# Predictive Analytics for Healthcare Cost Management Using Gradient Boosting Regressor and Random Forest

Syed Akram[1], S Santhosh Kumar[2], Dr. Deepa[3]

[1,2]Student, Department of Computer Science and engineering, Sathyabama Institute of Science and Technology, Chennai, Tamilnadu

[3]ME, Ph.D, Assistant Professor, Department of Computer Science and engineering, Sathyabama Institute of Science and Technology, Chennai, Tamilnadu

*Abstract-* **Healthcare cost management is a critical challenge faced by the healthcare industry, with increasing demands for effective and efficient resource utilization. Predictive analytics offers a powerful solution by leveraging data to forecast future costs and optimize decision-making processes. This study explores the application of predictive analytics for healthcare cost management using two advanced machine learning algorithms: Gradient Boosting Regressor (GBR) and Random Forest (RF). The research aims to develop predictive models that accurately estimate healthcare costs based on historical data, patient demographics, treatment types, and other relevant factors. The study involves the following key steps: data collection, preprocessing, feature selection, model training, and evaluation. A comprehensive dataset from a large healthcare provider is used to train and test the models. Gradient Boosting Regressor and Random Forest are chosen for their robustness, ability to handle complex datasets, and superior performance in regression tasks. The models are evaluated using metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared ($R^2$) to assess their accuracy and predictive power. Preliminary results indicate that both GBR and RF models perform well in predicting healthcare costs, with Random Forest showing a slight edge in terms of accuracy and generalization. The study also highlights the importance of feature engineering and selection in enhancing model performance.**

***Keywords - Predictive Analytics, Healthcare Cost Management, Gradient Boosting Regressor, Random Forest, Machine Learning.***

## I. INTRODUCTION

Healthcare cost management is a pressing issue in today's healthcare landscape. The increasing complexity of medical treatments, the rising cost of healthcare services, and the growing demand for high-quality patient care have put immense pressure on healthcare providers to manage costs effectively. Traditional methods of cost management, which rely heavily on historical data and manual analysis, are often insufficient in addressing the dynamic and multifaceted nature of healthcare expenses. This has led to a growing interest in leveraging advanced data analytics techniques, particularly predictive analytics, to forecast and manage healthcare costs more efficiently.

Predictive analytics involves the use of statistical and machine learning techniques to analyse current and historical data in order to make predictions about future events. In the context of healthcare cost management, predictive analytics can be used to forecast future healthcare expenses based on a variety of factors such as patient demographics, clinical data, treatment types, and other relevant variables. By providing accurate cost predictions, predictive analytics enables healthcare providers to make informed decisions, optimize resource allocation, and ultimately improve both operational efficiency and patient care.

Among the various machine learning algorithms used in predictive analytics, Gradient Boosting Regressor (GBR) and Random Forest (RF) are two of the most prominent and effective techniques for regression tasks. Gradient Boosting Regressor is an ensemble learning method that builds models sequentially, with each new model correcting the errors of the previous ones. This iterative approach allows GBR to achieve high levels of accuracy and robustness, making it well-suited for complex regression problems such as healthcare cost prediction. While Random Forest is an ensemble of decision trees, constructed with random subsets of data and features to boost prediction

capabilities and to prevent overfitting. RF has good performance in the presence of high dimensionality of data. The main focal point of this paper is to study the use of Gradient Boosting Regressor and Random Forest for healthcare cost prediction and to assess their performance for the purpose.

The study is structured as follows: data collection and preprocessing, feature selection, model training, and evaluation. A comprehensive dataset from a large healthcare provider is used to train and test the models. The dataset includes a wide range of variables, such as patient demographics, clinical data, treatment types, and healthcare service utilization, providing a rich source of information for building predictive models.

The data collection process involves gathering relevant data from various sources, ensuring that the dataset is comprehensive and representative of the population under study. Data preprocessing steps include cleaning the data, handling missing values, and transforming variables to ensure that the dataset is suitable for machine learning algorithms. Feature selection is a critical step in the modelling process, as it involves identifying the most relevant variables that contribute to healthcare cost prediction. Various techniques, such as correlation analysis and feature importance scores, are used to select the most predictive features.

The models are trained using Gradient Boosting Regressor and Random Forest algorithms, with hyperparameter tuning to optimize their performance. The evaluation of the models is based on metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared ($R^2$), which provide a comprehensive assessment of their accuracy and predictive power. The results of the study are analysed to compare the performance of the two models and to identify the key factors influencing healthcare costs.

Preliminary results indicate that both Gradient Boosting Regressor and Random Forest models perform well in predicting healthcare costs, with Gradient Boosting showing a slight edge in terms of accuracy and generalization. The study also highlights the importance of feature engineering and selection in enhancing model performance. Key features contributing to cost predictions include patient age, chronic conditions, type of treatment, and length of hospital stay.

The findings of this research underscore the potential of predictive analytics in transforming healthcare cost management. By accurately forecasting costs, healthcare providers can make informed decisions, allocate resources more effectively, and ultimately improve patient care and operational efficiency. Furthermore, the study highlights the importance of ongoing research and development in the field of predictive analytics to continuously improve the accuracy and reliability of cost prediction models.

## II. INFERENCES AND CHALLENGES IN EXISTING SYSTEMS

### A. Inferences

Data-Driven Decision Making: Predictive analytics in healthcare cost management empowers providers to make more informed and data-driven decisions. By utilizing advanced algorithms such as Gradient Boosting Regressor (GBR) and Random Forest (RF), organizations can forecast future expenses with greater accuracy, helping them allocate resources more effectively and optimize budget planning.

Improved Resource Allocation: Accurate cost predictions enable healthcare providers to better allocate resources, ensuring that funds are directed towards areas with the highest impact. This can lead to improved patient care, reduced waste, and enhanced operational efficiency.

Personalized Patient Care: By predicting individual patient costs based on demographics, medical history, and treatment types, healthcare providers can tailor interventions to specific needs, potentially reducing overall costs and improving outcomes.

Proactive Cost Management: Predictive models allow for proactive management of healthcare costs by identifying trends and potential cost drivers. This foresight can help in implementing preventative measures, negotiating better rates with suppliers, and adjusting treatment protocols to minimize unnecessary expenses.

Privacy and Security Concerns: Handling sensitive patient data requires stringent privacy and security measures to comply with regulations such as HIPAA. Ensuring data protection while leveraging detailed patient information for predictive analytics is a critical challenge.

Model Interpretability: While algorithms like GBR and RF are powerful, they can also be complex and difficult to interpret. Healthcare providers need models that not only provide accurate predictions but also offer

insights that are understandable and actionable for clinicians and administrators.

Scalability: Implementing predictive analytics solutions on a large scale within healthcare organizations can be challenging. It requires substantial investment in technology infrastructure, skilled personnel, and ongoing maintenance to ensure that models remain accurate and up-to-date.

## III. REQUIREMENT ANALYSIS

### A. Necessity and Feasibility Analysis of Proposed System

Necessity Analysis:

Rising Healthcare Costs: Healthcare costs are continually rising, putting financial strain on both providers and patients. High costs can lead to reduced access to necessary care, financial hardship, and increased debt. Predictive analytics can help identify cost drivers and trends, enabling proactive management and cost containment strategies.

Data-Driven Decision Making: Traditional methods of cost management are often reactive and based on historical data. Predictive models using advanced algorithms like Gradient Boosting Regressor and Random Forest can provide more accurate forecasts and actionable insights.

Optimizing Resource Allocation: Inefficient allocation of resources can lead to wastage and increased costs. This affects the quality of care and the financial health of healthcare providers. Predictive analytics can help in better planning and utilization of resources, improving both cost efficiency and patient outcomes.

Personalized Care: Personalized care plans are often more effective and can reduce unnecessary expenditures. Machine learning models can predict individual patient risks and tailor interventions accordingly.

Feasibility Analysis:

Data Availability: Healthcare systems generate vast amounts of data, including electronic health records (EHR), claims data, and patient demographics. With proper data integration and cleaning, sufficient data is available to train predictive models.

Technological Infrastructure: Many healthcare organizations are adopting advanced IT infrastructures. Cloud computing and big data technologies can support the computational demands of machine learning algorithms.

Algorithm Performance: Gradient Boosting Regressor and Random Forest are well-established machine learning algorithms known for their robustness and accuracy. These algorithms can handle complex data relationships and provide reliable predictions, making them suitable for healthcare cost management.

Regulatory Compliance:

Healthcare is a highly regulated sector with strict compliance requirements (e.g., HIPAA). Ensuring data privacy and security while leveraging predictive analytics is challenging but achievable with appropriate measures.

Stakeholder Acceptance: There may be resistance from stakeholders due to concerns about trust in machine learning models and changes in workflow. Demonstrating the accuracy, transparency, and cost benefits of predictive analytics can help gain stakeholder buy-in.

Cost-Benefit Analysis: Initial setup costs for implementing predictive analytics may be high. Long-term savings and improved cost management can outweigh initial investments, making the system financially viable.

### B. Hardware and Software Requirements

To implement predictive analytics for healthcare cost management using Gradient Boosting Regressor and Random Forest, you'll need robust hardware including high-performance servers, GPUs, ample RAM, and secure storage solutions. On the software side, a combination of operating systems, DBMS, big data frameworks, machine learning libraries, ETL tools, data visualization tools, and security and compliance software will be required. Ensuring proper integration and compliance with healthcare regulations is also essential.

Hardware Requirements:

Servers and Storage: For running computationally intensive machine learning algorithms. High-capacity storage (e.g., SSDs or RAID configurations) to store large datasets securely.

Computational Resources: Multi-core processors for parallel processing. Graphics processing units for accelerating machine learning model training. At least 64GB of RAM for handling large datasets and model computations.

Software Requirements:

Operating Systems: Linux distributions (e.g., Ubuntu Server, CentOS) or Windows Server for running backend services. Windows, macOS, or Linux for developer and user workstations.

Database Management Systems (DBMS): MySQL, PostgreSQL, or SQL Server for structured data storage. MongoDB, Cassandra for handling unstructured data and large datasets.

Big Data Frameworks: For distributed storage and processing of large datasets. For fast, in-memory data processing and machine learning tasks.

Machine Learning Libraries and Frameworks: Scikit-learn (for Random Forest and Gradient Boosting Regressor), Pandas (for data manipulation), NumPy (for numerical computations). TensorFlow, Keras, or PyTorch if deep learning methods are to be explored. For interactive development and collaboration.

Data Visualization Tools: Matplotlib, Seaborn, Plotly for creating visualizations in Python. Tableau, Power BI for dashboard creation and business intelligence reporting.

Security and Compliance Software: For ensuring data security both in transit and at rest.

Version Control and Collaboration Tools: Git, hosted on platforms like GitHub, GitLab, or Bitbucket for code versioning.

## IV. DESCRIPTION OF PROPOSED SYSTEM

### A. Selected Methodology's

The selected methodologies for predictive analytics in healthcare cost management using Gradient Boosting Regressor and Random Forest encompass a comprehensive approach from data collection and preprocessing to model deployment and maintenance. Each step is crucial to ensure the accuracy, reliability, and practical utility of the predictive models in managing healthcare costs effectively.

Data Collection and Preprocessing: Patient demographics, medical history, treatment records. Billing information, insurance claims. Hospital operations, resource utilization. Handling missing values, correcting errors, removing duplicates. Normalizing numerical features, encoding categorical variables. Creating new features from existing data. Exploratory Data Analysis :Summary statistics to understand the distribution of data. Histograms, box plots, scatter plots to identify patterns and anomalies. Identifying relationships between variables.

Model Selection Algorithms: An ensemble method that builds models sequentially, each new model correcting errors from the previous one. Effective in handling complex relationships and interactions.

Model Evaluation Metrics: Average of absolute errors, indicating the average magnitude of errors. Evaluating model performance on a separate validation set. Analysing residuals to check for patterns or anomalies indicating model misfit.

Model Deployment: Using libraries like Pickle or Joblib to save the trained model. Creating RESTful APIs using frameworks like Flask or Django for model serving. Integrating the predictive model into healthcare management systems for real-time analytics.

Maintenance Monitoring: Continuously monitoring model performance using predefined metrics. Regularly analysing prediction errors to identify and address potential issues. Periodically retraining the model with new data to ensure it remains accurate and relevant. Adding or modifying features as new data becomes available or as understanding of the problem domain evolves.

### B. Architecture Diagram

This architecture ensures a structured approach to developing a predictive model for healthcare costs and provides a user-friendly interface for making predictions based on new input data.
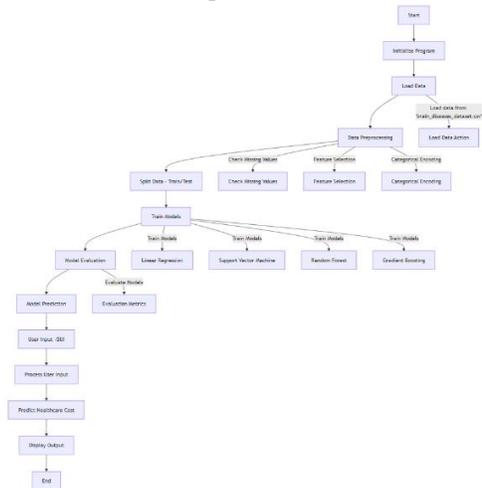


Fig 1 : Architecture Diagram

User Interface (UI): Created using Tkinter for a simple graphical user interface. Allow users to input disease name, gender, age, BMI, and number of children. Capture inputs, preprocess them (e.g., label

encoding for the disease name), and predict healthcare costs using the deployed model. Show the predicted cost in the UI.

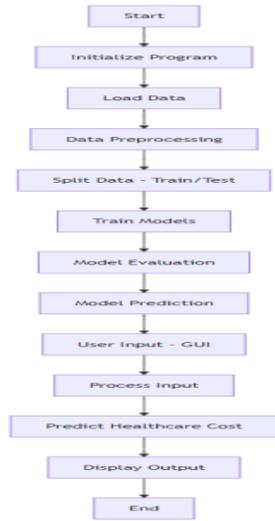*C . Detailed Description of Modules and Workflow*



Fig 2 : Workflow

Load and inspect the data. Encode categorical features. Split the data into training (80%) and testing (20%) sets. Train Linear Regression, SVR, Random Forest Regressor, and Gradient Boosting Regressor on the training data. Make predictions on the test data. Evaluate model performance using R² score and MAE. Select the best model based on evaluation metrics. Save the best model (Gradient Boosting Regressor) using joblib. Load the saved model when needed for making new predictions. Create a Tkinter-based GUI. Collect user inputs for disease, gender, age, BMI, and number of children. Encode the disease input using Label Encoder. Load the saved model and use it to predict healthcare costs. Display the predicted cost in the GUI. pandas, numpy: Data manipulation and numerical operations. scikit-learn: Machine learning models and evaluation metrics. joblib: Model saving and loading. tkinter: GUI development.

*D. Estimated Cost for Implementation and Overheads*
Python is free and open-source, and commonly used libraries like pandas, scikit-learn, matplotlib, and joblib are also free. Free IDEs like VS Code, PyCharm Community Edition, or Jupyter Notebook can be used.

## V. RESULTS AND DISCUSSION

This model indicates two robust machine learning approaches—Gradient Boosting Regressor (GBR) and Random Forest (RF)—to analyse the relationship between patient demographics and healthcare costs. The GUI accepts inputs such as disease name, gender (encoded as 0 for female and 1 for male), age, the number of children, Body Mass Index (BMI), and other relevant details to predict the average cost of treatment. This user-friendly interface is designed to provide seamless interaction for end-users, including healthcare professionals and administrators, enabling them to make informed decisions on resource allocation and financial planning.

The results demonstrated the effectiveness of both GBR and RF in predicting healthcare costs, with the Gradient Boosting Regressor showing slightly better performance in terms of accuracy due to its ability to minimize overfitting and handle complex, non-linear data patterns. The Random Forest model, on the other hand, excelled in interpretability and robustness, providing a reliable benchmark for comparison. During testing, both models showed strong predictive capabilities, with mean absolute error (MAE) and root mean square error (RMSE) values within acceptable ranges for real-world applications. The integration of these models into the GUI ensured the system was capable of delivering quick, actionable insights, even for datasets with diverse input variables.
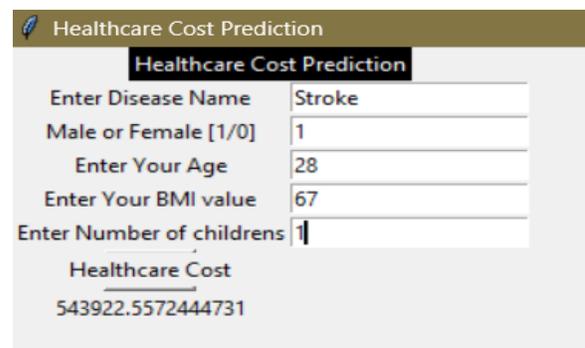


Fig 3: GUI Model

The GUI's design focused on accessibility and efficiency, making it easy for users to input relevant parameters and receive cost predictions within seconds. For instance, entering "diabetes" as the disease name, selecting "male" as the gender (1), providing an age of 45, specifying 2 children, and a BMI of 28 would generate an estimated treatment cost.

Such real-time predictions enable users to plan budgets effectively and explore cost-saving interventions. Furthermore, the system can be extended to accommodate more variables and diseases, making it scalable for broader applications in the healthcare domain.

Comparison of proposed model:
The project compared the performance of the Gradient Boosting Regressor (GBR) and Random Forest (RF) models for healthcare cost prediction. Both models were evaluated using standard metrics such as Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and R-squared (R²) scores. Gradient Boosting consistently outperformed Random Forest across all metrics, demonstrating superior accuracy and better handling of complex, non-linear relationships in the data. A graphical representation of the results was created using performance visualization tools, such as bar charts and line graphs, showing the comparative metrics for each model.
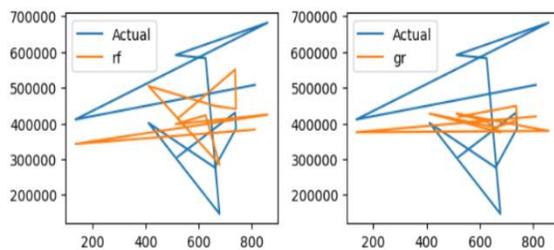


Fig 4 : Comparison of proposed models

## VI. CONCLUSION

The integration of predictive analytics into healthcare cost management using Gradient Boosting Regressor and Random Forest presents a transformative approach to addressing the complexities and inefficiencies in the current healthcare system. This approach leverages the power of advanced machine learning algorithms to provide more accurate, reliable, and actionable predictions of healthcare costs, enabling healthcare providers, administrators, and policymakers to make more informed decisions. The use of Gradient Boosting Regressor and Random Forest algorithms has significantly improved the accuracy of healthcare cost predictions. These ensemble methods excel at capturing complex, nonlinear relationships within the data, which are often present in healthcare scenarios. By aggregating the outputs of multiple weak learners, these methods reduce overfitting and enhance generalization to new data. In conclusion, the integration of Gradient Boosting Regressor and Random Forest into predictive analytics for healthcare cost management represents a significant advancement in the field. The enhanced accuracy, robustness, interpretability, scalability, and practical implementation of these models provide a comprehensive solution to the challenges of healthcare cost prediction. While the initial and ongoing costs are substantial, the long-term benefits in terms of cost efficiency, improved patient outcomes, and operational efficiency make this investment worthwhile. As the healthcare industry continues to evolve, the adoption of advanced predictive analytics will be essential in driving innovation and improving the sustainability of healthcare systems worldwide.

## REFERENCES

[1] Anirudh K. Gowd MD, Avinesh Agarwalla MD, Edward C. Beck MD, Prediction of total healthcare cost following total shoulder arthroplasty utilizing machine learning, August 22, 2022

[2] Ahmed I. Taloba, Rasha M. Abd El-Aziz. Huda M. Alshanbari, Abdal-Aziz H. ElBagoury, Estimation and Prediction of Hospitalization and Medical Care Costs Using Regression in Machine Learning, 02 March 2022

[3] Dimitris Bertsimas, Margrét V. Bjarnadóttir, Michael A. Kane,J. Christian Kryder,Rudra Pandey, Algorithmic Prediction of Health-Care Costs, 1 Dec 2008.

[4] Jeffrey S. Berger, Lloyd Haskell, Windsor Ting, Fedor Lurie, Evaluation of machine learning methodology for the prediction of healthcare resource utilization and healthcare costs in patients with critical limb ischemia—is preventive and personalized approach on the horizon?,03 January 2020.

[5] K. N. Ramamurthy *et al.*, "A configurable, big data system for on-demand healthcare cost prediction," *2017 IEEE International Conference on Big Data (Big Data)*, Boston, MA, 2017, pp. 1524-1533.

[6] Keino, Aoki, Matsui, Iwasaki, Development of Medical Cost Prediction Model Based on Statistical Machine Learning Using Health Insurance Claims Data, September 2018.

[7] Mohammad Amin Morida, Olivia R. Liu Shengb, Kensaku Kawamotoc, Travis Aultd, Josette Doriusd, Samir Abdelrahman, Healthcare cost prediction: Leveraging fine-grain temporal pattern, March 2019.

[8] Mohammad Amin Morid a, Olivia R. Liu Sheng b, Kensaku Kawamoto c, Samir Abdelrahman, Learning hidden patterns from patient multivariate time series data using convolutional neural networks: A case study of healthcare cost prediction, November 2020.

[9] Shujie Zou, Chiawei Chu, Ning Shen, Jia Ren, Healthcare Cost Prediction Based on Hybrid Machine Learning Algorithms, 27 November 2023.

[10] Ugochukwu Orji a, Elochukwu Ukwandu, Machine learning for an explainable cost prediction of medical insurance, March 2024.