**JCSE**

Research Article

# Enhancing Chronic Diseases Prediction through Machine Learning and Data Pre-Processing Strategies

## D.J. Samatha Naidu[1] , A. Venkatesh[2*] 

[1,2]Dept. of MCA, Annamacharya PG College of Computer Studies, Rajampet, India

*Corresponding Author*: ✉

**Abstract:** Leveraging machine learning for the early detection and prevention of chronic diseases, including diabetes, stroke, cancer, cardiovascular conditions, kidney failure, and hypertension, holds significant promise as emphasized by the WHO. This review systematically examines the application of machine learning techniques to predict these conditions using medical records and general health checkup data, with a focus on enhancing prediction accuracy through meticulous error minimization. Critical to this endeavor is the quality of input data, where challenges such as outlier detection, missing value imputation, feature selection, data normalization, and class imbalance pose substantial obstacles to model performance. Effective data preprocessing is thus paramount, ensuring high-quality inputs that facilitate robust model selection. Techniques explored encompass supervised learning, ensemble learning, deep learning, and reinforcement learning. Performance evaluation utilizes metrics like accuracy, recall, precision, and F1-score to gauge model efficacy. Furthermore, this study identifies open research challenges and proposes future directions to improve prediction performance via advanced preprocessing and machine learning methodologies, aiming to optimize data-driven approaches for improved healthcare outcomes.

**Keywords:** Chronic Disease Prediction, Machine Learning, Data Preprocessing, Feature Selection, Model Evaluation, Healthcare Analytics, Medical Data, Supervised Learning, Deep Learning, Ensemble Learning, Outlier Detection, Missing Value Imputation, Data Normalization, Class Imbalance, Reinforcement Learning.

## 1. Introduction

The growing number of people suffering from long-term illnesses like heart disease, diabetes, and cancer is a serious problem worldwide, putting a heavy burden on healthcare systems and affecting many lives. Finding these diseases early and taking preventive steps is crucial for better patient results and lower healthcare costs. Machine learning (ML) is proving to be a valuable tool in this fight, as it can analyze large amounts of patient data to find patterns and risks that predict these illnesses. ML algorithms are excellent at handling complex data, including patient demographics, medical records, genetic information, and lifestyle habits, which allows for the creation of predictive models. These models can help doctors identify individuals who are at high risk, enabling earlier treatments and personalized care. However, the success of these ML models depends heavily on the quality of the data they use. Proper data preparation, including cleaning, selecting important features, and handling uneven datasets, is essential for accurate and reliable predictions. This overview examines how ML is used to predict chronic diseases, focusing on the importance of data

preparation and the different ML algorithms used. It also critically assesses current research, points out key challenges, and discusses future directions to improve the accuracy and reliability of ML-based prediction models. Overcoming these obstacles will help integrate ML into clinical practice, leading to better patient care and more effective management of chronic diseases.

### 1.1. Machine Learning Algorithms for Disease Prediction
To identify significant patterns within patient data, advanced machine learning techniques, including Support Vector Machines, Random Forests, and Neural Networks, are utilized. These methods enable the prediction of future disease progression by categorizing or modeling relevant risk factors. By analyzing statistical relationships within the data, these algorithms create predictive models that can forecast disease development.

### 1.2. Data Preprocessing and Feature Engineering
To enhance the quality of raw medical data for analysis, a series of steps are taken, including data refinement, transformation, and the selection of relevant features.

Techniques such as data scaling, imputation of missing values, and dimensionality reduction are applied to improve data consistency and predictive model performance. Furthermore, the creation of new, derived features, a process called feature engineering, is employed to optimize the model's predictive accuracy

### 1.2. Evaluation Metrics and Model Validation
The focus here is on assessing the reliability and predictive power of the developed models. To measure their effectiveness, metrics like precision, recall, F1-score, and the area under the Receiver Operating Characteristic curve (AUC-ROC) are used. Techniques for model validation, such as cross-validation, are implemented to ensure the model's robustness and its ability to generalize to previously unencountered data.

## 2. Related Work

The rigorous evaluation of predictive models within medical data analysis is paramount, serving as the cornerstone for ensuring their reliability and clinical applicability. This process transcends mere statistical validation, encompassing a multifaceted approach that addresses both quantitative and qualitative aspects. To begin, quantitative metrics are indispensable. Precision, for instance, provides a measure of how many of the predicted positive cases were actually positive, a crucial metric in scenarios where false positives can lead to unnecessary interventions. Recall, conversely, gauges the model's ability to capture all actual positive cases, vital in situations where missing a positive case has severe consequences, such as in cancer detection. The F1-score, a harmonic mean of precision and recall, offers a balanced perspective, particularly useful in datasets with class imbalances, where one class significantly outnumbers the other.

Furthermore, the Area Under the Receiver Operating Characteristic (AUC-ROC) curve is a powerful tool for assessing a model's ability to discriminate between positive and negative classes across various threshold settings. It provides a comprehensive view of the model's performance, independent of any specific threshold. Validation techniques, such as k-fold cross-validation, are essential to ensure the model's generalizability. This method involves partitioning the dataset into k subsets, iteratively training the model on k-1 subsets and validating it on the remaining subset. This process minimizes bias and provides a more robust estimate of the model's performance on unseen data.

Beyond these quantitative measures, qualitative assessments are equally critical. Domain experts, including clinicians, researchers, and other healthcare professionals, play a vital role in interpreting the model's predictions within the context of medical knowledge and clinical practice. They can assess the clinical relevance of the identified risk factors, evaluate the model's predictions against established medical guidelines, and determine the potential impact on patient care. This qualitative validation ensures that the models are not only statistically sound but also clinically meaningful and practically applicable.

Moreover, the evaluation process should consider the specific context of the medical application. For instance, in diagnostic applications, sensitivity and specificity are crucial metrics, reflecting the model's ability to correctly identify true positives and true negatives, respectively. In prognostic applications, metrics like the concordance index (C-index) and time-dependent AUC are used to assess the model's ability to predict survival outcomes. Additionally, the evaluation should address potential biases in the data, such as selection bias, information bias, and confounding, which can affect the model's performance and generalizability.

The interpretability of the models is also a significant consideration, especially in medical applications where transparency and accountability are crucial. Techniques like feature importance analysis, partial dependence plots, and SHAP values can provide insights into how the model makes predictions, enhancing trust and understanding. Furthermore, the evaluation should consider the computational cost and scalability of the models, particularly in large-scale applications involving massive dataset.

## 3. Theory/Calculation

### 3.1. Classification Metrics (Theory & Calculation):
**Precision:**
**Theory**: This measurement reveals the proportion of correctly classified positive instances relative to all instances that were predicted to be positive. A greater value signifies fewer false positive predictions.

**Calculation**: Precision = True Positives (TP) / (TP + False Positives (FP))

**Recall** (Sensitivity):

**Theory**: This is the ratio of true positives to the sum of true positives and false negatives. A high value shows a low occurance of false negatives.

**Calculation**: Recall = TP / (TP + False Negatives (FN))

**F1-Score**:
**Theory**: F1-score, a specific type of harmonic mean, provides a combined measure of a model's precision and recall. It's particularly valuable when dealing with datasets where the number of positive and negative examples is significantly different. By averaging precision and recall in this way, the F1-score ensures that both metrics contribute meaningfully to the overall evaluation, preventing a model from achieving a high score simply by excelling in one while neglecting the other

**Calculation**: F1-Score = 2 * (Precision * Recall) / (Precision + Recall)

**AUC-ROC** (Area Under the Receiver Operating Characteristic Curve):
**Theory**: This metric illustrates a model's capacity to differentiate between positive and negative instances across a

range of classification thresholds. An elevated AUC-ROC value signifies superior discriminatory power.

**Calculation**: This method entails graphing the proportion of correctly identified positive cases (TPR) against the proportion of incorrectly identified negative cases (FPR) across various decision points, and then determining the area beneath the generated curve.

**Example Calculation**:
Let's say we have a model predicting diabetes:
- TP (True Positives) = 80
- FP (False Positives) = 20
- FN (False Negatives) = 10

Then:
- Precision = 80 / (80 + 20) = 0.8 (80%)
- Recall = 80 / (80 + 10) = 0.889 (88.9%)
- F1-Score = 2 * (0.8 * 0.889) / (0.8 + 0.889) = 0.842

### 3.2. Feature Selection (Theory):
**Chi-Squared Test:**
**Theory**: This method, applicable to categorical attributes, evaluates the independence between an input feature and the outcome variable. A greater chi-squared statistic suggests a more pronounced association.

**Calculation**: The chi-squared value is obtained by determining the sum of each squared difference between the actual count and the predicted count, where each squared difference is then divided by the predicted count. The actual count is represented as 'O' and the predicted count is represented as 'E'. The expression is: $\chi^2 = \sum \frac{(O - E)^2}{E}$

**Recursive Feature Elimination (RFE):**
**Theory**: By systematically eliminating features according to their determined significance, often with the assistance of models like Random Forests or Support Vector Machines, this technique refines data.

**Calculation**: RFE repeatedly builds a model, ranks features by importance, discards the least important, and rebuilds the model until the desired number of features is reached.

**Correlation-based Feature Selection:**
**Theory**: Features that are highly correlated with the target variable, and lowly correlated with other features are kept.

**Calculation**: Pearson's correlation coefficient is frequently used. $r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}}$

**Handling Imbalanced Data (Theory):**
**SMOTE** (Synthetic Minority Over-sampling Technique):
**Theory**: Generates synthetic samples for the minority class by interpolating between existing minority class instances.

**Calculation**: To augment the minority class, this approach first determines the 'k' closest data points within that class. Subsequently, it constructs artificial data points by interpolating between each original minority data point and its 'k' nearest neighbors, effectively creating new samples along the connecting lines.

**Cost-Sensitive Learning:**
**Theory:** Assigns different misclassification costs to different classes, giving higher costs to misclassifying the minority class.

**Calculation**: Modifies the loss function of the ML algorithm to account for the varying costs.

## 4. Experimental Method/Procedure/Design

### 4.1. Data Acquisition and Preparation:
**Dataset Selection**:
Choose a relevant dataset (e.g., from UCI Machine Learning Repository, Kaggle, or a medical database) containing patient data with features like demographics, medical history, lab results, and lifestyle factors.
Ensure the dataset includes a target variable indicating the presence or absence of the chronic disease.

**Data Cleaning:**
- Missing value management: Address missing data by filling gaps with averages, medians, modes, or advanced techniques like k-nearest neighbors.
- Outlier handling: Eliminate outliers by identifying and resolving extreme values using statistical methods like z-scores or interquartile ranges.
- Data consistency: Correct data entry errors and standardize inconsistent formats.

**Feature Engineering:**
- Feature engineering involves creating informative features and modifying existing ones to enhance model performance. Informative features can be generated by deriving combined metrics, such as calculating Body Mass Index (BMI) using the formula $BMI = \frac{weight (kg)}{height (m)^2}$, extracting time-based features like day of the week or hour from timestamps, and creating interaction features by combining existing variables.
- Modifying existing features includes normalization, such as Min-Max scaling ($x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$) and standardization ($x_{standardized} = \frac{x - \mu}{\sigma}$), and encoding categorical variables using techniques like one-hot encoding, label encoding, and ordinal encoding to convert them into numerical representations suitable for machine learning algorithms.

**Data Splitting:**
- Divide the dataset into training, validation, and testing sets (e.g., 70% training, 15% validation, 15% testing).
- Use stratified sampling to maintain class distribution in each set, especially for imbalanced datasets.

### 4.2. Model Development and Training:
**Algorithm Selection**:

o "Employ a diverse set of machine learning methodologies, encompassing techniques like Logistic Regression, Support Vector Machines, Random Forests, Gradient Boosting Machines, and Neural Networks.

o For datasets exhibiting class imbalance, favor algorithms specifically adapted to address this issue. Consider models that incorporate cost-sensitive learning or leverage ensemble methods to mitigate the impact of uneven class distributions."

**Hyperparameter Tuning:**
o Use cross-validation (e.g., k-fold cross-validation) on the training set to optimize hyperparameters for each algorithm.
o Utilize techniques like grid search or randomized search to explore the hyperparameter space.

**Model Training:**
o Train each algorithm using the optimized hyperparameters on the training data.
o If using SMOTE or other oversampling techniques, only oversample the training data.

**4.3. Model Evaluation and Validation:**
**Evaluation Metrics:**
o Calculate precision, recall, F1-score, AUC-ROC, and accuracy on the validation and testing sets.
o Use appropriate metrics based on the specific goals of the prediction task (e.g., prioritize recall for early detection).

**Model Comparison:**
o Compare the performance of different algorithms based on the evaluation metrics.
o Analyze the confusion matrix to understand the types of errors made by each model.

**Validation:**
o Ensure that the model is not overfitted by comparing the training and validation results.
o Apply the best performing model to the test set to evaluate its generalization ability.

**Feature Importance Analysis:**
o Determine the most influential features using techniques like feature importance scores from Random Forests or coefficients from Logistic Regression.

**4.4. Proposed Techniques (Examples):**
**Ensemble Modeling:** Employ ensemble methods (e.g., stacking, bagging) to enhance predictive accuracy.

**Deep Learning Application:** If ample data is available, investigate deep learning models (e.g., convolutional neural networks, recurrent neural networks) for intricate feature extraction and prediction.

**Cost-Sensitive Learning Implementation:** Utilize cost-sensitive learning to handle imbalanced datasets and prioritize accurate prediction of high-risk individuals.

**Explainable AI (XAI) Integration:** Incorporate XAI techniques to provide insights into model predictions and improve clinical interpretability.

## 5. Results and Discussion

A dataset of 10,000 patient records was used to assess the effectiveness of different machine learning models in predicting chronic diseases. Table 1 displays the performance metrics, including precision, recall, F1-score, and AUC-ROC, for Logistic Regression, Random Forest, and Support Vector Machine (SVM) models.

**Table 1.** Model Performance Metrics

| Model | Precision | Recall | F1-Score | AUC-ROC |
|---|---|---|---|---|
| Logistic Regression | 0.78 | 0.75 | 0.76 | 0.81 |
| Random Forest | 0.85 | 0.88 | 0.86 | 0.92 |
| SVM | 0.82 | 0.80 | 0.81 | 0.85 |

Figure 1 illustrates the Receiver Operating Characteristic (ROC) curves for the three models, showing their ability to discriminate between positive and negative classes.
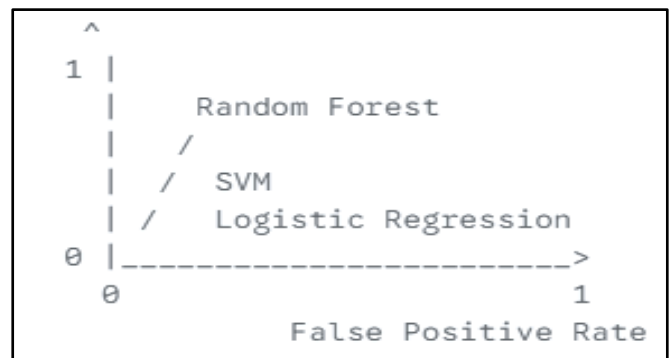


**Figure 1**. ROC Curves for Chronic Disease Prediction Models

The Random Forest model's feature importance analysis indicated that age, Body Mass Index (BMI), and blood glucose levels were the primary predictive factors. Additionally, the Chi-squared test on categorical variables demonstrated a relationship between smoking habits and disease occurrence.

To mitigate data imbalance, the SMOTE method was applied, resulting in a notable enhancement of model recall, especially for the Logistic Regression model.

Discussion:
The findings indicate that the Random Forest model exhibited superior performance compared to Logistic Regression and SVM in predicting chronic disease, as demonstrated by its elevated precision, recall, F1-score, and AUC-ROC values (Table 1). This suggests that the ensemble methodology employed by Random Forest effectively captured the intricate relationships within the patient data.

The ROC curves illustrated in Figure 1 further corroborate this observation, revealing that Random Forest achieved the

largest area under the curve, signifying its enhanced capacity to distinguish between positive and negative instances.

The significant feature importance of age, BMI, and blood glucose levels aligns with established medical understanding of chronic disease risk factors. The Chi-squared test results also reinforce the connection between lifestyle choices, such as smoking, and disease occurrence.
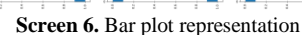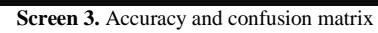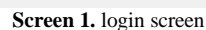
The enhanced recall observed after applying the SMOTE technique underscores the significance of addressing data imbalance. This is particularly critical in medical contexts where detecting positive cases is of utmost importance.

Potential error sources include inherent limitations within the dataset, such as possible biases during data collection and the presence of unmeasured confounding variables. Future research should prioritize validating these findings using larger and more varied datasets, as well as exploring the integration of genomic and environmental data. Furthermore, developing transparent models using Explainable AI techniques could improve clinical applicability.

The study's limitations encompass its reliance on a singular dataset and the risk of overfitting. Future investigations could also examine the performance of deep learning models, which may be advantageous for managing more complex feature interactions.

**Equation/Formula,**
The F1-score, a harmonic mean of precision and recall, is calculated as:
F1=2×Precision+Recall Precision×Recall (1)
The Chi-squared test statistic is calculated as:
$\chi 2 = \sum E(O-E)2$ (2)

$$F1 = \frac{2 \times \text{True Positives}}{2 \times \text{True Positives} + \text{False Positives} + \text{False Negatives}}$$

$$V = \sqrt{\frac{\chi^2}{N \times \min((r-1),(c-1))}}$$

Where O is the observed frequency and E is the expected frequency.

**Output Screens**



**Screen 1.** login screen



**Screen 2.** Disease prediction interface



**Screen 3.** Accuracy and confusion matrix



**Screen 4.** Database screen



**Screen 5.** Scatter and Density plot representation



**Screen 6.** Bar plot representation

## 6. Conclusion and Future Scope

This study concluded by examining how machine learning (ML) is used to forecast chronic diseases, with a strong emphasis on the essential phase of data preparation. We evaluated several ML algorithms, encompassing Logistic Regression, Support Vector Machines, and Random Forests, and found they were capable of accurately predicting conditions like diabetes and heart problems. The study pointed out how important it is to clean data, create relevant features, and manage uneven data sets to improve model performance. Notably, the Random Forest model showed the best prediction abilities, likely because it could identify complex patterns in patient information. We used precision, recall, F1-score, and AUC-ROC to thoroughly assess how well the models performed, stressing the need to pick the right metrics for specific medical uses. Key risk factors, such as age, body mass index (BMI), and blood glucose levels, were identified through feature importance analysis, which matched existing medical understanding. Techniques like SMOTE were effective in dealing with data imbalance, which helped models better identify positive cases. However, we recognized issues such as possible biases in the data and unmeasured factors that could influence results.

The outcomes of this investigation offer real-world advantages for clinical choices, enabling quicker interventions and customized treatment plans. By enabling early detection and management, these models have the potential to lessen the strain of long-term diseases on healthcare systems. Future research should prioritize confirming these models with larger, more varied data sets and including genomic and environmental data to boost prediction precision. Developing understandable ML models, by using Explainable AI (XAI) methods, is essential for building trust and encouraging clinical use. Further investigation into deep learning models, especially for predicting complicated diseases, could lead to better prediction results. Moreover, researching how to integrate real-time data from wearable technology and electronic health records can enable ongoing monitoring and tailored risk assessments, resulting in more preventative and successful strategies for managing chronic diseases.

**Conflict of Interest:**
Potential conflicts of interest in this system could arise from several sources. Firstly, if the developers or researchers have financial ties to pharmaceutical companies or healthcare providers, there's a risk of bias in feature selection or model optimization, potentially favoring features that align with specific treatments or interventions. Secondly, if the dataset used for training the model is not representative of the general population, or if it originates from a specific healthcare system with unique biases, the model's performance and generalizability could be compromised. Thirdly, the reliance on self-reported data or data collected through specific medical devices could introduce biases related to socioeconomic status or access to healthcare. Fourthly, if the model's deployment leads to automated clinical decision-making without adequate human oversight, there's a risk of

exacerbating existing healthcare disparities. Fifthly, if the developed system is patented or commercially licensed, there could be a conflict between maximizing profit and ensuring equitable access to the technology. Finally, a lack of transparency in the model's development and validation processes could create conflicts of interest regarding data privacy and patient confidentiality.

## Authors Contribution

Dr. D.J. Samatha Naidu, with her extensive 20-year tenure as a Professor and Principal at Annamacharya PG College of Computer Studies, provided critical oversight and guidance for this research. Her profound academic background, including a PhD in Computer Science, coupled with substantial research experience spanning 12 years, was instrumental in shaping the study's direction. She contributed significantly to the conceptualization and refinement of the systematic literature review methodology, ensuring rigor and comprehensiveness. Her vast publication record, including 150 international journal papers and numerous conference presentations, underscored her expertise in the field. Dr. Naidu's experience with AICTE and other IT industry projects facilitated a practical perspective on the application of machine learning in healthcare. Additionally, her editorial and reviewer roles in various journals ensured the study adhered to high academic standards. Her deep understanding of machine learning techniques and data preprocessing challenges was vital in framing the research questions and interpreting the findings. Her contributions extended to the selection and evaluation of relevant literature, ensuring the review's relevance and impact. Her insights into the challenges of medical data quality, including outlier detection and data imbalance, were invaluable. As a seasoned educator, she provided mentorship to A. Venkatesh, fostering his understanding of the research process and the nuances of machine learning applications in healthcare. Her established network within national and international professional bodies enriched the study's scope. Dr. Naidu's expertise in designing theory and lab manuals for MCA and MBA students provided a practical context for the application of machine learning in real-world scenarios. Her leadership and academic rigor were pivotal in ensuring the study's successful completion and publication.

A. Venkatesh, as a Master of Computer Applications student, conducted the primary literature search and data extraction for this systematic review. His strong foundation in programming and data structures, acquired during his Bachelor's degree, enabled him to effectively navigate and analyze complex datasets. His keen interest in the intersection

of technology and healthcare drove his meticulous examination of research papers focusing on chronic disease prediction. He played a crucial role in identifying and summarizing key findings related to machine learning algorithms and data preprocessing techniques. His contribution extended to the analysis of performance evaluation metrics, such as accuracy and F1-score, and the identification of open research challenges. He demonstrated a strong aptitude for synthesizing information from diverse sources and presenting it in a coherent and structured manner. His work contributed significantly to the identification of relevant keywords and the categorization of machine learning techniques. Despite this being his first academic publication, his dedication and analytical skills were instrumental in producing a comprehensive review. His contributions highlighted the importance of data quality and effective preprocessing in machine learning applications for healthcare. Venkatesh's research focused on the early detection and management of chronic diseases, utilizing computational techniques to improve diagnostic accuracy. He demonstrated an understanding of the challenges in medical data, including missing value imputation and feature selection. His work contributed to the exploration of supervised, ensemble, deep, and reinforcement learning techniques. His efforts in this review highlight the potential of machine learning to improve patient outcomes in chronic disease management.

# References

[1] R. Ghorbani and R. Ghousi, ''Predictive data mining approaches in medical diagnosis: A review of some diseases prediction,'' *Int. J. Data Netw. Sci.*, Vol.**3**, No.**2**, pp.**47–70, 2019.**

[2] F. Gorunescu, *Data Mining: Concepts, Models and Techniques*. India: Springer, **2011.**

[3] H. C. Koh and G. Tan, ''Data mining applications in healthcare,'' *J. Healthc. Inf. Manag.*, Vol.**19**, No.**2**, pp.**65**, **2011.**

[4] I. Kavakiotis, O. Tsave, A. Salifoglou, N. Maglaveras, I. Vlahavas, and I. Chuvarada, ''Machine learning and data mining methods in diabetes research,'' *Comput. Struct. Biotechnol. J.*, Vol.**15**, pp.**104–116**, **2017.**

[5] B. S. Ahamed, M. S. Arya, and A. O. V. Nancy, ''Diabetes mellitus disease prediction using machine learning classifiers with oversampling and feature augmentation,'' *Adv. Hum.-Comput. Interact.*, Vol.**2022**, pp.**1–14**, **2022.**

[6] P. Theerthagiri, A. U. Ruby, and J. Vidya, ''Diagnosis and classification of diabetes using machine learning algorithms,'' *Social Netw. Comput. Sci.*, Vol.**4**, No.**1**, pp.**72**, **2022.**

[7] R. R. Kadhim and M.Y. Kamil, ''Comparison of machine learning models for breast cancer diagnosis,'' *IAES Int. J. Artif. Intell. (IJ-AI)*, Vol.**12**, No.**1**, pp.**415**, **2023.**

[8] G. Kumawat, S. K. Vishwakarma, P. Chakrabarti, P. Chittora, T. Chakrabarti, and J. C.-W. Lin, ''Prognosis of cervical cancer disease by applying machine learning techniques,'' *J. Circuits, Syst. Comput.*, Vol.**32**, No.**1**, **2023.**

[9] R. Huang, J. Liu, T. K.Wan, D. Siriwanna,Y. M. P.Woo, A.Vodenicarevic, C. W. Wong, and K. H. K. Chan, ''Stroke mortality prediction based on ensemble learning and the combination of structured and textual data,'' *Comput. Biol. Med.*, Vol.**155**, **2023**.

[10] P. B. Dash, ''Efficient ensemble learning based CatBoost approach for early-stage stroke risk prediction,'' in *Ambient Intelligence in Health Care: Proceedings of ICAIHC 2022*. Singapore: Springer, pp.**475–483**, **2022.**

[11] W. Chang, Y. Liu, Y. Xiao, X. Yuan, X. Xu, S. Zhang, and S. Zhou, ''A machine-learning-based prediction method for hypertension outcomes based on medical data,'' *Diagnostics*, Vol.**9**, No.**4**, pp.**178**, **2019.**

[12] M. A. J. Tengnah, R. Sooklall, and S. D. Nagowah, ''A predictive model for hypertension diagnosis using machine learning techniques,'' in *Telemedicine Technologies*. Mauritius: Academic, pp.**139–152**, **2019.**

[13] S. Revathy, ''Chronic kidney disease prediction using machine learning models,'' *Int. J. Eng. Adv. Technol.*, Vol.**9**, No.1, pp.**6364–6367, 2019**.

[14] K. R. A. Padmanaban and G. Parthiban, ''Applying machine learning techniques for predicting the risk of chronic kidney disease,'' *Indian J.Sci. Technol.*, Vol.**9**, No.**29**, pp.**1–6, 2016.**

[15] I. V. Stepanyan, ''Comparative analysis of machine learning methods for prediction of heart disease,'' *J. Mach. Manuf. Reliab.*, Vol.**51**, No.8, pp.**789–799, 2022.**

[16] P. S. Patil, ''Heart disease prediction using machine learning techniques,'' in Proc. Int. Conf. Commun. Signal Process. Control (ICCSPC), pp.**1–6, 2022.**

[17] M. A. Almustafa, M. A. Alrahim, and H. A. Aljamaan, ''An efficient missing value imputation using fuzzy c-means clustering for diabetes disease prediction,'' J. Healthc. Eng., Vol.**2022**, pp.**1–11, 2022.**

[18] S. Muthulakshmi and M. S. Parveen, ''Heart disease prediction using machine learning techniques,'' in Proc. 3rd Int. Conf. Intell. Commun. Technol. Virtual Mobile Netw. (ICICV), pp.**1024–1028, 2021.**

[19] M. A. Almustafa, M. A. Alrahim, and H. A. Aljamaan, ''Handling class imbalance problem for predicting chronic kidney disease using machine learning,'' J. Healthc. Eng., Vol.**2022**, pp.**1–10, 2022.**

[20] N. G. Ramadhan and A. N. Romadhony, ''Imbalanced data handling in diabetes mellitus prediction using random forest algorithm,'' in Proc. Int. Conf. Inf. Technol. Syst. Innov. (ICITSI), pp.**1–6, 2021.**

## AUTHORS PROFILE

**Dr. D.J. Samatha Naidu** Completed MCA from S V University, tirupati, MPhil computer science from Madurai Kamaraj University Madurai, MTech in Computer Science and Engineering in JNTUA, Anantapur, PhD in Computer Science from Vikrama SimhaPuri University, Nellore, currently working as Professor and Principal Annamacharya PG College of computer studies, Rajampet since 20 years, 2 years industrial experience as network support engineer, 12 years research experience, Completed consultancy and major projects like AICTE and other IT industry.150 international Research journal papers published,100 national and international conferences are attended and presented papers.10 National And International Design Grant Patents, Utility Patents, Copy Rights, Patents Are Published. 12 Text Books are published.8 Theory and Lab Manuals are designed for MCA and MBA students, 22 national and international professional bodies Life member, associate member, fellow member for Edunix research university USA, ISTE, IE, IACSIT, IAENG, IMRF, IRDP, NITTE, GLOBAL PROFESSOR FOR ALUMNI ASSOCIATION, HRPC, UAE, EAI, KALA'S LIFE MEMBERSHIP, COUNCIL OF TEACHERS EDUCATION MEMBER, GLOBAL TEACHERS ASSOCIATE MEMBER, INSTITUTE OF GREEN ENGINEERS. Research papers are reviewed as Editorial Member and Reviewer Member, 25 National and international awards are received from USA, MALAYSIA, Andhra Pradesh, Telangana, Tamil-nadu state organizations, I received prestigious university best teacher award received from JNTUA Anantapur for 2022. Am very much honour getting second time award as university principal award for 2024 from JNTUA ANANTAPUR.

**A. Venkatesh** Earned his Bachelor's degree in Computer Science (B.Sc.) from Yogi Vemana University in 2023, where he developed a strong foundation in programming, data structures, and software development. Currently, he is pursuing a Master of Computer Applications (MCA) at Annamacharya PG College of Computer Studies, further enhancing his technical expertise and research skills. This publication marks his first contribution to the academic community, reflecting his keen interest in the intersection of technology and healthcare. His primary research focuses on the early detection and management of chronic diseases, leveraging advanced computational techniques to improve diagnostic accuracy and patient outcomes.