

PREDICTION OF STROKE USING LOGISTIC REGRESSION WITH A MACHINE LEARNING APPROACH

Ishiq Rana Aphrodita^{1*}, Ika Nur Fajri¹, Agung Nugroho¹

¹Information Systems, Amikom University Yogyakarta

email : *ishiqarana@students.amikom.ac.id

Abstract: Stroke is one of the leading causes of death and disability in various parts of the world, including in Indonesia. Along with the development of digital technology, the use of Machine Learning in the health sector is growing, one of which is in an effort to predict the occurrence of stroke. This study aims to implement the Logistic Regression algorithm in predicting the likelihood of a person having a stroke based on data from the Brain Stroke dataset. The research process includes data preprocessing (categorical encoding using One-Hot Encoding, normalization using StandardScaler, and class balancing using SMOTE), dividing the data into 80% training data and 20% test data, as well as model training. The model was then evaluated using several measures such as accuracy, precision, recall, F1-score, and ROC-AUC, as well as a confusion matrix. The model achieved an accuracy of 74.8%, precision of 14%, recall of 80%, F1-score of 24%, and ROC-AUC of 0.84. Then, the model is integrated into applications that use Streamlit, so it can be used interactively to predict stroke risk in new data. The results of this study show that the combination of Machine Learning and web-based applications has the potential to support efforts to detect early stroke risk.

Keywords: logistic regression; machine learning; prediction; streamlit; stroke.

Abstrak: Stroke adalah salah satu penyebab utama kematian dan kecacatan di berbagai belahan dunia, termasuk di Indonesia. Seiring perkembangan teknologi digital, penggunaan Machine Learning dalam bidang kesehatan semakin berkembang, salah satunya dalam upaya memprediksi terjadinya penyakit stroke. Penelitian ini bertujuan untuk mengimplementasikan algoritma Logistic Regression dalam memprediksi kemungkinan seseorang mengalami stroke berdasarkan data dari dataset Brain Stroke. Proses penelitian meliputi preprocessing data (Encoding Kategorikal menggunakan One-Hot Encoding, normalisasi menggunakan StandardScaler, dan penyeimbangan kelas menggunakan SMOTE), membagi data menjadi 80% data latih dan 20% data uji, serta pelatihan model. Model kemudian dievaluasi menggunakan beberapa ukuran seperti akurasi, precision, recall, F1-score, dan ROC-AUC, serta confusion matrix. Model kemudian dievaluasi menggunakan beberapa ukuran seperti akurasi, precision, recall, F1-score, dan ROC-AUC, serta confusion matrix. Model ini mencapai akurasi 74,8%, presisi 14%, recall 80%, F1-Score 24%, dan ROC-AUC 0,84. Kemudian, model tersebut diintegrasikan ke dalam aplikasi yang menggunakan Streamlit, sehingga dapat digunakan secara interaktif untuk memprediksi risiko stroke pada data baru. Hasil penelitian ini menunjukkan bahwa kombinasi Machine Learning dan aplikasi berbasis web berpotensi mendukung upaya deteksi dini risiko stroke.

Kata kunci: logistic regression; machine learning; prediksi; streamlit; stroke.

INTRODUCTION

Stroke is one of the leading causes of death and long-term disability worldwide. In Indonesia, according to the *Burden of Stroke* report, stroke accounts for around 15.4% of total deaths and remains the non-communicable disease with the highest incidence rate, reaching 8,3 per 1,000 population. [1] Globally, more than 12 million people suffer from stroke annually, with around 6.5 million deaths each year. Common risk factors include hypertension, smoking, high cholesterol, and advanced age. Early detection of stroke risk can significantly reduce mortality and disability rates; however, in developing countries such as Indonesia, challenges persist due to limited healthcare access, low awareness of stroke symptoms, and uneven distribution of medical specialists.

The advancement of digital technology, particularly *Machine Learning* (ML), provides new opportunities for predicting disease risk based on demographic, clinical, and lifestyle data. Among various ML algorithms, *Logistic Regression* is widely applied for binary classification problems, such as predicting whether a person is at risk of stroke. Prior studies have demonstrated its high performance after incorporating preprocessing steps such as handling missing values, normalization, and class balancing using SMOTE. For instance, the study “*Optimizing Accuracy of Stroke Prediction Using Logistic Regression*” achieved an accuracy of approximately 86% after applying SMOTE and feature selection [2] while “*Prediction of Patient’s Stroke Vulnerability Status Using Logistic Regression*” obtained 81% accuracy and an AUC-ROC value of around 90% [3] Another study also reported that the integration of preprocessing tech-

niques with *Logistic Regression* improved the AUC value up to 0.90 [4]

Based on these findings, this research focuses on developing a stroke prediction model using the Brain Stroke dataset and the *Logistic Regression* algorithm. The model construction involves several stages, including preprocessing (handling missing values, category encoding, normalization, and class balancing), training, and evaluation using metrics such as accuracy, precision, recall, F1-score, and ROC-AUC. Furthermore, the trained model is implemented into an interactive web-based application using Streamlit, allowing users to predict stroke risk more efficiently. The novelty of this study lies in integrating *Logistic Regression* with a practical, interactive web-based system, enabling both analytical accuracy and real-world usability in early stroke risk detection.[5]

METHOD

The research stages include data collection, preprocessing (imputation of missing values, encoding of categorical variables, normalization), data sharing with a stratified split of 80:20, training in the Logistic Regression model, evaluation using metrics (accuracy, precision, recall, F1-score, ROC-AUC), and implementation of Streamlit-based applications. This stage is in line with the best practices used in similar research.

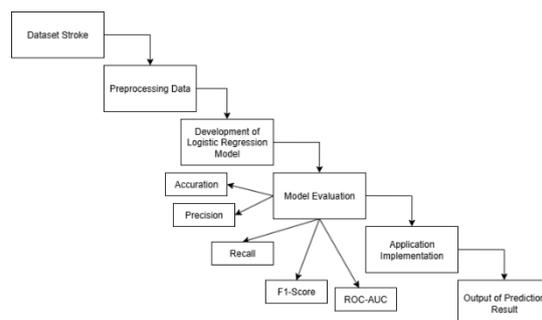


Figure 1. Research Flow

This research was carried out with several main stages, namely collecting datasets, preprocessing data, building models using the Logistic Regression method, evaluating the capabilities of the model, and implementing the model into web-based applications.

Dataset

The data used in this study is the *Brain Stroke Dataset*, which includes demographic attributes (age, gender), clinical factors (hypertension, heart disease, marital status), lifestyle indicators (smoking status), and physiological features such as BMI and average glucose level. The dataset is labeled with binary values: “1” for patients at risk of stroke and “0” for those not at risk. This dataset has also been widely utilized in previous studies to develop *machine learning*-based medical classification models. Here is the source of the data [Brain Stroke Dataset](#). [3]

Preprocessing Data

Preprocessing was conducted to enhance data quality before training the model. The steps include:

The dataset did not contain missing values in the BMI variable; therefore, no imputation was required. [1]

Encoding Categorical Variables, Variables such as *gender*, *ever_married*, and *smoking_status* were transformed using One-Hot Encoding to be processed numerically by the Logistic Regression algorithm. [2]

Numerical attributes (age, *avg_glucose_level*, *bmi*) were standardized using *StandardScaler* to ensure a mean of 0 and standard deviation of 1.

This step ensures consistent feature weighting during training. [2]

$$Z = \frac{X - \mu}{\sigma} \quad (1)$$

Where X represents the original value, μ is the mean of the feature, and σ is the standard deviation. This process ensures that all numerical features contribute proportionally during model training and prevents features with larger scales from dominating the learning process.

Due to class imbalance, SMOTE was applied to balance the stroke class and enhance model performance.

Model Development

The *Logistic Regression* algorithm was chosen because it performs effectively for binary classification problems in healthcare and produces interpretable outputs. The logistic regression model estimates the probability of stroke using the equation:

$$P(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}} \quad (2)$$

Where $P(y = 1|x)$ represents the probability of stroke, β_0 is an intercept, β_1 cohesion regression, x_1 is a predictor variable. Model parameters were optimized by minimizing the log-loss function as applied in previous medical prediction research. [3]

Model Evaluation

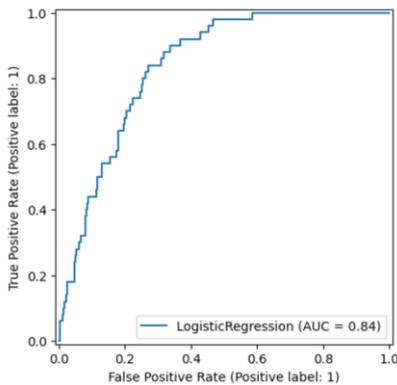


Figure 2. ROC Curve

The evaluation was carried out using train-test split with a ratio of 80:20. The model is measured using the following metrics:

Accuracy is the comparison of the number of correct predictions with the total predictions, Precision is the level of accuracy of positive predictions, Recall (Sensitivity) is the model's ability to detect positive cases (stroke), F-1 Score is the harmonization between precision and recall, ROC-AUC is the model's ability to distinguish between positive and negative classes.

The evaluation metrics table can be displayed as follows:

Table 1. Evaluation Metric

Metric	Formula
Accuracy	$\frac{TP + TN}{TP + TN + FP + FN}$ (3)
Precision	$\frac{TP}{TP + FP}$ (4)
Recall	$\frac{TP}{TP + FN}$ (5)
F1-Score	$\frac{2 \times Precision \times Recall}{Precision + Recall}$ (6)
ROC-AUC	The area of the ROC curve comparing True Positive Rate vs False Positive Rate

These metrics comprehensively assess the predictive performance of the

model and are widely applied in disease classification studies using *machine learning*. [2]

Application Implementation

After obtaining the optimal model, the final stage involves integrating it into a web-based application using Streamlit. The application allows users to input parameters such as age, hypertension history, and smoking habits, and instantly receive stroke risk predictions. Streamlit was chosen due to its simplicity, Python integration, and user-friendly interface, which facilitates accessibility for medical practitioners and non-technical users alike. A similar deployment approach has been adopted in other health informatics studies.

RESULT AND DISCUSSION

Dataset Stroke

The dataset used in this study is the [Brain Stroke Dataset](#). from Kaggle, the dataset consists of 4,981 patient records, with predictor variables including age, gender, hypertension, heart disease, marital status, work type, residence type, average glucose level, BMI, and smoking status. The target variable is binary, with “1” indicating stroke and “0” indicating no stroke. Initial exploration shows that patients’ ages range from 0–100 years (mean 45.2), *avg_glucose_level* ranges from 55.1–272.0 mg/dl (mean 106.8), and BMI from 10.3–97.6 (median 28.4). The gender distribution is nearly balanced (Male 51.2%, Female 48.8%), while stroke cases only represent 4.8% of the total population, confirming a significant class imbalance [6], [7]

Table 2. Target Class Distribution Table

Category	Sum	Presentase
No Stroke	4736	95%
Stroke	245	5%

This imbalance poses a major challenge in stroke prediction, as also reported in prior studies emphasizing the need for balancing techniques to avoid bias in model learning.

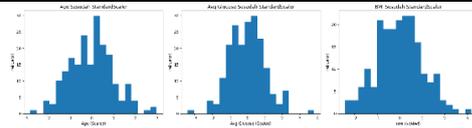


Figure 4. Numerical Variable Distribution After Normalization

Preprocessing Data

The BMI variable in the dataset used in this study did not contain missing values. Therefore, no imputation process was required. Data preprocessing focused on encoding categorical variables, normalization using StandardScaler, and handling class imbalance using SMOTE.

Data imbalance was addressed using the *Synthetic Minority Oversampling Technique (SMOTE)*, resulting in a more balanced class distribution and improved sensitivity to minority cases. [8]

Table 3 Missing Values Before and After Imputation

Variabel	Missing Value (Before)	Missing Value (After)
BMI	0	0

Development of Logistic Regression Model

The preprocessed dataset was divided into training and testing sets in an 80:20 stratified split to maintain class balance. The *Logistic Regression* model was trained with parameters summarized in Table 4:[2]

Categorical variables (*gender, ever_married, work_type, and smoking_status*) were encoded numerically to enable processing by the model, while numerical features (*age, avg_glucose_level, bmi*) were normalized using the *StandardScaler*, ensuring a mean of 0 and a standard deviation of 1[6]

Table 4. Parameter Model Logistic Regression

Parameter	Value
max_iter	1000
solver	lbfgs
random_state	42

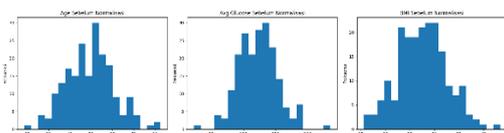


Figure 3. Distribution of Numerical Variables Before Normalization

The resulting regression coefficients (β) indicate that age, hypertension, heart disease, and average glucose level are the most influential predictors of stroke risk, while gender contributes insignificantly. These findings align with previous studies that identified similar key variables influencing stroke prediction using Logistic Regression.[9]

Model Evaluation

The model’s performance was evaluated on test data using metrics in-

cluding Accuracy, Precision, Recall, F1-Score, and ROC-AUC.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (7)$$

$$Precision = \frac{TP}{TP+FP} \quad (8)$$

$$Recall (sensitivity) = \frac{TP}{FP+FN} \quad (9)$$

$$F1 - Score = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (10)$$

The ROC-AUC is calculated by measuring the area under the ROC curve (plots TPR to FPR at various thresholds).

Table 5. Results of Logistic Regression Model Evaluation

Metrik	Nilai
Accuracy	74,8%
Precision	14%
Recall	80%
F1-Score	24%
ROC	0,84

Table 6. Confusion Matrix Logistic Regression

	Positive Actuals	Negative Actuals
Positive Predictions	706	241
Negative Predictions	10	40

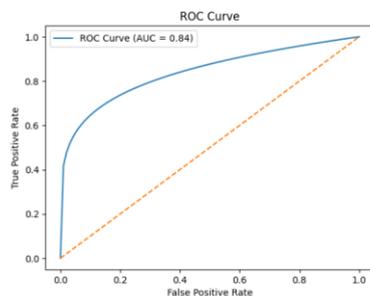


Figure 5. Kurva ROC Logistic Regression

The model shows high recall (80%), meaning it is capable of detecting most stroke cases. However, the precision is relatively low (14%), indicating a high number of false positive predictions. [10]

These results are consistent with similar studies achieving AUC values between 0.80–0.90, reinforcing the reliability of Logistic Regression for disease risk prediction.[11] Other research also highlights that Logistic Regression remains competitive against more complex algorithms like Random Forest and XGBoost when applied to structured medical datasets. This trade-off between recall and precision is common in imbalanced medical datasets, where improving sensitivity often increases false positive predictions. [12]

Application Implementation

The trained model is integrated into Streamlit-based applications. The app allows users (doctors, patients, or researchers) to enter patient data through a simple interface form, then the system will provide a prediction of stroke risk.

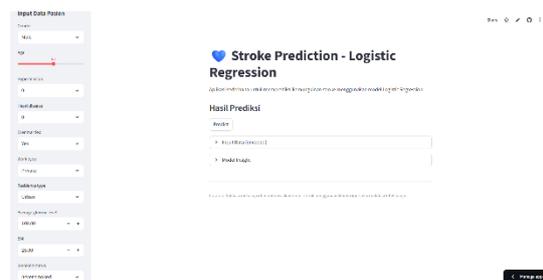


Figure 6. Application Implementation

The application features include Form Input Data, where users fill in data such as age, gender, hypertension, heart disease, marital status, type of occupation, type of residence, average glucose level, BMI, and smoking status. Predic-

tion process where data is processed according to preprocessing (encoding, normalization) and processed by the Logistic Regression model. Prediction results are the system displays results in the form of whether the patient is indicated for a stroke or not, along with the probability.[13]

This implementation demonstrates how predictive analytics can support decision-making in healthcare by providing a user-friendly and accessible platform.

Output of Prediction Results

To test the application, input is made on the test case as follows:

Table 7. Application Testing

Attribution	Information
Gender	Male
Age	67
Hypertension	0
Heart Disease	1
Ever Married	Yes
Work Type	Private
Residence Type	Urban
Avg Glucose Level	228.69
BMI	36.6
Smoking Status	Formerly Smoked

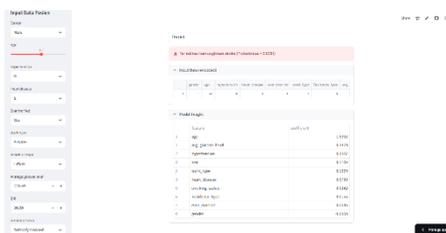


Figure 7. Application Test Results

From the test, a prediction of stroke was indicated with a probability of 0.8039. Other studies have shown that displaying predictive results in the form of probabilities other than binary classifi-

cations can help medical professionals better understand patients' risk levels.[14]

CONCLUSION

This study demonstrates that the Logistic Regression algorithm can predict stroke risk based on patient health data. The model achieved an accuracy of 74.8%, precision of 14%, recall of 80%, F1-score of 24%, and ROC-AUC of 0.84, indicating good sensitivity in identifying patients at risk of stroke despite the presence of false positive predictions. The data preprocessing stage—including handling missing values, categorical encoding, normalization, and class balancing—played a crucial role in improving model stability. Furthermore, the Streamlit-based application implementation highlights practical usability, providing medical professionals and patients with an interactive, accessible platform for early stroke detection and prevention.

Despite these results, this study is limited by the dataset size and lack of population diversity. Therefore, future research is suggested to employ larger, more comprehensive, and real-time datasets to improve the model’s accuracy, adaptability, and clinical relevance in real-world healthcare environments.

BIBLIOGRAPHY

[1] E. S. Darmawan, S. R. Hasibuan, V. Y. Permanasari, and D. Kusuma, “Disparities in health services and outcomes by National Health Insurance membership type for ischemic heart disease and stroke in Indonesia: analysis of claims, 2017–2022,” *Glob Health Res Pol*

- icy, vol. 10, no. 1, Dec. 2025, doi: 10.1186/s41256-025-00432-y.
- [2] M. Guhdar, A. Ismail Melhum, and A. Luqman Ibrahim, “Optimizing Accuracy of Stroke Prediction Using Logistic Regression,” *Journal of Technology and Informatics (JoTI)*, vol. 4, no. 2, pp. 41–47, Jan. 2023, doi: 10.37802/joti.v4i2.278.
- [3] O. A. Okwori, M. A. Agana, A. Ofem, and O. I. Ofem, “Article no.ABAARJ.1348 Original Research Article Okwori et al,” 2024.
- [4] A. G. Moelyo, M. N. Sitaresmi, and M. Julia, “Growth faltering or deceleration toward target height: Linear growth interpretation using WHO growth standard 2006 for Indonesian children,” *PLoS One*, vol. 20, no. 4 April, Apr. 2025, doi: 10.1371/journal.pone.0290053.
- [5] A. Hassan, S. Gulzar Ahmad, E. Ullah Munir, I. Ali Khan, and N. Ramzan, “Predictive modelling and identification of key risk factors for stroke using machine learning,” *Sci Rep*, vol. 14, no. 1, Dec. 2024, doi: 10.1038/s41598-024-61665-4.
- [6] Y. Fu, “A machine learning approach for predictings stroke,” *Medical Data Mining*, vol. 7, no. 3, Sep. 2024, doi: 10.53388/MDM202407015.
- [7] S. Li, “The prediction of stroke and feature importance analysis based on multiple machine learning algorithms,” *Applied and Computational Engineering*, vol. 18, no. 1, pp. 37–41, Oct. 2023, doi: 10.54254/2755-2721/18/20230961.
- [8] K. Swain *et al.*, “Enhancing Stroke Prediction Using LightGBM With SMOTE-ENN and Fine-Tuning: A Comprehensive Analysis,” *Cureus Journal of Computer Science*, Dec. 2024, doi: 10.7759/s44389-024-02268-y.
- [9] L. Li, “Stroke Prediction Base on Logistic Regression Model,” 2024.
- [10] T. Tazin, M. N. Alam, N. N. Dola, M. S. Bari, S. Bourouis, and M. Monirujjaman Khan, “Stroke Disease Detection and Prediction Using Robust Learning Approaches,” *J Healthc Eng*, vol. 2021, 2021, doi: 10.1155/2021/7633381.
- [11] R. Mitra and T. Rajendran, “EFFICIENT PREDICTION OF STROKE PATIENTS USING LOGISTIC REGRESSION ALGORITHM IN COMPARISON TO DECISION TREE ALGORITHM.”
- [12] A. Tashkova, S. Eftimov, B. Ristov, and S. Kalajdziski, “Comparative Analysis of Stroke Prediction Models Using Machine Learning,” May 2025, [Online]. Available: <http://arxiv.org/abs/2505.09812>
- [13] X. Huang *et al.*, “Novel Insights on Establishing Machine Learning-Based Stroke Prediction Models Among Hypertensive Adults,” *Front Cardiovasc Med*, vol. 9, May 2022, doi: 10.3389/fcvm.2022.901240.
- [14] M. Masuda *et al.*, “Recurrent cardiac arrests caused by Kounis syndrome without typical allergic symptoms,” *J Cardiol Cases*, vol. 27, no. 2, pp. 47–51, Feb. 2023, doi: 10.1016/j.jccase.2022.10.004.