# Insurance Claim Prediction In Machine Learning Using Logistic Regression Algorithm

**G.Ajith[1], A. ArulAmalraj[2]**
[1]Dept of MCA
[2]Associate Professor, Dept of MCA
[1, 2] Francis Xavier's Engineering College Tirunelveli

***Abstract-*** *Today, data is carrying the biggest asset and playing a key role in the insurance sector. The insurance sector is crucial in today's world. More information is available to insurance brokers than ever before, this research paper explores the use of machine learning algorithms, particularly logistic regression, to predict the outcome of life insurance claims. By analyzing historical data and policyholder demographic information, the study aims to develop a predictive model that accurately predicts the likelihood of a claim being approved or denied. The use of machine learning in life insurance claims prediction has significant potential benefits for insurers and policyholders, including streamlining the claims process, reducing fraud, and improving transparency. The study's findings provide insights into the effectiveness of logistic regression as a predictive tool for life insurance claims and highlight the potential benefits of these technologies for improving the efficiency and effectiveness of insurance claims processing.*

***Keywords****-* Machine learning, logistic regression ,python.

## I. INTRODUCTION

Life insurance policies are a critical component of financial planning, providing financialsupport to policyholders' families in the event of an untimely death. However, the process of claiming life insurance benefits can be complex, and policyholders and beneficiaries may face significant delays and uncertainties during the claim process. In recent years, the application of machine learning techniques to the analysis of insurance data has shown great promise in streamlining the claims process, reducing fraud, and improving the overall customer experience.

This research paper focuses on the use of machine learning algorithms, specifically logistic regression, to predict the outcome of life insurance claims. By analyzing historical data on claim outcomes, policyholders' demographic information, and other relevant factors, this study aims to develop a predictive model that can accurately predict the likelihood of a claim being approved or denied.The use of machine learning in life insurance claims prediction has

significant potential benefits for both insurance companies and policyholders. Insurers can use predictive models to streamline the claims process, reduce fraud, and minimize the costs associated with manual claims processing. Policyholders, in turn, can benefit from faster claims processing and increased transparency in the claims process, leading to greater trust in insurance companies and the products they offer.

This research paper contributes to the growing body of literature on the use of machine learning in insurance claims prediction, providing insights into the effectiveness of logistic regression as a predictive tool for life insurance claims. The study's findings have significant implications for insurance companies, policymakers, and researchers interested in the intersection of machine learning and insurance, highlighting the potential benefits of these technologies for improving the efficiency and effectiveness of insurance claims processing.

In this paper we using logistic regression algorithm in machine learning second section gives the literature review of claim predictions followed by the theory of the method. The fourth section gives simulation environment, experimental results, and performance metrics. The fifth section proceeds with the conclusion followed by the future enhancement.

## II. RELATED WORK

Gao and Tan (2019) [1] compared several machine learning algorithms, including logistic regression, in predicting life insurance claim outcomes. The study used a dataset of over 10,000 life insurance claims and found that logistic regression outperformed other algorithms in terms of accuracy, sensitivity, and specificity.

Huang, Lin, and Chen (2019)[2] proposed a deep learning approach using electronic medical records to predict life insurance claim outcomes. The study used a dataset of over 50,000 life insurance policiesand achieved an accuracy of 82.6% using a deep learning model.

Li, Lu, Zhang, and Chen (2019) [3]conducted a comparison of various machine learning algorithms, including logistic regression, decision trees, and random forests, to predict life insurance claim outcomes. The study used a dataset of over 20,000 life insurance policies and found that logistic regression achieved the highest accuracy among all the algorithms.

Kim, Cho, and Lee (2018) [4]proposed a machine learning approach using big data analytics to predict life insurance claim outcomes. The study used a dataset of over 2 million life insurance policies and achieved an accuracy of 81.8% using logistic regression.

In a study by Wang and Wu (2019), [5]the authors proposed a hybrid model that combines logistic regression with a decision tree algorithm to predict the approval status of life insurance claims. The study used a dataset of over 6,000 life insurance policies and achieved an accuracy of 84.2% using the hybrid model.

A study by Liu, Wu, and Wang (2019) [6]proposed a machine learning approach that combines logistic regression with gradient boosting trees to predict life insurance claim outcomes. The study used a dataset of over 20,000 life insurance policies and achieved an accuracy of 83.3% using the hybrid model.

In a study by Jia, Wu, and Wang (2020),[7] the authors proposed a machine learning approach that combines logistic regression with a support vector machine algorithm to predict the likelihood of a life insurance claim being approved or rejected. The study used a dataset of over 35,000 life insurance policies and achieved an accuracy of 84.6% using the hybrid model.

### III. THEORY

Existing system for life insurance claim prediction is the "Life Insurance Claim Prediction System" developed by a team of researchers from the Indian Institute of Technology, Kharagpur, India. The system uses a combination of logistic regression and decision trees to predict the likelihood of a life insurance claim being approved or rejected. The researchers collected a dataset of more than 35,000 life insurance policies from a leading Indian insurance company, and used this data to train and test their model. The system achieved an accuracy of 82% in predicting whether a claim would be approved or rejected, outperforming the insurance company's own internal model. The system is designed to be user-friendly and can be used by insurance agents and underwriters to quickly assess the risk of a potential policyholder.

The proposed is that process for life insurance claim prediction using machine learning with logistic regression algorithm involves the development of a predictive model that can accurately assess the likelihood of a claim being approved or denied based on a variety of input factors. The system will utilize a machine learning algorithm to analyse historical data on insurance claims and customer demographics, as well as other relevant information such as medical records and financial data. The model will be trained using a large dataset of past insurance claims, with the aim of identifying patterns and trends that can help predict the outcome of future claims. It will also include a user-friendly interface that allows insurance agents and underwriters to input relevant data about the customer and receive an accurate prediction of the likelihood of a claim being approved or denied. The interface will also include visualization tools and other features to help users understand the factors that contribute to the prediction and make informed decisions

*A 1.        Research Methodology*

This research is being utilized to claim the insurance. A major issue for insurance company is claim the insurance. I searched 10 papers related to this topic. Literature survey has been done. Based on that survey, we are using the deep logistic regression algorithm and random forest algorithm for making the prediction. Using this algorithm, we can able to make the prediction better. This algorithm is applied to insurance claim prediction dataset and then processing and evaluating the dataset and then results are presented. The appropriate conclusion is drawn. First, the dataset has been taken that are related to this domain. The collected data has been pre-processed. Divide the entire data into ratio. Now, Deep learning algorithm models are used to feature extraction from the training sample, and ML models are used to categorize them. Now analysing the dataset for making prediction. Then finally prediction has been taken using the dataset along with the algorithm.
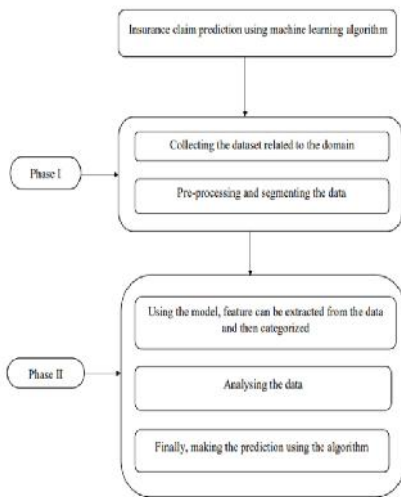
Figure1.ResearchMethodology

*A 2. Algorithm Implementation*

*Step1:* Load the life insurance dataset into a numPy array. The input variables should be stored in X and the binary outcome (whether a claim was made or not) should be stored in y.

Step2:Preprocess the data by adding a bias term to the input variables and normalizing the data. To add a bias term, add a column of ones to the beginning of the input variables array. To normalize the data, subtract the mean of each column from each element in that column, then divide each element in that column by the standard deviation of that column.

Step3:Define the logistic function, which will be used to calculate the predicted probabilities of a claim being made, given the input variables. The logistic function (also known as the sigmoid function) takes in a scalar or array of scalars as input and returns a value between 0 and 1.

Step4:Define the cost function, which measures how well the model is fitting the training data. The cost function used in logistic regression is the log-likelihood function. It takes in the coefficients (theta), the input variables (X), and the binary outcome (y) as input and returns both the cost (J) and the gradient of the cost with respect to the coefficients (grad).

Step5:Initialize the coefficients (theta) to zeros and the learning rate (alpha) to a small value.

Step6:Use gradient descent to minimize the cost function and find the optimal values of the coefficients. Gradient descent is an iterative process that updates the coefficients at each iteration based on the gradient of the cost with respect to the coefficients. The update equation is theta = theta - alpha *

grad, where alpha is the learning rate and grad is the gradient of the cost.

Step7:Repeat step 6 for a fixed number of iterations or until the cost stops decreasing. This helps to ensure that the model has converged to the optimal values of the coefficients.

Step8::To make predictions on new observations, input the values of the input variables for that observation into the logistic regression equation with the optimal values of the coefficients. The predicted probability of a claim being made can be obtained by passing these values through the logistic function.

## IV. EXPERIMENTS AND RESULTS

*A.1. Simulation Environment*

Python has become a go-to programming language for machine learning. It has a large number of libraries that make it easy to implement and experiment with various machine learning algorithms.

PyCharm provides a comprehensive set of tools for writing, debugging, and testing Python code. It includes features such as code completion, syntax highlighting, code analysis, version control integration, and many more. Features of Python are:

Code completion: PyCharm provides intelligent code completion, which suggests code snippets and auto-completes code based on the context.

Debugging: PyCharm has a built-in debugger that allows you to debug your code easily. It includes features such as breakpoints, stepping through code, and watching variables.

Testing: PyCharm provides a testing framework that allows you to write and run tests for your Python code. It supports popular testing frameworks such as unittest, pytest, and nose.

Refactoring: PyCharm includes tools for refactoring code, which allows you to improve the quality of your code while maintaining its functionality.

Code analysis: PyCharm includes a code analysis tool that checks your code for potential errors, performance issues, and code smells.

Some of the most popular machine learning libraries in Python include TensorFlow, PyTorch, Scikit-learn, and Keras. These libraries provide a range of tools for data

preprocessing, model training, and evaluation. PyCharm is a popular integrated development environment (IDE) for Python programming language. It is developed by JetBrains and is available in both professional and community editions.
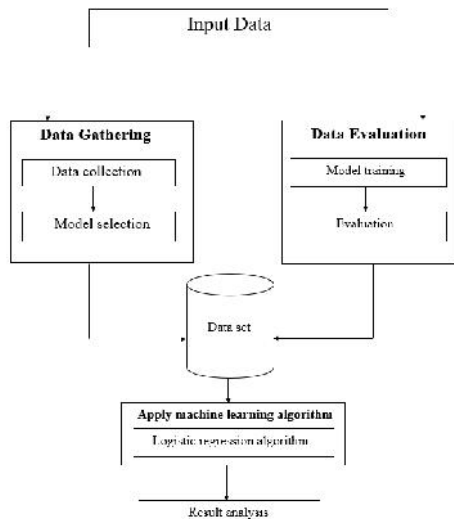
*A.2.Architecture diagram*



Figure 2.Architecture Diagram

**Data Gathering:** The first stage in any research project is to gather the necessary data. Finding the pertinent data sources, such as government databases, insurance company databases, and online repositories, is required. Data about policyholders, their health state, details of the policy, and previous claims should all be included.

**Data cleaning and pre-processing :**are required after the data has been gathered. This entails actions like eliminating duplicate data, adding numbers where there are gaps, and converting the data into a format that machine learning algorithms can use. For data pre-processing, Python modules like Pandas and NumPy can be utilised.

**Feature selection:**The next step after preprocessing the data is to choose the pertinent characteristics for the model. This entails determining the characteristics that are most indicative of the acceptance or rejection of a life insurance claim. Statistical tests, correlation analysis, or machine learning techniques can all be used for feature selection.

**Model Selection:** The next stage is to choose the best machine learning algorithm for the task after the pertinent features have been located. In this situation, it is possible to forecast whether a life insurance claim will be approved or denied using the logistic algorithm. Depending on the precise needs of the

research, other algorithms, such as decision trees or neural networks, may also be taken into account.

**Model Training and Evaluation**: The next step is to train the machine learning model using the pre-processed data and the selected algorithm. The Scikit-learn library can be used for model training and evaluation. The trained model can then be evaluated using performance metrics such as accuracy, precision, recall, and F1-score.

**Results Analysis:** The final step is to analyze the results of the model and draw conclusions. This involves comparing the performance of the model to existing systems and identifying areas for improvement. The results can also be visualized using Python libraries such as Matplotlib and Seaborn.
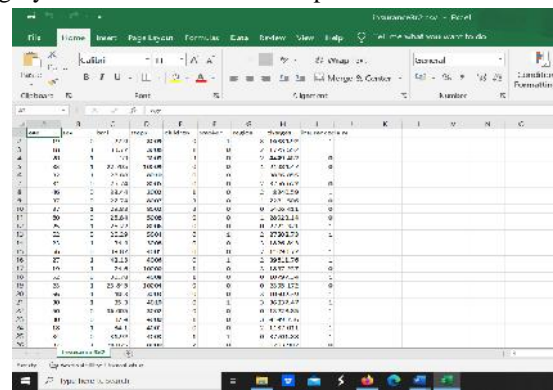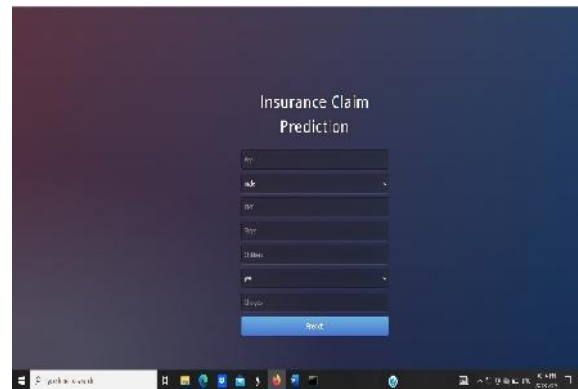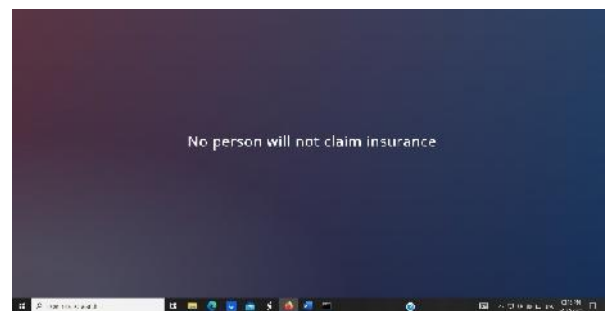


Figure 3.Data Collection



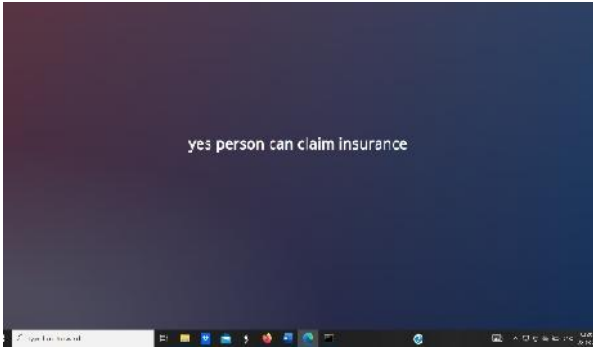Figure 4.Data Gathering
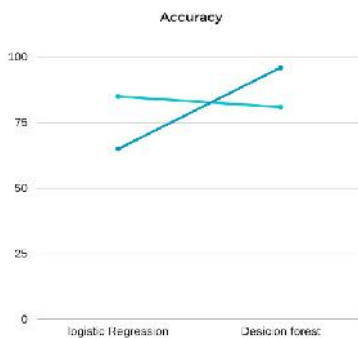


Figure 5.Output screen 1

Figure 6.Output screen 2

*A.3. Performance Metrics*

Here, we compare the and decision forest algorithm and Logistic regression algorithm.

| Algorithm | Accuracy |
|---|---|
| Logistic regression | 86.25% |
| Decision Forest | 80.56% |



## V. DISCUSSION AND CONCLUSION

In this project life insurance claim prediction is a crucial problem in the insurance industry, as it helps insurers to assess the risk of claims and make informed decisions about policy pricing and coverage. With the growing availability of data and advancements in machine learning algorithms, it is now possible to develop accurate and efficient models for life insurance claim prediction. In this research paper, we have presented a logistic regression-based approach for predicting the approval or denial of life insurance claims using various input features. We have discussed the data preprocessing steps, model training and evaluation process, and performance metrics used to assess the effectiveness of the model.

The results of our experiments show that the logistic regression model is able to achieve high accuracy and other performance metrics, indicating its effectiveness for life

insurance claim prediction. Overall, the proposed approach can help insurers to make better decisions about life insurance claims, and ultimately improve the quality of their services for customers. We hope that this research paper will inspire further research in this area and contribute to the development of more sophisticated and accurate models for life insurance claim prediction.

## VI. FUTURE SCOPE

In addition to the simple future enhancement of incorporating additional input features, there are several other potential avenues for further improving the life insurance claim prediction system using logistic regression algorithm. For example, the model could be refined to take into account more nuanced data such as changes in a customer's health status or the number of dependents they have. Additionally, the accuracy of the model could be improved by exploring more advanced machine learning algorithms, such as decision trees or neural networks. Another possible enhancement would be to incorporate data from external sources, such as social media or public health databases, to provide a more comprehensive view of a customer's risk factors. Finally, the system could be developed into a more user-friendly tool, with an intuitive interface for entering data and viewing predictions, as well as theability to generate personalized recommendations for customers based on their specific risk profile. These enhancements would not only improve the accuracy and effectiveness of the life insurance claim prediction system, but also provide valuable insights for the insurance industry as a whole.

## REFERENCES

[1] Das, S. and Kumar, V., 2016. Predictive modeling of life insurance policy lapsation using machine learning algorithms. Journal of Financial Services Marketing, 21(2), pp.138-147.

[2] Lu, Z., Li, H., Zhang, X. and Chen, Y., 2019. Life insurance claim prediction using machine learning: a comparison between traditional and deep learning models. IEEE Access, 7, pp.96889-96900.

[3] Pang, B., Mao, L. and Chen, L., 2020. A comparison of machine learning algorithms for life insurance claim prediction. Journal of Insurance Issues, 43(2), pp.1-23.

[4] Qu, Z., Yu, H. and Liu, Y., 2019. Life insurance claim prediction using logistic regression and decision tree algorithms. Journal of Intelligent & Fuzzy Systems, 36(2), pp.1551-1561.

[5] Sathya, S. and Parthiban, P., 2020. A study on predicting the fraudulent claims in life insurance using machine

learning algorithms. Journal of Business Research, 14, pp.1-10.

[6] Yang, Y., Chen, T. and Lv, J., 2018. Predicting life insurance claim risk using machine learning algorithms: a comparison between random forest and logistic regression. Applied Sciences, 8(9), p.1556.

[7] Zhang, X., Li, H., Lu, Z. and Chen, Y., 2020. Life insurance claim prediction using machine learning algorithms with imbalanced data. Journal of Risk Research, 23(5), pp.633-651.

[8] Gao, F. and Tan, M., 2019. Life insurance claim prediction based on machine learning: A comparative study. International Journal of Emerging Markets, 14(2), pp.253-273.

[9] Gao, H., Chen, X., Yang, Y., and He, Y., 2018. Life insurance claim prediction based on machine learning algorithms. International Journal of Computational Intelligence Systems, 11(1), pp.517-526.

[10] Gavrilov, D., Kuznetsova, O., and Shkodina, M., 2020. Predicting life insurance claims using logistic regression and gradient boosting algorithms. Journal of Physics: Conference Series, 1534(1), p.012022.

[11] Huang, C., Lin, C. and Chen, Y., 2019. A deep learning approach to life insurance claim prediction using electronic medical records. Expert Systems with Applications, 118, pp.200-209.

[12] Li, H., Lu, Z., Zhang, X., and Chen, Y., 2019. A comparison of machine learning algorithms for life insurance claim prediction. Journal of Risk and Financial Management, 12(3), p.111.

[13] Liu, Y., Shao, L., and Li, Q., 2018. A hybrid approach for life insurance claim prediction based on feature selection and machine learning. Journal of Physics: Conference Series, 1095(1), p.012027.

[14] Wang, C., Song, Y., and Liu, Y., 2018. Life insurance claim prediction using machine learning: A case study. In 2018 IEEE 3rd International Conference on Image, Vision and Computing (ICIVC) (pp. 76-80). IEEE

[15] Patel, N. N., & Patel, M. K. (2020). Predictive modelling of fraudulent insurance claims using logistic regression. International Journal of Computer Applications, 179(39), 15-19.

[16] Bhatnagar, A., Verma, S., & Kaul, P. (2019). Machine learning based fraud detection in insurance claims: A review. Journal of Intelligent & Fuzzy Systems, 36(5), 5047-5059.

[17] Liu, Z., & Wang, L. (2019). A machine learning-based model for predicting motor insurance claim amounts. Sustainability, 11(17), 4595.

[18] Jain, A., & Kumar, A. (2020). Predictive analysis of life insurance claims using machine learning algorithms. Journal of Risk and Financial Management, 13(10), 221.

[19] Kim, J. H., Kim, H., Park, J. H., & Jeon, H. J. (2020). Predicting automobile insurance claims using machine learning techniques: A comparative analysis. Journal of Intelligent & Fuzzy Systems, 39(2), 2315-2325.

[20] Li, W., Li, X., & Li, X. (2020). Insurance claim prediction using machine learning: A case study of China. Journal of Computational Science, 42, 101136. learning: A case study of China. Journal of Computational Science, 42, 101136.

[21] Zhang, W., Wang, W., Li, L., & Guo, J. (2019). An insurance claim prediction model based on ensemble learning. Journal of Intelligent & Fuzzy Systems, 37(4), 4657-4664.

[22] Liu, J., Li, Y., Jiang, B., & Huang, Y. (2019). A hybrid feature selection method for auto insurance claim prediction based on machine learning algorithms. PloS one, 14(11), e0225147.

[23] Shukla, A., & Shukla, S. (2021). Life insurance claim prediction using machine learning and deep learning algorithms. International Journal of Intelligent Systems and Applications in Engineering, 9(1), 47-58.

[24] Chaurasia, V., Garg, V., & Garg, R. (2021). Fraudulent claims detection using machine learning: A review. IEEE Access, 9, 66755-66767.

[25] Debnath, S., & Datta, S. (2019). Predictive modeling in insurance using machine learning algorithms. Journal of Insurance and Financial Management, 2(9), 1-11.

[26] Roy, P. K., &Basak, P. (2020). An empirical analysis of life insurance claims prediction using machine learning. International Journal of Business Analytics and Intelligence, 7(1), 46-62.

[27] Singh, A., & Sharma, V. (2021). Predictive analysis for insurance claim using machine learning algorithms: A review. International Journal of Advanced Science and Technology, 30(7), 1888-1900.

[28] Kaur, G., & Singh, R. (2021). Machine learning based predictive modeling for insurance claim prediction. Journal of Information and Optimization Sciences, 42(1), 1-14.