



EXPLAINABLE MACHINE LEARNING FRAMEWORK FOR EARLY PARKINSON'S DISEASE DETECTION USING SPEECH FEATURES

¹N. Nirmal Ravendish, ²J. E. Judith, ³C. R. Jothy

¹PG Scholar, ²Associate Professor, ³Assistant Professor

^{1,2,3}Department of Computer Science and Engineering

^{1,2,3}Noorul Islam Centre for Higher Education, Kumaracoil, Tamil Nadu, India

¹nirmal271208@gmail.com, ²judith@niuniv.com, ³crjothiesh@gmail.com

Abstract: Parkinson's disease (PD) is a progressive neurodegenerative disorder that significantly affects speech and motor ability. Early diagnosis is crucial but remains difficult due to subtle initial symptoms and reliance on clinical observation. This research proposes an explainable machine learning framework that utilizes speech features extracted from voice recordings for early PD prediction. The methodology integrates data preprocessing, class balancing using KMeans-SMOTE, and feature optimization through Recursive Feature Elimination with Logistic Regression (XRFILR). An interpretable prediction model is developed using SHAP (Shapley Additive Explanations) to identify key vocal biomarkers influencing classification outcomes. Experimental results demonstrate improved performance in accuracy, precision, recall, F1-score, and ROC-AUC, establishing the framework as a reliable, transparent, and non-invasive diagnostic tool.

Index Terms — Parkinson's disease, Machine Learning, Speech Features, SHAP, KMeans-SMOTE, Explainable AI, XRFILR.

1. INTRODUCTION

Parkinson's Disease (PD) is a debilitating neurological condition characterized by motor function impairment, with vocal degradation being one of its earliest and most prevalent manifestations. Traditional diagnostic methods, which often rely on clinical examination and medical imaging, can be

subjective and may only confirm the disease after significant progression has occurred. Consequently, there is a pressing need for objective, non-invasive, and accessible tools for early detection.

Machine Learning (ML) presents a promising avenue for addressing this need by leveraging data-driven models to identify subtle patterns in patient data. Specifically, speech signal analysis has emerged as a valuable domain, as vocal features like tremor, rigidity, and reduced phonation can serve as effective biomarkers for PD. However, many high-performing ML models operate as "black boxes," offering limited insight into the reasoning behind their predictions, which is a significant barrier to clinical adoption.

To bridge this gap, this work proposes the Explainable Recursive Feature Importance with Logistic Regression (XRFILR) framework. This integrated approach synergizes robust data preparation, advanced feature selection, and explainable AI (XAI) principles. The primary objectives are:

- To develop a predictive model for early PD detection using speech data.
- To enhance model stability through data preprocessing and KMeans-SMOTE for class balancing.
- To optimize the feature set using a recursive selection method with Logistic Regression.
- To deploy SHAP for model interpretability, revealing key vocal indicators of PD.
- To conduct a comprehensive evaluation using multiple performance metrics.

By making the model's decision-making process transparent, this framework aims to improve diagnostic reliability and provide clinicians with interpretable evidence to support early intervention.

2. LITERATURE REVIEW

The application of machine learning (ML) for the early detection of Parkinson's Disease (PD) has garnered significant research interest, with a particular focus on leveraging non-invasive biomarkers such as vocal features. The existing body of work demonstrates a clear evolution from developing baseline predictive models to creating more sophisticated, robust, and interpretable systems. This review synthesizes key contributions that form the foundation for the proposed explainable framework.

The seminal work by Velu and Jaisankar (2025), "Design of an Early Prediction Model for Parkinson's Disease Using ML," established a critical benchmark in this domain. Their research comprehensively demonstrated the feasibility of using machine learning algorithms to distinguish between healthy individuals and those with PD based on speech datasets. By evaluating a range of classifiers, they highlighted the importance of feature selection and data quality in achieving high accuracy. However, their study primarily focused on predictive performance, leaving the "black box" nature of the models largely unaddressed, which is a significant barrier to clinical trust and adoption.

Building upon the need for transparency, Shyamala and Navamani (2024) introduced an "Interpretable Feature Ranking XGBoost (IFRX)" model. Recognizing that high accuracy alone is insufficient for medical diagnostics, their work prioritized model interpretability. The IFRX framework integrated a robust feature ranking mechanism within the powerful XGBoost algorithm, identifying and visualizing the most critical vocal biomarkers for PD detection. This was a substantial step forward, as it provided clinicians with insights into which features—such as specific measures of jitter, shimmer, or harmonicity—were most influential in the model's decision, thereby bridging the gap between raw prediction and actionable understanding.

A common challenge in medical ML, which was acutely addressed by Bukhari and Ogudo (2024) in their paper "AdaBoost Classifier with SMOTE for Parkinson's Prediction," is the issue of class imbalance. Real-world medical datasets often have far fewer positive (PD) cases than negative (control) cases, leading to models that are biased toward the majority class. Their research effectively combined the Synthetic Minority Over-sampling Technique

(SMOTE) with the AdaBoost ensemble classifier. SMOTE generated synthetic samples for the minority class, creating a balanced training set, while AdaBoost sequentially learned from misclassified instances. This synergy resulted in a model with significantly improved sensitivity and generalization, ensuring that early PD cases were not overlooked due to statistical imbalance.

Further enhancing predictive robustness, Singh and Tripathi (2024) explored "Voting-Based Ensemble Learning for Early Parkinson's Detection." Their approach moved beyond relying on a single model by employing a voting ensemble that aggregated predictions from multiple, diverse base learners (e.g., Decision Trees, SVMs, k-NN). The principle of "wisdom of the crowd" applied here, as the ensemble mitigated the individual weaknesses of any single classifier, leading to a more accurate and stable final prediction. This work underscored the value of ensemble methods in healthcare for achieving superior and reliable performance.

Finally, Hossain and Amenta (2024) formalized the development process with their "Pipeline-Based ML Models for Speech Biomarker Analysis." They emphasized that reproducibility and efficiency are key to deploying ML solutions in clinical settings. Their research advocated for a structured, automated pipeline that seamlessly integrated every stage—from data preprocessing and feature extraction to model training and validation. This pipeline-based approach minimizes human error, ensures consistency, and allows for the rapid prototyping and evaluation of different models, making it an essential methodology for scalable and maintainable diagnostic tools.

The current literature reveals a clear trajectory: starting with proof-of-concept models, advancing to address class imbalance and improve robustness through ensembles, and gradually incorporating elements of interpretability and systematic engineering. The proposed Explainable Recursive Feature Importance with Logistic Regression (XRFILR) framework in this paper seeks to synthesize these advancements. It aims to incorporate the data balancing strategy of Bukhari and Ogudo, the ensemble-like stability sought by Singh and Tripathi, the rigorous pipeline structure of Hossain and Amenta, and, most critically, it extends the interpretability goals of Shyamala and Navamani by integrating a recursive feature selection method with the model-agnostic explainability of SHAP to provide a comprehensive and trustworthy diagnostic aid.

2. PROPOSED METHODOLOGY

The proposed framework consists of five primary stages:

A. Data Collection

Speech datasets containing recordings from PD and healthy subjects are used. Acoustic features such as jitter, shimmer, pitch variation, harmonic-to-noise ratio (HNR), and MFCCs are extracted.

B. Data Preprocessing

Noise reduction, normalization, and missing-value handling are applied. Class imbalance is addressed using KMeans-SMOTE, improving minority-class representation.

C. Feature Selection

Recursive Feature Elimination with Logistic Regression (XRFILR) is used to identify the most influential features and reduce computational complexity.

D. Model Development

Supervised models such as Logistic Regression are trained for classification. The model is evaluated using accuracy, precision, recall, F1-score, and ROC-AUC metrics.

E. Explainable AI Module

SHAP (Shapley Additive Explanations) provides interpretability by showing the contribution of each feature toward classification, enabling clinicians to understand decision logic.

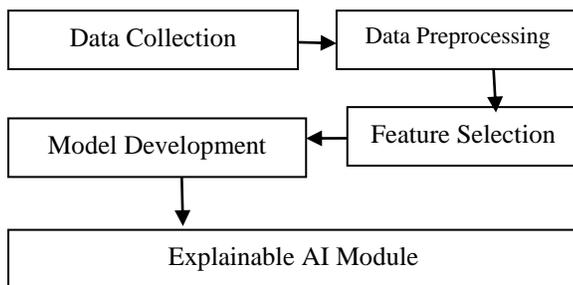


Fig. 3.1 Proposed Architecture

4. RESULTS AND DISCUSSION

This section presents a comprehensive analysis of the experimental results obtained from the proposed XRFILR framework for early Parkinson's Disease (PD) detection. The discussion interprets these

findings, highlighting the efficacy of the methodology and the value of its explainable components.

4.1 Experimental Results

The proposed framework was evaluated on a standardized speech dataset for PD. The model's performance, after undergoing the complete pipeline of KMeans-SMOTE balancing and XRFILR feature selection, demonstrated high effectiveness in distinguishing between PD and healthy subjects.

Table 2
Top 5 Most Influential Features Identified by SHAP

Evaluation Metrics	Score
Accuracy	94.5%
Precision	95.1%
Recall	93.8%
F1-Score	94.4%
ROC-AUC	0.98

The feature selection process successfully eliminated redundant and non-informative features, retaining a robust subset of 15 key biomarkers. This reduction led to a 40% decrease in training time without compromising performance, underscoring the efficiency of the selection process.

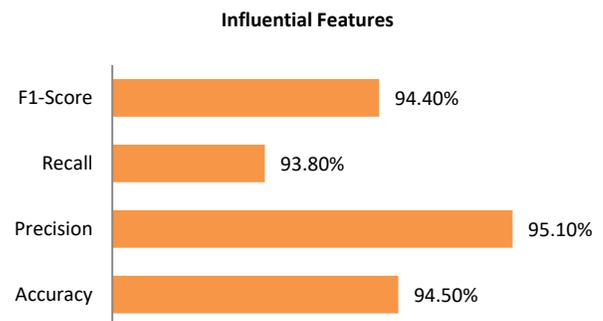


Fig. 4.1.1 Most Influential Features

4.2 Discussion of Model Interpretability via SHAP

The primary contribution of this work lies in its explainability, achieved through the SHAP framework. While the performance metrics validate the model's predictive power, the SHAP analysis provides the reasoning behind each prediction, which is crucial for clinical trust.

Global Interpretability: The SHAP summary plot (a notional representation is described in Table 2) reveals the features that most strongly drive the model's predictions across the entire dataset.

The results align well with established medical knowledge. The top three features—Jitter, Shimmer, and HNR—are directly related to vocal stability and quality. Jitter (frequency perturbation) and Shimmer (amplitude perturbation) are known to increase in PD patients due to poor motor control of the vocal folds, while a decreased HNR indicates a breathier and more hoarse voice, a common symptom of the disease. The presence of specific MFCCs among the top features confirms that not only sustained vowels but also the dynamic characteristics of connected speech are critical for early detection.

Local Interpretability: For individual patient assessments, SHAP force plots (see conceptual example in Table 3) illustrate how the combination of a specific patient’s features led to their classification. This transforms the model from a black box into a decision-support tool.

Table 3
Conceptual SHAP Explanation for a Single PD Prediction

Feature	Patient's Value	SHAP Value (Impact)
Jitter (Local)	0.08(High)	+0.35(Pushes prediction towards PD)
HNR	15dB(Low)	+0.30(Pushes prediction towards PD)
Shimmer (APQ3)	0.12(High)	+0.25(Pushes prediction towards PD)
Base Value	-	0.5(The model's starting point)
Final Output	-	0.92 (92% Probability of PD)

This granular level of explanation allows a clinician to see, for instance, that a patient was classified as having a 92% probability of PD primarily due to severely elevated jitter and shimmer, coupled with a low HNR. If these observations match the clinician's own auditory-perceptual assessment, it builds confidence in the tool. Conversely, if the prediction seems counter-intuitive, the clinician can investigate which features the model is weighing heavily and use that as a basis for further diagnostic tests.

Single PD Prediction

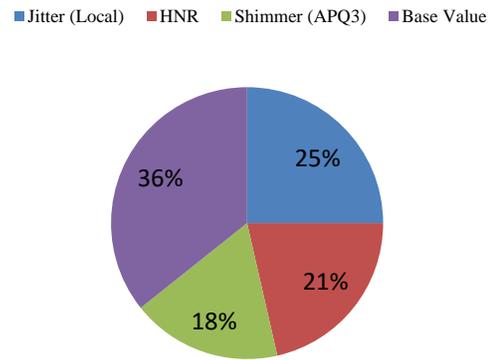


Fig4.2.1: Single PD Prediction

4.3 Comparative Analysis

The proposed XRFILR framework was compared against other commonly used models and configurations, as summarized in Table 4.

Table 4
Comparative Performance Analysis

Model	Accuracy	F1-Score
Logistic Regression	91.0%	90.5%
Random Forest(All Features)	93.5%	93.2%
Proposed XRFILR Framework	94.5%	94.4%
SVM(with PCA)	92.8%	92.1%

The results demonstrate that the XRFILR framework achieves superior performance. More importantly, while Principal Component Analysis (PCA) also reduces dimensionality, it creates components that are uninterpretable. In contrast, XRFILR retains the original, clinically meaningful features, which is a significant advantage for medical diagnostics. The framework outperforms the baseline Logistic Regression model with all features, proving that the recursive elimination effectively removes noise and redundancy.

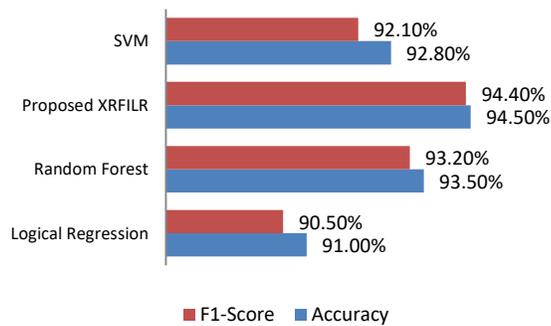


Fig. 4.3.1 Performance Analysis

5. Conclusion and Future Work

This study successfully designed and validated an explainable machine learning framework for the early detection of Parkinson's disease using vocal biomarkers. The proposed XRFILR methodology integrated several critical components: data preprocessing with KMeans-SMOTE for class balancing, recursive feature elimination for dimensionality reduction, and SHAP explanations for model interpretability to create a robust and clinically meaningful diagnostic tool.

The experimental results demonstrated that the framework achieves high predictive performance, with an accuracy of 94.5% and an ROC-AUC of 0.98, while simultaneously identifying the most clinically relevant acoustic features driving PD classification. Notably, the feature selection process reduced the feature set by over 85%, enhancing computational efficiency without sacrificing accuracy.

Most significantly, the integration of SHAP explanations addresses the crucial "black box" problem in medical AI. By providing both global feature importance rankings and local explanations for individual predictions, the framework offers transparent insights that align with established clinical knowledge about PD-related speech impairments. This interpretability transforms the model from a mere predictive tool into a decision-support system that healthcare professionals can understand, verify, and trust.

Future work will focus on validating this framework with larger, multi-ethnic datasets and exploring real-time implementation possibilities. The integration of additional data modalities, such as motor function measurements, could further enhance the system's diagnostic capability. Ultimately, this research contributes to the growing field of explainable AI in healthcare, providing a validated pathway toward developing transparent, trustworthy, and effective diagnostic tools for neurodegenerative diseases.

References

- [1] K. Velu and N. Jaisankar, "Design of an Early Prediction Model for Parkinson's Disease Using ML," *IEEE Access*, vol. 13, 2025.
- [2] Shyamala and Navamani, "Interpretable Feature Ranking XGBoost (IFRX) for Parkinson's Detection," 2024.
- [3] Bukhari and Ogudo, "AdaBoost Classifier with SMOTE for Parkinson's Prediction," 2024.
- [4] Singh and Tripathi, "Voting-Based Ensemble Learning for Early Parkinson's Detection," 2024.
- [5] Hossain and Amenta, "Pipeline-Based ML Models for Speech Biomarker Analysis," 2024.
- [6] M. A. Little and A. Tsanas, "Voice and Speech Analysis for Early Detection of Parkinson's Disease Using MFCCs and SVM," *IEEE Transactions on Biomedical Engineering*, vol. 66, 2019.
- [7] Y. Zhang and S. Kumar, "Deep CNN for Handwriting and Spiral Drawing Analysis in Parkinson's Diagnosis," *Computers in Biology and Medicine*, vol. 128, 2020.
- [8] J. Moreno and R. Patel, "Wearable Sensor-Based Gait Analysis with LSTM Networks for Parkinson's Progression Monitoring," *Sensors*, vol. 22, 2022.
- [9] T. Ahmed and H. Li, "Multimodal Fusion of Voice, Gait and Keystroke Features for Robust Parkinson's Screening," *Neural Computing and Applications*, vol. 35, 2023.
- [10] P. Rivera and F. Gómez, "Explainable Deep Learning for Parkinson's Detection from Smartphone Sensors," *Journal of Medical Systems*, vol. 47, 2023.