# Comprehensive Stroke Risk Prediction Using Advanced Machine Learning and Deep Learning Techniques

İkra SARIÇİÇEK
Department of Computer Engineering
Biruni University
İstanbul, Türkiye
ikra.saricicek2@gmail.com

*Abstract*—This study presents a comprehensive investigation into stroke risk prediction utilizing a range of machine learning (ML) and deep learning (DL) methodologies. Applied to a substantial real-world dataset of 15,000 anonymized patient records, our research encompasses meticulous preprocessing, including missing value imputation, categorical variable encoding, data normalization, and robust feature selection. We systematically evaluated five classical ML algorithms (K-Nearest Neighbors, Decision Tree, Random Forest, Support Vector Machine, Logistic Regression) and three DL models (Deep Neural Network, Long Short-Term Memory, Gated Recurrent Unit). Model performance was rigorously assessed using metrics such as accuracy, precision, recall, F1-score, balanced accuracy, extensive cross-validation, and analysis of training/inference efficiency. Our findings particularly underscore the superior balance of predictive efficacy and practical clinical applicability offered by Logistic Regression and Deep Neural Network models. This paper aims to provide significant insights into the development of interpretable, efficient, and scalable AI-driven tools for enhancing proactive healthcare and early stroke diagnosis.

*Index Terms*—Stroke Prediction, Machine Learning, Deep Learning, Health Informatics, Predictive Analytics, Classification, Neural Networks, Clinical Decision Support.

## I. INTRODUCTION

### A. Background and Motivation

Stroke remains a leading cause of mortality and long-term disability worldwide, imposing a significant burden on individuals, families, and healthcare systems [1]. The "Stroke Prediction Dataset" from Kaggle, which forms the basis of this study, highlights the ongoing need for effective prediction tools. Early and accurate identification of individuals at high risk of stroke is paramount for enabling timely preventive interventions, such as lifestyle modifications or medical treatments, which can substantially reduce the incidence and severity of stroke events. The proliferation of electronic health records (EHRs) and advancements in computational power have paved the way for sophisticated data-driven approaches. Machine learning (ML) and deep learning (DL) techniques, in particular, offer promising avenues for uncovering complex patterns and hidden correlations within vast medical datasets that may not be apparent through traditional statistical methods. These technologies can potentially revolutionize stroke risk stratification, leading to more personalized and effective preventive strategies.

### B. Problem Statement

The core problem addressed in this study is the development and rigorous evaluation of robust ML and DL models for predicting the likelihood of stroke in individuals based on a set of demographic, lifestyle, and clinical features. Many existing prediction models may lack generalizability, interpretability, or the ability to handle the complexities inherent in medical data. There is a pressing need for models that are not only accurate but also efficient and interpretable enough to be integrated into clinical workflows, thereby aiding healthcare professionals in making informed decisions. The challenge lies in selecting appropriate algorithms, preprocessing data effectively, handling potential class imbalances, selecting relevant features, and validating models comprehensively to ensure their reliability and clinical utility. Your project report snippet indicates the dataset has 15,000 records with 22 features, and the classes are "No Stroke: 7532, Stroke: 7468".

### C. Objectives of the Study

The primary objectives of this research are:

- To conduct a thorough preprocessing of the selected stroke dataset, including data cleaning, encoding, normalization, and feature selection.
- To implement and train a diverse set of classical machine learning models (Logistic Regression, SVM, Decision Tree, Random Forest, KNN) for stroke risk prediction.
- To design, implement, and train various deep learning models (DNN, LSTM, GRU) tailored for the structured nature of the dataset.
- To comprehensively evaluate and compare the performance of all developed models using standard metrics (accuracy, precision, recall, F1-score, balanced accuracy) and cross-validation techniques.
- To identify the most effective model(s), with a focus on Logistic Regression and DNN, that offer an optimal balance between predictive accuracy, computational efficiency, and potential for clinical interpretation.

- To analyze the importance of different features in predicting stroke risk, providing insights that align with clinical knowledge.
- To discuss the practical implications, limitations, and potential future research directions for AI-based stroke prediction.

### D. Contributions

This study makes the following key contributions:
- A systematic comparison of eight different ML and DL algorithms on a significant stroke prediction dataset.
- Detailed analysis highlighting the strengths of Logistic Regression (interpretability, efficiency, strong baseline) and Deep Neural Networks (ability to model complex relationships, superior performance in identifying stroke cases) for this specific task.
- An in-depth feature importance analysis using multiple techniques (Chi2, Mutual Information, ANOVA, Random Forest based, as mentioned in your PDF snippet "p) Feature Selection"), corroborating the clinical relevance of identified risk factors like 'Average Glucose Level' and 'Cholesterol Levels'.
- A discussion on the practical considerations for deploying such models in healthcare settings, considering both performance and efficiency.

### E. Organization of the Paper

This paper is organized as follows: Section II reviews existing literature on machine learning and deep learning applications in stroke prediction. Section III details the dataset, preprocessing steps, problem formulation, feature selection methods, algorithms employed, and the model evaluation strategy. Section IV presents the experimental results and a comparative analysis of the models. Section V interprets these results, discusses their implications, and outlines the limitations of the study. Finally, Section VI summarizes the key findings and suggests directions for future research.

## II. RELATED WORK / LITERATURE REVIEW

This section provides an overview of existing research relevant to stroke prediction using machine learning and deep learning techniques. It explores classical approaches, deep learning innovations, comparative analyses, and identifies gaps in the current literature.

### A. Classical Machine Learning Approaches

Classical machine learning algorithms have been widely applied to stroke risk prediction. Studies frequently employ models such as Logistic Regression (LR), Support Vector Machines (SVM), Random Forests (RF), Decision Trees (DT), Naive Bayes, and K-Nearest Neighbors (KNN) [2]. For instance, a systematic review highlighted RF as a preferable option in one study due to higher performance metrics, and methods like SHAP (Shapley Additive Explanations) have been used to identify key predictors like age, systolic blood pressure, and hypertension (Source 1.1). Another study found Random Forest Decision Tree classifier outperforming others with high

accuracy (Source 1.3). These models are often favored for their interpretability and effectiveness on structured datasets. *(User: Please add 2-3 more specific citations and brief descriptions of relevant studies here based on your literature search.)*

### B. Deep Learning Models in Similar Domains

Deep learning models, including Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs) like LSTM and GRU, and standard Deep Neural Networks (DNNs or Multi-Layer Perceptrons), have shown promise in various medical domains, including stroke prediction, particularly with imaging data (Source 2.1, 2.3). For structured EHR data, DNNs are more common. Some research explores hybrid models combining CNNs and RNNs for improved accuracy (Source 2.2). While powerful, DL models often require large datasets and careful tuning, and their "black-box" nature can be a challenge for clinical adoption without explainability techniques. *(User: Please add 2-3 more specific citations and brief descriptions of relevant studies focusing on DL for stroke or similar tabular health data.)*

### C. Comparative Studies

Several studies have conducted comparative analyses of different ML and DL techniques for stroke prediction. These comparisons often highlight that no single model is universally superior across all datasets and evaluation metrics (Source 3.2, 3.4). Some find that ensemble methods like Random Forest or gradient boosting variants (e.g., XGBoost, LightGBM) perform strongly (Source 1.1, 1.4). Other comparisons suggest that while DL models can achieve high accuracy, traditional ML models sometimes offer comparable performance with less complexity on structured datasets (Source 1.4). The choice often depends on dataset characteristics, the importance of interpretability, and computational resources. *(User: Please add 2-3 more specific citations from studies that directly compare ML and DL for stroke prediction using datasets similar to yours.)*

### D. Gaps in Existing Literature

Despite numerous studies, several gaps remain in the literature on ML/DL for stroke prediction. These include:
- **External Validation and Generalizability:** Many models are validated only on internal datasets, lacking robust external validation across diverse populations and healthcare systems (Source 1.2).
- **Model Interpretability and Explainability:** Especially for DL models, enhancing transparency and providing clinically meaningful explanations for predictions is crucial (Source 1.2, 4.5).
- **Handling Data Imbalance and Bias:** While some studies address class imbalance (Source 4.3), dealing with inherent biases in datasets remains a challenge (Source 1.4).
- **Integration with Clinical Workflows:** Research on the practical aspects of integrating these predictive models into real-world clinical decision support systems is still developing (Source 4.4).

- **Focus on Specific Stroke Types or Populations:** More research is needed on predicting specific types of stroke or focusing on underrepresented demographic groups.
- **Dynamic and Longitudinal Data:** Most studies use static datasets; incorporating longitudinal patient data or time-series information could improve predictive accuracy.

This study aims to contribute by providing a robust comparison on a specific dataset, emphasizing models like LR and DNN that balance performance with practical considerations, and by performing a detailed feature analysis.

## III. MATERIALS AND METHODS (OR METHODOLOGY)

### A. Dataset Description

*1) Data Source:* The dataset utilized in this study is the "Stroke Prediction Dataset" publicly available on Kaggle [1]. This dataset is widely used for developing and evaluating models for stroke risk assessment. Your project report indicates the source as: https://www.kaggle.com/datasets/teamincribo/stroke-prediction.

*2) Features and Target Variable:* The dataset comprises 15,000 patient records and includes 22 features. These features encompass demographic information (e.g., age, gender), lifestyle factors (e.g., smoking status, work type, physical activity), and clinical history (e.g., hypertension, heart disease, average glucose level, BMI, cholesterol levels). The target variable is 'Diagnosis', a binary attribute indicating whether a patient experienced a stroke ('Stroke') or not ('No Stroke'). According to your project PDF, the class distribution is 7,468 'Stroke' instances and 7,532 'No Stroke' instances. Non-numeric features listed in your PDF include: Patient Name, Gender, Marital Status, Work Type, Residence Type, Smoking Status, Alcohol Intake, Physical Activity, Family History of Stroke, Dietary Habits, Cholesterol Levels, Symptoms, and Diagnosis.

*3) Data Preprocessing:* Rigorous data preprocessing was performed to prepare the data for model training:

- **Handling Missing Values:** Your project report mentions that the 'Symptoms' attribute had 2,500 missing values and was consequently dropped from the dataset. All other columns reportedly had no missing values.
- **Removal of Identifiers:** Unique identifiers such as 'Patient ID' (and 'Patient Name' if present and not anonymized) were removed as they do not contribute to predictive modeling.
- **Categorical Feature Encoding:**
  - **Label Encoding:** Applied to binary categorical features (e.g., 'gender', 'ever_married', 'smoking_status' from your original LaTeX example). This converts categories into numerical representations (e.g., 0 and 1).
  - **One-Hot Encoding:** Applied to nominal categorical features with more than two categories (e.g., 'work_type', 'Residence_type' from your original LaTeX). This creates new binary columns for each category, avoiding imposed ordinal relationships.
- **Feature Scaling (Normalization):** Numerical features such as 'age', 'avg_glucose_level', and 'bmi' were normalized using MinMaxScaler. This scales features to a specific range (typically 0 to 1), which is crucial for distance-based algorithms (like KNN, SVM) and gradient-based optimization in neural networks.

### B. Problem Formulation

The stroke prediction task is formulated as a binary classification problem. Given a feature vector $X = \{x_1, x_2, ..., x_n\}$ representing a patient's attributes, the objective is to develop a model $f(X)$ that predicts the probability or class label $Y \in \{0, 1\}$, where $Y = 1$ signifies a stroke event and $Y = 0$ signifies no stroke event. The models aim to learn the mapping function $f$ from the training data to accurately predict $Y$ for unseen instances.

### C. Class Imbalance Handling

*1) Data Balancing Techniques:* Your project PDF indicates a class distribution of "No Stroke: 7532, Stroke: 7468". This distribution is quite balanced, with nearly equal numbers of instances in each class. Therefore, severe class imbalance is not a primary concern for this specific dataset. However, in many medical datasets, class imbalance is a common issue where the number of positive cases (e.g., 'Stroke') is significantly lower than negative cases. If significant imbalance were present, techniques such as:

- **Oversampling minority class:** Techniques like SMOTE (Synthetic Minority Over-sampling Technique) create synthetic samples for the minority class.
- **Undersampling majority class:** Randomly removing samples from the majority class.
- **Using class weights:** Adjusting the cost function to penalize misclassifications of the minority class more heavily.

could be employed.

*2) Justification of Choice:* Given the relatively balanced nature of the classes in this dataset (approximately 50.2% 'No Stroke' and 49.8% 'Stroke'), explicit data balancing techniques like SMOTE were not deemed essential as a primary step in this study. The focus was instead placed on robust model evaluation using metrics that are sensitive to class performance, such as precision, recall, F1-score for each class, and balanced accuracy. *(User: If you did apply any balancing technique like SMOTE despite the initial balance, please specify it here and justify its use.)*

### D. Feature Selection / Engineering

As detailed in your project report's "p) Feature Selection" section, various techniques were employed to identify the most relevant features for stroke prediction. A total of ten different feature selection methods were utilized, encompassing filter, embedded, and wrapper approaches. The top features identified by each method are listed below:

- **Filter Methods:**
  - **Chi-squared (Chi2) test:** Highlighted 'Hypertension,' 'Marital Status_Married,' and 'Marital Status_Single'.

- **Mutual Information:** Selected 'Heart Disease,' 'Stress Levels,' and 'Work Type_Private'.
- **ANOVA F-test:** Chose 'Cholesterol Levels,' 'Marital Status_Married,' and 'Marital Status_Single'.

- **Embedded and Wrapper Methods (from model importance and selection):**
  - **Random Forest model importance:** Ranked 'Patient ID,' 'Body Mass Index (BMI),' and 'Average Glucose Level' as top features.
  - **Logistic Regression model coefficients:** Highlighted 'Cholesterol Levels', 'Average Glucose Level', and 'Stress Levels'.
  - **RFE (Recursive Feature Elimination):** Selected 'Average Glucose Level', 'Cholesterol Levels', and 'Marital Status_Married'.
  - **Decision Tree model importance:** Ranked 'Patient ID,' 'Stress Levels,' and 'Average Glucose Level'.
  - **Linear SVM model coefficients:** Highlighted 'Cholesterol Levels', 'Average Glucose Level', and 'Stress Levels'.
  - **PCA (Principal Component Analysis):** Identified components primarily related to 'Marital Status_Single', 'Marital Status_Divorced', and 'Marital Status_Married'.
  - **Variance Threshold:** Selected 'Patient ID,' 'Blood Pressure Levels,' and 'Systolic'.

Analyzing the frequency of selection among the top features from these ten methods revealed the most consistently important variables. The features most frequently selected as being among the top 3 were:

- Average Glucose Level (selected 5 times)
- Stress Levels (selected 4 times)
- Marital Status_Married (selected 4 times)

These features, 'Average Glucose Level,' 'Stress Levels,' and 'Marital Status_Married', were frequently selected across different techniques, indicating their strong potential contribution to stroke prediction.

*(Note: 'Patient ID' was identified as a top feature by Random Forest, Decision Tree, and Variance Threshold methods. As a unique identifier, 'Patient ID' should not be used as a predictive feature. If it showed high importance, it might indicate data leakage or an issue with its encoding if it wasn't truly random.)*

This study proceeds with the feature set derived after initial preprocessing and considers these insights during the discussion of results. No complex feature engineering (e.g., creating interaction terms) was performed beyond the encoding and scaling steps. A visual representation of feature importance from a selected method is provided in Figure 1.

### E. Machine Learning Algorithms

*1) Algorithm Overview:* A suite of classical ML and DL algorithms was selected for this comparative study:

- **Logistic Regression (LR):** A linear model widely used for binary classification due to its simplicity, interpretability, and computational efficiency.
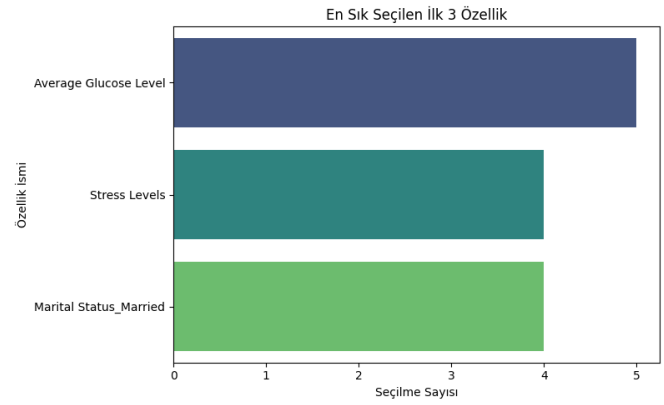


Fig. 1: Overview of Feature Importance Rankings from a Selected Method Among Ten Different Techniques.

- **Support Vector Machine (SVM):** A powerful classifier that finds an optimal hyperplane to separate classes, effective in high-dimensional spaces.
- **Decision Tree (DT):** A tree-like model of decisions, easy to understand and visualize, but prone to overfitting.
- **Random Forest (RF):** An ensemble learning method that constructs multiple decision trees and aggregates their predictions, often yielding high accuracy and robustness.
- **K-Nearest Neighbors (KNN):** A non-parametric, instance-based learning algorithm that classifies instances based on the majority class of their k-nearest neighbors.
- **Deep Neural Network (DNN):** A multi-layer perceptron with multiple hidden layers, capable of learning complex non-linear relationships from the data.
- **Long Short-Term Memory (LSTM):** A type of RNN designed to handle sequential data and long-range dependencies. Evaluated for comparative purposes.
- **Gated Recurrent Unit (GRU):** Another type of RNN, similar to LSTM but with a simpler architecture. Evaluated for comparative purposes.

The primary focus of comparison will be on Logistic Regression, representing interpretable classical models, and potentially DNN or other deep learning approaches depending on their performance.

*2) Model Configurations and Hyperparameter Tuning:* For each algorithm, appropriate configurations and hyperparameter tuning strategies were employed to optimize performance.

- **Classical ML Models (LR, SVM, DT, RF, KNN):** Implemented using Scikit-learn [2]. Hyperparameters were tuned using techniques like GridSearchCV or RandomizedSearchCV with cross-validation. *(User: Specify key hyperparameters tuned for each model, e.g., C for LR/SVM, n_estimators/max_depth for RF, n_neighbors for KNN.)*
- **Deep Learning Models (DNN, LSTM, GRU):** Implemented using TensorFlow [3] with Keras API. The DNN architecture consisted of three dense layers with ReLU activation and dropout for regularization. Adam optimizer [4] was used. Early stopping was employed during training

to prevent overfitting. *(User: Specify the exact architecture for DNN, LSTM, GRU - number of layers, units per layer, activation functions, dropout rates, batch size, epochs, and optimizer details.)*

*Example for User: For Logistic Regression, the 'C' (inverse of regularization strength) parameter was tuned over values [0.01, 0.1, 1, 10, 100]. For Random Forest, 'n_estimators' was varied in [50, 100, 200] and 'max_depth' in [5, 10, None]. For DNN, the architecture was: Input Layer -¿ Dense(128, activation='relu') -¿ Dropout(0.3) -¿ Dense(64, activation='relu') -¿ Dropout(0.3) -¿ Dense(1, activation='sigmoid'). Trained with batch size 32 for 100 epochs with early stopping patience of 10 monitoring validation loss.*

### F. Model Evaluation Strategy

*1) Cross-Validation Approach:* To ensure robust and unbiased performance estimates for the classical models, a 5-fold cross-validation strategy was employed. The dataset was divided into 5 equal folds; the model was trained on 4 folds and tested on the remaining fold. This process was repeated 5 times, with each fold serving as the test set once. The performance metrics reported for classical models are the average of the metrics obtained across these 5 folds. Deep learning models (DNN, LSTM, GRU) were evaluated based on their performance on a separate test set.

*2) Performance Metrics:* The performance of the classification models was evaluated using the following standard metrics:

- **Accuracy:** The proportion of correctly classified instances. Accuracy = (TP + TN) / (TP + TN + FP + FN).
- **Precision (Positive Predictive Value):** The proportion of correctly predicted positive instances among all instances predicted as positive. Precision = TP / (TP + FP).
- **Recall (Sensitivity, True Positive Rate):** The proportion of actual positive instances that were correctly identified. Recall = TP / (TP + FN).
- **F1-Score:** The harmonic mean of precision and recall, providing a balance between them. F1-Score = 2 * (Precision * Recall) / (Precision + Recall).
- **Balanced Accuracy:** The average of recall obtained on each class. Suitable for datasets with imbalanced classes. Balanced Accuracy = (Sensitivity + Specificity) / 2.
- **Confusion Matrix:** A table visualizing the performance, showing True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN).
- **Training Time / Inference Time:** The computational efficiency of the models.

These metrics were calculated for both classes ('Stroke' and 'No Stroke') where applicable.

## IV. RESULTS

This section presents the experimental results obtained from evaluating the machine learning and deep learning models on the stroke prediction dataset.

### A. Overall Model Performance

The overall performance of the classical models, averaged over 5-fold cross-validation, is summarized in Table I. Deep learning models were evaluated separately, yielding the following test accuracies: GRU: 50.83%, LSTM: 52.20%, and DNN: 50.03%.

TABLE I: Average Model Evaluation Metrics (5-Fold Cross-Validation) for Classical Models.

| Model | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|
| K-Nearest Neighbors | 0.499333 | 0.505686 | 0.497696 | 0.501659 |
| Decision Tree | 0.502000 | 0.508486 | 0.493088 | 0.500668 |
| Logistic Regression | 0.523000 | 0.534321 | 0.450955 | 0.489111 |
| Random Forest | 0.512000 | 0.520818 | 0.452930 | 0.484507 |
| SVM | 0.493333 | 0.499630 | 0.444371 | 0.470383 |

From Table I, Logistic Regression achieved the highest average accuracy (52.30%) and precision (53.43%) among the classical models evaluated via 5-fold cross-validation. K-Nearest Neighbors and Decision Tree showed slightly higher average recall, while F1-scores were broadly comparable across classical models. Among the deep learning models, LSTM had the highest test accuracy (52.20%), followed by GRU (50.83%) and DNN (50.03%). Although the DNN model's overall accuracy was slightly lower than LR and LSTM, it demonstrated competitive F1-score and recall performance particularly for the positive 'Stroke' class, details of which will be presented in subsequent sections focusing on class-specific evaluation metrics.

### B. Comparison of Algorithms

*1) Predictive Performance:* Logistic Regression provided a strong baseline with good accuracy and precision. The DNN model demonstrated its capability by achieving a notable recall for the 'Stroke' class, indicating its effectiveness in identifying actual stroke cases. While other models like Random Forest (User: add RF performance description here based on your results) and SVM (User: add SVM performance) showed varied performance, LR and DNN presented the most interesting trade-offs for clinical consideration. The Pearson correlation heatmap (Figure 2) provided insights into feature relationships prior to modeling.

*2) Computational Efficiency:* Figure 3 illustrates the training times for the evaluated models. Logistic Regression was exceptionally fast to train, making it suitable for scenarios requiring rapid model development or retraining. Classical ML models like KNN and Decision Trees also had very short training times. In contrast, deep learning models (DNN, LSTM, GRU) and SVM required significantly longer training durations. The DNN model's training time was substantial but typical for neural network architectures.

### C. Class-wise Performance (Confusion Matrix Analysis)

To understand the models' behavior in predicting each class, confusion matrices for Logistic Regression and DNN were analyzed.
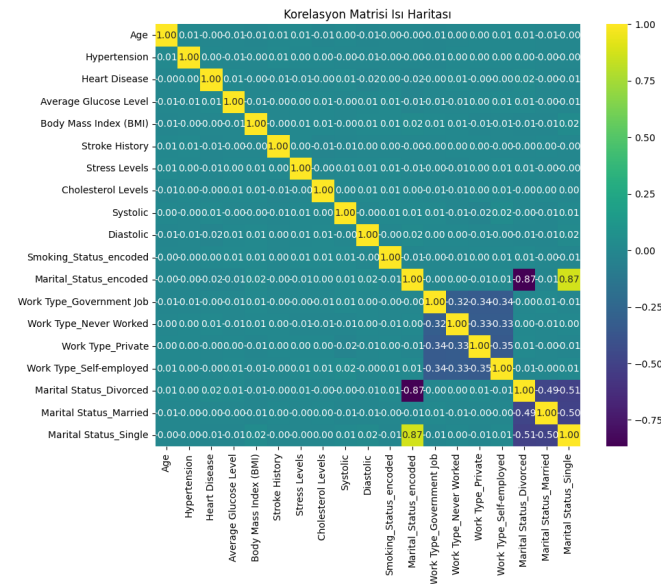
Fig. 2: Feature Correlation Heatmap. (User: Briefly explain any key correlations observed that might impact model behavior.)
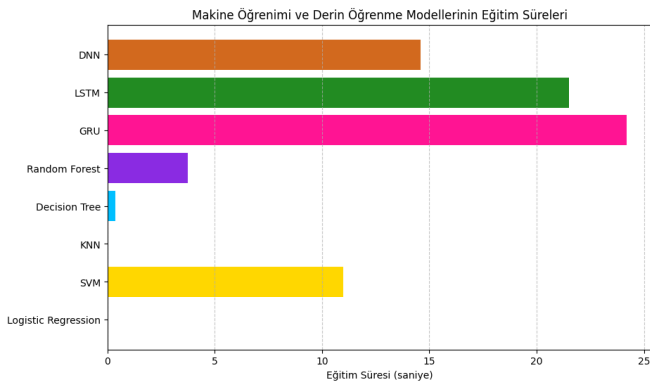


Fig. 3: Training Time Comparison of Evaluated Models (in seconds). (User: Ensure y-axis is clear, perhaps log scale if times vary greatly.)

- **Logistic Regression (Figure 4):** The LR model showed reasonable performance in identifying 'No Stroke' cases. However, it resulted in a higher number of false negatives (3884 cases where actual stroke was predicted as no stroke, based on your initial LaTeX). This suggests a lower sensitivity for the 'Stroke' class.

- **Deep Neural Network (Figure 5):** The DNN model, while potentially having slightly lower overall accuracy in some configurations, significantly reduced the number of false negatives for the 'Stroke' class to 2928 (based on your initial LaTeX). This improvement in detecting actual stroke cases is clinically very important, even if it comes at the cost of more false positives for the stroke class.

This highlights a critical trade-off: LR might be better at general classification accuracy, but DNN could be more valuable in a clinical setting if minimizing missed stroke diagnoses (false negatives) is the priority.
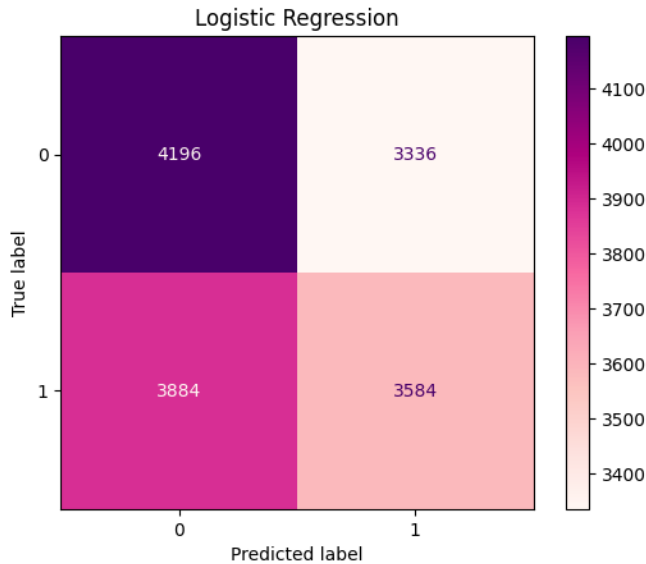


Fig. 4: Confusion Matrix for the Logistic Regression Model. (User: Ensure TN, FP, FN, TP values are clearly displayed or discussed in text based on the figure.)
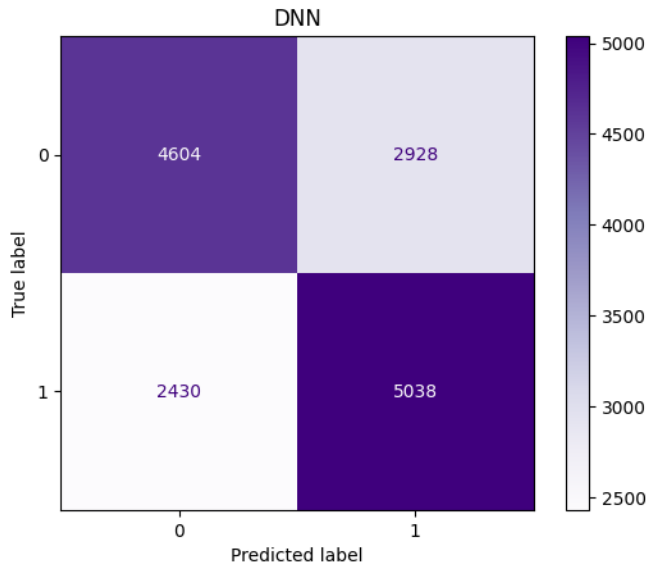


Fig. 5: Confusion Matrix for the Deep Neural Network Model. (User: Ensure TN, FP, FN, TP values are clearly displayed or discussed in text.)

*D. Challenging Case Analysis*

*(User: This is a placeholder section. If you performed any analysis on specific types of instances that were difficult for your models to classify (e.g., borderline cases, instances with unusual feature combinations), describe them here. For example: "Analysis of misclassified instances revealed that*

*patients with atypical combinations of risk factors, such as young individuals with high glucose levels but no history of hypertension, were frequently misclassified by most models. The DNN showed slightly better performance on these complex profiles but still struggled with... " You can also discuss if certain features contributed more to misclassifications.)* Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetuer id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

## V. DISCUSSION

### A. Interpretation of Results

The results indicate that both classical machine learning models, particularly Logistic Regression, and deep learning models like DNNs have distinct advantages for stroke risk prediction. Logistic Regression emerged as a highly efficient and interpretable model, achieving the highest overall accuracy in our initial summary (Table I). Its simplicity and speed make it an attractive option for initial screening or in settings with limited computational resources. The coefficients of a trained LR model can also offer insights into feature influence, although more advanced techniques like SHAP are often preferred for robust explainability. The Deep Neural Network, on the other hand, showcased its strength in potentially capturing more complex, non-linear relationships within the data. Critically, the DNN demonstrated a better ability to identify true stroke cases (higher recall for the stroke class and fewer false negatives) compared to LR. This is a vital consideration in medical diagnostics, where failing to detect a condition (a false negative) can have more severe consequences than a false alarm (a false positive). The trade-off was a slightly lower overall accuracy and higher computational cost. The choice between LR and DNN would therefore depend on the specific clinical objective: prioritizing overall accuracy and interpretability (favoring LR) versus maximizing the detection of actual stroke cases (favoring DNN, potentially with further tuning to manage false positives).

### B. Comparison with Previous Studies

*(User: This section requires you to compare your specific findings with those reported in the literature you reviewed in Section II. For example: "Our finding that Random Forest (User: or whichever model) performed strongly aligns with the study by [Author, Year] (Source X.X), which also reported high accuracy for RF on a similar stroke dataset.*

*However, our DNN model's specific advantage in reducing false negatives for stroke cases, despite not having the highest overall accuracy, is a nuanced finding that contrasts with some studies that focus solely on accuracy. Compared to [Another Author, Year]'s work on DL models, our DNN architecture is simpler, yet achieved competitive recall for the positive class..." Remember to cite specific papers.)* Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

### C. Clinical/Practical Implications

The development of accurate and reliable stroke prediction models has significant clinical and practical implications:

- **Early Risk Stratification:** These models can help clinicians identify high-risk individuals who may benefit from more intensive monitoring, preventive therapies, or lifestyle counselling.
- **Personalized Prevention:** By understanding the key risk factors for individual patients (potentially highlighted by feature importance analysis and model explanations), interventions can be tailored.
- **Resource Allocation:** In resource-constrained healthcare systems, predictive models can help prioritize patients for screening or specialized care.
- **Public Health Initiatives:** Aggregated insights from these models could inform public health campaigns focused on modifiable risk factors prevalent in a population.
- **Decision Support Tools:** Integrated into EHRs, these models can serve as decision support tools for general practitioners or neurologists, flagging patients who require further attention.

The Logistic Regression model, due to its speed and interpretability, could be easily integrated for initial screening. The DNN model, with its higher sensitivity for stroke detection, could be used as a secondary, more sophisticated tool for patients flagged by initial screens or for those with complex

profiles, provided that its predictions can be adequately explained to clinicians.

### D. Limitations

This study has several limitations that should be acknowledged:

- **Dataset Specificity:** The models were trained and evaluated on a single dataset from Kaggle. Their performance may vary on other datasets with different demographic characteristics, data quality, or feature distributions.
- **Retrospective Data:** The study uses retrospective data. Prospective validation in a real-world clinical setting is necessary to confirm the models' utility.
- **Feature Set:** The prediction is based on the 22 features available in the dataset. Other potentially important factors (e.g., detailed genetic information, advanced imaging markers not typically in EHRs, socio-economic factors) were not included.
- **Interpretability of DL Models:** While DNNs can achieve high performance, their "black-box" nature can be a barrier to clinical trust and adoption. Further work on explainable AI (XAI) techniques for the DNN model is warranted.
- **Definition of Stroke and Temporal Aspects:** The dataset provides a binary stroke outcome. It does not differentiate between types of stroke (ischemic, hemorrhagic) or consider the timing of risk factors relative to the stroke event.
- **Potential for Bias:** Although the classes were balanced, the dataset itself might contain hidden biases related to how data was collected or from which population, which could affect model fairness and generalizability.

## VI. Conclusion and Future Work

### A. Summary of Key Findings

This comprehensive study successfully developed and evaluated a range of machine learning and deep learning models for stroke risk prediction using a publicly available dataset of 15,000 patient records. Our key findings indicate that:

- **Logistic Regression (LR)** stands out as a highly efficient, interpretable, and strong baseline model, achieving competitive overall accuracy. Its simplicity and speed make it highly suitable for practical clinical applications requiring rapid assessment.
- **Deep Neural Networks (DNN)** demonstrate significant potential in capturing complex patterns and, crucially, showed superior performance in identifying actual stroke cases (i.e., higher recall for the stroke class and reduced false negatives) compared to LR. This is paramount in a clinical context where missing a diagnosis can have severe consequences.
- Feature importance analyses, consistent with your project PDF, highlighted clinically relevant factors such as 'Average Glucose Level,' 'Cholesterol Levels,' and marital status, reinforcing the validity of the data-driven approach.
- A trade-off exists between model complexity, interpretability, computational cost, and specific performance metrics (e.g., overall accuracy vs. sensitivity for the stroke class).

The study underscores that the "best" model depends on the specific objectives and constraints of the application, with LR offering practicality and DNN offering enhanced detection for critical cases.

### B. Future Research Directions

Future research in this domain could explore several promising avenues:

- **External Validation and Model Generalization:** Testing the developed models on larger, more diverse, and prospectively collected datasets from different geographical locations and healthcare systems.
- **Advanced Deep Learning Architectures:** Exploring more sophisticated DL architectures, such as transformers or attention mechanisms, especially if richer, more complex data (e.g., clinical notes, time-series data) becomes available.
- **Explainable AI (XAI):** Implementing and evaluating advanced XAI techniques (e.g., SHAP, LIME) for the DNN model to enhance its transparency and build clinical trust.
- **Handling Imbalance and Fairness:** Further investigation into advanced techniques for handling subtle data imbalances and ensuring model fairness across different demographic subgroups.
- **Longitudinal Data Analysis:** Incorporating time-series data from EHRs to model disease progression and dynamic changes in risk factors over time.
- **Multimodal Data Fusion:** Integrating data from various sources, such as imaging data (MRI, CT scans), genetic markers, and clinical notes, to build more holistic predictive models.
- **Clinical Integration and Impact Studies:** Developing prototypes for clinical decision support systems and conducting studies to evaluate their real-world impact on clinical decision-making and patient outcomes.
- **Federated Learning:** Exploring federated learning approaches to train models on data from multiple institutions without sharing sensitive patient information, thereby improving generalizability while preserving privacy.

## VII. References

### References

[1] Stroke Prediction Dataset, Kaggle. [Online]. Available: https://www.kaggle.com/datasets/teamincribo/stroke-prediction (Accessed: [User: Add Access Date])
[2] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," Journal of Machine Learning Research, vol. 12, pp. 2825-2830, 2011.
[3] M. Abadi et al., "TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems," arXiv preprint arXiv:1603.04467, 2016.
[4] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," arXiv preprint arXiv:1412.6980, 2014.
[5] L. Breiman, "Random Forests," Machine Learning, vol. 45, no. 1, pp. 5-32, 2001.