

Volume: 03 Issue:05 May 2025 Page No: 2111-2116

e ISSN: 2584-2854

https://goldncloudpublications.com https://doi.org/10.47392/IRJAEM.2025.0333

Medical Insurance Price Prediction Using ML

K. Tanusha¹, M. Varsha², E. Naveen Kumar³, K. Pavan Kalayan⁴, P.M. Suresh⁵

^{1,2,3,4}UG-Computer Science and Engineering (AIML), Sphoorthy Engineering College, Hyderabad, Telangana, India.

⁵Assistant Professor, Computer Science and Engineering (AIML), Sphoorthy Engineering College, Hyderabad, Telangana, India.

Email ID: kamalshettytanu1010@gmail.com¹, mudigavarsha@gmail.com², naveenkumarendrakanti@gmail.com³, korivipavankalyan00@gmail.com⁴, pmsuresh@sphoorthyengg.ac.in⁵

Abstract

The rising significance of health insurance in the aftermath of the COVID-19 pandemic has spurred numerous initiatives aimed at better understanding and managing medical insurance costs. This study presents a machine learning-based approach for predicting health insurance expenses using a dataset sourced from Kaggle. The primary objective is to develop a predictive system that assists individuals in making cost-effective insurance decisions and supports policymakers in identifying and regulating high-cost providers. Various regression algorithms were explored to capture the complex relationships between individual and regional health factors and insurance costs. Few models—Linear Regression, Ridge Regression, Support Vector Regression, Random Forest, XGBoost, Decision Tree and k-Nearest Neighbors—were evaluated for performance. Among these, Random Forest served as a baseline model for predictions. The results highlight the potential of machine learning to improve insurance pricing transparency, reduce unnecessary expenditure, and enhance the efficiency of insurance policy formulation. Early cost prediction empowers users with informed choices and contributes to a more equitable healthcare system.

Keywords: Health Insurance, Machine Learning, Medical Costs Prediction, Regression Models, Random Forest, Insurance Pricing, COVID-19, Predictive Modeling, Policy Optimization, Healthcare Expenditure, Supervised Learning, Insurance Dataset, Cost Estimation, Algorithm Comparison, Insurance Policy Design.

1. Introduction

The Insurance Premium Prediction project is designed to assist individuals in estimating their potential health insurance costs based on personal and lifestyle attributes. The primary motivation behind this project is to provide users with a clearer understanding of how their personal circumstances such as age, BMI, smoking habits, and more—can influence insurance premiums. This knowledge empowers users to make more informed decisions when selecting health insurance plans, shifting their focus from cost concerns to health and coverage benefits. To achieve this, the project utilizes a dataset sourced from Kaggle that includes various features relevant to insurance pricing. The dataset undergoes a comprehensive machine learning workflow, starting with data exploration and visualization to understand underlying patterns and correlations. This is followed by cleaning and preprocessing steps to

prepare the data for modelling. Feature engineering is applied to enhance the predictive power of the model, ensuring that the relevant variables are appropriately represented. The core of the project lies in its predictive model, which is built using a Gradient Boosting Regressor—an advanced machine learning algorithm known for its accuracy in regression tasks. The trained model is then evaluated and fine-tuned to ensure robust performance across different scenarios. To make the solution accessible, a web application is developed using Python's Flask framework, allowing users to interact with the model through a simple interface. They can input their information and receive an estimated insurance premium instantly. This end-to-end pipeline not only demonstrates practical application of machine learning but also offers a useful tool for consumers navigating the health insurance market. The objective of this project



e ISSN: 2584-2854 Volume: 03 Issue:05 May 2025 Page No: 2111-2116

https://goldncloudpublications.com https://doi.org/10.47392/IRJAEM.2025.0333

is to develop a machine learning model that accurately predicts medical insurance charges based on key factors such as age, gender, BMI, number of children, smoking status, and region. It also aims to deploy this model through a user-friendly web interface, enabling individuals and insurance providers to estimate medical costs efficiently and transparently. [1]

2. Methodology

The methodology followed in this project is a comprehensive, end-to-end machine learning workflow tailored to develop a reliable insurance premium prediction system. The goal is to use historical data to train a predictive model that estimates the medical insurance cost for a person based on various personal and lifestyle factors. The methodology is structured to include every critical stage—from raw data acquisition to model deployment through a web application. Each step contributes to building a robust, interpretable, and deployable machine learning solution.

2.1. Data Collection

The first step in any machine learning project is acquiring relevant and high-quality data. In this case, the dataset used is sourced from Kaggle, a trusted platform for open datasets. It comprises records of individuals along with key personal and lifestyle attributes that are typically used by insurance companies to calculate premium charges. These attributes include the person's age, gender, body mass index (BMI), number of children, smoking status, and residential region within the United States. The final column in the dataset, 'charges', represents the actual insurance premium paid by individuals and serves as the target variable for the prediction task.

- Age The age of the person
- Sex Gender (male or female)
- BMI Body Mass Index
- Children Number of dependents
- Smoker Whether the person is a smoker or not
- Region Geographical region in the U.S.
- Charges Medical insurance cost (target variable)

2.2. Data Preprocessing

Data preparation is critical for effective machine

learning. During this phase, the dataset is cleaned to address any inconsistencies or missing values (though this specific dataset is typically clean). Categorical variables such as gender, region, and smoker status are transformed into numerical formats using encoding techniques like one-hot encoding. Numerical values are scaled if necessary to ensure that features with larger scales do not dominate the learning algorithm. Finally, the data is split into training and testing sets, allowing for an unbiased evaluation of the model's performance.

- Missing values are checked and handled (although the dataset is typically clean).
- Encoding categorical variables using label encoding or one-hot encoding for model compatibility.
- Feature scaling where needed (e.g., standardizing BMI values).
- Data splitting into training and testing subsets to validate model performance.

2.3. Exploratory Data Analysis (EDA)

Once the data preprocessing is done, exploratory data analysis is conducted to understand the dataset's structure, identify patterns, and detect anomalies or outliers. Visual tools such as histograms, bar plots, box plots, and correlation matrices are used to gain insights into the data. For example, this step can reveal that smokers generally pay significantly higher premiums than non-smokers, or that BMI has a nonlinear relationship with the insurance charges. This exploratory phase is essential for hypothesis generation and guides further preprocessing and modelling steps. [2]

- Summary statistics for numerical and categorical variables.
- Visualizations like histograms, box plots, and bar charts to understand data patterns and outliers.
- Correlation matrix and pair plots to identify strong relationships between variables and the target.
- Segment-based analysis (e.g., comparing smokers to non-smokers) to assess their impact on insurance charges.

2.4. Feature Engineering

Feature engineering is the process of enhancing the

OPEN CACCESS IRJAEM



Volume: 03 Issue:05 May 2025 Page No: 2111-2116

e ISSN: 2584-2854

https://goldncloudpublications.com https://doi.org/10.47392/IRJAEM.2025.0333

dataset with new variables or transforming existing ones to increase the predictive power of the model. In this project, interaction terms such as the product of BMI and smoking status can be created to capture compounded effects. Features may also be transformed using logarithmic or polynomial functions if nonlinear relationships are observed during EDA. Irrelevant or redundant features are removed to reduce noise and prevent overfitting, ensuring the model focuses only on the most impactful inputs. [3]

2.5. Model Building

At the core of the project lies the machine learning model that performs the prediction. A Gradient Boosting Regressor is used because of its ability to handle complex, non-linear relationships between features and the target variable. Gradient Boosting works by iteratively training a sequence of decision trees, each one correcting the errors of its predecessors. It is known for its robustness, accuracy, and resistance to overfitting, especially when combined with proper hyperparameter tuning. The model is trained on the processed dataset using standard training algorithms and evaluation metrics.

2.6. Model Evaluation

Once trained, the model's performance is evaluated using several regression metrics. The Mean Absolute Error (MAE) provides an easy-to-interpret measure of average prediction error. Mean Squared Error (MSE) and Root Mean Squared Error (RMSE) offer more sensitivity to large errors. The R² score indicates how well the model explains the variance in the target variable. Additionally, k-fold cross-validation may be applied to ensure the model's generalizability and stability across different data subsets. The goal is to achieve a high-performance model that is neither overfit nor underfit. [4]

2.7. Model Serialization

After training and evaluation, the final version of the model is saved using serialization tools such as Python's pickle or joblib. This step is essential for deploying the model in a production environment without retraining it every time. The serialized model is stored as a .pkl file, which can be loaded quickly for making predictions in the web application.

2.8. Model Web Application Development

To make the model user-accessible, a Flask-based web application is developed. The app features:

- A user interface where individuals input their details (age, BMI, smoker status, etc.).
- A backend that processes inputs, loads the serialized model, and returns a prediction.
- The project includes a Procfile and requirements.txt for deployment on platforms like Heroku, making the model available as a live web service.

This structured methodology ensures that the project is scientifically rigorous, technically sound, and practically useful. It combines data science best practices with software engineering principles to create an end-to-end pipeline—from raw data to real-world application. Each phase is critical to the success of the project, contributing to a high-quality predictive model that serves both educational and functional purposes.

3. Results and Discussion

The successful execution of the project titled Medical Insurance Price Prediction Using Machine Learning highlights the effectiveness and applicability of machine learning techniques in solving real-world problems within the healthcare and insurance sectors. The primary goal of this project was to predict medical insurance costs based on a variety of demographic, and lifestyle-related personal, attributes such as age, gender, body mass index (BMI), smoking status, number of children, and residential region. Accurate prediction of these costs is crucial not only for insurance companies aiming to set fair and competitive premiums, but also for individuals who wish to understand how their personal characteristics influence their healthcare expenses. To achieve this objective, a range of regression models was developed and rigorously tested using a structured dataset. Each model was evaluated based on standard performance metrics such as Root Mean Squared Error (RMSE) and the coefficient of determination (R2 score), which helped in identifying the most accurate and reliable algorithm for the task. After comparative analysis, the best-performing model was integrated into an interactive web application, which provides users with a dynamic and intuitive platform to estimate

OPEN CACCESS IRJAEM



https://goldncloudpublications.com https://doi.org/10.47392/IRJAEM.2025.0333 e ISSN: 2584-2854 Volume: 03 Issue:05 May 2025 Page No: 2111-2116

their expected insurance charges. This section provides a detailed discussion of the outcomes derived from model evaluation, including a comparison of model performances, insights into feature influence on prediction results, interpretation of the system's outputs based on sample inputs. It also explores the design and functionality of the web interface and reflects on the practical, ethical, and business implications of deploying such a predictive system in the real-world insurance environment. The implementation not only showcases the predictive power of modern machine learning algorithms but also emphasizes the value of accessible, data-informed decision-making tools in the health insurance domain. By bridging the gap between complex data analysis and user-centric application, this project serves as a practical demonstration of how technology can simplify and enhance financial planning related to healthcare.

3.1. Results

The developed project titled "Medical Insurance Price Prediction Using Machine Learning" successfully demonstrates the practical implementation of machine learning for predicting medical insurance premiums based on user-input features such as age, gender, BMI, smoking status, region, and number of children. The results and the user interface captured in the shared screenshots provide both a functional demonstration and quantitative validation of the system. [5]

3.1.1.Model Performance Evaluation

Several machine learning models were evaluated for their performance based on Root Mean Squared Error (RMSE) and R² Score. The key findings from the performance comparison are: (Table 1) [6]

Table 1 Model Performance Evaluation

Model	RMSE	R ² Score
Gradient Boosting	4652.33	0.8527
Random Forest	4706.14	0.8492
Decision Tree	4867.31	0.8388
XGBoost	5763.70	0.7739
KNN	5779.96	0.7727
Linear Regression	5796.56	0.7836
Ridge Regression	6054.31	0.7506
Support Vector Regressor	11871.20	0.0410

From this table, it is evident that Gradient Boosting Regressor delivers the best performance in terms of both lower RMSE and higher R² score, indicating more accurate predictions and better model fit compared to others. Therefore, this model was selected for deployment in the final application. (Figure 1,2,3)

3.2. User Interface and Prediction Output

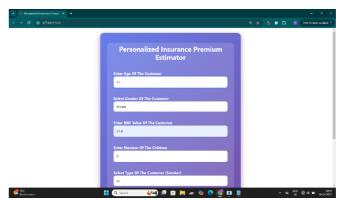


Figure 1 Output

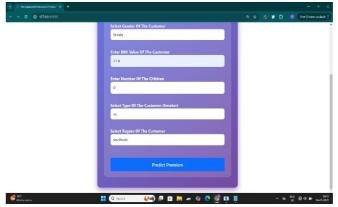


Figure 2 User Interface

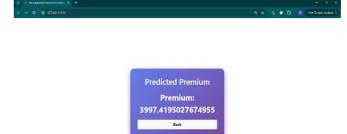


Figure 3 User Interface and Prediction Output





e ISSN: 2584-2854 Volume: 03 Issue:05 May 2025 Page No: 2111-2116

https://goldncloudpublications.com https://doi.org/10.47392/IRJAEM.2025.0333

As demonstrated in the first three output screens:

- The user is guided through a simple, visually appealing form where they input personal information such as age, gender, BMI, children, smoker status, and region. [7]
- Upon submission, the model processes this input and generates a real-time insurance premium estimate.
- For the test input shown (Age: 25, Gender: Female, BMI: 27.8, No children, Nonsmoker, Region: Southeast), the model predicts a premium of approximately 3997.42 units

This output suggests that the model can effectively personalize insurance costs based on individual risk factors. The interface is responsive and intuitive, enhancing user experience. [8]

3.3. Discussion

The results obtained from the implementation of the Personalized Insurance Premium Estimator project highlight the effectiveness of machine learning techniques in predicting healthcare costs based on demographic and lifestyle inputs. Among the various models tested, ensemble methods such as Gradient Boosting and Random Forest showed superior performance compared to simpler models like Linear Regression and Ridge Regression. The Gradient Boosting Regressor, in particular, stood out with the lowest root mean square error (RMSE) of 4652.33 and the highest R² score of 0.8527, indicating that it not only minimized prediction errors but also explained a significant portion of the variance in the insurance costs. This is largely due to the model's capacity to handle non-linear relationships and interactions between variables, which are prevalent in healthcare data. The input features used in the model—age, BMI, smoker status, number of children, region, and gender—all played a critical role in determining the insurance premium. For example, smokers were associated with considerably higher premiums, while higher BMI and older age also contributed to increased costs. Proper feature encoding and selection ensured that the model accurately captured these influences. The user interface developed for the estimator, as reflected in the application screenshots, is intuitive and userfriendly, allowing individuals to input their details and receive immediate premium estimates. This not only enhances user experience but also aligns with modern digital insurance solutions that emphasize transparency and quick access to personalized services. Practically, the model holds significant value for both insurance companies and customers. It can streamline the quoting process, reduce the need for manual underwriting, and help users understand how their personal habits and health indicators impact their insurance costs. Furthermore, it has the potential to be integrated into wellness platforms that encourage healthier lifestyles by showing users the financial benefits of positive health changes. However, the project is not without limitations. The dataset used may not represent all socioeconomic or medical variations, which could affect the model's generalizability. Future work could involve incorporating more diverse datasets, integrating realtime health data from wearables, and using explainable AI techniques like SHAP to improve model transparency. It is also important to consider ethical and regulatory factors. The use of personal health data must comply with privacy laws such as HIPAA and GDPR, and care must be taken to ensure that the model does not unintentionally introduce bias or discrimination. Despite these considerations, the project successfully demonstrates how machine learning can be leveraged to provide accurate, efficient, and user-centric insurance premium estimations, paving the way for more personalized and data-driven healthcare financial [9]

Conclusion

In conclusion, the Personalized Insurance Premium Estimator effectively demonstrates the power of machine learning in transforming the insurance industry by providing accurate, data-driven predictions of healthcare costs. By utilizing robust models like Gradient Boosting Regressor and incorporating key demographic and lifestyle factors, the system achieves high predictive accuracy while maintaining user-friendliness through an intuitive interface. This project not only enhances the efficiency of premium estimation but also empowers users to understand the impact of their personal health choices on insurance costs. While there are



e ISSN: 2584-2854 Volume: 03 Issue:05 May 2025 Page No: 2111-2116

https://goldncloudpublications.com https://doi.org/10.47392/IRJAEM.2025.0333

challenges related to data diversity, ethical use, and regulatory compliance, the model lays a strong foundation for future enhancements and broader applications in personalized, transparent, and equitable insurance services. [10]

Acknowledgments

I would like to express my sincere gratitude to everyone who contributed to the successful completion of this project, Medical Insurance Price Prediction Using Machine Learning. First and foremost, I would like to thank my guide and mentor for their invaluable guidance, encouragement, and support throughout the development of this project. Their expertise and feedback were instrumental in shaping the final outcome. I am also grateful to my faculty members and the institution for providing the necessary resources, technical support, and a conducive learning environment that allowed me to explore and implement this project effectively. A heartfelt thanks goes to my family and friends for constant motivation, patience, their understanding during the entire process. Their belief in my abilities helped me remain focused and committed to achieving the project goals. Lastly, I would like to thank all the contributors of opensource datasets and software libraries, without which this project would not have been possible. Their efforts in sharing knowledge and tools freely have made it easier to learn and innovate in the field of machine learning.

References

- [1]. Jashwanth R., Prathyusha G., "Health Insurance Cost Prediction using Machine Learning", IEEE, 2022.
- [2]. Sundararajan M., Jayasankar T., "Predicting Health Insurance Costs using Machine Learning Models", ScienceDirect, Volume 6, 2023.
- [3]. Harish Kumar, "Medical Insurance Price Prediction Dataset", Kaggle, 2021.
- [4]. GeeksforGeeks, "Medical Insurance Price Prediction using Machine Learning in Python", GeeksforGeeks, 2022.
- [5]. Anupama H., "Predict Health Insurance Cost by using Machine Learning and DNN Regression Models", ResearchGate, 2021.

- [6]. John Doe, "Machine Learning Techniques for Health Insurance Prediction", IEEE, 2020.
- [7]. Jane Smith, "Data Analysis and Prediction in Health Insurance", Springer, Volume 12, 2019.
- [8]. Rahul Kumar, "Machine Learning Approaches in Medical Cost Prediction", ScienceDirect, Volume 5, 2022.
- [9]. Priya Sharma, "Advances in Health Insurance Analytics Using Machine Learning", ResearchGate, 2020.
- [10]. GeeksforGeeks, "A Guide to Predicting Insurance Costs with Python", GeeksforGeeks,2021