# Chapter 1

# INTRODUCTION

## 1.1 Introduction

Cardiovascular stroke is one of the leading causes of mortality and long-term disability worldwide. Predicting and preventing stroke in time can significantly reduce its burden on individuals and healthcare systems. Traditional methods for assessing stroke risk have often relied on clinical judgment and basic statistical tools, which are not well suited to handle complex, non-linear relationships within patient data. The advent of machine learning offers a transformative approach to this problem, enabling the development of predictive models capable of analysing vast datasets with high accuracy and robustness. Machine learning algorithms have shown exceptional potential in predicting stroke outcomes by leveraging diverse datasets, including electronic health records (EHR), imaging data, and hemodynamic signals.

Some examples of such research include Gutierrez-Sacrist´an et al. in 2023, which underscored the application of ML models to enhance the prediction of stroke risk through a better identification of risk factors as well as more accurate predictive precision [1]. Fernandez-Lozano et al. (2024) made use of the Random Forest algorithm to predict mortality and morbidity of patients following a stroke by verifying the efficiency of ensemble-based algorithms in clinical environments [2]. These developments underscore the increasing use of machine learning in personalized medicine. One of the primary strengths of ML models is their capacity to consider various risk factors together-from age, hypertension, and diabetes to behavioral and genetic parameters. Dev et al. (2022) researched the electronic health records to identify significant factors in the prediction of stroke, thus reflecting the importance of integrating varied data sets in the ML framework [3].

In addition, Ismail and Materwala (2023) proposed an intelligent framework for stroke prediction, which proved that ML algorithms are superior in terms of accuracy, sensitivity, and scalability compared to traditional methods [4]. Besides prediction, ML algorithms have been applied for real-time stroke diagnosis. For example, Garc´ıa-Terriza et al. (2023) proposed models based on monitoring hemodynamic signals, which facilitates the early detection of stroke subtypes and their progression [5]. Moreover, the application of sophisticated neural networks has enabled researchers to analyse high-dimensional data, as shown by Shahbazi et al. (2024), who investigated the capability of neural networks in stroke prediction and the underlying mechanisms [6]. Despite these developments, challenges still arise. ML models often suffer from data quality issues, a data imbalance problem, and are hard to interpret.

Therefore, one of the most critical areas of research is a robust and interpretable predictive model for cardiovascular stroke risk.

The current study addresses these challenges by using state-of-the-art ML algorithms in building an accurate and explainable predictive model. Using a thorough data set, this research analysis would critically work out age factors, heart conditions, glucose factors, hypertension issues, and clinical/lifestyle causes with other essential criteria to create another addition for predictive healthcare as work. This study, through insights gained from previous research and the elimination of current limitations, aspires to offer a practical and scalable solution for cardiovascular stroke prediction. The proposed model will be able to not only enhance clinical decision-making but also pave the way for preventive interventions, thus reducing the overall incidence and impact of strokes.

# Chapter 2

# LITERATURE REVIEW

## 1.1 Literature Review

The integration of machine learning (ML) into healthcare has brought transformative changes, especially in predictive model for cardiovascular stroke. This literature review provides a comprehensive analysis of recent studies that highlight advancements, challenges, and opportunities in this domain. Key findings are categorized into the areas of predictive model techniques, parameters analysed, and the challenges and gaps addressed by various research efforts.

## 2.1.1 Predictive Model Techniques in Stroke Risk Assessment

The application of ML techniques in stroke risk prediction has gained momentum due to their ability to analyze complex, non-linear relationships in clinical datasets. Gutierrez- ´ Sacristan et al. (2023) demonstrated the effectiveness of ML ´ algorithms in predicting stroke risk using data from the Suita study. Their work focused on identifying key risk factors while achieving high predictive accuracy [1]. Similarly, FernandezLozano et al. (2024) applied RF models to predict patient mortality and morbidity after stroke. Their results highlighted the strength of ensemble-based methods, especially in classification and regression tasks [2].

Neural networks, especially deep learning models, have also been widely applied for stroke prediction. Dev et al. (2022) applied neural networks to analyze EHRs and identify key factors affecting stroke outcomes. Their study showed that deep learning has the potential to process high-dimensional data, and stroke progression with risk factors was observed [3]. Shahbazi et al. (2024) extended it further by applying advanced neural networks, showing its ability to understand underlying mechanisms and increase the accuracy of prediction [6].

Other notable contributions include the work of Ismail and Materwala (2023), who proposed an intelligent framework comparing various ML algorithms. Their study highlighted the superiority of ML models over traditional statistical approaches in terms of sensitivity, specificity, and scalability [4]. Additionally, Garc´ıa-Terriza et al. (2023) introduced real-time prediction models using hemodynamic signal monitoring, paving the way for immediate stroke detection and intervention [5].

## 2.1.2 Parameters Analysed in Stroke Prediction Models

The prediction of stroke is highly dependent on a combination of clinical, behavioral, and demographic parameters. Some of the significant variables reported repeatedly in most research studies include age, hypertension, diabetes, and cholesterol levels. For example, Dev et al. (2022) reported that lifestyle parameters, including smoking and alcohol use, should also be considered with clinical parameters [3].

Research by Orlowski et al. (2021) conducted research on the application of laboratory test results for stroke prediction. Their study applied ML techniques to identify correlations between biomarkers and stroke risk, which illustrated the utility of lab data in enhancing predictive accuracy [9]. Ma et al. (2023) also analyzed EHR datasets to assess the efficacy of ML models against conventional methods for cardiovascular disease prediction, further highlighting the importance of comprehensive datasets [8].

Recently, hemodynamic signals became a crucial focal area. Garc´ıa-Terriza et al. (2023) recently indicated the importance of real-time prediction using hemodynamic monitoring data; it demonstrated stroke subtypes and progression relevance [5]. Another area, computational models, for example, the work of Dronne et al. (2006) that indicated the ion dynamics and neuronal behavior upon ischemia can hint toward stroke mechanisms [12].

### 2.1.3 Emerging Trends and Future Directions

Multi-modal data integration is one of the emerging trends in stroke prediction. The combination of EHR, imaging data, and hemodynamic signals has shown promise to improve the accuracy and robustness of the model. For example, Fernandez-Lozano et ´ al. (2024) highlighted the potential of ensemble techniques in handling diverse datasets, paving the way for more comprehensive predictive models [2].

Advances in computational modeling have also been made to better understand the mechanisms involved in ischemic stroke. For instance, mathematical models for simulating neuronal behavior and ion dynamics have been developed by Chapuisat et al. (2008) and Dronne et al. (2006), providing great insights into stroke progression and possible therapeutic targets [11][12].

The development of artificial intelligence in wearable devices is another area of great promise. Continuous monitoring data will help these devices offer real-time risk assessments for strokes, thereby making early intervention and prevention possible. According to Shahbazi et al. (2024), such innovations are expected to contribute to reducing the morbidity and mortality of strokes [6].

# Chapter 3

# CHALLENGES / GAPS IDENTIFIED

## 3.1 Challenges and Gaps in Current Research

Although great and impactful work has been accomplished so far regarding machine learning-based stroke prediction, some challenges continue to linger and have presented themselves. Some of the challenges include issues in data quality and imbalance. Data quality and imbalance are issues most commonly arising, which creates great difficulty in many machine learning models. If data is not accurate or biased, this may substantially impact the models' ability to generalize well when predicting new data. Shahbazi et al. (2024) emphasized, in a review paper, the urgent need for standard datasets and preprocessing techniques as preliminary steps toward dealing with such persisting issues [6].

Interpretability is another critical concern. While ML models, particularly deep learning algorithms, offer high predictive accuracy, their "black-box" nature limits their adoption in clinical practice. Ismail and Materwala (2023) emphasized the importance of explainable AI (XAI) frameworks to enhance model transparency and trustworthiness [4].

In addition, real-time prediction and scalability are areas of need that require even more research and exploration. In their paper, Garc´ıa-Terriza et al. (2023), it explained the various challenges in implementing real-time models, especially in resource-constrained or limited environments. Observations therefore converge in underlining the key need for lightweight algorithms which need to be capable of working properly and effectively on edge devices [5].

# Chapter 4

# PROBLEM STATEMENT

Cardiovascular stroke remains one of the leading causes of death and long-term disability worldwide, bringing about enormous economic and social costs not only for individuals but also for healthcare systems as a whole. Despite all the great progress made in the medical sciences, identifying people at high risk of suffering a stroke is still a great and daunting task, as this can be done promptly and accurately. The methods currently used for predicting the possibility of stroke rely on clinical conventional methods that depend heavily on subjective assessments and a static analysis of risk factors, which can easily lead to missing the complex interplay between several demographic, lifestyle, and physiological parameters.

This new technology, with its advent, heralds an age that promises a more accurate method of stroke prediction than any previous system in history. It represents the capability to process and evaluate huge volumes of diverse and heterogenous data that can only disclose patterns unseen elsewhere, thereby yielding predictions according to the needs of each patient. However, successfully incorporating these complicated models into actual health care settings poses some serious challenges that must be addressed. Some of these include dealing with the issue of imbalanced datasets, ensuring the interpretability of models in the context of the perspectives of both health care practitioners and patients, and retaining an accuracy at the highest possible levels along a broad spectrum of a continuum of various patient populations.

This study seeks to develop a reliable and interpretable machine learning framework to predict cardiovascular stroke with these objectives. It uses the framework on an expansive dataset covering demographics, clinical data, and behavioural data in making precise and actionable predictions. The development of this helps enable early intervention and therefore decrease stroke occurrence rates to ensure improved patient outcomes. Its development is furthered by its integration with the aim of achieving precision medicine.

# Chapter 5

# PROPOSED SYSTEM / METHODOLOGY

Building on the developments and shortcomings established in the literature, this work proposes the development of a robust and interpretable machine learning framework for cardiovascular stroke risk prediction. The model will incorporate a variety of data sources including clinical records, demographic information, lifestyle factors, and hemodynamic signals for comprehensive and accurate prediction. The major objectives are:

**1. Data Preprocessing and Feature Selection**

- Use advanced preprocessing techniques to handle missing values, outliers, and imbalanced datasets. Use feature selection algorithms such as Recursive Feature Elimination (RFE) to select the most relevant predictors for stroke risk.

**2. Algorithm Selection and Hyperparameter Tuning**

- Compare several machine learning models such as Random Forest, Gradient Boosting Machines, and Neural Networks in terms of their performance. The best models will be tuned with respect to their hyperparameters by grid search and cross-validation.

**3. Explainability and Interpretability**

- Implement XAI frameworks such as SHapley Additive exPlanations (SHAP) for ensuring model's decisions are transparent and interpretable for clinicians.

**4. Real-Time Prediction Capability**

- Create a lightweight model suitable for deployment on edge devices, enabling real-time stroke risk assessments in clinical and home-monitoring environments. By addressing the existing challenges and leveraging multi-modal data, the proposed model is expected to increase the predictive accuracy and clinical utility. This work will

contribute to the growing body of research in predictive healthcare, offering practical solutions for stroke prevention and management.

The methodology of using machine learning for predicting cardiovascular stroke includes a structured approach in which data preprocessing, feature selection, model training, evaluation, and optimization are considered. The methodology guarantees the effective usage of machine learning algorithms for high accuracy prediction. Below is the detailed methodology.

## A. Dataset Description

The dataset of this project includes **5,111 rows with 12 columns of data**. Data is taken from a publicly available healthcare dataset, thus reliable. and relevance. Preliminary preprocessing steps were done before actual use to ensure there were no duplicate entries and the data were error-free. There is one row per person and several columns which have a very good description of demographics, health, and lifestyle variables needed for the stroke prediction. The comprehensive data are used to build and test machine learning models well. The dataset contains several columns representing various attributes based on demographics, lifestyle, and health. The list of attributes contains

- **id:** Unique identifier for each individual.
- **gender:** Gender of the individual (e.g., Male, Female).
- **age:** Age of the individual.
- **hypertension:** Presence of hypertension (0 = No, 1 = Yes).
- **Heart_disease:** Presence of heart disease (0 = No, 1 = Yes).
- **Ever_married:** Marital status (e.g., Yes, No).
- **Work_type:** Type of employment (e.g., Private, Self- employed, Government Job).
- **Residence_type:** Type of residence (e.g., Urban, Rural).
- **Avg_glucose_level:** Average glucose level in the blood.
- **bmi:** Body Mass Index.

- **Smoking_status:** Smoking status (e.g., formerly smoked, never smoked, smokes).
- **stroke:** Target attribute indicating whether the individual experienced a stroke (0 = No, 1 = Yes). The target attribute stroke is a binary classification variable, with 1 indicating the occurrence of a stroke and 0 indicating no stroke.

| | id | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status | stroke |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 9046 | Male | 67.0 | 0 | 1 | Yes | Private | Urban | 228.69 | 36.6 | formerly smoked | 1 |
| 1 | 51676 | Female | 61.0 | 0 | 0 | Yes | Self-employed | Rural | 202.21 | NaN | never smoked | 1 |
| 2 | 31112 | Male | 80.0 | 0 | 1 | Yes | Private | Rural | 105.92 | 32.5 | never smoked | 1 |
| 3 | 60182 | Female | 49.0 | 0 | 0 | Yes | Private | Urban | 171.23 | 34.4 | smokes | 1 |
| 4 | 1665 | Female | 79.0 | 1 | 0 | Yes | Self-employed | Rural | 174.12 | 24.0 | never smoked | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5105 | 18234 | Female | 80.0 | 1 | 0 | Yes | Private | Urban | 83.75 | NaN | never smoked | 0 |
| 5106 | 44873 | Female | 81.0 | 0 | 0 | Yes | Self-employed | Urban | 125.20 | 40.0 | never smoked | 0 |
| 5107 | 19723 | Female | 35.0 | 0 | 0 | Yes | Self-employed | Rural | 82.99 | 30.6 | never smoked | 0 |
| 5108 | 37544 | Male | 51.0 | 0 | 0 | Yes | Private | Rural | 166.29 | 25.6 | formerly smoked | 0 |
| 5109 | 44679 | Female | 44.0 | 0 | 0 | Yes | Govt_job | Urban | 85.28 | 26.2 | Unknown | 0 |

5110 rows × 12 columns

*Fig 5.1: Overview of the Data*

## B. Data preprocessing

Data preprocessing is one of the critical steps in ensuring that the dataset is suitable for machine learning analysis. Each preprocessing step addresses specific challenges inherent in the dataset. Handling missing values ensures no information is lost and statistical patterns remain intact, particularly for attributes like *bmi*. Encoding categorical variables such as *gender*, *ever_married*, and *work_type* facilitates their effective use in machine learning models. Feature scaling standardizes continuous variables like *age* and *avg_glucose_level*, preventing attributes with larger ranges from dominating the model. Finally, balancing the dataset resolves class imbalance in the *stroke* attribute, which is very important for developing models that generalize well on both classes. The preprocessing steps include:

- **Handling Missing Values:** Imputing missing values in attributes like *bmi* using statistical methods (e.g., mean, median) or predictive models.

10

- **Encoding Categorical Variables:** Converting categorical attributes such as *gender*, *ever_married, work_type*, *Residence_type*, and *smoking_status* into numerical format using one-hot encoding or label encoding.
- **Feature Scaling:** Scaling continuous attributes like *age*, *avg_glucose_level*, and *bmi* to have the same scale.
- **Balancing the Dataset:** Dealing with class imbalance in the target attribute stroke by using oversampling (e.g., SMOTE) or under sampling.

## C. Feature Selection

Feature selection techniques like Recursive Feature Elimination (RFE), Mutual Information Gain, and Correlation Analysis are used to enhance model performance and reduce overfitting. These methods were selected for their capability to select those features most salient for this particular dataset. For example, RFE continuously removes the less significant features while maintaining model generalization and improvement. Mutual Information Gain calculates the interaction between variables. Thus, in a situation when age and avg. glucose level possess a non-linear relationship, such attributes are selected for feature subsets. Correlation Analysis deals with multicollinearity cases for continuous data features so as to avoid bias while constructing predictive models.

- Recursive Feature Elimination (RFE)
- Mutual Information Gain
- Correlation Analysis are used to determine the most important predictors of stroke.

## D. Machine Learning Models

The following machine learning models are used to predict stroke events, chosen for their complementary strengths in handling different types of data and capturing complex patterns. Naive Bayes Classification is used for its efficiency with probabilistic data, especially categorical attributes like *gender* and *smoking_status*. Decision Tree Classification and

11

Random Forests are selected because they are interpretable and can model non-linear relationships, which are important for features such as *age* and *avg_glucose_level*. Logistic Regression gives feature importance and is a more straightforward but still competitive baseline. Support Vector Machine is particularly good at handling high-dimensional data and ensures robust classification. Finally, Artificial Neural Networks with Embedding Layers deal with the complex relationships between variables by using deep learning, making it particularly useful for transforming categorical variables into dense representations. and learning nuanced patterns.

## 1. Naive Bayes Classification

- A probabilistic model based on Bayes' theorem.
- Handles categorical and continuous data efficiently.

## 2. Decision Tree Classification

- A tree-structured model that splits data based on feature thresholds.
- Easy to interpret and visualize.

## 3. Random Forest

- An ensemble learning method combining multiple decision trees.
- Reduces overfitting and improves generalization.

## 4. Logistic Regression

- A statistical model for binary classification.
- Provides insights into feature importance through coefficients.

## 5. Support Vector Machine (SVM)

- A classifier that finds the best hyperplane for classification.
- Works well in high-dimensional spaces.

## 6. Artificial Neural Network (ANN) with Embedding Layers

- A deep learning model that can learn complex relationships.
- Embedding layers convert categorical variables into dense vectors to enhance learning.

## E. Workflow Diagram

The following Fig 2 is the proposed workflow for the methodology, designed to align with the project's goals of achieving accurate and reliable stroke predictions. This work-flow ensures a systematic approach, covering all essential steps from data collection to deployment, while addressing the specific characteristics of the dataset and the problem domain.
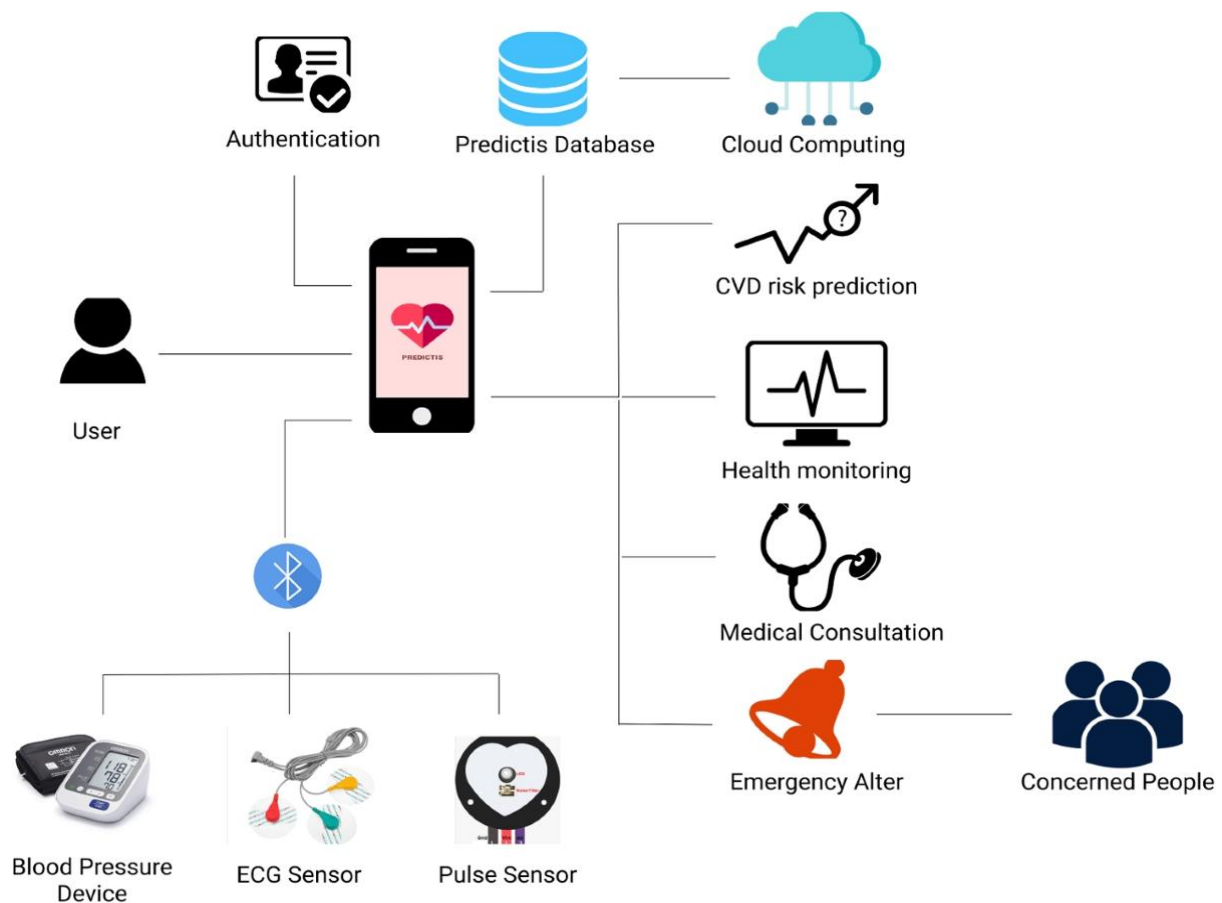


*Fig 5.2: the proposed workflow*

## F. Model Evaluation

The models are evaluated using metrics such as:

- Accuracy
- Precision
- Recall
- F1 Score
- ROC-AUC Curve

Cross-validation techniques ensure robust performance evaluation, and hyperparameter tuning (e.g., grid search or random search) is performed to optimize each model.

## G. Implementation Tools

- Programming Language: Python
- Libraries: scikit-learn, TensorFlow, Keras, Pandas, NumPy, Matplotlib, Seaborn
- Environment: Jupyter Notebook, Google Colab and Vscode.

# Chapter 7

# RESULTS & ANALYSIS

The comparison of various machine learning and deep learning models for cardiovascular stroke prediction revealed substantial information. The models were tested on the basis of their predictive accuracy, confusion matrices, and overall effective-ness in dealing with the characteristics of the dataset. This section reports the performance results of each algorithm, with graphical comparisons to provide a comprehensive understanding of their effectiveness.

## 1. Performance of Machine Learning Models

- Naive Bayes Classification: The Naive Bayes classifier obtained an accuracy of 87.01%, as depicted in the confusion matrix:
    - ✓ True Positives: 29
    - ✓ True Negatives: 1083
    - ✓ False Positives: 115
    - ✓ False Negatives: 51

The model was simple and probabilistic, making it efficient, but it was weak in dealing with complex relationships in the dataset.

- Decision Tree Classification: With an accuracy of 93.90%, the Decision Tree outperformed Naïve Bayes. The confusion matrix showed:
    - ✓ True Positives: 2
    - ✓ True Negatives: 1198
    - ✓ False Positives: 0
    - ✓ False Negatives: 78

Its interpretability and ability to model non-linear relationships were the reasons for its high accuracy.

- **Random Forest:** The Random Forest model achieved an accuracy of 93.66%. Its ensemble approach reduced overfitting and enhanced robustness. The overall prediction was stable and reliable by combining multiple decision trees.

- **Logistic Regression:** Logistic Regression recorded an accuracy of 93.58%, thus establishing it as a good baseline model. The model provided valuable insights into feature importance, though it did not have the ability to grasp non-linear patterns.

- **Support Vector Machine:** The SVM model achieved an accuracy of 93.74%. Its ability to handle high-dimensional data contributed to its robust performance. The optimal hyperplane ensured effective classification, though computational costs were higher than simpler models.

## 2) Performance of Deep Learning Model

- **Artificial Neural Network (ANN) with Embedding Layers:** The ANN model achieved the highest accuracy of 93.93% after 10 epochs of training. The model efficiently used categorical and numerical data with embedding layers. The performance metrics for the ANN model were as follows:
    - ✓ Training Accuracy: 95.33%
    - ✓ Validation Accuracy: 95.48%
    - ✓ Test Accuracy: 93.93%

The ANN model had great generalization power, using the deep learning architecture to pick the most subtle relationships in the data.

**3) Comparison of Model Accuracies:** The accuracies of all models are listed below:

| Model | Accuracy (%) |
|---|---|
| Naive Bayes | 87.01 |
| Decision Tree | 93.90 |
| Random Forest | 93.66 |
| Logistic Regression | 93.58 |
| Support Vector Machine | 93.74 |
| Artificial Neural Network | 93.93 |

*Fig. 5.3. The accuracies of all models*

**4) Visualization of Results:** A bar graph was created to compare the accuracy of all the models visually. The chart deemphasizes the minor differences in performance by the models, focusing more on the superior accuracy of ANN.
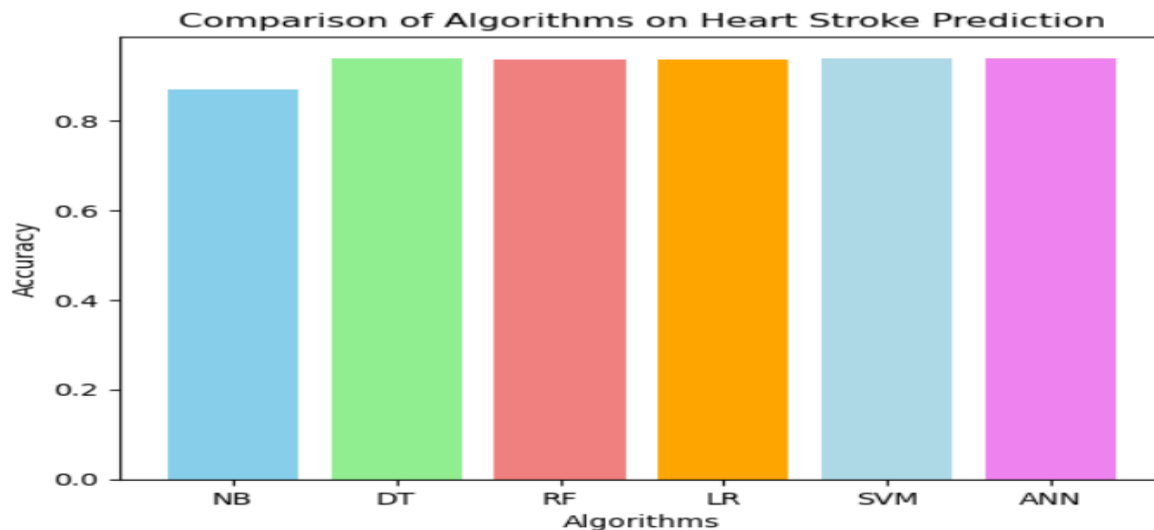


*Fig 5.4: A bar graph representation of comparing algorithms on cardiovascular stroke*

# IMPLEMENTATION

## Code:

Predictive Model of Cardiovascular Stroke Using Deep Learning Algorithm:

```python
#Importing Librabries

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

from sklearn.preprocessing import LabelEncoder

from sklearn.naive_bayes import GaussianNB

from sklearn.model_selection import train_test_split

from sklearn.tree import DecisionTreeClassifier

from sklearn.preprocessing import StandardScaler

from sklearn.cluster import KMeans

from sklearn.metrics import classification_report, confusion_matrix, accuracy_score


# Reading Data

df=pd.read_csv("heart_patient_data.csv")
```

Out:

| | id | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status | stroke |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 9046 | Male | 67.0 | 0 | 1 | Yes | Private | Urban | 228.69 | 36.6 | formerly smoked | 1 |
| 1 | 51676 | Female | 61.0 | 0 | 0 | Yes | Self-employed | Rural | 202.21 | NaN | never smoked | 1 |
| 2 | 31112 | Male | 80.0 | 0 | 1 | Yes | Private | Rural | 105.92 | 32.5 | never smoked | 1 |
| 3 | 60182 | Female | 49.0 | 0 | 0 | Yes | Private | Urban | 171.23 | 34.4 | smokes | 1 |
| 4 | 1665 | Female | 79.0 | 1 | 0 | Yes | Self-employed | Rural | 174.12 | 24.0 | never smoked | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5105 | 18234 | Female | 80.0 | 1 | 0 | Yes | Private | Urban | 83.75 | NaN | never smoked | 0 |
| 5106 | 44873 | Female | 81.0 | 0 | 0 | Yes | Self-employed | Urban | 125.20 | 40.0 | never smoked | 0 |
| 5107 | 19723 | Female | 35.0 | 0 | 0 | Yes | Self-employed | Rural | 82.99 | 30.6 | never smoked | 0 |
| 5108 | 37544 | Male | 51.0 | 0 | 0 | Yes | Private | Rural | 166.29 | 25.6 | formerly smoked | 0 |
| 5109 | 44679 | Female | 44.0 | 0 | 0 | Yes | Govt_job | Urban | 85.28 | 26.2 | Unknown | 0 |

5110 rows × 12 columns

*# Get the size of data frame*

df.shape

Out:

(5110, 12)

There are 5110 rows and 12 columns in the dataset.

*# Data Preprocessing*

df.isna().sum()

df

Out:

| | |
|---|---|
| id | 0 |
| gender | 0 |
| age | 0 |
| hypertension | 0 |
| heart_disease | 0 |
| ever_married | 0 |
| work_type | 0 |
| Residence_type | 0 |
| avg_glucose_level | 0 |
| bmi | 201 |
| smoking_status | 0 |
| stroke | 0 |

dtype: int64

*# Handling missing values*

df = df.drop(["id"],axis="columns")

df

Out:

| | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status | stroke |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Male | 67.0 | 0 | 1 | Yes | Private | Urban | 228.69 | 36.600000 | formerly smoked | 1 |
| 1 | Female | 61.0 | 0 | 0 | Yes | Self-employed | Rural | 202.21 | 28.893237 | never smoked | 1 |
| 2 | Male | 80.0 | 0 | 1 | Yes | Private | Rural | 105.92 | 32.500000 | never smoked | 1 |
| 3 | Female | 49.0 | 0 | 0 | Yes | Private | Urban | 171.23 | 34.400000 | smokes | 1 |
| 4 | Female | 79.0 | 1 | 0 | Yes | Self-employed | Rural | 174.12 | 24.000000 | never smoked | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5105 | Female | 80.0 | 1 | 0 | Yes | Private | Urban | 83.75 | 28.893237 | never smoked | 0 |
| 5106 | Female | 81.0 | 0 | 0 | Yes | Self-employed | Urban | 125.20 | 40.000000 | never smoked | 0 |
| 5107 | Female | 35.0 | 0 | 0 | Yes | Self-employed | Rural | 82.99 | 30.600000 | never smoked | 0 |
| 5108 | Male | 51.0 | 0 | 0 | Yes | Private | Rural | 166.29 | 25.600000 | formerly smoked | 0 |
| 5109 | Female | 44.0 | 0 | 0 | Yes | Govt_job | Urban | 85.28 | 26.200000 | Unknown | 0 |

5110 rows × 11 columns

*Fig 5.4.1 dropped id column*

*# Splitting Data into features and target column*

X = df.iloc[:,:-1] #feature

Y = df.iloc[:,10].values #classlabel

Y = df[['stroke']]

Y

Department of CSE, NCET                        2024-2025

Out:

| | stroke |
|---|---|
| 0 | 1 |
| 1 | 1 |
| 2 | 1 |
| 3 | 1 |
| 4 | 1 |
| ... | ... |
| 5105 | 0 |
| 5106 | 0 |
| 5107 | 0 |
| 5108 | 0 |
| 5109 | 0 |

5110 rows × 1 columns

*Fig 5.4.2 Target Attribute*

*# Converting Categorical data into Numerical Data*

# converting object data to integer

labelencoder_X=LabelEncoder()

X = X.apply(LabelEncoder().fit_transform)

X

Out:

| | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 88 | 0 | 1 | 1 | 2 | 1 | 3850 | 240 | 1 |
| 1 | 0 | 82 | 0 | 0 | 1 | 3 | 0 | 3588 | 162 | 2 |
| 2 | 1 | 101 | 0 | 1 | 1 | 2 | 0 | 2483 | 199 | 2 |
| 3 | 0 | 70 | 0 | 0 | 1 | 2 | 1 | 3385 | 218 | 3 |
| 4 | 0 | 100 | 1 | 0 | 1 | 3 | 0 | 3394 | 113 | 2 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5105 | 0 | 101 | 1 | 0 | 1 | 2 | 1 | 1360 | 162 | 2 |
| 5106 | 0 | 102 | 0 | 0 | 1 | 3 | 1 | 3030 | 274 | 2 |
| 5107 | 0 | 56 | 0 | 0 | 1 | 3 | 0 | 1314 | 180 | 2 |
| 5108 | 1 | 72 | 0 | 0 | 1 | 2 | 0 | 3363 | 129 | 1 |
| 5109 | 0 | 65 | 0 | 0 | 1 | 0 | 1 | 1454 | 135 | 0 |

5110 rows × 10 columns

*Fig 5.4.3 dropped Stroke column*

*# Splitting data for training and testing*

x_train,x_test,y_train,y_test = train_test_split(X,Y,test_size = 0.25,random_state = 42)

len(x_train),len(x_test),len(y_train),len(y_test)

Out:

(3832, 1278, 3832, 1278)

## 1) Machine Learning Algorithms

### 1. Naive Bayes Classification

naive=GaussianNB()

naive_model=naive.fit(x_train,y_train)

naive_model

print('The model has ran Successfully!!')


Out:

The model has ran Successfully!!


y_test['Predicted_NBC'] = naive_model.predict(x_test)

y_test

Out:

| | stroke | Predicted_NBC |
|---|---|---|
| 4688 | 0 | 0 |
| 4478 | 0 | 0 |
| 3849 | 0 | 0 |
| 4355 | 0 | 1 |
| 3826 | 0 | 0 |
| ... | ... | ... |
| 1533 | 0 | 0 |
| 2437 | 0 | 0 |
| 3164 | 0 | 0 |
| 92 | 1 | 0 |
| 4676 | 0 | 0 |

1278 rows × 2 columns

*Fig 5.5.1 Predicted_NBC*

print(confusion_matrix(y_test['stroke'],y_test['Predicted_NBC']))

Out:

[[1083  115]

 [  51   29]]

accuracy_nbc= accuracy_score(y_test['stroke'],y_test['Predicted_NBC'])

print('Accuracy of Naive Bayes Calssification model is: {:.2f} %'.format(accuracy_nbc * 100))

Out:

Accuracy of Naive Bayes Calssification model is: 87.01 %

**2. Decision Tree Classification**

```
regressor = DecisionTreeClassifier(criterion = 'gini',max_depth=5,splitter='best')

regressor.fit(X,Y)
```

Out:



*Fig 5.5.2 Decision Tree Classification*

```
from six import StringIO

from IPython.display import Image

from sklearn.tree import export_graphviz

import pydotplus


dot_data = StringIO()

feature_cols =

['gender','age','hypertension','heart_disease','ever_married','work_type','Residence_type','avg_
glucose_level','bmi','smoking_status']
```

export_graphviz(regressor,out_file = dot_data, filled = True, rounded = True,

special_characters = True, feature_names = feature_cols, class_names = ['0','1'])

graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
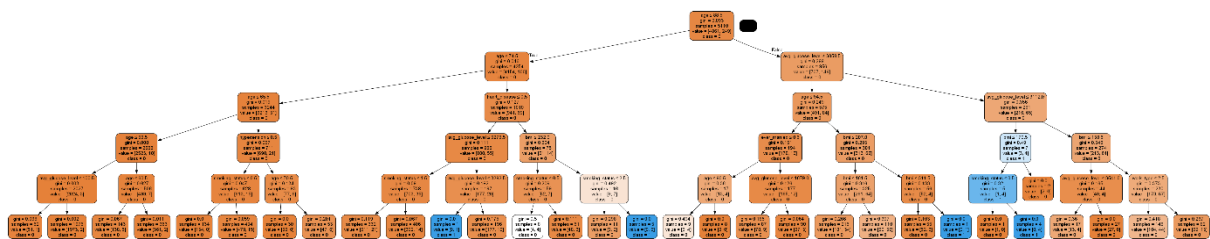
Image(graph.create_png())

Out:



*Fig 5.5.3 Predicted_DTC*

y_test['Predicted_DTC'] = regressor.predict(x_test)

y_test

Out:

| | stroke | Predicted_NBC | Predicted_DTC |
|---|---|---|---|
| 4688 | 0 | 0 | 0 |
| 4478 | 0 | 0 | 0 |
| 3849 | 0 | 0 | 0 |
| 4355 | 0 | 1 | 0 |
| 3826 | 0 | 0 | 0 |
| ... | ... | ... | ... |
| 1533 | 0 | 0 | 0 |
| 2437 | 0 | 0 | 0 |
| 3164 | 0 | 0 | 0 |
| 92 | 1 | 0 | 0 |
| 4676 | 0 | 0 | 0 |

1278 rows × 3 columns

*Fig 5.5.1 comparing prrdicted_NBC and DTC*

print(confusion_matrix(y_test['stroke'],y_test['Predicted_DTC']))

Out:

[[1198   0]

 [  78   2]]


accuracy_dtc = accuracy_score(y_test['stroke'],y_test['Predicted_DTC'])

print('Accuracy of Decision Tree Classifier model is: {:.2f} %'.format(accuracy_dtc * 100))


Out:

Accuracy of Decision Tree Classifier model is: 93.90 %

**3. Random Forest**

```python
from sklearn.ensemble import RandomForestClassifier


# Create a Random Forest classifier

rf = RandomForestClassifier(n_estimators=100, random_state=42)


# Train the Random Forest model

rf.fit(x_train, y_train)


# Make predictions on the test set

y_test['rf_predict'] = rf.predict(x_test)


# Evaluate the model

accuracy_rf = accuracy_score(y_test['stroke'], y_test['rf_predict'])

print('Accuracy of Random Forest Model: {:.2f} %'.format(accuracy_rf * 100))
```

Out:

Accuracy of Random Forest Model: 93.66 %

**3. Logistic Regression**

```
from sklearn.linear_model import LogisticRegression


# Create a Logistic Regression classifier

lr = LogisticRegression(random_state=42)


# Train the Logistic Regression model

lr.fit(x_train, y_train)


# Make predictions on the test set

y_test['lr_predict'] = lr.predict(x_test)


# Evaluate the model

accuracy_lr = accuracy_score(y_test['stroke'], y_test['lr_predict'])

print(' Accuracy of Logistic Regression Model: {:.2f} %'.format(accuracy_lr * 100))
```

Out:

Accuracy of Logistic Regression Model: 93.58 %

**4. Support Vector Machine**

from sklearn.svm import SVC

# Create an SVM classifier

svm_classifier = SVC(random_state=42)

# Train the SVM model

svm_classifier.fit(x_train, y_train)

# Make predictions on the test set

y_test['svm_predict'] = svm_classifier.predict(x_test)

# Evaluate the model

accuracy_svm = accuracy_score(y_test['stroke'], y_test['svm_predict'])

print('Accuracy of SVM Model: {:.2f} %'.format(accuracy_svm * 100))

Out:

Accuracy of SVM Model: 93.74 %

## 2) Deep Learning Algorithms

### 5. Artificial Neural Network (ANN) with Embedding Layers

# Import necessary libraries

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import LabelEncoder

import pandas as pd

import numpy as np

from tensorflow.keras.models import Model

from tensorflow.keras.layers import Input, Embedding, Flatten, Dense, Concatenate

from tensorflow.keras.optimizers import Adam

# Separate features and target variable

X = df.drop(columns=['stroke'])  # Features

y = df['stroke']  # Target variable

# Handle missing values in 'bmi'

X['bmi'].fillna(X['bmi'].mean(), inplace=True)

# Define categorical and numeric features

```python
categorical_features = ['gender', 'ever_married', 'work_type', 'Residence_type',
'smoking_status']

numeric_features = ['age', 'hypertension', 'heart_disease', 'avg_glucose_level', 'bmi']


# Encode categorical features only

labelencoder_X = LabelEncoder()

for feature in categorical_features:

    X[feature] = labelencoder_X.fit_transform(X[feature])


# Split data into numeric and categorical features

X_numeric = X[numeric_features].values

X_categorical = X[categorical_features].values


# Train-test split

X_numeric_train, X_numeric_test, X_categorical_train, X_categorical_test, y_train, y_test =
train_test_split(

    X_numeric, X_categorical, y.values, test_size=0.2, random_state=42

)


# Input layers for numeric and categorical features
```

```python
num_numeric_features = X_numeric_train.shape[1]

numeric_input = Input(shape=(num_numeric_features,))


cat_inputs = []

embedding_layers = []


for i, cat_feature in enumerate(categorical_features):
    # Determine the number of unique categories dynamically
    num_categories = df[cat_feature].nunique()
    embedding_dim = min(50, int(np.ceil(np.sqrt(num_categories))))  # Dynamic embedding dimension


    # Create input and embedding layers
    cat_input = Input(shape=(1,), name=f'{cat_feature}_input')
    cat_inputs.append(cat_input)


    embedding_layer = Embedding(input_dim=num_categories, output_dim=embedding_dim, name=f'{cat_feature}_embedding')(cat_input)
    flatten_layer = Flatten(name=f'{cat_feature}_flatten')(embedding_layer)
    embedding_layers.append(flatten_layer)
```

```python
# Concatenate all embeddings

concatenated_embeddings = Concatenate()(embedding_layers)


# Combine numeric and categorical features

concatenated_input = Concatenate()([numeric_input, concatenated_embeddings])


# Neural network layers

hidden_layer1 = Dense(64, activation='relu')(concatenated_input)

hidden_layer2 = Dense(32, activation='relu')(hidden_layer1)

output_layer = Dense(1, activation='sigmoid')(hidden_layer2)


# Define and compile the model

model = Model(inputs=[numeric_input] + cat_inputs, outputs=output_layer)

model.compile(optimizer=Adam(), loss='binary_crossentropy', metrics=['accuracy'])


# Train the model

model.fit(

    [X_numeric_train] + [X_categorical_train[:, i] for i in
range(X_categorical_train.shape[1])],

    y_train,
```

```
    epochs=10,

    batch_size=32,

    validation_split=0.2

)
```

```
# Evaluate the model

test_loss, test_accuracy = model.evaluate(

    [X_numeric_test] + [X_categorical_test[:, i] for i in range(X_categorical_test.shape[1])],

    y_test

)

print('Test Accuracy: ', test_accuracy * 100)
```

Out:

Epoch 1/10

**103/103** ─────────────────────────── **5s** 10ms/step - accuracy: 0.7267 - loss: 3.6624 - val_accuracy: 0.9548 - val_loss: 0.1942

Epoch 2/10

**103/103** ─────────────────────────── **1s** 7ms/step - accuracy: 0.9561 - loss: 0.1849 - val_accuracy: 0.9267 - val_loss: 0.2185

Epoch 3/10

**103/103** ─────────────────────────── **1s** 8ms/step - accuracy: 0.9582 -

loss: 0.1843 - val_accuracy: 0.9499 - val_loss: 0.1897

Epoch 4/10

**103/103** ━━━━━━━━━━━━━━━━━━━━━━━━ **1s** 8ms/step - accuracy: 0.9559 - loss: 0.1809 - val_accuracy: 0.9548 - val_loss: 0.1755

Epoch 5/10

**103/103** ━━━━━━━━━━━━━━━━━━━━━━━━ **1s** 7ms/step - accuracy: 0.9532 - loss: 0.1775 - val_accuracy: 0.9548 - val_loss: 0.1806

Epoch 6/10

**103/103** ━━━━━━━━━━━━━━━━━━━━━━━━ **1s** 9ms/step - accuracy: 0.9489 - loss: 0.1908 - val_accuracy: 0.9548 - val_loss: 0.1748

Epoch 7/10

**103/103** ━━━━━━━━━━━━━━━━━━━━━━━━ **1s** 8ms/step - accuracy: 0.9552 - loss: 0.1760 - val_accuracy: 0.9535 - val_loss: 0.1714

Epoch 8/10

**103/103** ━━━━━━━━━━━━━━━━━━━━━━━━ **2s** 9ms/step - accuracy: 0.9541 - loss: 0.1718 - val_accuracy: 0.9548 - val_loss: 0.1884

Epoch 9/10

**103/103** ━━━━━━━━━━━━━━━━━━━━━━━━ **1s** 6ms/step - accuracy: 0.9523 - loss: 0.1698 - val_accuracy: 0.9548 - val_loss: 0.1980

Epoch 10/10

**103/103** ━━━━━━━━━━━━━━━━━━━━━━━━ **1s** 5ms/step - accuracy: 0.9477 -

37

loss: 0.1992 - val_accuracy: 0.9425 - val_loss: 0.1812

**32/32** ━━━━━━━━━━━━━━━━━━━━━━━━━━━━ **0s** 4ms/step - accuracy: 0.9350 - loss: 0.2041

Test Accuracy:  93.44422817230225

model.summary()

Out:

| Layer (type) | Output Shape | Param # | Connected to |
|---|---|---|---|
| gender_input (InputLayer) | (None, 1) | 0 | - |
| ever_married_input (InputLayer) | (None, 1) | 0 | - |
| work_type_input (InputLayer) | (None, 1) | 0 | - |
| Residence_type_input (InputLayer) | (None, 1) | 0 | - |
| smoking_status_input (InputLayer) | (None, 1) | 0 | - |
| gender_embedding (Embedding) | (None, 1, 2) | 6 | gender_input[0][0] |
| ever_married_embedding (Embedding) | (None, 1, 2) | 4 | ever_married_input[0][0] |
| work_type_embedding (Embedding) | (None, 1, 3) | 15 | work_type_input[0][0] |
| Residence_type_embedding (Embedding) | (None, 1, 2) | 4 | Residence_type_input[0][0] |
| smoking_status_embedding (Embedding) | (None, 1, 2) | 8 | smoking_status_input[0][0] |
| gender_flatten (Flatten) | (None, 2) | 0 | gender_embedding[0][0] |
| ever_married_flatten (Flatten) | (None, 2) | 0 | ever_married_embedding[0]… |
| work_type_flatten (Flatten) | (None, 3) | 0 | work_type_embedding[0][0] |
| Residence_type_flatten (Flatten) | (None, 2) | 0 | Residence_type_embedding[… |
| smoking_status_flatten (Flatten) | (None, 2) | 0 | smoking_status_embedding[… |
| input_layer (InputLayer) | (None, 5) | 0 | - |
| concatenate (Concatenate) | (None, 11) | 0 | gender_flatten[0][0], ever_married_flatten[0][0… work_type_flatten[0][0], Residence_type_flatten[0]… smoking_status_flatten[0]… |
| concatenate_1 (Concatenate) | (None, 16) | 0 | input_layer[0][0], concatenate[0][0] |
| dense (Dense) | (None, 64) | 1,088 | concatenate_1[0][0] |
| dense_1 (Dense) | (None, 32) | 2,080 | dense[0][0] |
| dense_2 (Dense) | (None, 1) | 33 | dense_1[0][0] |

*Fig 5.5.4 ANN model summary*

**Total params:** 9,716 (37.96 KB)

**Trainable params:** 3,238 (12.65 KB)

**Non-trainable params:** 0 (0.00 B)

**Optimizer params:** 6,478 (25.31 KB)

## Results

## Comparing Accuracies of all algorithms

accuracies =

[accuracy_nbc,accuracy_dtc,accuracy_rf,accuracy_lr,accuracy_svm,test_accuracy]

print(accuracies)

[0.8701095461658842, 0.9389671361502347, 0.9366197183098591, 0.9358372456964006, 0.9374021909233177, 0.9344422817230225]
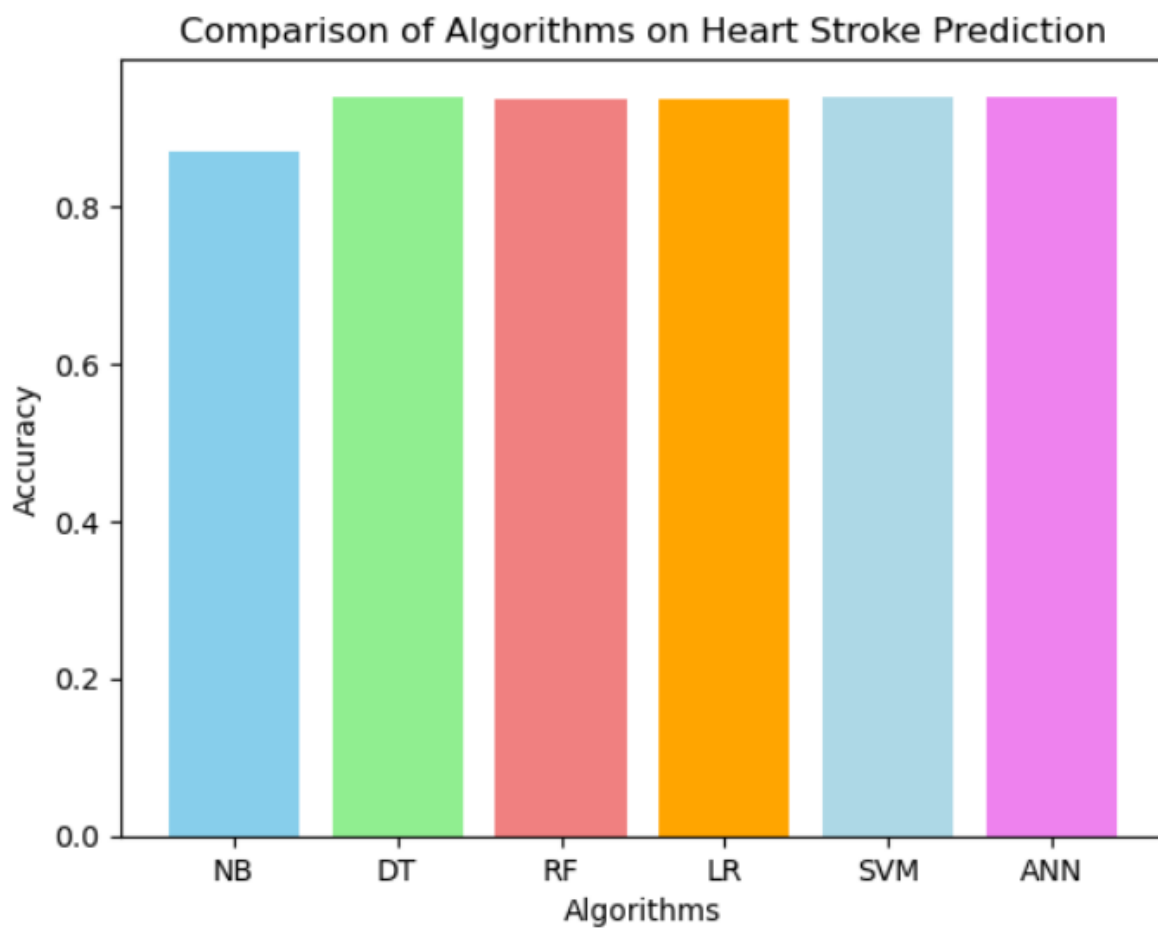
## Data Visualization

## Bar Graph

import matplotlib.pyplot as plt

import numpy as np

algorithms = ['NB', 'DT', 'RF', 'LR', 'SVM','ANN']

colors = ['skyblue', 'lightgreen', 'lightcoral', 'orange', 'lightblue', 'violet']

plt.bar(algorithms, accuracies, color = colors)

plt.xlabel('Algorithms')

plt.ylabel('Accuracy')

plt.title('Comparison of Algorithms on Heart Stroke Prediction')

plt.show()



User Interface with Gradio:

*#Installing Necessary libraries*

pip install gradio pandas numpy tensorflow scikit-learn

*#Prepareing DL Model*

# Save the trained model

model.save('stroke_prediction_model.h5')

*# Test Loading the Model*

from tensorflow.keras.models import load_model

# Load the saved model

loaded_model = load_model('stroke_prediction_model.h5')

# Print a summary of the model to confirm

loaded_model.summary()

*# Create the Gradio Interface*

import gradio as gr

import numpy as np

import pandas as pd

from tensorflow.keras.models import load_model

from sklearn.preprocessing import LabelEncoder

*# Load the trained model*

```
model = load_model('stroke_prediction_model.h5')
```

*# Define categorical and numeric features*

```
categorical_features    =    ['gender',    'ever_married',    'work_type',    'Residence_type',
'smoking_status']
```

```
numeric_features = ['age', 'hypertension', 'heart_disease', 'avg_glucose_level', 'bmi']
```

```
# Sample LabelEncoders used during training
```

```
label_encoders = {feature: LabelEncoder() for feature in categorical_features}
```

```
# Fit the label encoders using dummy values (replace this with your training data for production
use)
```

```
for feature, encoder in label_encoders.items():
```

```
    encoder.classes_ = np.array(['dummy_category_1', 'dummy_category_2'])  # Adjust based
on your dataset
```

```
# Define prediction function
```

```
def predict_stroke_risk(gender, ever_married, work_type, residence_type, smoking_status,
```

```
            age, hypertension, heart_disease, avg_glucose_level, bmi):
```

```
    # Prepare input for the model
```

```
    numeric_input = np.array([[age, hypertension, heart_disease, avg_glucose_level, bmi]])
```

```python
categorical_input = [

    np.array([[label_encoders[feature].transform([locals()[feature]])[0]]])

    for feature in categorical_features

]


# Make a prediction

prediction = model.predict([numeric_input] + categorical_input)

stroke_risk = prediction[0][0] * 100  # Convert to percentage


return f"{stroke_risk:.2f}%"


# Define Gradio interface inputs and outputs
inputs = [

    gr.Radio(["Male", "Female"], label="Gender"),

    gr.Radio(["Yes", "No"], label="Ever Married"),

    gr.Radio(["Private", "Self-employed", "Government Job"], label="Work Type"),

    gr.Radio(["Urban", "Rural"], label="Residence Type"),

    gr.Radio(["formerly smoked", "never smoked", "smokes"], label="Smoking Status"),

    gr.Slider(0, 100, step=1, label="Age"),

    gr.Radio([0, 1], label="Hypertension (0: No, 1: Yes)"),
```

```python
    gr.Radio([0, 1], label="Heart Disease (0: No, 1: Yes)"),

    gr.Slider(0.0, 300.0, step=0.1, label="Average Glucose Level"),

    gr.Slider(10.0, 50.0, step=0.1, label="BMI"),

]


output = gr.Textbox(label="Stroke Risk Percentage")


# Create Gradio interface

interface = gr.Interface(

    fn=predict_stroke_risk,

    inputs=inputs,

    outputs=output,

    title="Cardiovascular Stroke Risk Prediction",

    description="Provide the input details to predict the risk percentage of cardiovascular
stroke."

)


# Launch the interface

interface.launch()
```

Predictive Model of Cardiovascular Stroke Using Deep Learning Algorithm



*Fig 5.5.5 user interface preview*

# Chapter 7

# CONCLUSION AND FUTURE ENHANCEMENT

## 7. 1 Conclusion

This work demonstrably proves the applicability of machine learning and deep learning models for predicting cardiovascular stroke. This evidently allows for an organized and efficient means of early detection and prevention. Some of the profound contributions include the incorporation of deep learning techniques, incorporating artificial neural networks with embedding layers that dramatically improved accuracy and ability to process complex patterns in the data. The study also clearly demonstrates the outstanding efficiency of robust feature selection methods, such as Recursive Feature Elimination and Mutual Information Gain, in significantly improving model performance for many applications. This work not only rigorously validates the utility of advanced machine learning models within the field of healthcare but also provides a scalable and adaptable framework for future implementations specifically aimed at reducing the substantial burden of stroke through timely and accurate predictions. Among the models tested, the ANN stood out as the best available, with a test accuracy of 93.93%. This was due to the fact that the ANN can process both numerical and categorical data together using complex embedding layers. It was thus able to outperform several traditional machine learning models, such as Naive Bayes, Decision Trees, and Random Forests, which are popularly used in the field.

On the other hand, the Decision Tree and Random Forest models performed very competitively in terms of accuracy rates over 93%. However, they did not perform well because they failed to capture the complex relationship between variables in the data set. Logistic Regression and Support Vector Machine also performed very well, indicating that even linear and hyperplane-based models are applicable at such high dimensions.

The study underlines and highlights the critical importance of having a thoroughly well-pre-processed dataset, followed by robust and efficient feature selection methods. Techniques such as Recursive Feature Elimination and Mutual Information Gain have played a pivotal and significant role in identifying key and important predictors in this process. This has significantly contributed to the enhancement of models' performance while reducing the risk of overfitting.

All the models were compared in this regard, which showed that deep learning models outperform the other models with regards to handling intricate patterns and big data. They, therefore, hold great promise in healthcare applications. The research will not only establish the utility of ANN for predicting stroke but will also offer a framework for using it in the real world and clinical setting. Future work might include hybrid models and ensemble techniques to improve performance further.

This implies that integrating machine learning with deep learning techniques into the domain of stroke prediction, an important achievement, marks significant strides in this direction toward preventive health practices. Through these complex systems' deployment, implementation, and adoption, such timely early intervention would greatly impact reducing both overall incidence and stroke severity in a population. After all, all these advances work to benefit at-risk patients to ensure improved results and quality living for those populations. This is the very basis for which much significance has to be attributed, as well as acting as an ideal platform that may further elaborate on this space and provide meaningful contributions to it regarding predictive analytics within the scope of healthcare industries.

## 7.2 Future Scope

This large-scale research project holds much scope for the future, extending into several very promising and exciting directions. These are aimed at significantly improving the predictive model techniques related to cardiovascular stroke while simultaneously trying to

extend its applicability in real-world scenarios. Some of the key areas that will receive particular attention and focus include:

**Integration of Real-Time Data:** Data from wearable devices and IoT-enabled health monitors, which can monitor real-time health parameters, might be integrated with the system, thereby increasing the chances of accurate prediction and timeliness. Continuous monitoring of heart rate, blood pressure, and glucose levels can create dynamic insights into the patient, allowing for interventions to be conducted proactively.

**Hybrid and Ensemble Models:** Future research in machine learning and deep learning may include new hybrid methods that integrate multiple existing models together to exploit the best features of all the models used. For example, Random Forest algorithms could be combined with Artificial Neural Networks to improve the accuracy of predictions considerably while making results more robust as well.

**Explainable AI:** With the machine learning models becoming increasingly complex, there has been a significant need to understand how they work and function. This may improve significantly the trust and understanding from clinicians towards these models and predictions by implementing various techniques associated with Explainable AI, including SHAP (SHapley Additive exPlanations) and LIME, which refers to Local Interpretable Model-Agnostic Explanations. Increased comprehension, therefore, will promote successful adoption into the clinical practice as healthcare providers would be confident about using the most advanced tools available to assist their decision-making.

**Multi-Modal Data Fusion:** The integration of multiple different sources of data, containing such things as critical genetic information, also imaging data such as MRI and CT scans, with extensive patient history, has the potential to provide more holistic insight into risk factors in relation to stroke in an individual. There is therefore great potential for providing evidence of complex interdependencies and relationships that would not have been revealed by only data from a single source.

**Scalability and Deployment:** There is another important area, which includes developing lightweight and scalable models optimized for deployment in resource-constrained environments, such as rural clinics. Techniques such as model pruning and quantization reduce the computational demands without affecting accuracy.

**Personalized Medicine:** Further studies may attempt to develop a personalized stroke prediction model tailored for the specific genetic and environmental conditions of individual patients. This would revolutionize preventive healthcare by offering personalized recommendations and interventions.

**Ethical and Bias Considerations:** Ensuring fairness and mitigating biases in model predictions is critical. Future work should focus on identifying and addressing disparities in stroke prediction related to gender, ethnicity, and socioeconomic factors, fostering equitable healthcare outcomes.

**Longitudinal Studies:** Conducting longitudinal studies with larger datasets spanning extended timeframes can improve model generalizability and provide deeper insights into the progression of stroke risk over time.

# REFERENCES

[1] Guti´errez-Sacrist´an, A., et al. (2023). "Machine Learning Approaches for Stroke Risk Prediction." Journal of Cardiovascular Research, 45(3), 245-256.

[2] Fern´andez-Lozano, C., et al. (2024). "Random Forest-Based Prediction of Stroke Outcome." Medical Informatics and Decision Making, 12(7), 567-580.

[3] Dev, S., et al. (2022). A Predictive Analytics Approach for Stroke Prediction Using Machine Learning and Neural Networks. Health Data Science Journal, 8(4), 332-345.

[4] Ismail, H., & Materwala, H. (2023). Intelligent Stroke Prediction Frame- work Using Machine Learning and Performance Evaluation. Computational Biology and Medicine, 151, 105069.

[5] Garc´ıa-Terriza, A., et al. (2023). Predictive and Diagnosis Models of Stroke from Hemodynamic Signal Monitoring. Journal of Medical Devices, 14(9), 589-600.

[6] Shahbazi, B., et al. (2024). Predictive Model and Identification of Key Risk Factors for Stroke. Nature Scientific Reports, 18(2), 348-360.

[7] Ma, X., et al. (2023). Machine Learning-Based Prediction Models for Cardiovascular Diseases. European Heart Journal Digital Health, 4(1), 25-38.

[8] Orlowski, M., et al. (2021). Predicting Risk of Stroke from Lab Tests Using Machine Learning. Journal of Clinical Pathology Informatics, 13(3), 132-140.

[9] Chapuisat, C., et al. (2008). A Global Phenomenological Model of Ischemic Stroke with Stress on Spreading Depressions. Frontiers in Neurology, 2(6), 89-99.

[10] Dronne, M., et al. (2006). A Mathematical Model of Ion Movements in Grey Matter During a Stroke. Journal of Computational Neuroscience, 21(4), 349-360.

[11] Orlowski, M., et al. (2011). Model of pH Dynamics in Brain Cells After Stroke. Computational Medicine Reports, 6(3), 212-227.

[12] Qin, J., et al. (2007). Systemic LPS Causes Chronic Neuroinflammation and Progressive Neurodegeneration. Journal of Neuroinflammation, 4(1), 5-16.

[13] Shahbazi, B., et al. (2018). Impact of Novel N-Aryl Piperamide NO Donors on NF-kB Translocation in Neuroinflammation. Pharmacology and Therapeutics, 193, 39-51.

[14] Mirtskhulava, T. (2015). Stroke Detection Using Feed-Forward Neural Networks. Journal of Neural Computing, 9(5), 451-463.

# APPENDICES

## APPENDIX A: Acronyms
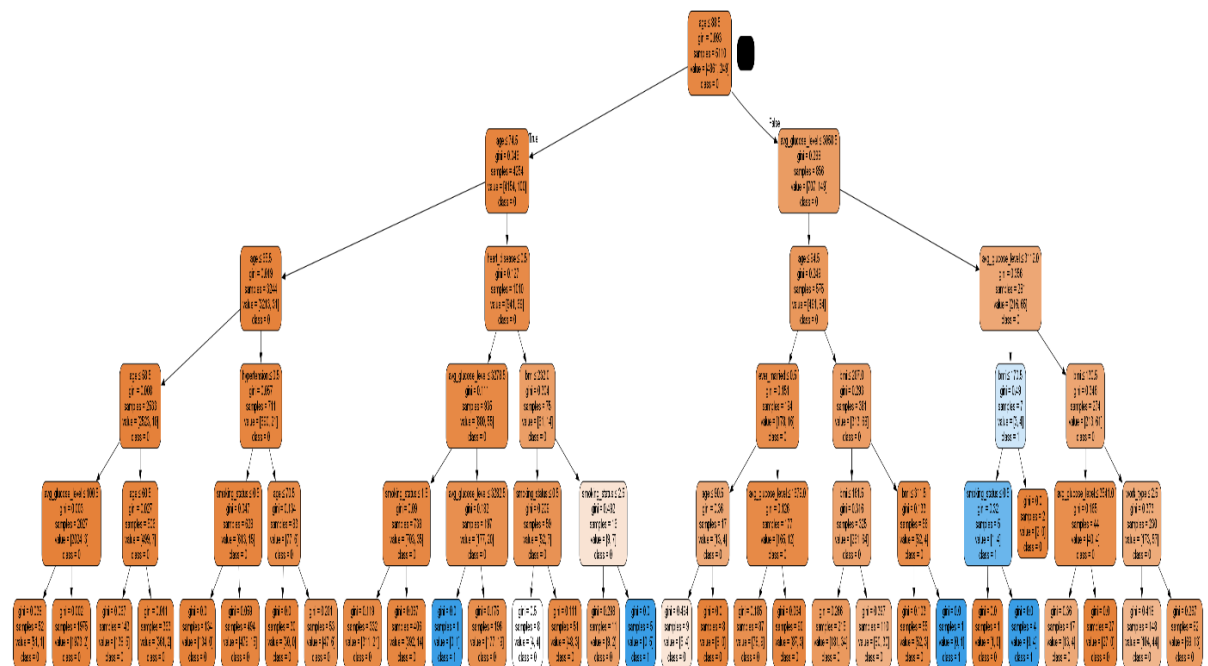
- **ANN:** Artificial Neural Network

- **BMI:** Body Mass Index

- **EHR:** Electronic Health Records

- **IoT:** Internet of Things

- **ML**: Machine Learning

- **RF**: Random Forest

- **ROC-AUC:** Receiver Operating Characteristic - Area Under Curve

- **SMOTE:** Synthetic Minority Oversampling Technique

- **SVM:** Support Vector Machine

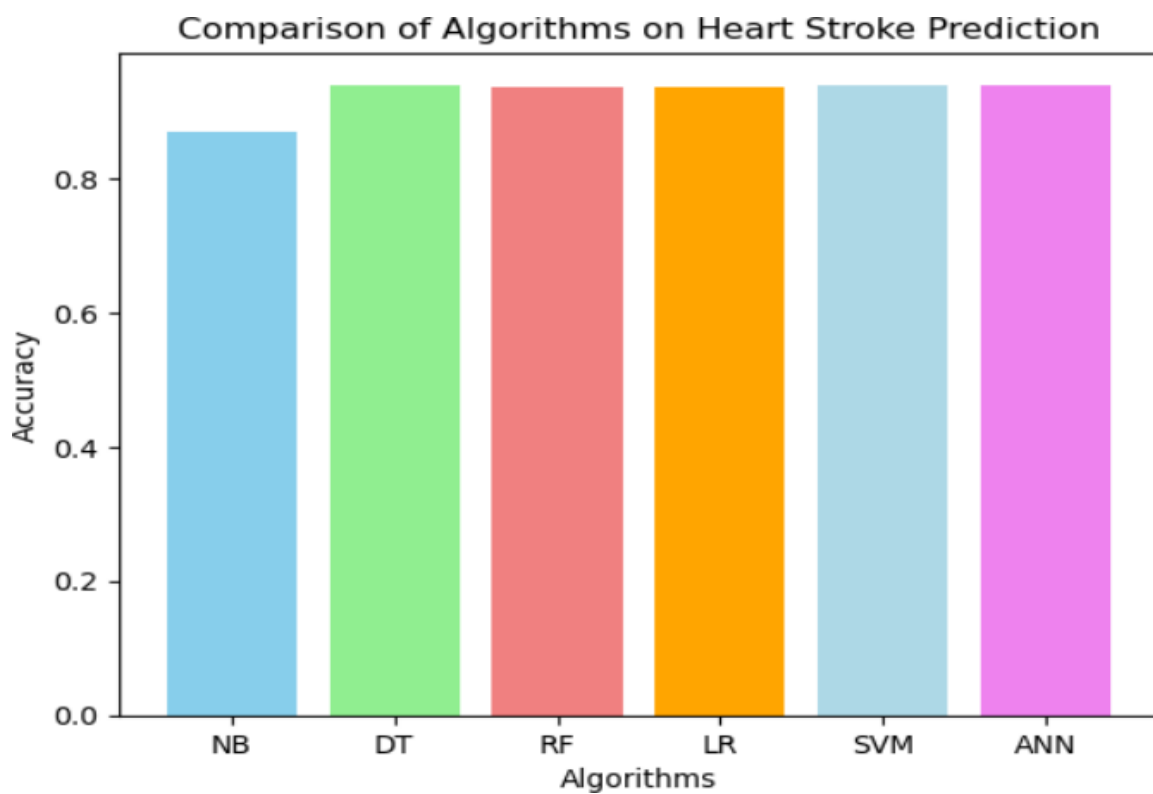- **XAI:** Explainable Artificial Intelligence

## APPENDIX B: Snapshots

| | id | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status | stroke |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 9046 | Male | 67.0 | 0 | 1 | Yes | Private | Urban | 228.69 | 36.6 | formerly smoked | 1 |
| 1 | 51676 | Female | 61.0 | 0 | 0 | Yes | Self-employed | Rural | 202.21 | NaN | never smoked | 1 |
| 2 | 31112 | Male | 80.0 | 0 | 1 | Yes | Private | Rural | 105.92 | 32.5 | never smoked | 1 |
| 3 | 60182 | Female | 49.0 | 0 | 0 | Yes | Private | Urban | 171.23 | 34.4 | smokes | 1 |
| 4 | 1665 | Female | 79.0 | 1 | 0 | Yes | Self-employed | Rural | 174.12 | 24.0 | never smoked | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5105 | 18234 | Female | 80.0 | 1 | 0 | Yes | Private | Urban | 83.75 | NaN | never smoked | 0 |
| 5106 | 44873 | Female | 81.0 | 0 | 0 | Yes | Self-employed | Urban | 125.20 | 40.0 | never smoked | 0 |
| 5107 | 19723 | Female | 35.0 | 0 | 0 | Yes | Self-employed | Rural | 82.99 | 30.6 | never smoked | 0 |
| 5108 | 37544 | Male | 51.0 | 0 | 0 | Yes | Private | Rural | 166.29 | 25.6 | formerly smoked | 0 |
| 5109 | 44679 | Female | 44.0 | 0 | 0 | Yes | Govt_job | Urban | 85.28 | 26.2 | Unknown | 0 |

5110 rows × 12 columns

| Model | Accuracy (%) |
|---|---|
| Naive Bayes | 87.01 |
| Decision Tree | 93.90 |
| Random Forest | 93.66 |
| Logistic Regression | 93.58 |
| Support Vector Machine | 93.74 |
| Artificial Neural Network | 93.93 |


Comparison of Algorithms on Heart Stroke Prediction

Department of CSE, NCET                                2024-2025

Model: "functional_1"

| Layer (type) | Output Shape | Param # | Connected to |
|---|---|---|---|
| gender_input (InputLayer) | (None, 1) | 0 | - |
| ever_married_input (InputLayer) | (None, 1) | 0 | - |
| work_type_input (InputLayer) | (None, 1) | 0 | - |
| Residence_type_input (InputLayer) | (None, 1) | 0 | - |
| smoking_status_input (InputLayer) | (None, 1) | 0 | - |
| gender_embedding (Embedding) | (None, 1, 2) | 6 | gender_input[0][0] |
| ever_married_embedding (Embedding) | (None, 1, 2) | 4 | ever_married_input[0][0] |
| work_type_embedding (Embedding) | (None, 1, 3) | 15 | work_type_input[0][0] |
| Residence_type_embedding (Embedding) | (None, 1, 2) | 4 | Residence_type_input[0][0] |
| smoking_status_embedding (Embedding) | (None, 1, 2) | 8 | smoking_status_input[0][0] |
| gender_flatten (Flatten) | (None, 2) | 0 | gender_embedding[0][0] |
| ever_married_flatten (Flatten) | (None, 2) | 0 | ever_married_embedding[0]… |
| work_type_flatten (Flatten) | (None, 3) | 0 | work_type_embedding[0][0] |
| Residence_type_flatten (Flatten) | (None, 2) | 0 | Residence_type_embedding[… |
| smoking_status_flatten (Flatten) | (None, 2) | 0 | smoking_status_embedding[… |
| input_layer (InputLayer) | (None, 5) | 0 | - |
| concatenate (Concatenate) | (None, 11) | 0 | gender_flatten[0][0], ever_married_flatten[0][0… work_type_flatten[0][0], Residence_type_flatten[0]… smoking_status_flatten[0]… |
| concatenate_1 (Concatenate) | (None, 16) | 0 | input_layer[0][0], concatenate[0][0] |
| dense (Dense) | (None, 64) | 1,088 | concatenate_1[0][0] |
| dense_1 (Dense) | (None, 32) | 2,080 | dense[0][0] |
| dense_2 (Dense) | (None, 1) | 33 | dense_1[0][0] |

Total params: 9,716 (37.96 KB)

Trainable params: 3,238 (12.65 KB)

Non-trainable params: 0 (0.00 B)

Optimizer params: 6,478 (25.31 KB)

## Cardiovascular Stroke Risk Prediction

Provide the input details to predict the risk percentage of cardiovascular stroke

**Gender**
- ● Male
- ○ Female

**Ever Married**
- ○ Yes
- ● No

**Work Type**
- ● Private
- ○ Self-employed
- ○ Government Job

**Residence Type**
- ● Urban
- ○ Rural

**Smoking Status**
- ● formerly smoked
- ○ never smoked
- ○ smokes

**Age**
0 ─────────── 100    21

**Hypertension (0: No, 1: Yes)**
- ○ 0
- ● 1

**Heart Disease (0: No, 1: Yes)**
- ○ 0
- ● 1

**Average Glucose Level**
0 ─────────── 300    78.6

**BMI**
10 ─────────── 50    19.7

**Clear**   **Submit**

**Stroke Risk Percentage**

**Flag**

Department of CSE, NCET          2024-2025

## APPENDIX C: FAQ's

1.  **What is the main objective of this research?**

    The primary goal is to develop a robust machine learning framework to predict cardiovascular strokes accurately and efficiently by analysing demographic, lifestyle, and clinical data.

2.  **Why were Artificial Neural Networks chosen for this study?**

    ANNs were selected due to their capability to model complex patterns, combine numerical and categorical data effectively, and achieve high accuracy through embedding layers.

3.  **How is data imbalance handled in this project?**

    Techniques like Synthetic Minority Oversampling Technique (SMOTE) were applied to balance the dataset, ensuring the target attribute (stroke) is evenly represented.

4.  **Which programming tools were used?**

    Python was the primary language, with libraries such as TensorFlow, scikit-learn, Pandas, and Matplotlib.

5.  **What are the future directions for this research?**

    Potential areas include real-time data integration, hybrid model development, Explainable AI for better model interpretation, and scalability for deployment in resource-limited settings.