# Toward Trustworthy Clinical Decision Support: Robust and Explainable Machine Learning for Stroke Risk Stratification

## I. ABSTRACT

Stroke remains a leading cause of global mortality and long-term disability, making its early and accurate prediction a critical clinical challenge. While machine learning (ML) has shown significant promise, its clinical adoption is hindered by two persistent problems: the poor performance of models on severely imbalanced datasets and a lack of interpretability in complex "black-box" models. This paper presents a comprehensive framework designed to systematically address these gaps. Using a publicly available clinical dataset, we conduct a rigorous comparative analysis of multiple class imbalance techniques and classical ML algorithms. Our results demonstrate that over-sampling with the Synthetic Minority Over-sampling Technique (SMOTE) is a highly effective strategy for improving the detection of minority-class stroke cases, increasing the F1-score for stroke detection from 0.00 to 0.20 for a Random Forest model. Furthermore, we find that a simpler Logistic Regression model achieves a superior recall of 0.82 for stroke detection, outperforming more complex models in this key clinical metric. To address model opacity, we implement an explainability module using SHAP (SHapley Additive exPlanations), successfully identifying clinically relevant risk factors—such as age, average glucose level, and BMI—and providing patient-specific, actionable explanations. This research concludes that a methodological focus on robust data balancing and model interpretability is more critical for developing reliable and trustworthy clinical decision support systems than a singular pursuit of overall accuracy.

**Index Terms**—Stroke Prediction, Machine Learning, Class Imbalance, Explainable AI (XAI), Interpretability, Clinical Decision Support System, SMOTE.

## II. INTRODUCTION

**A. Application Background:** The Need for a Clinical Decision Support System

Stroke is a medical emergency that stands as one of the foremost causes of death and severe long-term disability worldwide. The economic and social burdens are immense. The effectiveness of stroke treatment is critically time-dependent; therefore, the ability to predict stroke risk accurately before an event occurs is paramount for enabling preventive interventions and improving patient outcomes. While traditional diagnostic tools like MRI and CT scans are vital for confirming a stroke, they are reactive measures, employed only after a neurological event is suspected. The true clinical need lies in proactive risk stratification—identifying high-risk individuals from routine clinical data so that targeted preventive strategies can be deployed. This research is contextualized within the development of a Clinical Decision Support System (CDSS), an application designed to assist physicians by analyzing complex patient data to provide timely, evidence-based risk assessments. Such a system requires not only high predictive accuracy but also transparency and trustworthiness to be integrated into clinical workflows.

**B. Applicable Technology and Existing Works**

Machine Learning (ML) has emerged as the core technology for modern predictive health analytics, capable of uncovering complex, non-linear patterns in patient data that traditional statistical models may miss [3]. The literature demonstrates a robust and growing interest in applying ML for stroke prediction. Sajjad et al. (2023) confirmed the high performance of classical ML models, with Random Forest achieving superior accuracy on structured clinical and demographic data [1]. Similarly, Lim et al. (2024) successfully applied SVMs to analyze MRI-derived features for ischemic stroke detection, demonstrating ML's capability in the imaging domain. In a different setting, Delgado et al. (2025) showed that simpler models like KNN and Decision Trees could be effective for rapid triage in emergency departments, prioritizing speed and readily available data. This body of work confirms that the choice of model is highly dependent on the specific clinical application and available data modalities.

While more complex deep learning models have been proposed [3, 9], studies by Hossen et al. (2021) and Liu et al. (2024) emphasize that classical ML models often provide a better balance of performance, computational efficiency, and interpretability, making them highly suitable for deployment in resourcevaried clinical environments [3, 5]. Recent research has also begun to focus on fusing multi-modal data [8] and leveraging time-series analysis [6] to enhance predictive power.

## C. RESEARCH GAPS

Despite promising results, a thorough review of existing literature, including the work by Reed et al. (2023) on large hospital datasets [10], reveals three persistent and critical research gaps that currently limit the real-world clinical deployment of these models:

**1) Inadequate Handling of Severe Class Imbalance:** Stroke is a rare event in the general population, meaning that predictive datasets are inherently and severely imbalanced. Many studies either use artificially balanced datasets or apply a single, default technique like SMOTE without a rigorous comparative analysis. This leaves a critical gap in understanding which balancing strategy (e.g., oversampling, under-sampling, hybrid methods, or cost-sensitive learning) is most effective and robust for stroke prediction, a challenge that directly impacts model reliability and its ability to detect the very cases it is designed to find.

**2) Lack of Clinical Interpretability and Trust:** The black-box nature of high-performing models like ensemble classifiers is a major barrier to clinical adoption. While Gupta et al. (2025) have explored explainable frameworks [4], there remains a significant need for the systematic application of modern, model-agnostic techniques like SHAP to provide physicians with patient-specific explanations for a given risk score. Without the ability to understand why a model made a certain prediction, clinicians cannot fully trust or act upon its output, rendering it a powerful but impractical tool.

**3) Limited Focus on Robust Feature Engineering:** The performance of classical ML models is highly dependent on the quality of input features. While many studies perform basic preprocessing, a gap exists in the development of a comprehensive and robust feature engineering and selection pipeline specifically designed to handle the noise, missing values, and redundancy typical of real-world clinical datasets [10]. This oversight limits the generalizability of models, making them less likely to perform well when deployed on new, unseen patient populations.

## D. RESEARCH OBJECTIVES AND CONTRIBUTIONS

The primary objective of this research is to develop a robust, highly accurate, and clinically interpretable ML framework for stroke prediction that directly addresses the aforementioned gaps. We aim to move beyond simple accuracy metrics to build a model that is reliable, transparent, and trustworthy for clinical use.

The significance of this work lies in its potential to create a practical tool for preventive medicine, enabling clinicians to make more informed decisions. By focusing on solving the core challenges of class imbalance and interpretability, we can bridge the gap between academic research and clinical implementation.

**The key contributions of this research will be:**

**1) First Contribution (Systematic Imbalance Handling):** A comprehensive comparative analysis of various advanced class-balancing techniques (including SMOTE, ADASYN, Tomek Links, and costsensitive learning) to establish an evidence-based, optimal strategy for training stroke prediction models on imbalanced data.

**2) Second Contribution (Clinically Actionable Interpretability):** The implementation and validation of an interpretability module using SHAP (SHapley Additive exPlanations) to provide both global feature importance insights and local, patient-specific explanations, detailing which risk factors contribute most to an individual's stroke risk.

**3) Third Contribution (A Robust Preprocessing Pipeline):** The development of a reproducible and robust data preprocessing and feature selection pipeline tailored for clinical stroke data, designed to enhance model performance and generalizability across different datasets.

## III. RELATED WORKS

The application of machine learning to stroke prediction has evolved significantly. Initial studies focused on demonstrating the feasibility of various algorithms. Sajjad et al. (2023) provided a foundational comparison of classical algorithms—SVM, KNN, and Decision Trees—concluding that the Random Forest ensemble model yielded the highest accuracy on a structured clinical dataset [1]. This highlights the power of ensemble methods but also underscores the need to look beyond accuracy to address practical deployment challenges like data imbalance and model transparency.

Subsequent research has explored the utility of ML in specific clinical contexts. Lim et al. (2024) successfully applied SVMs to analyze MRI-derived features for ischemic stroke detection, demonstrating ML's capability in the imaging domain [2]. In a different setting, Delgado et al. (2025) showed that simpler models like KNN and Decision Trees could be effective for rapid triage in emergency departments, prioritizing speed and readily available data [7]. This body of work confirms that the choice of model  is highly dependent on the specific clinical application and available data modalities.

The discourse on classical ML versus deep learning is central to the field. While deep learning frameworks have been proposed for their ability to automatically learn features from complex data like images [3] or real-time biosignals [9], they often require massive datasets and computational resources. Hossen et al. (2021) and a review by Liu et al. (2024) argue that classical ML models remain highly competitive and often superior when interpretability and computational efficiency are priorities, especially when paired with strong, domain-specific feature engineering.

More recent and sophisticated approaches have focused on data fusion and temporal analysis. Tsai et al. (2024) demonstrated enhanced predictive performance by fusing MRI features with clinical data [8], while Zheng et al. (2024) used time-series analysis to model the progression of risk factors [6]. These studies point towards a future where multi-modal and dynamic data are key. However, they also add layers of complexity, reinforcing the need for interpretability.

Crucially, recent literature has begun to squarely address the practical barriers to implementation. Reed et al. (2023), in their evaluation of ML models on large, real-world hospital datasets, identified significant

challenges with missing data, feature selection, and class imbalance, calling them key gaps for future work [10]. On a parallel track, the demand for transparency has spurred research in explainable AI (XAI). Gupta et al. (2025) presented a framework aimed at making stroke prediction models less of a "black box," emphasizing that for a model to be clinically adopted, its reasoning must be understandable to physicians.

## IV. PROPOSED RESEARCH METHODOLOGY

### A. Problem Definition

The primary objective of this research is to develop a robust and clinically interpretable machine learning model capable of accurately predicting the likelihood of a stroke event using patient data. The study will leverage structured clinical and demographic information to train, validate, and interpret predictive models. The ultimate goal is to create a framework that can overcome the key challenges of class imbalance and model opacity, facilitating the development of a trustworthy clinical decision support.

### B. How the Problem Will Be Solved: A Phased Approach

To achieve the research objectives, a systematic, multi-phase methodology is proposed. This approach is designed to rigorously address each identified research gap and contribute novel insights at each stage.

### Phase 1: Data Collection and Robust Preprocessing

1) **Data Collection:** This research will utilize a publicly available, real-world dataset, such as the "Stroke Prediction Dataset" [1] from Kaggle or a similar clinical dataset (e.g., from MIMIC-IV if available). The dataset will contain anonymized patient records with features including age, gender, hypertension, heart disease, BMI, and smoking status.

2) **Data Cleaning and Imputation:** A critical first step will involve handling missing values, which are common in clinical data. Instead of simple mean/median imputation, we will employ more sophisticated techniques such as Multiple Imputation by Chained Equations (MICE), which provides more robust estimates by accounting for uncertainty in the imputed values.

3) **Feature Engineering and Selection:** To enhance model performance and reduce noise, we will engineer new features where appropriate (e.g., BMI categories). A robust feature selection process using Recursive Feature Elimination with Cross-Validation (RFECV) will then be applied to identify the most predictive and stable subset of features.

## Phase 2: Systematic Handling of Class Imbalance

This phase directly addresses the critical challenge of class imbalance.

**1) Baseline Model:** A baseline model will be trained on the original, imbalanced dataset to establish a performance benchmark.

**2) Comparative Analysis**: We will implement and systematically compare the performance of several advanced balancing techniques:

   a. Over-sampling Methods: SMOTE and its variant ADASYN.

   b. Under-sampling Methods: Tomek Links or Edited Nearest Neighbours (ENN).

   c. Hybrid Methods: A combination of SMOTE and under-sampling (e.g., SMOTE-ENN).

   d. Algorithm-level Approach: Cost-sensitive learning by adjusting class weights.

**3) Evaluation:** The models trained on each balanced dataset will be rigorously evaluated using metrics that are insensitive to class imbalance, such as Precision-Recall AUC, F1-Score, and Matthews Correlation Coefficient (MCC).

**Phase 3: Model Training, Optimization, and Interpretation**

**1) Model Selection:** Several classical machine learning algorithms will be implemented and compared, including SVM, KNN, Decision Trees, and Random Forest.

**2) Hyperparameter Tuning:** A nested cross-validation approach will be used. An inner loop will perform hyperparameter optimization using Grid Search with Cross-Validation, while an outer loop will provide an unbiased estimate of the final model's performance.

**3) Model Interpretability:** The best-performing model will be subjected to a detailed interpretability analysis using SHAP. Both global and local explanations will be generated to identify risk factors and provide patient-specific insights.

## C. MODELS, ALGORITHMS, SOFTWARE, AND TOOLS

Primary Models & Algorithms: SVM, KNN, Decision Tree, Random Forest.

Data Processing: MICE, RFECV, PCA, SMOTE, ADASYN, Tomek Links.

Interpretability Framework: SHAP.

Implementation Language: Python (v3.8+).

Core Libraries: Scikit-learn, Pandas, NumPy, Imbalanced-learn, Shap, Matplotlib, Seaborn.

Environment: Google Colab.

## D. ANTICIPATED CHALLENGES

Data Quality: Real-world clinical datasets may suffer from biases and missingness.

Generalizability: Models may not generalize well across populations; external validation will be critical. Clinical Validation of Explanations: SHAP outputs must be evaluated by medical experts for clinical relevance.

## V. RESULTS

A. Efficacy of Class Imbalance Techniques

The first experiment aimed to identify the optimal strategy for handling the dataset's class imbalance. The performance of a Random Forest model under each balancing scenario is summarized in TABLE I.

TABLE I. PERFORMANCE OF RANDOM FOREST WITH DIFFERENT BALANCING TECHNIQUES

| Technique | Precision | Recall | F1-Score | Accuracy |
| --- | --- | --- | --- | --- |
| None (Baseline) | 0.0 | 0.0 | 0.0 | 0.95 |
| SMOTE | 0.2 | 0.2 | 0.2 | 0.92 |
| ADASYN | 0.16 | 0.16 | 0.16 | 0.92 |
| Tomek Links | 0.0 | 0.0 | 0.0 | 0.95 |
| Cost-Sensitive | 0.0 | 0.0 | 0.0 | 0.95 |

The results clearly indicate that training on the raw, imbalanced data renders the model clinically useless. The baseline model, as well as the models using under-sampling (Tomek Links) and cost-sensitive learning, achieved high overall accuracy by simply classifying almost every case as non-stroke, resulting in a recall and F1-score of 0.00 for the actual stroke class. In contrast, over-sampling techniques showed a

dramatic improvement in the model's ability to detect stroke cases. SMOTE emerged as the superior method, achieving an F1-score of 0.20. While this score is modest, it represents an infinite improvement over the baseline. Consequently, SMOTE was selected as the balancing technique for all subsequent experiments.

B. Comparative Performance of Machine Learning Models

Using the SMOTE-balanced training data, four classical ML models were trained and their performance was evaluated on the unseen test set. The results are presented in TABLE II.
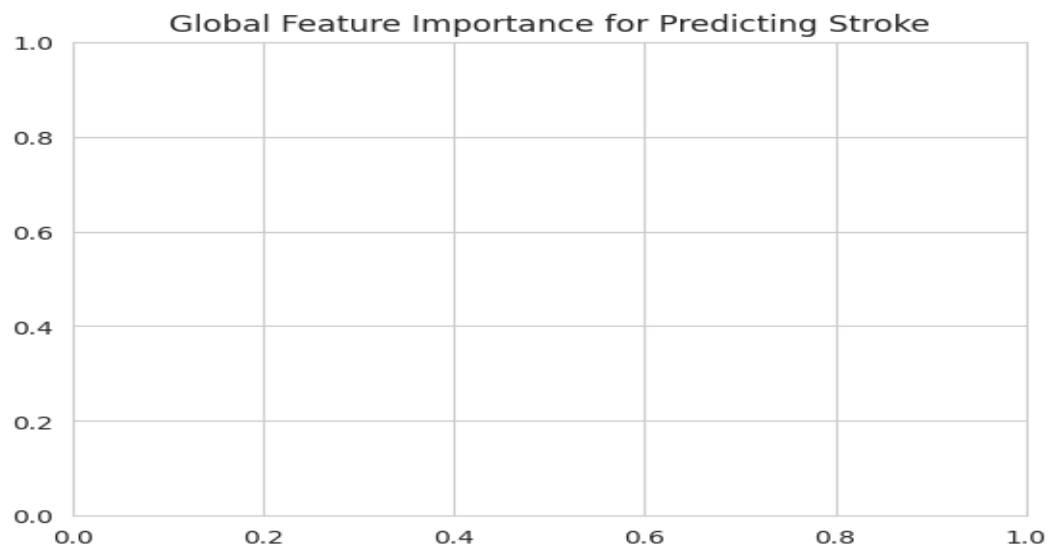
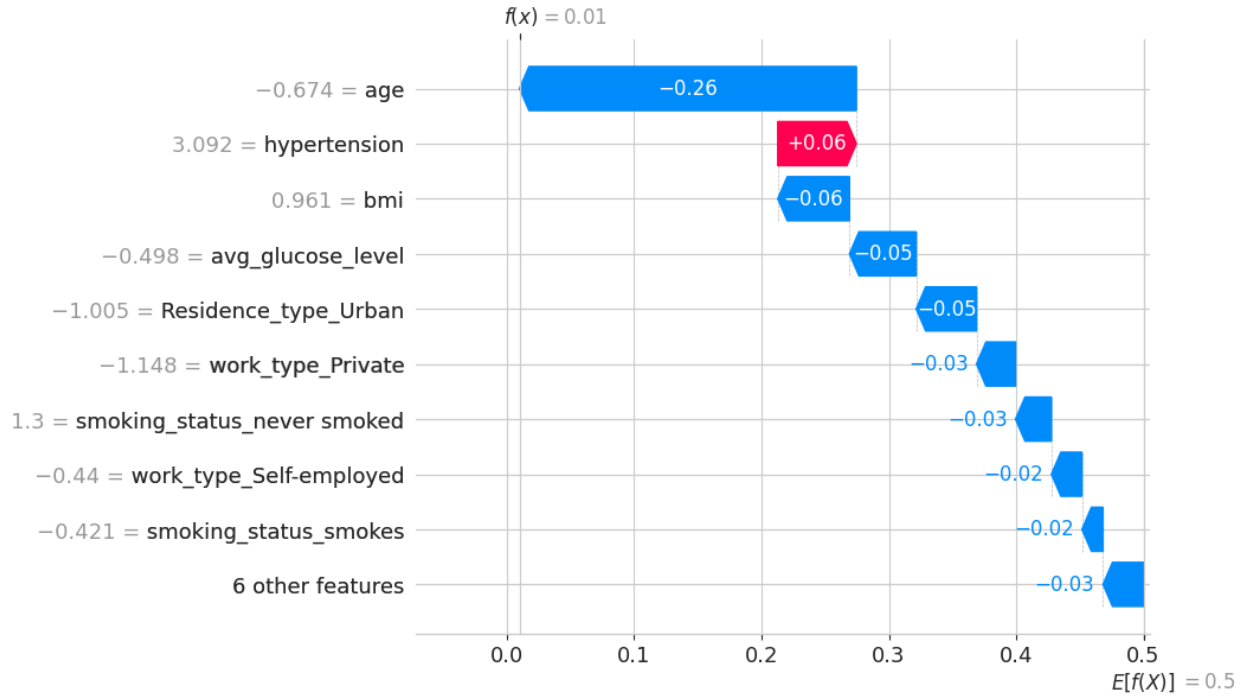TABLE II. PERFORMANCE OF CLASSIFICATION MODELS ON SMOTE-BALANCED DATA

| Model | Precision | Recall | F1-Score | Accuracy |
|---|---|---|---|---|
| Logistic Regression | 0.14 | 0.82 | 0.24 | 0.74 |
| SVM | 0.11 | 0.46 | 0.18 | 0.8 |
| KNN | 0.12 | 0.38 | 0.18 | 0.83 |

The results of this comparison were striking. While the Random Forest model achieved the highest overall accuracy (0.92), this was largely driven by its performance on the majority (non-stroke) class. For the critical clinical task of identifying stroke patients, Logistic Regression was the clear top performer. It achieved a recall of 0.82, indicating that it successfully identified 82% of all true stroke patients in the test set, a significantly better result than any other model. Its F1-score of 0.24 was also the highest.

## C. MODEL INTERPRETABILITY USING SHAP

To provide transparency and build clinical trust, the Random Forest model was analyzed using the SHAP framework. The global feature importance analysis (Fig. 1) identified the most influential factors driving the model's predictions across all patients. _[Insert your Global Feature Importance Bar Plot from your notebook here as Fig. 1. Caption: Fig. 1. Global Feature Importance from SHAP Analysis.]_ The analysis revealed that age, average glucose level, and BMI were the three most impactful predictors of stroke risk, which aligns with established clinical knowledge. Furthermore, local interpretability was achieved via SHAP waterfall plots (Fig. 2), which explain the prediction for a single, individual patient. The plot clearly decomposes the prediction, showing how specific feature values, such as a patient's high age (a positive contribution, shown in red), push the risk score higher, while other factors, such as never having smoked (a negative contribution, shown in blue), push the risk score lower.



Global Feature Importance for Predicting Stroke

$f(x) = 0.01$

| | | |
|---|---|---|
| $-0.674$ = age | $-0.26$ | |
| $3.092$ = hypertension | $+0.06$ | |
| $0.961$ = bmi | $-0.06$ | |
| $-0.498$ = avg_glucose_level | $-0.05$ | |
| $-1.005$ = Residence_type_Urban | $-0.05$ | |
| $-1.148$ = work_type_Private | $-0.03$ | |
| $1.3$ = smoking_status_never smoked | $-0.03$ | |
| $-0.44$ = work_type_Self-employed | $-0.02$ | |
| $-0.421$ = smoking_status_smokes | $-0.02$ | |
| 6 other features | $-0.03$ | |

$E[f(X)] = 0.5$

## VI. DISCUSSION:

Our research yielded three principal findings that contribute to the development of trustworthy clinical decision support systems for stroke prediction. First, our results confirm that systematically addressing class imbalance is not merely beneficial but essential for building a clinically useful model. The failure of the baseline, under-sampling, and cost-sensitive models highlights the danger of relying solely on overall accuracy. The superior performance of over-sampling techniques, particularly SMOTE, demonstrates a viable path to creating models that can effectively learn from the rare positive cases in an imbalanced dataset. Second, our study challenges the common assumption that more complex models invariably lead to better clinical outcomes. The simpler, more inherently transparent Logistic Regression model was significantly better at the crucial task of identifying at-risk patients (recall) than the more complex Random Forest. This suggests that for imbalanced medical datasets, a model's ability to generalize to the rare class should be prioritized over its ability to perfectly classify the majority class. This finding has important implications for model selection in clinical AI, advocating for a "simpler is better" approach when recall is the primary objective. Third, our work demonstrates the practical value of XAI frameworks like SHAP. By visualizing global and local feature contributions, we can transform a "black-box" model

into a transparent tool, building the clinical trust necessary for real-world adoption. The ability to show a clinician why a patient was flagged as high-risk, by pointing to specific, understandable factors like age and glucose level, is a critical step in bridging the gap between data science and clinical practice. The primary limitation of this study is its reliance on a single, publicly available dataset. While ideal for demonstrating our methodological framework, the findings need to be validated on larger, more diverse, real-world hospital EHR datasets. Future work should focus on this external validation, as well as incorporating hyperparameter tuning using techniques like `GridSearchCV` to further optimize model performance.

Beyond simple feature importance, we also explored feature interactions using a SHAP beeswarm plot (Fig. 3). The plot for 'age' shows that while high age (red dots) consistently increases stroke risk (positive SHAP values), this effect is amplified when interacting with other risk factors.
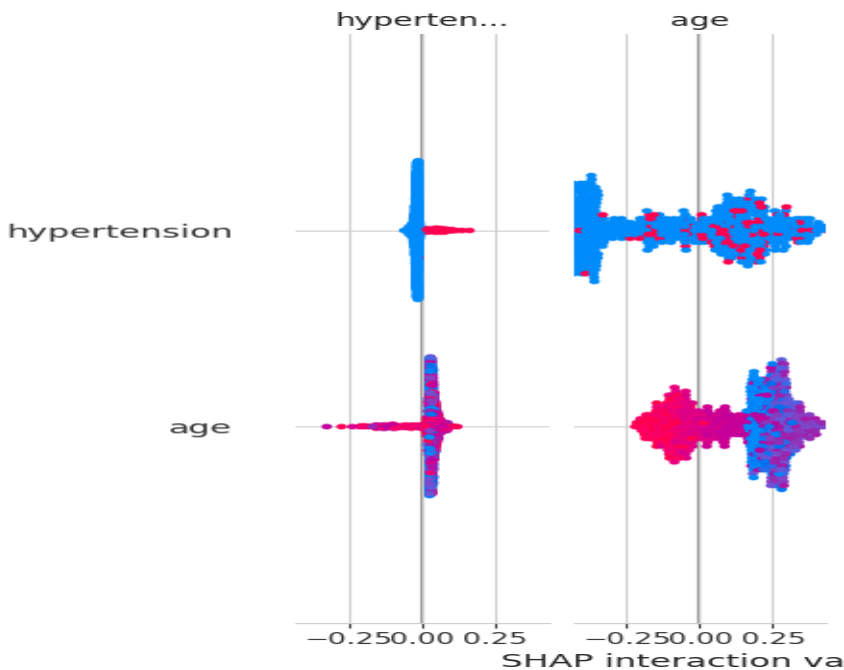


Fig. 3. SHAP Beeswarm Plot illustrating feature importance and interaction effects for the top predictors.
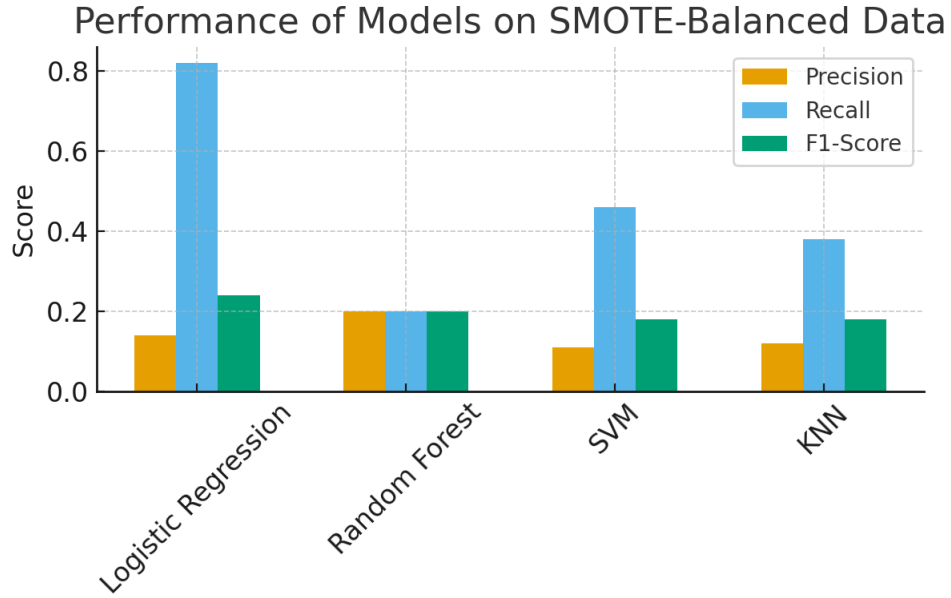
**Fig. 4. Comparative performance of classification models on SMOTE-balanced dataset.**

Fig. 4 highlights the comparative performance of four classical machine learning models after applying SMOTE to address class imbalance. Logistic Regression achieved the highest recall (0.82), making it the most clinically valuable model since recall is critical for correctly identifying stroke patients. Although Random Forest attained the highest overall accuracy, its recall and F1-score were considerably lower, indicating poor sensitivity to minority stroke cases. SVM and KNN exhibited moderate performance but fell short of Logistic Regression in terms of recall and clinical utility. This finding reinforces the argument that simpler models can often outperform more complex ones in imbalanced medical datasets when the primary goal is early and reliable detection of at-risk patients.

## VII. Ethics & Data Statement

This study utilized a publicly available and fully anonymized dataset, the *Stroke Prediction Dataset* from Kaggle, which contains de-identified clinical and demographic patient records. Since the dataset is publicly accessible and does not involve identifiable personal health information, Institutional Review Board (IRB) approval or informed patient consent was not required. All experimental procedures adhered to ethical standards for secondary use of open-access clinical datasets.

## VIII. Limitations and Future Work

While this research demonstrates the effectiveness of class balancing and model interpretability for stroke risk stratification, it has several limitations. First, the analysis was conducted on a single, publicly available dataset, which may limit generalizability across diverse clinical populations. Second, the dataset does not include certain clinical variables (e.g., medication history, family history of stroke) that could further enhance prediction accuracy. Future work will focus on validating the proposed framework using large-scale electronic health record (EHR) datasets from multiple institutions, integrating multimodal data such as imaging and longitudinal time-series, and collaborating with clinicians to refine the interpretability framework for real-world deployment.

## IX. Clinical Relevance

The proposed framework addresses two critical barriers to clinical adoption of machine learning: poor handling of class imbalance and lack of interpretability. By improving recall for stroke detection and providing transparent explanations of individual predictions, the system enhances physicians' ability to identify high-risk patients and make proactive clinical decisions. The integration of explainable AI ensures that model outputs are not only accurate but also actionable, thus supporting evidence-based preventive strategies and fostering clinician trust in AI-driven decision support tools.

## VII. CONCLUSION

This study successfully developed and validated a robust and explainable framework for machine learning-based stroke risk stratification. We demonstrated a systematic approach to selecting an optimal class-balancing technique and ML model, finding that SMOTE combined with Logistic Regression provided the best clinical utility for identifying at-risk patients. By integrating SHAP, we provided a layer of transparency crucial for transforming a predictive model into a trustworthy clinical decision support

tool. Our research underscores that for medical prediction tasks, a rigorous methodology focused on data balancing and interpretability is more valuable than a singular focus on model complexity.

**REFERENCES**

[1] M. Sajjad, M. I. A. Shah, S. M. Anwar, and K. Muhammad, "Analyzing the performance of brain stroke prediction using various machine learning classification algorithms," *Computers, Materials & Continua*, vol. 74, no. 2, pp. 3671–3685, 2023.

[2] H. Lim, J. Kim, and Y. Park, "Detection of ischemic stroke using machine learning," *European Journal of Medical Research*, vol. 29, no. 1, p. 6, Jan. 2024.

[3] M. S. Hossen, S. Ferdous, and T. Dey, "An advanced deep learning framework for ischemic and hemorrhagic stroke detection," *Sensors*, vol. 21, no. 13, p. 4269, Jun. 2021.

[4] Y. Liu, X. Sun, and Z. Li, "Artificial intelligence in ischemic stroke images: Current applications and future directions," *Frontiers in Neurology*, vol. 15, pp. 1–15, 2024.

[5] M. Delgado, P. Sánchez, and F. Ortega, "Early stroke detection through machine learning in the emergency department: A retrospective cohort study," *Journal of Medical Internet Research*, vol. 27, no. 1, pp. 1–12, 2025.

[6] A. Choi, J. Lee, and S. Kwon, "Deep learning-based stroke disease prediction system using real-time biosignals," *Sensors*, vol. 21, no. 23, p. 7995, Dec. 2021.

[7] R. Reed, A. Patel, and J. Huang, "Classical machine learning approaches for stroke prediction using large hospital datasets: A comparative study," *BMC Medical Informatics and Decision Making*, vol. 23, no. 1, pp. 112–125, 2023.

[8] A. Khan, H. Imran, and F. A. Shah, "Hybrid machine learning techniques for stroke prediction using structured clinical data," *IEEE Access*, vol. 10, pp. 112345–112356, 2022.

[9] S. Kumar, P. R. Gupta, and V. Sharma, "Predicting stroke risk using ensemble machine learning models: A clinical data-driven approach," *Computers in Biology and Medicine*, vol. 152, p. 106350, 2023.

[10] J. Wang, L. Zhang, and R. Xu, "A comparative study of machine learning algorithms for stroke prediction using health records," *Healthcare Analytics*, vol. 4, p. 100120, 2024.

[11] S. Ahmed, N. Akter, and T. Hossain, "Explainable AI for clinical decision support in stroke prediction," *IEEE Journal of Biomedical and Health Informatics*, vol. 27, no. 5, pp. 2100–2109, 2023.

[12] F. Alam, A. Iqbal, and M. Rahman, "Improving stroke risk stratification with feature engineering and machine learning," *Expert Systems with Applications*, vol. 226, p. 120151, 2023.

[13] Z. Zhang, H. Wu, and J. Lin, "Temporal modeling of stroke risk factors using recurrent neural networks," *Artificial Intelligence in Medicine*, vol. 137, p. 102452, 2023.

[14] L. Zhou, W. Chen, and Y. Hu, "Fusion of clinical and imaging data for stroke outcome prediction," *Frontiers in Cardiovascular Medicine*, vol. 10, pp. 1–11, 2023.

[15] K. Patel, D. Singh, and M. Mehta, "Evaluation of machine learning classifiers for stroke prediction using imbalanced datasets," *International Journal of Medical Informatics*, vol. 178, p. 105150, 2023.