080/01616412.2019.1609159</u>

[41] Shilo, S., Rossman, H. and Segal, E. (2020) Axes of a Revolution: Challenges and Promises of Big Data in Healthcare. *Nature Medicine*, 26, 29-38. <u>https://doi.org/10.1038/s41591-019-0727-5</u>