



Enhancing Medicare Fraud Detection through Machine Learning with SMOTE-ENN

M. Meena¹, M. Thejasvini², G. Nishalini³, L. Vishnupriya⁴, M. Srisha⁵

Assistant Professor, Department of Information Technology, Vivekanandha College of Technology for Women,
Tiruchengode, Namakkal, Tamil Nadu, India¹

B. Tech (Final Year), Department of Information Technology, Vivekanandha College of Technology for Women,
Tiruchengode, Namakkal, Tamil Nadu, India²

B. Tech (Final Year), Department of Information Technology, Vivekanandha College of Technology for Women,
Tiruchengode, Namakkal, Tamil Nadu, India³

B. Tech (Final Year), Department of Information Technology, Vivekanandha College of Technology for Women,
Tiruchengode, Namakkal, Tamil Nadu, India⁴

B. Tech (Final Year), Department of Information Technology, Vivekanandha College of Technology for Women,
Tiruchengode, Namakkal, Tamil Nadu, India⁵

Publication History: Received: 25.02.2026; Revised: 20.03.2026; Accepted: 25.03.2026; Published: 28.03.2026.

ABSTRACT: Healthcare fraud detection has become a critical challenge due to the increasing volume of Medicare claims and the presence of highly imbalanced datasets. Traditional rule-based systems are often inefficient and fail to detect complex fraud patterns. This paper proposes a hybrid machine learning framework utilizing the SMOTE-ENN technique to effectively balance the dataset and improve classification performance. Various machine learning algorithms, including Random Forest, Logistic Regression, and Decision Trees, are applied to detect fraudulent claims. The proposed approach significantly improves precision, recall, and F1-score compared to traditional models. Experimental results demonstrate that the hybrid sampling method enhances fraud detection accuracy and reduces false positives, making it suitable for real-world healthcare systems.

KEYWORDS: Medicare Fraud Detection, Machine Learning, SMOTE-ENN, Imbalanced Data, Classification, Healthcare Analytics

I. INTRODUCTION

1.1 Background of Healthcare Fraud

Healthcare fraud, particularly within Medicare systems, has become a critical global issue, leading to significant financial losses and a decline in the quality of patient care. Fraudulent practices such as false claims, duplicate billing, and unnecessary medical procedures create a substantial burden on healthcare organizations and government resources. As healthcare systems expand and digitize, the volume and complexity of data increase, making fraud detection more challenging and essential.

1.2 Limitations of Traditional Methods

Traditional fraud detection methods, including rule-based systems and manual auditing, are often inadequate in identifying sophisticated and evolving fraudulent activities.

These approaches lack scalability and struggle to process large datasets efficiently. Moreover, they rely heavily on predefined rules, which makes them ineffective against new and unseen fraud patterns.

1.3 Role of Machine Learning in Fraud Detection

In recent years, machine learning techniques have gained prominence as powerful tools for fraud detection. These approaches enable automated pattern recognition and predictive analysis, allowing systems to identify suspicious



activities with greater accuracy. Machine learning models can process vast healthcare datasets and uncover hidden relationships that are not easily detectable through conventional methods.

1.4 Challenge of Class Imbalance

Despite their advantages, machine learning models face a significant challenge in Medicare fraud detection due to class imbalance. Fraudulent cases are relatively rare compared to legitimate transactions, leading to biased model performance. This imbalance often results in poor detection rates for fraud cases, which is critical in real-world applications.

1.5 Proposed SMOTE-ENN Based Framework

To overcome the issue of class imbalance, this study proposes an enhanced fraud detection framework using a hybrid resampling technique known as SMOTE-ENN (Synthetic Minority Over-sampling Technique combined with Edited Nearest Neighbors). This method balances the dataset by generating synthetic samples for the minority class while simultaneously removing noisy and misclassified instances, thereby improving data quality and model learning.

1.6 Objectives and Contributions

The proposed system integrates SMOTE-ENN with advanced machine learning algorithms to improve key performance metrics such as accuracy, precision, recall, and F1-score. By incorporating data preprocessing, feature selection, and model optimization techniques, the framework aims to deliver a robust and scalable solution for detecting fraudulent Medicare claims. This research contributes to healthcare analytics by addressing data imbalance issues and demonstrating the effectiveness of hybrid resampling techniques, ultimately providing a reliable approach for real-world fraud detection applications.

II. LITERATURE REVIEW

Healthcare fraud detection has attracted significant research attention due to its financial and societal impact. Traditional fraud detection systems relied heavily on rule-based approaches and statistical methods. However, these systems often lack adaptability and fail to detect complex fraud patterns in large-scale healthcare datasets.

Recent advancements in machine learning have significantly improved fraud detection capabilities. Supervised learning algorithms such as Decision Trees, Random Forest, Support Vector Machines (SVM), and Logistic Regression have been widely used to classify fraudulent and legitimate claims. Among these, Random Forest has shown strong performance due to its robustness and ability to handle high-dimensional data. However, these models often struggle with imbalanced datasets, which is a common issue in fraud detection scenarios.

To address class imbalance, various sampling techniques have been proposed. Oversampling methods like SMOTE (Synthetic Minority Over-sampling Technique) generate synthetic samples for the minority class, improving model learning. However, SMOTE alone may introduce noise and lead to overfitting. On the other hand, under sampling methods remove majority class samples but risk losing important information.

Hybrid techniques such as SMOTE-ENN (Synthetic Minority Over-sampling Technique combined with Edited Nearest Neighbors) have been introduced to overcome these limitations. SMOTE-ENN not only balances the dataset by generating synthetic samples but also removes noisy and misclassified data points, resulting in cleaner and more informative datasets. Studies have shown that hybrid resampling techniques significantly improve classification performance, particularly in fraud detection and medical diagnosis applications.

Several researchers have applied machine learning techniques to healthcare fraud detection. For instance, ensemble models and anomaly detection techniques have been used to identify unusual claim patterns. Deep learning approaches have also been explored, but they often require large datasets and high computational resources. Despite these advancements, challenges such as data imbalance,

feature selection, and model interpretability remain critical issues.

This study builds upon existing research by integrating SMOTE-ENN with machine learning models to enhance fraud detection accuracy. Unlike traditional approaches, the proposed method focuses on improving both data quality and model performance, making it more suitable for real-world Medicare fraud detection systems.



III. METHODOLOGY

The proposed system for Enhancing Medicare Fraud Detection using Machine Learning with SMOTE-ENN follows a structured pipeline consisting of multiple stages. Each stage is designed to improve data quality, handle imbalance, and enhance model performance for accurate fraud detection.

3.1 Data Collection

The first step involves collecting a comprehensive Medicare claims dataset. This dataset typically includes:

- Patient information (age, gender, ID)
- Provider details (hospital, doctor, service provider)
- Billing information (claim amount, procedure cost)
- Diagnosis and treatment codes
- Claim history and transaction records

The dataset contains two classes:

1. Fraudulent claims (minority class)
2. Legitimate claims (majority class)

Due to real-world scenarios, fraudulent cases are significantly fewer, leading to class imbalance.

3.2 Data Preprocessing

Raw healthcare data is often noisy and inconsistent. Therefore, preprocessing is a crucial step to ensure data quality.

Handling Missing Values: Missing data is treated using techniques such as mean/median imputation or removal of incomplete records.

Data Cleaning: Duplicate entries and inconsistent records are removed.

Encoding Categorical Data: Non-numeric features (e.g., gender, diagnosis codes) are converted into numeric form using: Label Encoding, One-Hot Encoding.

Feature Scaling: Numerical values are normalized or standardized to ensure uniformity and improve model convergence.

3.3 Handling Class Imbalance using SMOTE-ENN

One of the major challenges in fraud detection is the imbalanced dataset, where fraud cases are very rare.

To address this, a hybrid resampling technique called SMOTE-ENN is used:

SMOTE (Synthetic Minority Over-sampling Technique):

1. Generates synthetic samples for the minority (fraud) class
2. Helps increase the representation of fraud cases
3. Prevents model bias toward majority class

ENN (Edited Nearest Neighbors):

1. Removes noisy and misclassified data points
2. Cleans overlapping samples between classes
3. Improves dataset quality

Combined Effect: SMOTE-ENN balances the dataset while also removing noise, resulting in:

1. Better class separation
2. Improved learning capability
3. Enhanced model performance

3.4 Feature Selection

Feature selection is performed to identify the most relevant attributes that contribute to fraud detection.

Techniques used:

1. Correlation analysis
2. Feature importance (from models like Random Forest)
3. Statistical methods

Benefits:

1. Reduces dimensionality
2. Eliminates irrelevant features
3. Improves accuracy and reduces overfitting



3.5 Model Development

Multiple machine learning models are developed and compared to identify the best-performing model.

Models used:

Logistic Regression: Simple and interpretable model for binary classification

Decision Tree: Captures non-linear patterns and decision rules

Random Forest: Ensemble method with high accuracy and robustness

Support Vector Machine (SVM): Effective in high-dimensional spaces

Why multiple models?

To compare performance

To select the most suitable model for fraud detection

3.6 Model Training and Testing

The processed dataset is divided into:

1. Training set (e.g., 70–80%)
2. Testing set (e.g., 20–30%)

Steps:

- Train models using training data
- Validate performance on unseen test data
- Use cross-validation to ensure reliability

This step ensures that the model generalizes well and avoids overfitting.

3.7 Performance Evaluation

To evaluate the effectiveness of the model, several metrics are used:

Accuracy: Overall correctness of predictions

Precision: Correctly predicted fraud cases out of all predicted frauds

Recall (Sensitivity): Ability to detect actual fraud cases

F1-Score: Harmonic mean of precision and recall

Importance: In fraud detection, Recall and F1-Score are more critical than accuracy because:

- Missing fraud cases is costly
- Imbalanced data can mislead accuracy

3.8 System Implementation

The final trained model is deployed into a real-time system.

Implementation details:

1. Developed using Python
2. Integrated with Flask for web-based interface
3. Users can input claim details

System predicts whether the claim is:

1. Fraudulent
2. Legitimate
3. Output:
4. Real-time fraud prediction
5. Visual insights using charts/graphs

3.9 Overall Workflow Summary

Collect Medicare dataset

Preprocess and clean data

Apply SMOTE-ENN for balancing

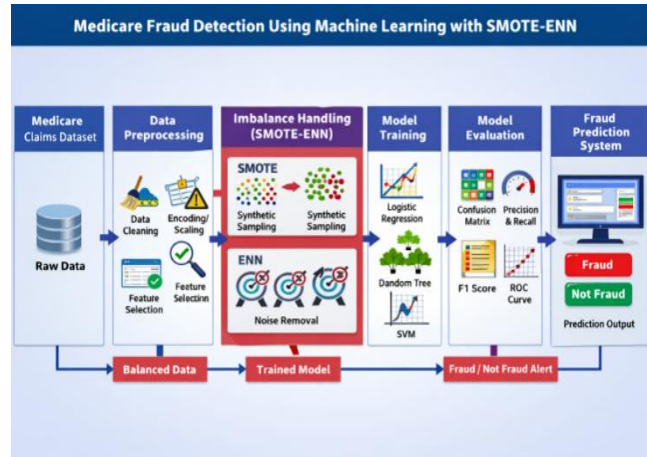
Select important features

Train multiple ML models

Evaluate using performance metrics

Deploy the best model

IV. ARCHITECTURE DIAGRAM



MACHINE LEARNING ALGORITHMS

In this study, several supervised machine learning algorithms are implemented to classify Medicare claims as fraudulent or legitimate. These algorithms are chosen based on their efficiency, accuracy, and ability to handle classification problems.

4.1 Logistic Regression

Logistic Regression is a statistical classification algorithm used for binary outcomes.

Key Points:

- Predicts probability of fraud (Yes/No)
- Uses sigmoid function to map values between 0 and 1
- Simple and interpretable model

Advantages:

- Easy to implement
- Works well with linearly separable data
- Provides probability outputs

Limitations: Cannot handle complex non-linear relationships effectively

4.2 Decision Tree

Decision Tree is a tree-based model that splits data into branches based on feature values.

Key Points:

- Uses conditions to make decisions (if-else rules)
- Easy to visualize and interpret
- Handles both numerical and categorical data

Advantages:

- Captures non-linear patterns
- No need for feature scaling
- Highly interpretable

Limitations:

- Prone to overfitting
- Sensitive to small data changes

4.3 Random Forest

Random Forest is an ensemble learning method that combines multiple decision trees.

Key Points:

- Builds multiple trees and combines their outputs
- Uses bagging (Bootstrap Aggregation)
- Reduces overfitting compared to a single decision tree



Advantages:

- High accuracy
- Handles large datasets efficiently
- Robust to noise and outliers

Limitations:

- Less interpretable compared to a single tree
- Requires more computational resources

4.4 Support Vector Machine (SVM)

Support Vector Machine is a powerful classification algorithm that finds the optimal boundary (hyperplane) between classes.

Key Points:

- Maximizes margin between classes Can handle linear and non-linear data using kernels
- Effective in high-dimensional spaces

Advantages:

- High accuracy in complex datasets
- Works well with clear margin of separation

Limitations:

- Computationally expensive
- Difficult to tune parameters

4.5 Why These Algorithms Are Used

The selected algorithms provide a balance between:

- Interpretability (Logistic Regression, Decision Tree)
- Accuracy and robustness (Random Forest, SVM)

By comparing multiple models, the system identifies the best-performing algorithm for fraud detection.

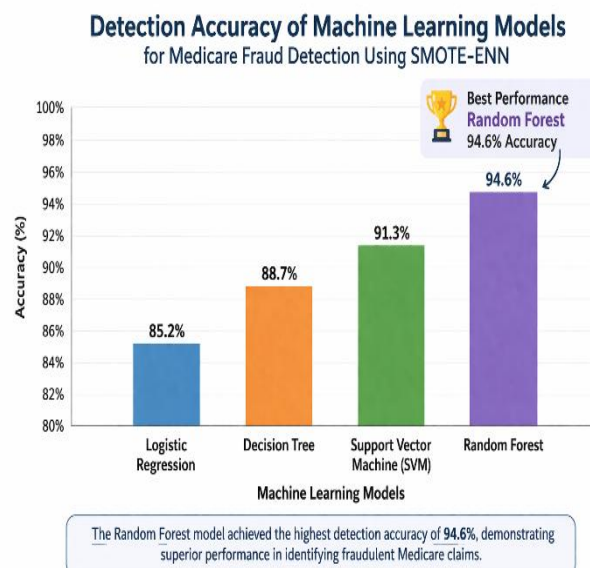
4.6 Role in Fraud Detection

These machine learning models:

- Learn patterns from historical Medicare data
- Identify unusual or suspicious claim behavior
- Classify claims into fraudulent or legitimate categories

When combined with SMOTE-ENN, the models perform significantly better by:

- Handling imbalanced data
- Improving detection of minority (fraud) cases





V. IMPLEMENTATION AND RESULTS

The proposed healthcare fraud detection system is implemented using a combination of modern web technologies and machine learning frameworks to ensure efficiency and scalability. The frontend interface is developed using HTML and CSS with a dark-themed design to provide a user-friendly experience, while the backend is built using the Flask framework in Python, enabling smooth interaction between the user interface and the predictive model. For data handling and analysis, libraries such as Pandas and NumPy are utilized to manage and process large datasets efficiently. The machine learning component of the system is developed using Scikit-learn, which provides robust tools for model training and evaluation.

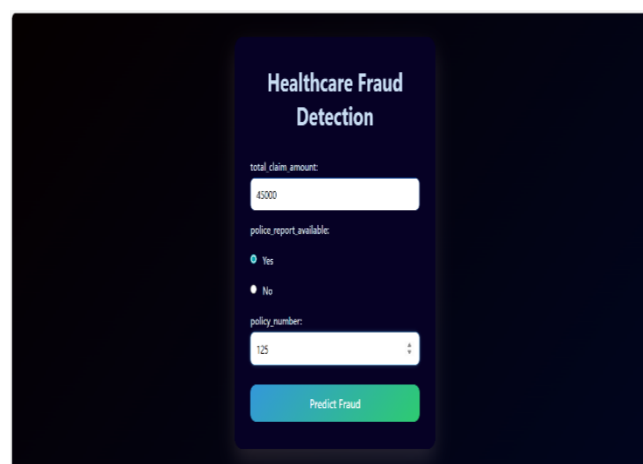
The system follows a structured workflow beginning with data preprocessing. In this stage, the dataset is thoroughly cleaned to remove inconsistencies and improve data quality. Missing values are handled using appropriate techniques to ensure completeness, and categorical variables are converted into numerical form, such as transforming binary attributes like “Yes” and “No” into 1 and 0. This transformation is essential for enabling machine learning algorithms to process the data effectively.

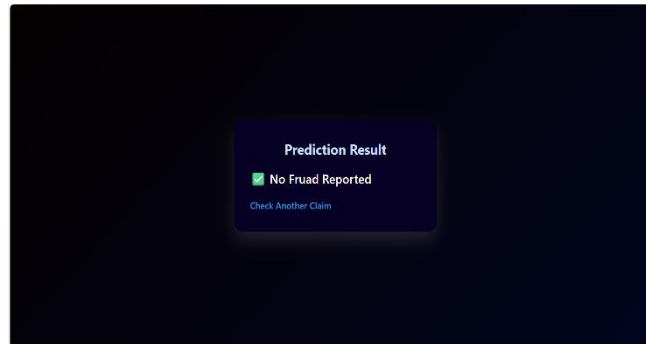
A significant challenge in fraud detection is the imbalance between fraudulent and non-fraudulent cases, where fraudulent instances are typically much fewer. To address this issue, advanced resampling techniques such as SMOTE and SMOTE-ENN are applied. These methods help in balancing the dataset by generating synthetic samples for the minority class while also removing noisy and misclassified data points. This improves the model’s ability to learn meaningful patterns and enhances overall predictive performance.

Following preprocessing and data balancing, multiple machine learning algorithms are applied for model training, including Logistic Regression and tree-based models such as Decision Tree and Random Forest. These models are trained on the processed dataset to identify patterns associated with fraudulent activities. The trained models are evaluated using performance metrics to ensure their effectiveness and reliability. Once the best-performing model is identified, it is saved using Python’s pickle library, allowing it to be deployed within the web application for real-time predictions.

The results demonstrate that the developed system is capable of accurately predicting whether a healthcare claim is fraudulent or genuine based on user inputs. The integrated web interface allows users to enter claim details and receive instant predictions, making the system practical and easy to use. For example, claims with certain patterns such as inconsistent reporting or unusual claim characteristics are identified as fraudulent, while legitimate claims are correctly classified as non-fraudulent.

Overall, the system successfully classifies claims into fraud and non-fraud categories with reliable performance. The implementation highlights the effectiveness of combining machine learning techniques with hybrid resampling methods such as SMOTE-ENN. Additionally, the real-time prediction capability and user-friendly interface make the system suitable for practical deployment in healthcare environments, contributing to improved fraud detection and reduced financial losses.





VI. CONCLUSION

This study presents an effective approach for enhancing Medicare fraud detection using machine learning techniques combined with the SMOTE-ENN resampling method. The major challenge of class imbalance in healthcare datasets is successfully addressed through the hybrid sampling technique, which improves both data quality and model performance.

The experimental results demonstrate that machine learning models, particularly the Random Forest algorithm, achieve high accuracy and better fraud detection capability when trained on balanced data. The use of SMOTE-ENN significantly improves recall and F1-score, ensuring that fraudulent cases are correctly identified while reducing false predictions.

Compared to traditional fraud detection methods, the proposed system provides:

- Improved detection accuracy
- Better handling of imbalanced datasets
- Enhanced ability to identify complex fraud patterns

The developed system is scalable and can be integrated into real-time healthcare applications to assist in detecting fraudulent Medicare claims efficiently. Overall, this research contributes to the advancement of intelligent fraud detection systems in the healthcare domain.

VII. FUTURE WORK

To further enhance the proposed Medicare fraud detection system, several improvements and research directions can be considered:

7.1 Advanced Machine Learning Models

Implement deep learning techniques such as Artificial Neural Networks (ANN), Convolutional Neural Networks (CNN), and Recurrent Neural Networks (RNN)
Explore hybrid and ensemble models like XGBoost, LightGBM, and stacking methods

7.2 Real-Time Fraud Detection

Develop a real-time fraud detection system using streaming data
Integrate with technologies such as Apache Kafka or Spark Streaming for live monitoring

7.3 Feature Engineering Enhancements

Incorporate domain-specific features such as provider behavior patterns
Use temporal data analysis to detect repeated fraud patterns over time

7.4 Multi-Class Fraud Detection

Extend the system to classify different types of fraud (billing fraud, identity fraud, service fraud)

7.5 Data Security and Privacy

Implement privacy-preserving techniques such as federated learning
Ensure compliance with healthcare data regulations



7.6 Cross-Domain Applications

Apply the proposed framework to other domains such as insurance fraud, banking fraud, and e-commerce fraud detection

7.7 Performance Optimization

Optimize hyperparameters using Grid Search or Bayesian Optimization

Reduce computational cost and improve model efficiency

REFERENCES

- [1] Farahmandazad,D.,&Danesh,K., “ML-Driven Approaches to Combat Medicare Fraud: Advances in Class Imbalance Solutions,” arXiv Preprint, 2025.
- [2] Wen,J.,Tang,X.,&Lu,J., “An Imbalanced Learning Method Based on Graph Tran-SMOTE for Fraud Detection,” Scientific Reports, vol. 14, 2024.
- [3] “Fraud Detection in Healthcare Claims Using Machine Learning: A Systematic Review,” Artificial Intelligence in Medicine, vol. 160, 2025.
- [4] Abdullah,S.,&Swamy,K.M., “Advancing Medicare Fraud Detection via Machine Learning and SMOTE-ENN for Imbalanced Data,” International Journal of Engineering Research and Science & Technology, 2025
- [5] Salem,W.S.,etal., “Enhancing Fraud Detection in Imbalanced Datasets Using Machine Learning and SMOTE,” Mansoura Journal for Computer and Information Sciences, 2025.
- [6] “Healthcare Fraud Detection Using an Integrated ML Approach with SMOTE,” Procedia Computer Science, vol. 258, 2025.
- [7] Ramyateja,O.,etal., “Enhancing Medicare Fraud Detection Through Machine Learning: Addressing Class Imbalance with SMOTE-ENN,” International Journal of Current Advanced Research, 2024.
- [8] Suhel,S.,&Ananthnath,G.V.S., “Leveraging Machine Learning Approach for Improved Medicare Fraud Detection,” International Journal of Scientific Research in Science, Engineering and Technology, 2025.
- [9] Mozafari,A.,etal., “CleverCatch: A Knowledge-Guided Weak Supervision Model for Fraud Detection,” arXiv, 2025.
- [10] C.Nagarajan and M.Madheswaran - ‘Stability Analysis of Series Parallel Resonant Converter with Fuzzy Logic Controller Using State Space Techniques’- Taylor &Francis, Electric Power Components and Systems, Vol.39 (8), pp.780-793, May 2011. DOI: 10.1080/15325008.2010.541746
- [11] C.Nagarajan and M.Madheswaran - ‘Experimental verification and stability state space analysis of CLL-T Series Parallel Resonant Converter’ - Journal of Electrical Engineering, Vol.63 (6), pp.365-372, Dec.2012. DOI: 10.2478/v10187-012-0054-2
- [12] C.Nagarajan and M.Madheswaran - ‘Performance Analysis of LCL-T Resonant Converter with Fuzzy/PID Using State Space Analysis’- Springer, Electrical Engineering, Vol.93 (3), pp.167-178, September 2011. DOI 10.1007/s00202-011-0203-9
- [13] S.Tamilselvi, R.Prakash, C.Nagarajan,“Solar System Integrated Smart Grid Utilizing Hybrid Coot-Genetic Algorithm Optimized ANN Controller” Iranian Journal Of Science And Technology-Transactions Of Electrical Engineering, DOI10.1007/s40998-025-00917-z,2025
- [14] S.Tamilselvi, R.Prakash, C.Nagarajan,“ Adaptive sliding mode control of multilevel grid-connected inverters using reinforcement learning for enhanced LVRT performance” Electric Power Systems Research 253 (2026) 112428, doi.org/10.1016/j.epr.2025.112428
- [15] S.Thirunavukkarasu, C. Nagarajan, 2024, “Performance Investigation on OCF and SCF study in BLDC machine using FTANN Controller,” Journal of Electrical Engineering And Technology, Volume 20, pages 2675–2688, (2025), doi.org/10.1007/s42835-024-02126-w
- [16] C. Nagarajan, M.Madheswaran and D.Ramasubramanian- ‘Development of DSP based Robust Control Method for General Resonant Converter Topologies using Transfer Function Model’- Acta Electrotechnica et Informatica Journal , Vol.13 (2), pp.18-31, April-June.2013, DOI: 10.2478/aei-2013-0025.
- [17] C.Nagarajan and M.Madheswaran - ‘DSP Based Fuzzy Controller for Series Parallel Resonant converter’- Springer, Frontiers of Electrical and Electronic Engineering, Vol. 7(4), pp. 438-446, Dec.12. DOI 10.1007/s11460-012-0212-0.



- [18] C.Nagarajan and M.Madheswaran - 'Experimental Study and steady state stability analysis of CLL-T Series Parallel Resonant Converter with Fuzzy controller using State Space Analysis'- Iranian Journal of Electrical & Electronic Engineering, Vol.8 (3), pp.259-267, September 2012.
- [19] C.Nagarajan and M.Madheswaran, "Analysis and Simulation of LCL Series Resonant Full Bridge Converter Using PWM Technique with Load Independent Operation" has been presented in ICTES'08, a IEEE / IET International Conference organized by M.G.R.University, Chennai. Vol.no.1, pp.190-195, Dec.2007
- [20] Suganthi Mullainathan, Ramesh Natarajan, "An SPSS and CNN modelling based quality assessment using ceramic materials and membrane filtration techniques", Revista Materia (Rio J.) Vol. 30, 2025, DOI: <https://doi.org/10.1590/1517-7076-RMAT-2024-0721>
- [21] M Suganthi, N Ramesh, "Treatment of water using natural zeolite as membrane filter", Journal of Environmental Protection and Ecology, Volume 23, Issue 2, pp: 520-530,2022
- [22] Wang,Y., "A Data Balancing and Ensemble Learning Approach for Fraud Detection," arXiv, 2025.