



BIOSIGNALS BASED SMOKING STATUS PREDICTION USING STANDARD AUTOENCODER AND ARTIFICIAL NEURAL NETWORK

Dr. N. Palanivel¹, S. Deivanai², G. Lakshmi Priya³ and B. Sindhuja⁴

¹ Associate Professor, Dept. of Computer Science and Engineering, Manakula Vinayagar Institute of Technology, Puducherry, India

^{2,3,4} UG Students, Dept. of Computer Science and Engineering, Manakula Vinayagar Institute of Technology, Puducherry, India

E-mail address: palanivelcse@mvit.edu.in, deivanaisubramanian6@gmail.com, lakshmiPriyaganesan02@gmail.com, sindhujacse29@gmail.com

Received ## Mon. 20##, Revised ## Mon. 20##, Accepted ## Mon. 20##, Published ## Mon. 20##

Abstract: Smoking is still a major global health concern since it causes a host of illnesses and early deaths around the globe. Utilizing biosignals to predict smoking status can yield insightful information for tailored interventions and smoking cessation programs. This work presents a novel method that combines an artificial neural network (ANN) and a regular autoencoder to predict smoking status based on biosignals. The proposed method involves preprocessing biosignal data to extract relevant features, which are then input into an autoencoder for dimensionality reduction. The output of an autoencoder is used as input for predicting smoking status using an ANN. The model is trained and evaluated using a dataset containing biosignal data from individuals with a known smoking status. The suggested strategy is effective, as seen by the experimental results, which show a high degree of prediction accuracy about smoking status. The model's performance is further validated through comparisons with existing methods, showing superior performance in terms of accuracy and robustness. The developed model is integrated into a user-friendly application aimed at promoting smoking cessation. In addition to specific online pages aimed at enlightening users about the negative consequences of smoking and the advantages of stopping, the program offers users individualized insights into their smoking status based on biosignals. Additionally, a menu-based chatbot is included to address user queries and provide support for smoking cessation efforts. The implemented deep learning model achieves the desired level of accuracy in predicting smoker status, and the user-friendly application offers a convenient platform for public health and personalized healthcare interventions.

Keywords: Biosignals, Smoking status, Autoencoder, Artificial Neural Network, Deep learning, Feature Extraction.

1. INTRODUCTION

Smoking remains one of the most pressing public fitness-traumatic situations globally, exacting a heavy toll on people's health and well-being and straining healthcare systems. In spite of worldwide reputation campaigns and crook responsibilities, smoking remains a wonderful contributor to avoidable deaths and ailments. Predicting smoking reputation correctly is essential for developing powerful smoking cessation strategies and personalized interventions aimed at decreasing smoking incidence and its associated fitness burdens.

In current years, there has been a developing hobby of leveraging biosignals to count on smoking fame. Biosignals, inclusive of physiological and biochemical

markers obtained from people, offer precious insights into their health popularity and behaviours, along with smoking behavior. Robust predictive models that reliably affirm a person's smoking fame may be advanced by using state-of-the-art computer techniques to investigate those biosignals.

This paper gives a completely unique approach for predicting smoking reputation based totally on biosignals using a fusion of autoencoder and artificial neural community (ANN) architectures. Autoencoders are unsupervised deep learning algorithms able to learn inexperienced representations of input records with the useful resource of encoding them right into a decrease-dimensional latent area and then reconstructing them. By leveraging the latent representations found by an

E-mail: palanivelcse@mvit.edu.in, deivanaisubramanian6@gmail.com, lakshmiPriyaganesan02@gmail.com, sindhujacse29@gmail.com

<http://journals.uob.edu.bh>



autoencoder, in conjunction with additional features extracted from biosignal data, an ANN version is educated to expect smoking recognition appropriately.

The proposed technique has several advantages over conventional strategies of smoking status prediction. Firstly, it allows the automatic extraction of applicable capabilities from biosignal records, putting off the need for manual characteristic engineering and doubtlessly taking pictures of diffused patterns indicative of smoking behavior. Secondly, the fusion of autoencoder and ANN architectures allows the model to investigate complicated relationships between biosignals and smoking, improving its predictive performance.

Extending a predictive version that can accurately determine smoking popularity based on biosignal facts is the main goal of this research. In order to gain from this, the actions listed below are taken:

- *Data Preprocessing and Feature Extraction:*

Biosignal records, along with physiological and biochemical parameters including age, weight, levels of cholesterol, and so forth., are pre-processed to ensure consistency and reliability. Relevant functions are extracted from the biosignal facts, together with the uncooked measurements and derived metrics.

- *Autoencoder is primarily based on Dimensionality Reduction:*

The pre-processed biosignal information is entered into an autoencoder, which learns a lower-dimensional representation of the entered facts. This latent instance captures the important capabilities of the biosignals while decreasing noise and, besides the factor information.

- *ANN is primarily based on Smoking Status Prediction:*

The latent representations observed through the usage of the autoencoder, collectively with additional features, are used as input for training an ANN model to assume smoking reputation. The ANN version is professional in the usage of supervised studying strategies, in which the goal variable is the binary classification of smoking repute (smoker or non-smoker).

- *Model Evaluation and Validation:*

The mean efficacy of the developed prediction model is assessed by a variety of criteria, including F1-score, accuracy, precision, and consideration. In order to assess the model's generalization functionality, flow validation techniques are also used during testing.

- *Integration into a user-friendly application:*

The advanced predictive model is included in consumer-pleasant software aimed at promoting smoking cessation. The application offers users custom-designed insights into their smoking recognition primarily based on biosignals, along with educational assets on the

detrimental consequences of smoking and the blessings of quitting. A menu-frequently-based chatbot is blanketed to deal with user queries and offer useful resources for smoking cessation efforts.

Overall, this research contributes to the improvement of predictive modelling techniques for smoking prediction based on biosignals. The advanced version gives a promising technique for identifying individuals who smoke and non-people who smoke appropriately, thereby facilitating focused interventions and public health responsibilities geared toward decreasing smoking incidence and improving population health consequences.

2. RELATED WORKS

Machine learning applications for predicting smoking behaviour conducted by Etter et al. [1] explores the domain of smoking cessation interventions, with a specific focus on predicting changes in smoking behavior six months after utilizing the Stop-Tabac smartphone app. Drawing upon data from a randomized trial, the researchers utilized machine learning techniques, particularly the Random Forest classification algorithm, to discern predictors associated with smoking cessation, reduction, and relapse. The study unveils the intricate interplay of factors influencing smoking behaviour, ranging from motivational factors to levels of dependence and the utilization of nicotine replacement products. Although the results do not explicitly declare the accuracy of the machine learning models, they do point to promising results when it comes to tailoring smoking cessation programs to the unique profiles of individual users. This study provides information that may improve the effectiveness of apps designed to help people quit smoking and, in turn, improve public health outcomes related to the battle against tobacco use. In the area of cessation applications, it represents a substantial breakthrough.

The study conducted by Issabakhsh et al. uses information from the Population Assessment of Tobacco and Health (PATH) survey [2] to investigate the use of machine learning techniques in forecasting adult US smoking cessation. The authors uncover factors that influence smoking cessation and create machine learning classification models targeted at adult smokers who are currently smokers through analysis of waves 1-3 of the PATH project. The study predicts smoking cessation with 72% accuracy by wave 2 and 70% accuracy by wave 3 using algorithms like random forest and gradient boosting machines. The use of e-cigarettes in the past, smoking cigarettes before stopping, age when smoking started, number of years of smoking, use of polytobacco, and BMI are among the important factors that were found. This research highlights the significance of machine



learning in understanding smoking cessation dynamics and provides insights for the development of effective cessation interventions.

Nuryunarsih et al.'s study demonstrates a significant stride in healthcare analytics by employing artificial neural network (ANN) machine learning [3] to forecast smoking behaviour among Indonesian health professionals. By employing ANN algorithms, the researchers have developed predictive models capable of accurately analysing various factors influencing smoking behaviour, including demographic information, knowledge about smoking-related diseases, and attitudes towards smoking cessation. Achieving an impressive accuracy rate of 81%, these models provide valuable insights into the determinants of smoking behaviour within the Indonesian healthcare workforce. The integration of ANN technology holds promise for enhancing smoking cessation interventions and advancing tobacco control initiatives in Indonesia, ultimately leading to improved patient and public health outcomes. Continued collaboration between healthcare professionals and AI systems is crucial for refining predictive models and ensuring their effectiveness in addressing public health challenges.

Frank et al.'s research paper [4] presents a fresh method to healthcare analytics by predicting smoking status through statistical analysis and machine learning algorithms. By using a range of classification algorithms and statistical approaches, the researchers expect to accurately predict smoking behaviour based on medical information acquired during hospital stays. Using a dataset of 40,000 patients from a community hospital in the Greater Pittsburgh Area, significant differences in blood test results between smokers and nonsmokers are discovered. The INR, HB, and HCT levels are where these variations are most noticeable. The study demonstrates that the logistic model performs better than the other models in terms of accuracy rates when it comes to predicting smoking status. The rates for precision, recall, and F-measure are 83%, 83.4%, 83.2%, and 83.44%, respectively. In conclusion, the current study shows how machine learning algorithms and statistical analysis can predict smoking behavior, improving healthcare analytics and assisting tobacco-reduction-focused public health efforts. In conclusion, the current study shows how machine learning algorithms and statistical analysis can predict smoking behavior, improving healthcare analytics and assisting tobacco-reduction-focused public health efforts. To improve smoking cessation efforts and refine predictive models, data scientists and healthcare practitioners must work together.

Thakur et al. provide a multi-class classification model based on machine learning for the real-time prediction of smoking behavior during activities of daily

living (ADLs) [5]. Using data collected from wrist-wearable devices equipped with IMU sensors, the study focused on activities such as walking, running, and smoking. Time-domain, frequency-domain, and descriptive features were retrieved from streaming sensor data using a sliding window method. The developed models achieved high predictive accuracy rates of up to 98.7% for identifying smoking activity, indicating potential for real-time detection and intervention. These findings offer implications for monitoring and motivating smokers to quit, improving healthcare interventions, and expanding into preventive healthcare applications.

Lai et al. concentrate on creating machine learning models to forecast the results of quitting smoking [6]. The research employs a range of machine learning algorithms, such as support vector machines (SVM), artificial neural networks (ANN), logistic regression (LoR), random forests (RF), k-nearest neighbors (KNN), naïve Bayes (NB) and classification and regression trees (CART), and is carried out on patients participating in a smoking cessation program in Northern Taiwan. With a sensitivity of 0.704, specificity of 0.567, accuracy of 0.640, and area under the receiver operating characteristic (ROC) curve of 0.660, the ANN model performed marginally better. By incorporating patient parameters such as nicotine dependence level, daily cigarette consumption, and previous quit attempts, the developed predictive model offers a personalized approach to smoking cessation treatment. These predictive models have the potential to assist healthcare providers in identifying smokers likely to succeed in quitting, thereby enhancing counselling and treatment strategies and ultimately reducing the burden of smoking-related diseases on public health.

Caccamisi et al. Examine how machine learning and natural language processing (NLP) approaches can be combined to automatically extract and classify patients' smoking status from unstructured electronic medical records (EMRs) [7]. Their research involves training 32 prediction models with different combinations of classifiers, tokenization strategies, and attribute selection procedures. The analysis of 85,000 classified sentences taken from EMR data serves as the basis for this work. The best-performing model outperforms rule-based methods with an astounding accuracy of 98.14% and an F-score of 0.981 using Support Vector Machine (SVM) with a combination of unigrams and bigrams as tokens. This research highlights the potential of automated classification systems to accurately assess smoking status from EMRs without manual intervention, thereby enhancing clinical research, healthcare delivery, and public health interventions targeting smoking cessation.

Wang et al. Examine the viability of using magnetic resonance imaging (MRI) with deep learning to



distinguish between smoking and non-smoking [8]. 127 participants, including smokers and non-smokers, provided head MRI 3D-T1WI images for the study since smoking is a serious public health concern. A convolutional neural network coupled with a recurrent neural network with long short-term memory architecture (ConvLSTM) and a deep 3D convolutional neural network (Conv3D) were the two deep learning models created. The Conv3D model achieved a promising accuracy of 80.6%, while the ConvLSTM model outperformed, attaining an accuracy of 93.5%. These findings point to the possibility of deep learning-based MRI to outperform conventional support vector machine (SVM) techniques in terms of smoking status prediction accuracy. This strategy has the potential to improve our knowledge of the structural alterations in the brain caused by smoking and to improve smoking cessation programs, which will lessen the impact of smoking-related illnesses on public health.

Haque et al. Determine the validity of smoking status ascertainment methods in population-based electronic health databases by conducting a comprehensive review and meta-analysis [9]. With 57 articles and 116 algorithms analyzed, predominantly based on electronic medical record (EMR) data, The study evaluates these algorithms' efficacy, a number of them make use of diagnosis codes for illnesses connected to smoking. Approximately 50% of the algorithms incorporate machine-learning models. Even though algorithm performance varied greatly, the pooled estimates show a positive predictive value of 0.843, a sensitivity of 0.672, and a specificity of 0.918. Meta-regression analysis identifies model-based algorithms and EMR data as factors associated with higher sensitivity. Some algorithms have limits in terms of sensitivity and positive predictive value; however, those that use several linked data sources and machine learning models show enhanced validity. The present analysis highlights the significance of strong algorithms in accurately determining smoking status in electronic health data, which is crucial for comprehending the dangers associated with smoking and guiding public health measures.

In their study, Vaghefi et al. explore the use of Convolutional Neural Networks (CNNs) to detect smoking status from retinal images, highlighting the impact of smoking on retinal vasculature and its association with cardiovascular disorders [10]. Utilizing a dataset of 165,104 retinal images labelled with self-reported smoking status, the researchers developed CNN models using two pre-processing methods: "contrast-enhanced" and "skeletonized". While the contrast-enhanced model demonstrated an accuracy of 88.88% and a specificity of 93.87%, the skeletonized model

showed lower performance. Attention maps generated by the contrast-enhanced model emphasized predictive features such as retinal vasculature, perivascular regions, and the fovea. Despite achieving high accuracy, the study suggests further research to improve sensitivity, potentially by considering smoking frequency, duration, and dosage. This research underscores the potential of CNNs in detecting smoking status from retinal images, providing valuable insights into the physiological changes associated with smoking.

In their study, Dr. N. Palanivel et al. explore the use of Deep Neural Networks (DNN) for the early detection of Alzheimer's Disease (AD) using a large MRI dataset comprising normal and diseased subjects [11]. AD is a progressive neurodegenerative disorder characterized by cognitive decline and memory loss, making early diagnosis crucial for effective management and intervention. The DNN algorithms developed in this research achieved a high accuracy of approximately 95.1% in predicting AD, surpassing the performance of CNN models. This study addresses the significant clinical, social, and economic need for early detection of AD, offering promising implications for improving patient outcomes and healthcare management.

In their study, Dr. N. Palanivel et al. introduce novel techniques for diabetes classification and prediction in healthcare applications, leveraging deep learning and advanced genetic optimization [12]. The vast amount of data in medical databases poses a challenge for accurate data classification, particularly in disease diagnosis. The study proposes three methods to accurately classify patient medical data by establishing optimal association rules: Advanced Genetic Optimization (AGO)-based deep neural connectivity classification (ABGO-DNAC), integrated gradient-based classification, and amplified classification. These approaches aim to enhance diabetes diagnosis outcomes by improving classification accuracy and reducing processing time. The AGO method specifically focuses on enhancing diabetes classification accuracy by utilizing an expanded evolutionary algorithm to optimize association rules and employing an amplified gradient derivative classifier. Attributes are organized based on IDGBC selection trees, further optimizing the classification process. This research contributes to the advancement of healthcare analytics, offering promising implications for improving diabetes diagnosis and patient care.

3. PROPOSED WORK

A. Overview

The proposed work introduces biosignal data for smoking status prediction through the fusion of standard autoencoders and artificial neural networks. On the other



hand, no relevant research has been done on the use of biosignal data to forecast smoking status or smoking cessation. A number of crucial elements are included in the suggested system for predicting smoker status using biosignals, including data collection, preprocessing, model construction, and application integration. The various machine learning models have been explored, but comparatively, the model developed using a fusion of an autoencoder and an artificial neural network (ANN) gives good performance. Overall, the proposed system aims to provide a comprehensive and accessible tool for predicting smoker status and promoting smoking cessation using biosignal data and deep learning techniques.

B. Dataset Description

1) *Demographic data:*

This category includes Age, Height, Weight, and Waist circumference.

2) *Clinical data:*

- Vision: Eyesight measurements (left and right)
- Hearing: Hearing measurements (left and right)

- Cardiovascular: Blood pressure (systolic and relaxation)
- Metabolic: Fasting blood sugar levels, Cholesterol levels, Triglyceride levels, Haemoglobin levels
- Renal: Urine protein levels, Serum creatinine levels
- Liver function: AST (Aspartate Aminotransferase) levels, ALT (Alanine Aminotransferase) levels, and GTP (Gamma-Glutamyl Transferase) levels
- Oral health: Dental caries status

C. Proposed Algorithm

The proposed algorithm aims to first extract meaningful features from the input data using an autoencoder and then utilize these features to train a neural network classifier for predicting smoker status. The use of an autoencoder helps in capturing the underlying patterns in the data, potentially improving the classification performance of the neural network

TABLE I. DATASET DESCRIPTION

Feature Category	Feature Name	Description
Demographic Information	Age (years)	User's age in years.
Anthropometric Measurements	Height (cm)	User's height in centimeters.
	Weight (kg)	User's weight in kilograms.
	Waist (cm)	User's waist circumference in centimeters.
Visual acuity	Eyesight (left/right)	It represents the eyesight health
Auditory Perception	Hearing (left/right)	It represents the hearing health
Vital Signs	Systolic Blood Pressure (mmHg)	Systolic blood pressure reading in millimeters of mercury.
	Relaxation Blood Pressure (mmHg)	Diastolic blood pressure reading in millimeters of mercury.
Biochemical Measures	Fasting Blood Sugar (mg/dL)	Blood sugar level measured after an overnight fast (mg/dL).
	Cholesterol (mg/dL)	Total cholesterol level in the blood (mg/dL).
	Triglyceride (mg/dL)	Level of triglycerides (a type of fat) in the blood (mg/dL).
	HDL (mg/dL)	Level of high-density lipoprotein (HDL) or "good" cholesterol in the blood (mg/dL).
	LDL (mg/dL)	Level of low-density lipoprotein (LDL) or "bad" cholesterol in the blood (mg/dL).
	Haemoglobin (g/dL)	Haemoglobin level in the blood (g/dL), which carries oxygen.
Renal	Urine protein (mg/dL)	Amount of protein found in the urine (mg/dL).
	Serum creatinine (mg/dL)	Level of creatinine (a waste product filtered by the kidneys) in the blood (mg/dL).
Liver function	AST (U/L)	Aspartate aminotransferase, an enzyme found in the liver, muscles, and red blood cells
	ALT (U/L)	Alanine aminotransferase, another enzyme found in the liver.
	GTP (U/L)	Gamma-glutamyl transferase, an enzyme found in the liver, bile ducts, and pancreas.
Target Variable	Smoker Status (categorical)	Binary value indicating smoker (1) or non-smoker (0).



classifier. The proposed algorithm construction is divided into two phases, namely the autoencoder architecture and building the neural network classifier.

Autoencoder Architecture:

- Defined an autoencoder with an input layer, two hidden layers for encoding (128 neurons, 64 neurons), and a mirrored architecture for decoding.
- Compile the autoencoder using the 'adam' optimizer and 'mean_squared_error' loss function.
- Fit the autoencoder to the standardized training data for 50 epochs with a batch size of 32.
- Validate the autoencoder on the standardized test data.
- Use the encoder part of the trained autoencoder to extract features from the standardized training and test data.

Building the Neural Network Classifier:

- Created a sequential model for classification.
- Add a dense layer with 128 neurons and 'relu' activation function, a dropout layer with 0.5 dropout rate, another dense layer with 64 neurons and 'relu' activation function, and another dropout layer.
- Add a dense output layer with 1 neuron and 'sigmoid' activation function for binary classification.
- Compile the model using the 'adam' optimizer, the 'binary_crossentropy' loss function, and the accuracy metric.
- Validate the encoded test features by training the model with the encoded features for 50 epochs at a batch size of 32.
- Finally, assessing the accuracy and loss measures of the model's performance on the test set.

D. Workflow of the Algorithm

1) *Data Acquisition and Preparation:* Obtaining bio-signal data for smokers and non-smokers. This data includes demographic data and clinical data. Ensure the data is properly labeled, indicating whether each data point belongs to a smoker or a non-smoker. Clean and preprocess the data, which involves handling missing values, outliers, or inconsistencies.

2) *Data Splitting:* Dividing labeled data into two sets: training and testing data. The training data (typically 70–80% of the total data) will be used to train the model. The testing data (remaining 20–30%) will be used to evaluate the model's performance on unseen data.

3) *Autoencoder Training:* Build and train the autoencoder model using the scaled training data. Train

the autoencoder to reconstruct the original data from the encoded features as accurately as possible. This process helps the autoencoder learn underlying patterns and relationships within the biosignals.

4) *Feature Extraction (Using Encoder):* Once trained, use the encoder part of the autoencoder to extract the encoded features from both the scaled training data and the original testing data. These encoded features represent compressed versions of the original bio-signals, hopefully capturing the most important information for smoker classification.

5) *Neural Network Classifier Training:* Build a neural network classifier model and train it on the extracted encoded features from the training data and the corresponding labels (smoker or non-smoker). The neural network classifier learns to associate specific patterns in the encoded features with the smoker and non-smoker labels.

6) *Model Evaluation:* Evaluate the performance of the trained neural network classifier on the encoded features from the testing data. Common metrics used for classification include accuracy, precision, recall, F1-score and AUC-ROC score.

4. METHODOLOGY

A. Data Source

The overview that outlines the specific measurements obtained for each biosignal data, including demographic information, vision and hearing assessments, physiological parameters, lipid panel and complete blood count tests, urine protein levels, serum creatinine levels, liver function tests, and dental caries status. Each of the biosignal data plays a crucial role in informing the predictive model, contributing essential insights into the individual's health status and smoking behaviour. Table II. discusses the dataset source and their tests for collecting the biasignal data.

B. Data Preprocessing

Data preprocessing encompasses various techniques that clean, transform, and format the data to make it suitable for model training. One essential step in this process is Normalization using techniques like StandardScaler.

Normalization with StandardScaler:

StandardScaler transforms each feature in the data (bio-signals) to have a mean of 0 and a standard deviation of 1. Prior to feeding the bio-signal

characteristics into the autoencoder and neural network classifier, StandardScaler is utilized to normalize them. As a result, the model's ability to distinguish between smokers and non-smokers may be enhanced and all bio-signals are guaranteed to contribute equally to its learning process. By bringing all features down to a common scale, this effectively prevents any one feature from controlling the learning process by having a significantly bigger or smaller scale than the rest.

StandardScaler divides the output by the feature's standard deviation after subtracting each feature's mean from each data point. Here, x_{scaled} represents the normalized value of the feature x .

In terms of math, for a feature x :

$$x_{scaled} = \frac{(x - mean(x))}{standard_deviation(x)}$$

There are several advantages to using StandardScaler before feeding data into a machine learning model:

1) *Improved Model Performance:* Numerous machine learning methods may be sensitive to the size of the features, particularly those that compute the distances between data points (such as certain classification algorithms). Normalization guarantees that every feature makes an equal contribution to the model's computations,

which may improve convergence and performance.

2) *Faster Training:* When features have different scales, the model might take longer to learn because it needs to adjust for these differences. Normalization can help the model converge faster as the features are presented on a similar scale.

3) *Prevents Feature Domination:* If a feature has a much larger scale compared to others, it can overshadow the influence of the other features on the model's learning process. Normalization prevents this by putting all features on a level playing field.

C. Deep Learning Architecture

1) *Autoencoder for Feature Extraction and Dimensionality Reduction:*

Standard autoencoders use Mean Square Error for reconstruction, which offers a general-purpose approach for dimensionality reduction and feature extraction. Standard autoencoder that doesn't specifically target any particular distributions or noise management techniques, instead learning a compressed representation of the biosignal data. This process not only helps in reducing the computational complexity of the subsequent modelling steps but also ensures that the most relevant features are retained for accurate prediction.

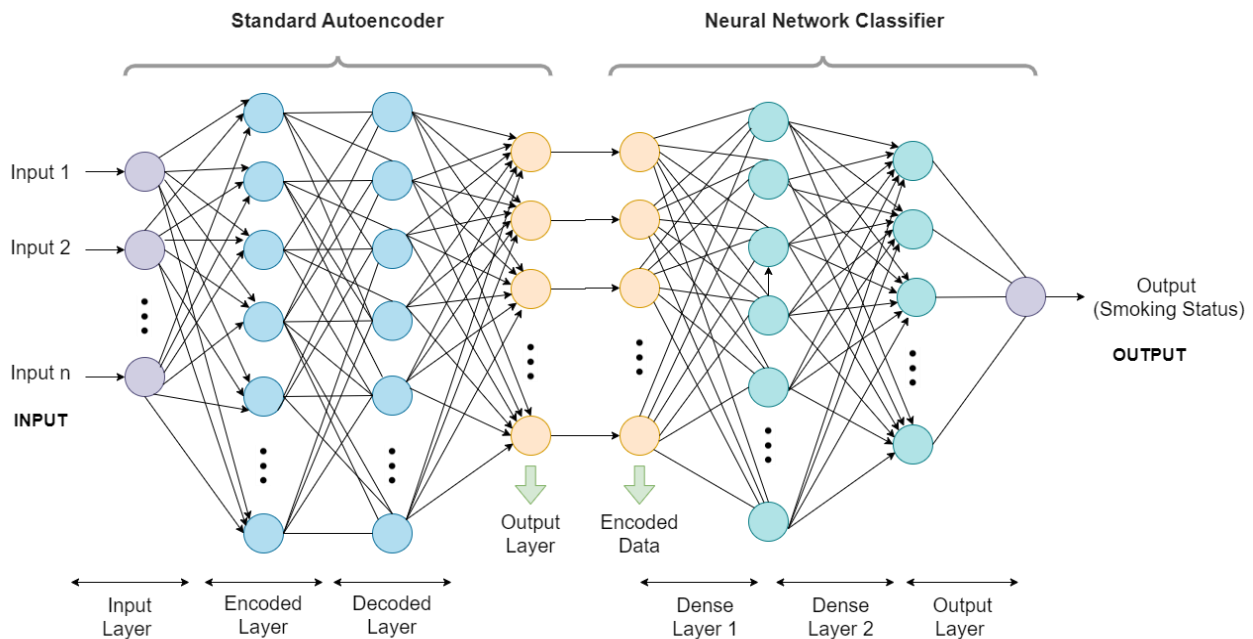


Figure 1. Proposed Algorithm Architecture



2) ANN for Prediction and Classification:

ANN is a broad term encompassing various neural network architectures, the network proposed is Dense Feedforward Neural Network. It accepts the encoded data representing bio-signals after processing by the autoencoder. Composed of densely connected layers with ReLU activation functions. These layers learn complex patterns within the bio-signal data that may correlate with smoking status. A single neuron that functions as a sigmoid activator. By converting the output from the hidden layers into a value between 0 and 1, it determines the likelihood that the individual smokes. Dense feedforward neural networks are well-suited for classification tasks on structured data, like the encoded biosignal features.

D. Flask Application Development

In the realm of smoker status prediction and smoking cessation initiatives, the development of a Flask application serves as a pivotal tool for delivering actionable insights and facilitating user engagement. By leveraging the versatility and scalability of Flask, Intuitive and user-friendly platform that has been created which empowers individuals to assess their smoker status, access resources for smoking cessation, and gain valuable insights into the benefits of quitting smoking.

E. Key features of the web application

1) Smoker Status Prediction:

The core functionality of the Flask application revolves around smoker status prediction based on the provided dataset. Users can input relevant health metrics such as age, weight, cholesterol levels, and more, and receive a personalized prediction of their smoker status. This prediction is powered by the deep learning model integrating autoencoder and ANN architectures, ensuring high accuracy and reliability.

2) Interactive Visualization:

To enhance user experience and facilitate data interpretation, the Flask application incorporates interactive visualization tools. Users can explore their health metrics through dynamic charts, graphs, and visualizations, allowing them to gain deeper insights into the factors influencing their smoker status.

3) Educational Resources:

The Flask application has parts devoted to offering thorough information on the dangers of smoking, the advantages of stopping, and evidence-based cessation techniques because it recognizes the critical role that education plays in smoking cessation efforts.

4) Why Quit Smoking Page:

A dedicated page within the application is designed to educate users about the compelling reasons to quit smoking. Through compelling narratives, testimonials, and statistics, this page reinforces the importance of smoking cessation and motivates users to take proactive steps towards quitting.

5) Benefits of Quitting Smoking Page:

Another integral component of the Flask application is the Benefits of Quitting Smoking page, which outlines the myriad health benefits associated with smoking cessation. Users can research how quitting smoking can enhance their overall health, reduce their chance of getting chronic illnesses, and increase their quality of life.

6) Menu-Based Chatbot:

To address user queries and provide personalized support, the Flask application integrates a menu-based chatbot. Users can interact with the chatbot to seek answers to common questions about smoking cessation, receive tips for quitting, and access additional resources for support.

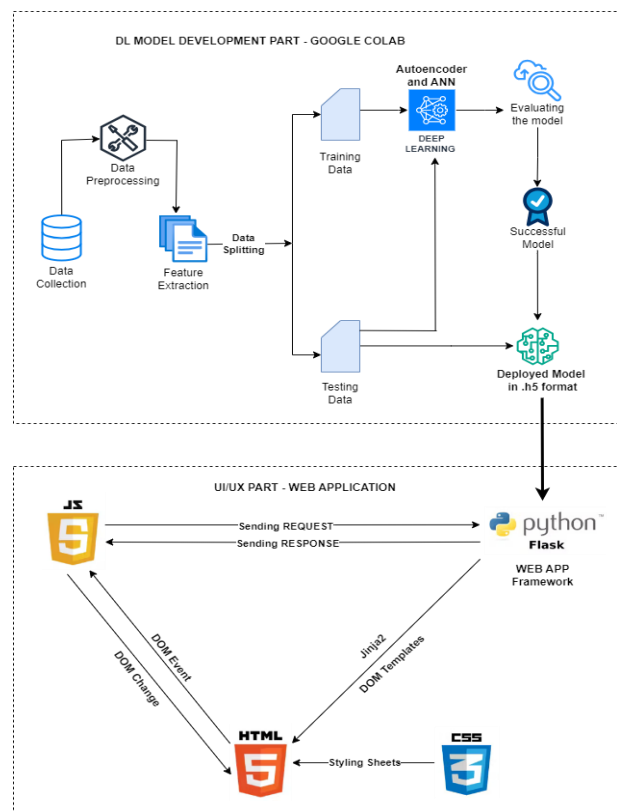


Figure 2. System Architecture Diagram



F. System Architecture

The system architecture in Figure 2. depicts a Flask web application for smoker status prediction with chatbot integration. The user interacts with the web application through a web browser, likely filling out a form with biosignal data. A Flask backend receives the data, utilizes a deep learning model (autoencoder and ANN) for prediction, and delivers the results along with the prediction value and various web pages for smoking cessation awareness through the web browser. Figure 2. is the pictorial representation of System Architecture.

5. RESULTS

Using the biosignal dataset, the effectiveness of several machine learning (ML) methods and the deep learning (DL) model for predicting smoker status was assessed. The outcomes show the advantages and disadvantages of each strategy, shedding light on whether DL is a better option than conventional ML methods.

Machine Learning Algorithm Performance:

Using the biosignal dataset, the effectiveness of several machine learning (ML) algorithms for predicting smoking status was assessed. The outcomes reveal that the Gradient Boosting Classifier attained the highest testing accuracy of 78%, trailed closely by the Adaboost Classifier and Random Forest Classifier, also achieving an accuracy of 77%. Table III. discusses the performance of various ML algorithms with evaluation metrics.

The machine learning algorithms were tested for accuracy, with Random Forest achieving the highest score (100%) but potentially overfitting the data, while Decision Tree shows similar performance but with a testing accuracy of 70%. Figure 3. is the graphical representation, which depicts the Training accuracy

versus Testing accuracy of various ML algorithms.

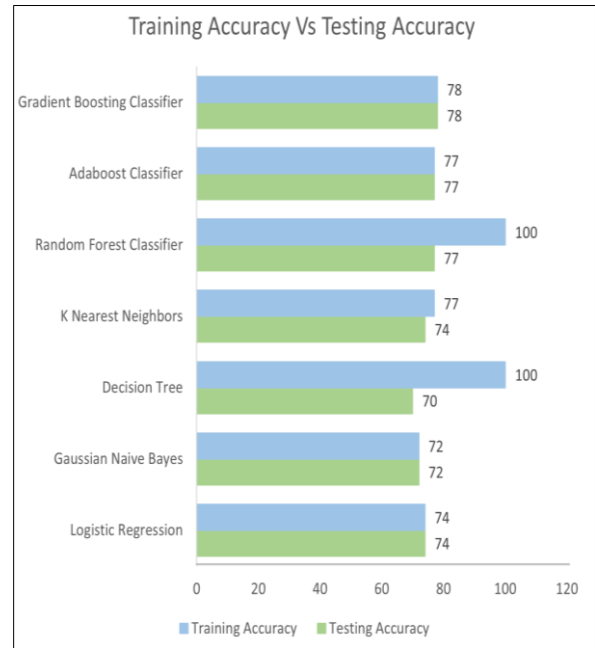


Figure 3. Training Accuracy Vs Testing Accuracy of various ML algorithm

In an evaluation of various machine learning algorithms on a biosignal dataset, three algorithms emerged as the top performers. Adaboost, Gradient Boosting, and Random Forest classifiers all achieved a precision score of 72%. When it comes to recall, the Gradient Boosting Classifier takes the crown with a score of 82%, followed by Random Forest at 80%. Combining precision and recall, the Gradient Boosting Classifier again secures the top spot with an F1-score of 77%, with Random Forest trailing closely behind at 76%. Overall,

TABLE III. PERFORMANCE OF VARIOUS ML ALGORITHM

Model	Training Accuracy	Testing Accuracy	Precision	Recall	F1-score	AUC-ROC score
Logistic Regression	74	74	71	71	71	74
Gaussian Naive Bayes	72	72	66	78	71	73
Decision Tree	100	70	66	65	65	69
K Nearest Neighbors	77	74	70	73	71	74
Random Forest Classifier	100	77	72	80	76	78
Adaboost Classifier	77	77	72	78	75	77
Gradient Boosting Classifier	78	78	72	82	77	78



these results suggest that Gradient Boosting and Random Forest are strong contenders for biosignal analysis tasks, with Adaboost also showing promise for accurate classification. Figure 4. is the graphical representation, which depicts the performance metrics such as precision, recall and F1 score of various ML algorithms.

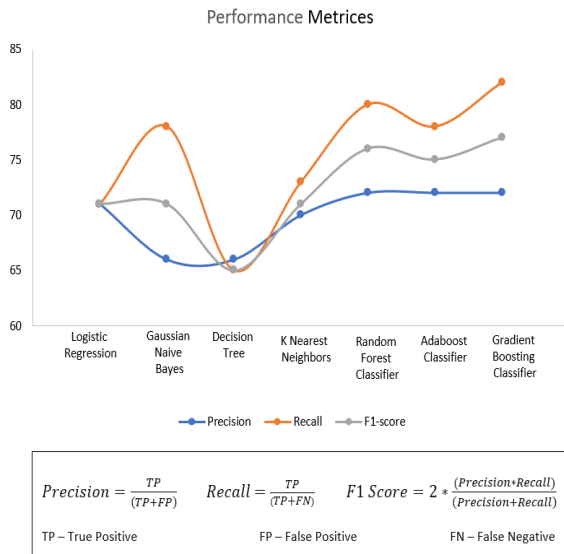


Figure 4. Performance metrics of various ML model

However, AUC-ROC offers a robust and informative metric for evaluating binary classification models, considering different thresholds, and comparing various models. The highest score is for Logistic Regression (78%), while the lowest score is for Decision Tree (69%). It's important to note that the overall

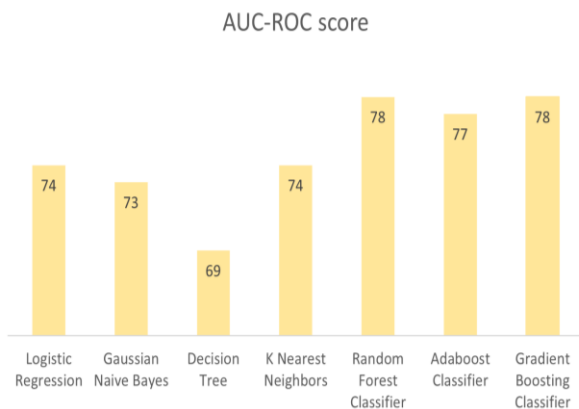


Figure 5. AUC ROC Curve of various ML model

performance of all models is relatively low, as none of the AUC-ROC scores are above 80%. Hence, the proposed approach goes on with a deep learning model. Figure 5. is the graphical representation, which depicts the AUC ROC performance metrics of various ML algorithms.

Deep Learning Algorithm Performance:

The evaluation metrics provided indicate that the model is performing reasonably well for predicting a person's smoking status for cessation using biosignals. Here's a brief analysis of each metric:

1) *Accuracy:* The accuracy of 0.80 indicates that the model is correctly predicting the smoking status of individuals approximately 80% of the time. While the accuracy score is commendable, it's crucial to contextualize it within the problem domain and consider the implications of false positives and false negatives.

2) *Precision:* Precision measures the proportion of positive predictions that are correct. A precision of 0.710 means that when a model predicts a person as a smoker, it is correct about 71% of the time. This indicates that the model has a relatively low rate of false positives.

3) *Recall:* Recall, often referred to as sensitivity, gauges the percentage of true smokers that the model correctly identifies. A recall of 0.835 means that the model is able to correctly identify about 83.5% of actual smokers. This indicates that the model has a relatively low rate of false negatives.

4) *F1-score:* The F1-score provides a balance between precision and recall by taking the harmonic mean of the two measures. The model appears to have an excellent balance between recall and precision, as indicated by its 77% F1-score.

5) *AUC-ROC score:* The AUC-ROC score evaluates the area under the receiver operating characteristic (ROC) curve, reflecting the model's capacity to differentiate between smokers and non-smokers. An AUC-ROC score of 86% suggests that the model demonstrates strong classification performance.

When comparing machine learning methods to the deep learning model, findings indicate that the deep learning model achieved an accuracy rate of 80%, comparable to the 77% accuracy rates of both the Gradient Boosting Classifier and Random Forest Classifier. But in terms of precision, recall, and

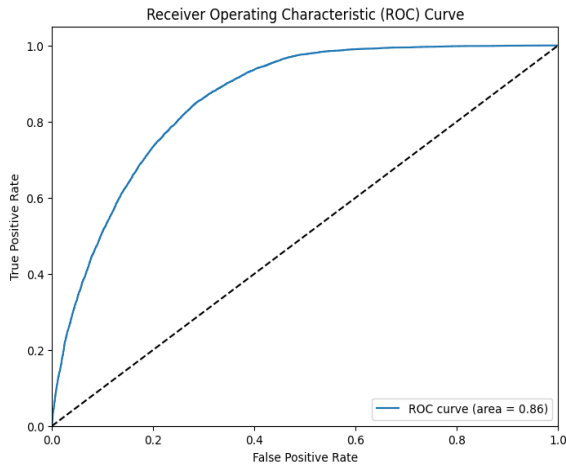


Figure 6. AUC ROC Curve of Deep Learning model

F1-score, deep learning - a hybrid of conventional autoencoder and artificial neural network (ANN) models - performed better than these approaches.

In contrast to the majority of machine learning algorithms, the deep learning model's precision of 71% shows that it can accurately identify smokers while reducing false positives. The recall of the DL model was 84%, indicating its effectiveness in capturing true positive instances of smokers. The F1-score of the DL model was 77%, which is higher than most ML algorithms, indicating its overall effectiveness in correctly classifying smokers while maintaining a balance between precision and recall.

6. DISCUSSION

Our research involved the development of a web application alongside the deep learning model for smoker status prediction. This user-friendly web application serves as a practical interface for interaction with the model, allowing users to easily input their Biosignal data and receive personalized predictions.

Benefits of the Web App: The web application offers several benefits in the context of smoker status prediction and potential behavior change:

- 1) *Accessibility*: The web application makes smoker status prediction readily accessible to a wider audience. Users can conveniently access the application from any web browser without requiring specialized software installation.
- 2) *Ease of Use*: The user interface is designed to be intuitive and user-friendly, allowing individuals with varying levels of technical expertise to easily interact with the model and obtain predictions.
- 3) *Personalized Feedback*: The application provides personalized predictions based on the user's unique Biosignal data. This can be more impactful than generic information about smoking risks.
- 4) *Educational Support*: By integrating educational content based on the predicted smoker status, the web app can potentially serve as a valuable tool for raising awareness about smoking risks and promoting smoking cessation efforts.

A. Application Home Page UI

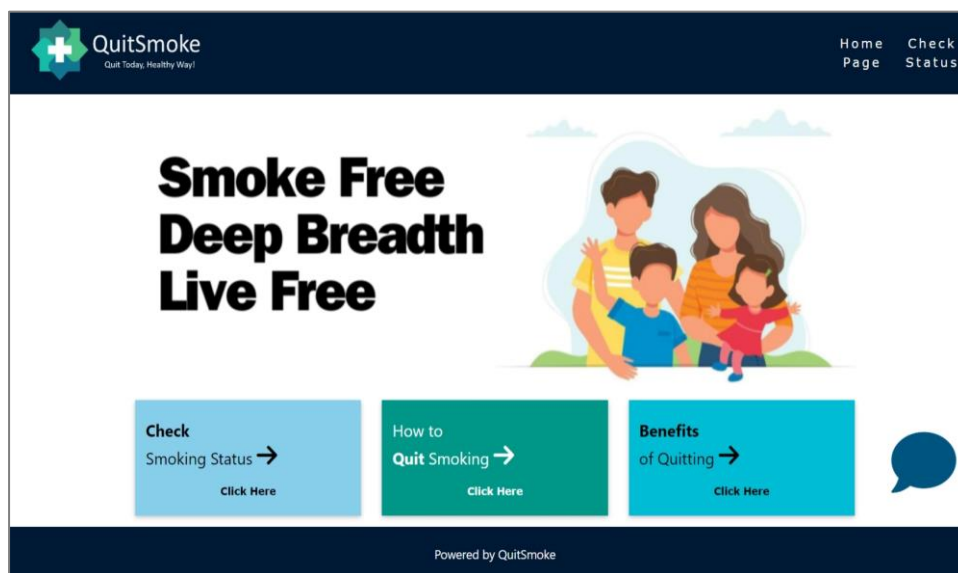


Figure 7. Home Page UI

B. Educational Content Pages UI



Figure 8. How to Quit Smoking Page UI

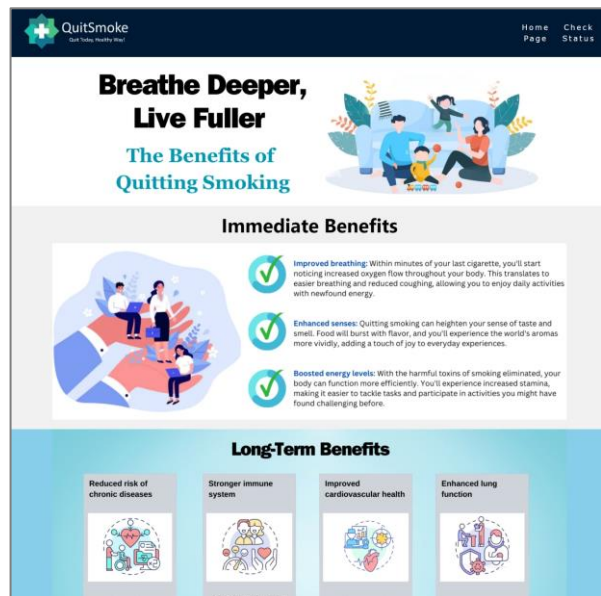


Figure 9. Benefits of Quitting Page UI

C. Smoking Status Prediction Pages UI

Figure 10. Data Collection Form for Prediction

Figure 11. Interactive health dashboard with Smoking Status



7. CONCLUSION

In conclusion, our research endeavours to revolutionize smoking status prediction and smoking cessation interventions by pioneering the integration of biosignal data and deep learning techniques. By leveraging a vast dataset of 1.5 lakhs biosignal data points, we have developed a robust predictive model that surpasses existing methodologies in accuracy and performance. Through the fusion of standard autoencoders and artificial neural networks, our model demonstrates its effectiveness in precisely determining smoker status with an amazing 80% accuracy and 86% AUC-ROC score. The deep learning models (autoencoder and ANN) outperformed traditional ML algorithms. Their higher recall and accuracy were a result of their capacity to recognize subtle patterns, process high-dimensional data, and understand complicated features.

Moreover, our comprehensive approach extends beyond model development, encompassing crucial stages such as data collection, preprocessing, model construction, and application integration. By addressing these fundamental elements, we ensure the reliability, scalability, and practicality of our proposed system. Furthermore, our user-friendly Flask frontend, featuring dedicated pages for Smoker Status Prediction, Why Quit Smoking, and Benefits of Quitting Smoking, functions as a valuable asset for individuals aiming to comprehend and address their smoking habits. Additionally, the integration of a menu-based chatbot enhances user engagement and accessibility, providing personalized support and information tailored to each user's needs. By empowering individuals with insights into the detrimental effects of smoking and the benefits of cessation, our system aims to catalyse positive behaviour change and promote healthier lifestyles.

REFERENCES

- [1] J.-F. Etter, G. Vera Cruz, and Y. Khazaal, "Predicting smoking cessation, reduction and relapse six months after using the Stop-Tabac app for smartphones: a machine learning analysis," *BMC Public Health*, vol. 23, no. 1, Article no. 1076, 2023.
- [2] M. Issabakhsh et al., "Machine learning application for predicting smoking cessation among US adults: An analysis of waves 1-3 of the PATH study," *PLoS One*, vol. 18, no. 6, Article e0286883, Jun. 2023. doi: 10.1371/journal.pone.0286883.
- [3] D. Nuryunarsih et al., "Artificial neural network machine learning prediction of the smoking behavior and health risks perception of Indonesian health professionals," *Environ Anal Health Toxicol*, vol. 38, no. 1, Article e2023003-0, Mar. 2023. doi: 10.5620/eaht.2023003.
- [4] C. Frank et al., "Predicting Smoking Status Using Machine Learning Algorithms and Statistical Analysis," *Adv. Sci. Technol. Eng. Syst. J.*, vol. 3, no. 2, pp. 184-189, 2018. DOI: 10.25046/aj030221.
- [5] S. S. Thakur et al., "Real-time prediction of smoking activity using machine learning based multi-class classification model," *Multimed Tools Appl.*, vol. 81, no. 10, pp. 14529-14551, 2022. DOI: 10.1007/s11042-022-12349-6.
- [6] C.-C. Lai et al., "Development of Machine Learning Models for Prediction of Smoking Cessation Outcome," *Int J Environ Res Public Health*, vol. 18, no. 5, p. 2584, Mar. 2021. DOI: 10.3390/ijerph18052584.
- [7] A. Caccamisi et al., "Natural language processing and machine learning to enable automatic extraction and classification of patients' smoking status from electronic medical records," *Scand J Public Health*, vol. 48, no. 3, pp. 316-324, Jul. 2020. DOI: 10.1080/03009734.2020.1792010.
- [8] S. Wang et al., "Discrimination of smoking status by MRI based on deep learning method," *Quant Imaging Med Surg*, vol. 8, no. 11, pp. 1113-1120, Dec. 2018. DOI: 10.21037/qims.2018.12.04.
- [9] M. A. Haque et al., "The validity of electronic health data for measuring smoking status: a systematic review and meta-analysis," *BMC Med. Inform. Decis. Mak.*, vol. 24, article no. 33, Feb. 2024. DOI: 10.1186/s12911-024-01822-1.
- [10] E. Vaghefi et al., "Detection of smoking status from retinal images: a Convolutional Neural Network study," *Sci. Rep.*, vol. 9, article no. 7180, May 2019. DOI: 10.1038/s41598-019-43702-3.
- [11] Indumathi, R., Palanivel, N., Kumar, V., Kannan, S., & Ahmed, J. M. (2022). Alzheimer's Disease Detection using Deep Neural Network. *TELEMATIQUE Vol 21, No 1*.
- [12] Selvi, V., Saraswathi, G., Indumathi, R., Palanivel, N., Reshma, S., & Subiksha, R. (2023). "Deep Learning Based Diabetes Classification and Prediction for Healthcare Applications Using Advanced Genetic Optimization", *International Conference on System, Computation, Automation and Networking (ICSCAN)* (pp. 1-6). Puducherry, India: IEEE. doi: 10.1109/ICSCAN58655.2023.10394940
- [13] Yupu Zhang, Jinhai Liu, Zhihang Zhang, and Junnan Huang, "Prediction of Daily Smoking Behavior Based on Decision Tree Machine Learning Algorithm," in *2019 IEEE 9th International Conference on Electronics Information and Emergency Communication (ICEIEC)*, Beijing, China, 12-14 July 2019. DOI: 10.1109/ICEIEC.2019.8784698.
- [14] Khishigsuren Davagdorj, Jong Seol Lee, Kwang Ho Park, and Keun Ho Ryu, "A machine-learning approach for predicting success in smoking cessation intervention," in *2019 IEEE 10th International Conference on Awareness Science and Technology (iCAST)*, Morioka, Japan, 23-25 October 2019. DOI: 10.1109/ICAwST.2019.8923252.
- [15] Pranta Roy, Fahad Hossain, and Nusrat Jahan, "Machine Learning Approach to Predict Influence of Smoking on Student Life," in *2021 12th International Conference on Computing Communication and Networking Technologies (ICCCNT)*, Kharagpur, India, 06-08 July 2021. DOI: 10.1109/ICCCNT51525.2021.9579775.
- [16] M. P. Milali, S. S. Kiware, N. J. Govella, F. Okumu, N. Bansal, S. Bozdag, J. D. Charlwood, M. F. Maia, S. B. Ogoma, F. E. Dowell, G. F. Corliss, M. T. Sikulu-Lord, & R. J. Povinelli, "An autoencoder and artificial neural network-based method to estimate parity status of wild mosquitoes from near-infrared spectra," *PLoS ONE*, vol. 15, no. 6, e0234557, Jun. 2020. DOI: 10.1371/journal.pone.0234557.



- [17] S. A. Ebiaredoh-Mienye, E. Esenogho, and T. G. Swart, "Integrating Enhanced Sparse Autoencoder-Based Artificial Neural Network Technique and Softmax Regression for Medical Diagnosis," *Electronics*, vol. 9, no. 11, p. 1963, Nov. 2020. DOI: 10.3390/electronics9111963.
- [18] N. Deshai, S. Ramya, Dr. B V D S Sekhar, and Dr. S. V. Ramana, "Prediction of Heart Disease with Autoencoder based ANN," in *ICRADL – 2021 (Volume 09 – Issue 05)*, Apr. 2021. DOI: 10.17577/IJERTCONV9IS05095.
- [19] A. B. W. Putra, R. Malani, B. Suprpty, and A. F. O. Gaffar, "A Deep Auto Encoder Semi Convolution Neural Network for Yearly Rainfall Prediction," in *2020 International Seminar on Intelligent Technology and Its Applications (ISITIA)*, Surabaya, Indonesia, 22-23 July 2020. DOI: 10.1109/ISITIA49792.2020.9163775.
- [20] Ramo, D. E., Thrul, J., Delucchi, K. L., Ling, P. M., Hall, S. M., & Prochaska, J. J. "The Tobacco Status Project (TSP): Study protocol for a randomized controlled trial of a Facebook smoking cessation intervention for young adults." *BMC Public Health*, vol. 15, 2015, article no. 897. doi: 10.1186/s12889-015-2227-8.
- [21] Wang, Z., Masoomi, A., Xu, Z., Boueiz, A., Lee, S., Zhao, T., Bowler, R., Cho, M., Silverman, E. K., Hersh, C., Dy, J., & Castaldi, P. J. (2021). Improved prediction of smoking status via isoform-aware RNA-seq deep learning models. *PLoS Comput Biol*, 17(10), e1009433. doi: 10.1371/journal.pcbi.1009433.



Dr. N. Palanivel received the Bachelor of Computer Applications (BCA) from CSJM University, Kanpur, in 2001 and the Master of Computer Application (MCA) from Madurai Kamarajar University, Madurai, in 2004. He has the Masters in Engineering (Computer Science and Engineering) from Anna University Chennai in 2009 and the Ph.D (Computer Science and Engineering) from Manonmaniam

Sundaranar University, Tirunelveli, Tamil Nadu, in 2016. He has been working as an Associate Professor with the Department of Computer Science and Engineering, with almost 15+ years of experience. He is also the coordinator of the Research and Development cell at Manakula Vinayagar Institute of Technology and the project coordinator for UG and PG students in order to guide them in bringing innovative ideas from their minds into successful projects. He consistently publishes his works in reputed top-indexed journals and presents papers at international conferences on emerging topics such as artificial intelligence, machine learning, neural networks, image processing, and IOT. He served as session chair at an IEEE sponsored international conference (ICSCAN'2020, ICSCAN'2021, and ICSCAN'2023). He has reviewed papers for IEEE, reputed journals and conferences.



S. Deivanai received the B.Tech. degree in Computer Science and Engineering from Manakula Vinayagar Institute of Technology, Pondicherry University, Puducherry, in 2024. She has done some industry projects that explore emerging technologies such as artificial intelligence, machine learning, neural networks, speech recognition, and computer vision, specifically in the healthcare and finance domains. She had published her innovative project work in journals and conference papers. She had participated in and contributed projects to the hackathons and ideathons conducted by the organizations.



G. Lakshmi Priya received the B.Tech. degree in Computer Science and Engineering from Manakula Vinayagar Institute of Technology, Pondicherry University, Puducherry, in 2024. She has published a paper on "Web Application for Cardiovascular Disease Diagnosis using Data Science" in the *International Journal of Innovative Research in Technology*. She also

presented two papers at IEEE conferences, focusing on cancer diagnosis using YoloV8 and real-time object detection using CNN.



B. Sindhuja received the B.Tech. degree in Computer Science and Engineering from Manakula Vinayagar Institute of Technology, Pondicherry University, Puducherry, in 2024. She presented two papers at IEEE conferences, focusing on cancer diagnosis using YoloV8 and real-time object detection using CNN. She also presented IOT patent project named "Smart Door Lock

Automation" at Manakula Vinayagar Institute of Technology on National Science Day celebration (SCIMIT '23). She had participated in the workshop on Augmented Reality (AR) and Virtual Reality (VR).