



International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com ISSN 2582-7421

Health Insurance Cost Prediction App

Gaurav Kumar¹, Nidhi Prajapati²

B. Tech Student, Information Technology, Buddha Institute of Technology, Gorakhpur, India
1srivastavagaurav7781@gmail.com, 2nidhiprajapati076@gmail.com

ABSTRACT –

Predicting health insurance costs is crucial for insurance providers to determine premiums, allocate resources, and make informed business decisions. In this paper, we present a study on health insurance cost prediction using machine learning techniques. We collected and preprocessed data on demographics, medical history, and lifestyle factors of a sample population. We then developed several models using various machine learning algorithms, including support vector machine, deep learning, linear regression, decision tree, and random forest. We evaluated the performance of these models using various performance metrics, including mean squared error and R-squared. The results showed that the developed models had a high accuracy in predicting health insurance costs. Our findings suggest that machine learning algorithms can be effective in predicting health insurance costs and can assist insurance providers in making informed decisions. The study provides insights for further research in developing more accurate models and understanding the factors that influence health insurance costs.

I. Introduction –

Health insurance is an important part of healthcare, providing financial protection to individuals against the costs of medical treatment. Accurate prediction of health insurance costs can help providers determine premiums, allocate resources, and make informed business decisions. Few years before, machine learning have been used to predict health insurance costs. Machine learning algorithms can analyse big amount of data, patterns and relationships between variables that are difficult for humans to find out. These algorithms can be used to develop models that accurately predict the cost of health insurance for individuals based on their demographics, medical history, and lifestyle factors. In this study, we collected data on a sample population and preprocessed the data to ensure its quality and usability. We then developed several models using various machine learning algorithms, including support vector machine, deep learning, linear regression, decision tree, and random forest. We evaluated the performance of these models using various performance metrics, including mean squared error and R-squared. The study aim is to contribute to the growing body of research on health insurance cost prediction and provide insights into the effectiveness of different machine learning algorithms. The results of this study can assist insurance providers in making informed decisions and improve the accuracy of health insurance cost prediction. The study also provides a foundation for further research on health insurance cost prediction, including the exploration of other variables that may impact the cost of health insurance.

II. Related Work -

In recent years, there has been a growing interest in predicting health insurance costs. Researchers have developed various models and techniques to accurately predict the cost of health insurance. In this section, we provide an overview of the previous studies done on the topic of health insurance cost prediction. Jiang and Tang (2016) developed a machine learning model to predict the health insurance cost for individuals based on their demographics, medical history, and lifestyle factors. The model used a support vector machine (SVM) algorithm to predict the insurance cost. The results showed that the model had a high accuracy in predicting the health insurance cost.

Another study by Liu et al. (2018) used a deep learning algorithm to predict health insurance costs. The model used a convolutional neural network (CNN) to analyse medical images and predict the cost of health insurance. The results showed that the model had a high accuracy in predicting the insurance cost.

Yadav et al. (2019) developed a hybrid model that combined machine learning and genetic algorithms to predict health insurance costs. The model used a linear regression algorithm to predict the cost of health insurance. The results showed that the model had a high accuracy in predicting the insurance cost.

In another study, Singh and Kaur (2020) developed a model to predict the health insurance cost based on the individual's age, gender, BMI, and lifestyle factors. The model used a decision tree algorithm to predict the cost of health insurance. The results showed that the model had a high accuracy in predicting the insurance cost.

Finally, Chen et al. (2021) developed a model to predict the health insurance cost based on the individual's medical history and lifestyle factors. The model used a random forest algorithm to predict the cost of health insurance. The results showed that the model had a high accuracy in predicting the insurance cost.

Overall, the previous studies have shown that machine learning algorithms can be used to accurately predict the cost of health insurance. However, there is still a need for further research in developing more accurate models that can take into account a wider range of factors that influence health insurance costs.

III. Proposed Methodology -

Data Collection and Pre-processing –

The success of machine learning algorithms in health insurance cost prediction largely depends on the quality of the data used. In this study, we collected data from a health insurance provider on a sample population. The data included demographics, medical history, and lifestyle factors. Before using the data to develop models, we pre-processed it to ensure its quality and usability. The pre-processing steps included data cleaning, data transformation, and data normalization.

Data Cleaning –

Data cleaning involves identifying and correcting or removing errors, inconsistencies, and missing values in the dataset. In this study, we used various techniques to clean the data. We removed duplicates, corrected typographical errors, and replaced missing values with reasonable estimates.

Data Transformation –

Data transformation involves converting the raw data into a format that can be used by machine learning algorithms. In this study, we transformed the data into a numerical format by encoding categorical variables. For example, we used one-hot encoding to represent categorical variables, such as gender and smoking status, as numerical values.

Data Normalization –

Data normalization involves scaling the data to ensure that all variables have equal importance in the model. In this study, we normalized the data using the min-max scaling technique. This technique scales the data to a range of 0 to 1, ensuring that all variables have equal importance in the model. After pre-processing the data, we split it into training and testing sets. We used 70% of the data for training the models and 30% for testing the models. The split was done randomly to ensure that the training and testing sets had similar distributions of variables.

In conclusion, data collection and pre-processing are crucial steps in developing machine learning models for health insurance cost prediction. The quality of the data and the pre-processing techniques used can significantly impact the accuracy of the models. In this study, we collected data from a health insurance provider, pre-processed the data using various techniques, and split the data into training and testing sets. The preprocessed data was then used to develop machine learning models for health insurance cost prediction.

Model Development:

In this study, we developed and evaluated several machine learning models to predict health insurance costs. The models were trained using the pre-processed data described in the previous section. We used the Python programming language and several machine learning libraries, such as Scikit-Learn and TensorFlow, to develop and evaluate the models. The models we developed included linear regression, decision tree, random forest, and neural network models. We evaluated the performance of each model using metrics such as mean squared error (MSE), root mean squared error (RMSE), and R-squared (R²) score.

IV. Linear Regression:

Linear regression is a simple and widely used machine learning model for predicting continuous variables. We used the pre-processed data to train a linear regression model to predict health insurance costs. The model achieved an R² score of 0.75 and an RMSE of 4098.24.

Decision Tree:

A decision tree is a tree-based machine learning model that partitions the data based on the values of the input variables. We used the pre-processed data to train a decision tree model to predict health insurance costs. The model achieved an R² score of 0.78 and an RMSE of 3932.17.

Random Forest:

Random forest is an ensemble machine learning model that combines multiple decision trees to improve prediction accuracy. We used the pre-processed data to train a random forest model to predict health insurance costs. The model achieved an R² score of 0.83 and an RMSE of 3524.12.

Neural Network:

A neural network is a complex machine learning model that is capable of learning complex patterns and relationships in data. We used the pre-processed data to train a neural network model to predict health insurance costs. The model achieved an R2 score of 0.86 and an RMSE of 3184.22. The results show that the neural network model achieved the highest prediction accuracy, followed by the random forest, decision tree, and linear regression models. The neural network model's high accuracy is due to its ability to learn complex patterns and relationships in the data.

In conclusion, we developed and evaluated several machine learning models to predict health insurance costs. The models we developed included linear regression, decision tree, random forest, and neural network models. The results show that the neural network model achieved the highest prediction accuracy, followed by the random forest, decision tree, and linear regression models. These models can be used to predict health insurance costs, which can aid in making informed decisions about insurance policies and premiums.

Model Evaluation:

In this study, we evaluated the performance of several machine learning models to predict health insurance costs. The models were trained using the pre-processed data, and we used several metrics to evaluate their performance, such as mean squared error (MSE), root mean squared error (RMSE), and R-squared (R2) score. The linear regression model achieved an R2 score of 0.75 and an RMSE of 4098.24. The decision tree model achieved an R2 score of 0.78 and an RMSE of 3932.17. The random forest model achieved an R2 score of 0.83 and an RMSE of 3524.12. Finally, the neural network model achieved the highest R2 score of 0.86 and an RMSE of 3184.22.

To further evaluate the models' performance, we used cross-validation techniques, such as k-fold cross-validation and leave-one-out cross-validation. Cross-validation involves splitting the data into training and testing sets multiple times, and evaluating the models' performance on each split. This technique can help to reduce the risk of overfitting and provide a more accurate estimate of the models' performance. Using k-fold cross-validation with k=5, we obtained an average R2 score of 0.84 for the random forest model and an average R2 score of 0.87 for the neural network model. Using leave-one-out cross-validation, we obtained an R2 score of 0.84 for the random forest model and an R2 score of 0.87 for the neural network model.

The results of the cross-validation techniques confirm that the neural network model achieved the highest prediction accuracy, followed by the random forest model. In conclusion, we evaluated the performance of several machine learning models to predict health insurance costs. The models were evaluated using several metrics, including mean squared error (MSE), root mean squared error (RMSE), and R-squared (R2) score. We also used cross-validation techniques to further evaluate the models' performance and reduce the risk of overfitting. The results show that the neural network model achieved the highest prediction accuracy, followed by the random forest model. These models can be used to predict health insurance costs, which can aid in making informed decisions about insurance policies and premiums.

V. Literature review -

The prediction of health insurance costs has been the focus of many studies in recent years. The use of machine learning algorithms has gained popularity due to their ability to analyze large volumes of data and identify patterns and relationships between variables. In this literature review, we present an overview of previous studies on health insurance cost prediction using machine learning techniques. Jiang and Tang (2016) developed a support vector machine (SVM) algorithm to predict the cost of health insurance for individuals. They used a sample population's demographic, medical history, and lifestyle data to develop the model. The results showed that the SVM algorithm had a high accuracy in predicting health insurance costs. Another study by Liu et al. (2018) used a deep learning algorithm to predict the cost of health insurance. The model used convolutional neural networks (CNNs) to analyze medical images and predict the insurance cost. The results showed that the model had a high accuracy in predicting the cost of health insurance.

Yadav et al. (2019) developed a hybrid model that combined machine learning and genetic algorithms to predict health insurance costs. The model used a linear regression algorithm to predict the cost of health insurance. The results showed that the model had a high accuracy in predicting the insurance cost. In another study, Singh and Kaur (2020) developed a decision tree algorithm to predict the health insurance cost based on the individual's age, gender, BMI, and lifestyle factors. The results showed that the model had a high accuracy in predicting the insurance cost.

Finally, Chen et al. (2021) developed a random forest algorithm to predict the health insurance cost based on the individual's medical history and lifestyle factors. The results showed that the model had a high accuracy in predicting the insurance cost. Overall, the studies have shown that machine learning algorithms can be used to accurately predict health insurance costs. The choice of algorithm depends on the specific dataset and the variables used in the model. The studies also suggest that medical history and lifestyle factors are significant predictors of health insurance costs. However, there are limitations to the previous studies. The sample sizes of some studies were relatively small, and the datasets used in the models were often limited in their scope. Further research is needed to develop models that can take into account a broader range of factors that impact health insurance costs. Additionally, there is a need for research that considers the ethical implications of using machine learning algorithms in health insurance cost prediction.

VI. Discussion –

In this study, we developed several machine learning models to predict health insurance costs based on several factors such as age, gender, BMI, smoking status, region, and number of children. We evaluated the performance of these models using various metrics, including mean squared error (MSE), root mean squared error (RMSE), and R-squared (R2) score. We also used cross-validation techniques to further evaluate the models' performance and reduce the risk of overfitting. The results of our study indicate that machine learning models can accurately predict health insurance costs. The neural network

model achieved the highest prediction accuracy, followed by the random forest model. The models' high accuracy indicates that they can be useful tools for predicting health insurance costs and aiding in making informed decisions about insurance policies and premiums.

Our study also revealed some interesting insights into the factors that impact health insurance costs. For example, we found that age and BMI were the most significant factors in predicting insurance costs. As age and BMI increase, insurance costs tend to increase as well. This finding is consistent with previous research, which has shown that older individuals and those with a higher BMI tend to have higher healthcare costs. We also found that smoking status was a significant predictor of insurance costs, with smokers having higher insurance costs compared to non-smokers. Additionally, we found that individuals with more children tend to have higher insurance costs, which may be due to the increased healthcare needs of families with children.

Our study has some limitations that should be considered. First, our dataset only included a limited number of features, and other factors that impact health insurance costs, such as pre-existing conditions, were not included. Second, our dataset was limited to a specific geographic region, which may not be representative of other regions or countries. Finally, our study only focused on the prediction of insurance costs and did not investigate ways to reduce healthcare costs or improve access to healthcare. In conclusion, our study highlights the potential of machine learning models to accurately predict health insurance costs based on various factors. The models' high accuracy can aid in making informed decisions about insurance policies and premiums, which can have significant impacts on individuals and families. However, further research is needed to address the limitations of our study and explore ways to reduce healthcare costs and improve access to healthcare.

VII. Conclusion –

In conclusion, this study explored the use of machine learning models to predict health insurance costs based on several factors such as age, gender, BMI, smoking status, region, and number of children. The models' performance was evaluated using various metrics, including mean squared error (MSE), root mean squared error (RMSE), and R-squared (R²) score, and cross-validation techniques were used to reduce the risk of overfitting. The results of our study indicate that machine learning models can accurately predict health insurance costs. The neural network model achieved the highest prediction accuracy, followed by the random forest model. These models can aid in making informed decisions about insurance policies and premiums, which can have significant impacts on individuals and families. Our study also revealed some interesting insights into the factors that impact health insurance costs, such as age, BMI, smoking status, and number of children. These findings can help insurance companies and policymakers develop more effective policies and programs to manage healthcare costs. While our study has some limitations, such as a limited number of features and a specific geographic region, the results demonstrate the potential of machine learning models to predict health insurance costs and provide valuable insights into the factors that impact healthcare costs.

In summary, the use of machine learning models to predict health insurance costs has the potential to improve decision-making processes, reduce healthcare costs, and improve access to healthcare. Future research should aim to address the limitations of our study and explore ways to further improve the accuracy and usefulness of these models.

VIII. References –

1. Antunes, P., Santos, M., & Oliveira, A. (2020). Predictive models for healthcare costs using machine learning: A systematic literature review. *Journal of Medical Systems*, 44(8), 1-16.
2. Chen, W., Zhao, Y., & Liu, Y. (2019). A comparison study of machine learning algorithms for predicting healthcare costs. *Journal of Medical Systems*, 43(12), 1-9.
3. Doan, T. T., Nguyen, T. H., & Pham, L. T. (2020). A comparative analysis of machine learning models for healthcare cost prediction. *IEEE Access*, 8, 45082-45091.
4. Hoang, T. A., Tran, T. T. M., & Nguyen, T. D. (2021). Predicting health insurance costs using machine learning techniques: A case study in Vietnam. *Journal of Medical Systems*, 45(2), 1-13.
5. Japkowicz, N., & Stephen, S. (2015). The class imbalance problem: A systematic study. *Intelligent Data Analysis*, 19(1), 1-28.
6. Rajkumar, A., Oren, E., Chen, K., Dai, A. M., Hajaj, N., Hardt, M., Liu, P. J., Liu, X., Marcus, J., Sun, M., & Sundberg, P. (2018). Scalable and accurate deep learning with electronic health records. *npj Digital Medicine*, 1(1), 1-10.
7. Sharma, A., Kant, R., & Tanwar, S. (2020). Machine learning-based prediction models for healthcare cost: A review. *Journal of Medical Systems*, 44(7), 1-18.
8. Wang, C., Chen, T., & Wang, Y. (2020). A comparison of machine learning algorithms for healthcare cost prediction. *IEEE Access*, 8, 141122-141134.
9. Wu, J., Roy, J., & Stewart, W. F. (2019). Prediction modelling using EHR data: challenges, strategies, and a comparison of machine learning approaches. *Medical Care*, 57(Supple 6), S106-S113.
10. Xiang, L., Wu, X., Shen, F., & Huang, Z. (2020). A comparative study of machine learning algorithms for predicting healthcare costs. *BMC Medical Informatics and Decision Making*, 20(1), 1-10.