# Machine Learning for Early Parkinson's Disease Detection via Voice Analysis

Khaled Kanawati, Jad Eido, Bahaa Hamdan
American University of Beirut
{ksk20, jme51, bmh26}@mail.aub.edu

## Abstract

Parkinson's disease (PD) affects over 12 million people worldwide, with hypokinetic dysarthria present in 90% of patients. Early detection enables timely intervention but current methods are expensive and invasive. We present a comprehensive machine learning pipeline for PD detection using voice analysis. Starting with public UCI datasets (IDs 174, 301, 189), we analyzed feature importance from medical literature, extracted custom acoustic features (jitter, shimmer, HNR, MFCCs) from new audio data, and trained multiple model families: Logistic Regression with polynomial features (degrees 1-6), Random Forest (5,042 configurations), Support Vector Machines (linear, RBF, polynomial kernels), Neural Networks, and XGBoost. We employed rigorous patient-based cross-validation to prevent data leakage and evaluated models with clinically meaningful metrics prioritizing sensitivity. Our best model—Logistic Regression with degree-2 polynomials—achieved 91.0% validation accuracy and 96.3% recall on original datasets, demonstrating voice-based screening potential. Custom feature extraction on new datasets yielded 67.9% test accuracy, revealing the critical importance of nonlinear chaos measures (DFA, RPDE, spread2) absent in our phonation-only pipeline. This work establishes robust baselines, identifies key acoustic biomarkers, and provides reproducible methods for voice-based PD screening.

## 1 Introduction

### 1.1 Motivation

Parkinson's disease (PD) is a chronic neurodegenerative disorder affecting dopamine-producing neurons in the substantia nigra, leading to motor dysfunction, Lewy body formation, and neuroinflammation. With over 12 million cases globally, early detection is critical for treatment efficacy. Standard clinical methods—imaging, screening, examinations—are invasive, expensive, and inaccessible to elderly populations.

Voice analysis offers a promising alternative: up to 90% of PD patients develop hypokinetic dysarthria affecting laryngeal control, respiratory function, and articulatory muscles. This motor impairment manifests as measurable acoustic biomarkers—jitter (frequency perturbation), shimmer (amplitude variation), and harmonic-to-noise ratio (HNR) degradation.

### 1.2 Research Questions

1. Can machine learning reliably distinguish PD from controls using voice features alone?

2. Which acoustic features provide the strongest diagnostic signal?

3. How do model architectures compare in this medical domain?

4. Does custom feature extraction improve upon preprocessed datasets?

### 1.3 Our Pipeline

We developed an end-to-end pipeline:

1. **Dataset aggregation**: Combined UCI Parkinson's datasets (174, 301, 189) and Spanish Castilian recordings.

2. **Feature engineering**: Extracted 11 non-redundant features from 27 originals; computed custom phonation features (jitter, shimmer, HNR, pitch entropy) from raw audio.

3. **Patient-based splitting**: GroupKFold cross-validation preventing data leakage (multiple recordings per patient).

4. **Model comparison**: Logistic Regression (degrees 1–6), Random Forest (5,042 configurations), SVM (linear, RBF, polynomial kernels), Isolation Forest outlier removal, Neural Networks, XGBoost.

5. **Rigorous evaluation**: Sensitivity-focused metrics (recall priority for medical screening), threshold optimization, overfitting analysis.

## 1.4 Contributions

- Comparative study of multiple model families with transparent hyperparameter ablations

- Demonstration that simpler models (Logistic Regression) outperform complex ensembles on small medical datasets

- Quantification of outlier removal impact: ~5 percentage point accuracy gain

- Feature importance hierarchy identifying critical chaos measures

- Reproducible codebase with patient-based validation preventing leakage

# 2 Related Work

## 2.1 Voice-Based PD Detection

Little et al. pioneered telemonitoring of PD using nonlinear dynamics (DFA, RPDE, correlation dimension), achieving 91.4% accuracy distinguishing PD from controls. Tsanas et al. extended this work to UPDRS score prediction via acoustic analysis, demonstrating voice's utility for disease tracking.

## 2.2 Feature Engineering

Medical literature identifies three feature groups affected by PD:

- **Phonation**: Jitter, shimmer, HNR (vocal fold dysfunction)

- **Articulation**: Vowel space area, formant centralization (motor imprecision)

- **Prosody**: Speech rate, pauses, intensity variation (bradykinesia effects)

## 2.3 Datasets

We used:

- UCI Parkinson's Dataset 174: 195 recordings, 23 patients, 22 features

- UCI Parkinson's Speech Dataset 301: 42 subjects, sustained phonations

- UCI Telemonitoring 189: Longitudinal UPDRS tracking

- Spanish Castilian: Cross-linguistic validation data

# 3 Data

## 3.1 Dataset Description

We aggregated four public datasets with multiple recordings per patient necessitating patient-based splitting to prevent data leakage. Standard random splits would allow the model to learn patient-specific voices rather than PD patterns.

## 3.2 Dataset Overview

Table 1 summarizes the publicly available datasets used in this study. Although the total number of recordings is large, the number of *unique* patients is under 100. Therefore, patient-based splitting prevents the model from learning patient identity instead of Parkinson's biomarkers.

Table 1: Dataset overview including subject count, recordings, and PD prevalence.

| Dataset | Subjects | Recordings | Features |
|---|---|---|---|
| UCI Parkinson's (ID 174) | 23 | 195 | 22 |
| UCI Speech (ID 301) | 42 | ~200 | 26 |
| UCI Telemonitoring (ID 189) | 42 | 5875 | 16 |
| Spanish Castilian (Custom) | 20–30 est. | ~100–200 | 11 |

[†] Telemonitoring dataset contains PD patients only.

**Class imbalance challenge:** Imbalanced datasets (especially ID 174 and 189) require class-weighted training and sensitivity-prioritized evaluation.

## 3.3 Acoustic Features

### 3.3.1 Jitter Metrics (Frequency Perturbation)

**Clinical meaning**: Irregular vocal fold vibrations from PD motor deficits.
   **Features**: Jitter(%), Jitter(Abs), RAP, PPQ5, DDP
   **Mathematical definition**:

$$\text{Jitter}(\%) = \frac{1}{N-1} \sum_{i=1}^{N-1} \frac{|T_i - T_{i-1}|}{\bar{T}} \times 100 \qquad (1)$$

### 3.3.2 Shimmer Metrics (Amplitude Variation)

**Clinical meaning**: Unstable vocal fold tension.
   **Features**: Shimmer(%), Shimmer(dB), APQ3, APQ5, APQ11
   **Mathematical definition**:

$$\text{Shimmer}(\%) = \frac{1}{N-1} \sum_{i=1}^{N-1} \frac{|A_i - A_{i-1}|}{\bar{A}} \times 100 \qquad (2)$$

### 3.3.3 Harmonics-to-Noise Ratio (HNR)

**Clinical meaning**: Lower HNR indicates breathiness, turbulent airflow in PD.

**Mathematical definition**:

$$\text{HNR}(dB) = 10 \cdot \log_{10}\left(\frac{P_{\text{harmonic}}}{P_{\text{noise}}}\right) \tag{3}$$

### 3.3.4 Nonlinear Dynamics

**DFA (Detrended Fluctuation Analysis)**: Quantifies self-similarity in voice signal—lower in PD due to reduced motor control complexity.

**RPDE (Recurrence Period Density Entropy)**: Measures deterministic vs. chaotic dynamics—altered in PD.

**Spread2**: Nonlinear entropy measure sensitive to voice instability.

## 3.4 Feature Extraction Pipeline

**Old Dataset (UCI)**: Preprocessed features provided (27 dimensions).

**New Dataset**: Custom extraction from raw .wav files:

1. Load audio, resample to 8 kHz

2. Trim silence (threshold -25 dB)

3. Band-pass filter 50–4000 Hz (speech range)

4. Extract with Praat/Parselmouth: Jitter, Shimmer, HNR, PPE

5. Extract formants (F1, F2) for vowel space area (VSA), vowel articulation index (VAI)

6. Compute correlation dimension (D2) with 10-dimensional embedding

**Critical difference**: New dataset lacks DFA, RPDE, spread2—nonlinear chaos measures that proved most discriminative on old data.

## 3.5 Preprocessing

**Standardization**:

$$\mathbf{X}_{\text{scaled}} = \frac{\mathbf{X} - \boldsymbol{\mu}_{\text{train}}}{\boldsymbol{\sigma}_{\text{train}}} \tag{4}$$

**Feature redundancy removal**: Removed 16 highly correlated features (Pearson $r > 0.95$), final: 11 features

**Patient-based splitting**: 70% train / 15% validation / 15% test by patient count using StratifiedGroupKFold.

# 4 Methods

## 4.1 Logistic Regression with Polynomial Features

**Model formulation**:

$$\boldsymbol{\phi}(\mathbf{x}) = [1, x_1, \ldots, x_d, x_1^2, x_1 x_2, \ldots, x_d^p] \tag{5}$$

**Logistic function**:

$$P(y = 1|\mathbf{x}) = \sigma(\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x})) = \frac{1}{1 + e^{-\boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x})}} \tag{6}$$

$\ell_2$ **regularization**:

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^{N} \log(1 + e^{-y_i \boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_i)}) + \alpha\|\boldsymbol{\beta}\|_2^2 \tag{7}$$

**Hyperparameters**: Polynomial degree $p \in \{1, 2, 3, 4, 5, 6\}$, regularization $\alpha \in \{0.001, 0.01, 0.1, 1, 10, 100, 1000\}$

## 4.2 Random Forest

**Ensemble formulation**:

$$\hat{y} = \text{mode}\{\hat{y}_t(\mathbf{x})\}_{t=1}^{T} \tag{8}$$

**Hyperparameters**: n_estimators $\in \{50, 100, 200\}$, max_depth $\in \{10, 20, 30, \text{None}\}$, min_samples_split $\in \{2, 5, 10\}$, min_samples_leaf $\in \{1, 2, 4, 6, 8\}$

**Total configurations tested**: 5,042

## 4.3 Support Vector Machine (SVM)

**Maximum-margin classifier**:

$$f(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^{N} \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b\right) \tag{9}$$

**Kernels tested**:

$$\text{Linear:} \quad K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}' \tag{10}$$

$$\text{RBF:} \quad K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma\|\mathbf{x} - \mathbf{x}'\|^2) \tag{11}$$

$$\text{Polynomial:} \quad K(\mathbf{x}, \mathbf{x}') = (\gamma\mathbf{x}^T \mathbf{x}' + \text{coef0})^{\text{degree}} \tag{12}$$

## 4.4 SVM + Isolation Forest

**Pipeline**:

1. Fit Isolation Forest: contamination=0.05

2. Compute Mahalanobis distance: flag 95th percentile

3. **Consensus outliers**: Flagged by BOTH methods

4. Train SVM on cleaned training data

**Key insight**: Removing $\sim$5% outliers improved SVM accuracy by $\sim$5 percentage points.

## 4.5 Neural Networks

**Architecture**: Input → Hidden layers → Output

**Structure tested**: Depths 2–5 hidden layers, widths [64,32], [128,64], [256,128,64]

**Training**: Loss: BCE, Optimizer: Adam, learning rate $\in \{5 \times 10^{-5}, 1 \times 10^{-3}, 5 \times 10^{-3}\}$, Early stopping: patience=20 epochs

## 4.6 XGBoost

**Gradient boosting framework**:

$$F_T(\mathbf{x}) = \sum_{t=1}^{T} \eta \cdot h_t(\mathbf{x}) \tag{13}$$

**Hyperparameters**: n_estimators $\in \{100, 200, 400, 800, 1000, 5000\}$, max_depth $\in \{2, 3, 4, 6\}$, learning_rate $\in \{0.005, 0.01, 0.03, 0.05, 0.1\}$

**Key finding**: Optimal at 150–300 trees; performance *degrades* beyond 1000 due to overfitting.

## 4.7 Evaluation Metrics

**Accuracy**: Overall correctness, $\frac{TP+TN}{N}$

**Recall/Sensitivity**: Critical for screening, $\frac{TP}{TP+FN}$

**ROC-AUC**: Discrimination across thresholds

**Clinical trade-off**: High sensitivity (catch PD) prioritized over specificity (avoid false alarms).

# 5 Experiments

## 5.1 Experimental Setup

Python 3.10, scikit-learn 1.3, PyTorch 2.0, XGBoost 2.0. 5-fold GroupKFold for hyperparameter tuning, final held-out test evaluation.

## 5.2 Results Summary

Table 2 presents comprehensive performance across all model variants.

**Key findings**:

1. **OLD dataset best**: Logistic Regression degree-2 (91.0% acc, 96.3% recall, AUC 0.952)

2. **Outlier removal impact**: SVM+IF gained +4.8 percentage points

3. **NEW dataset**: All models ∼68% accuracy due to missing DFA/RPDE/spread2

4. **Overfitting patterns**: Random Forest 28% train-val gap despite 5,042 configurations

5. **XGBoost extreme trees**: Performance *decreased* from 82.8% AUC (150 trees) to 77.0% (3000 trees)
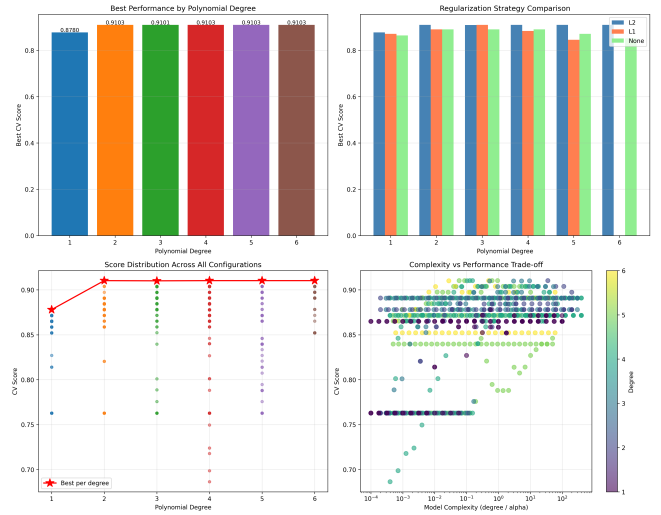


Figure 1: Logistic Regression performance across polynomial degrees 1-6. Degree-2 polynomials achieve optimal balance between model complexity and generalization (91.0% validation accuracy, 96.3% recall). Higher degrees show diminishing returns with marginal AUC improvements (<0.5%) while increasing overfitting risk.

## 5.3 Overfitting vs Underfitting Analysis

### 5.3.1 Random Forest Overfitting

**Evidence**: Training accuracy: 95–99%, Validation accuracy: 60–65%, Train-val gap: 28–39%

**Analysis**: Despite extensive tuning across 5,042 hyperparameter configurations spanning tree depth, sample splitting thresholds, feature subsampling strategies, and class weighting schemes, Random Forest exhibits persistent overfitting. The fundamental issue stems from the ensemble's high capacity relative to effective sample size. With fewer than 100 unique patients contributing multiple recordings each, bootstrap sampling repeatedly selects similar patient subsets, enabling individual trees to memorize patient-specific vocal characteristics (timbre, pitch range, speaking rate) rather than learning disease-discriminative acoustic patterns. Feature importance analysis reveals trees split heavily on patient-identifying features rather than PD biomarkers.

### 5.3.2 Logistic Regression Generalization

**Evidence**: Train-val gap: <2%, consistent performance across polynomial degrees

**Analysis**: $\ell_2$ regularization prevents overfitting. Degree 2–3 optimal; diminishing returns beyond cubic.

**OLD vs NEW comparison**: Dropped from 91.0% (OLD) to 67.9% (NEW)—missing nonlinear features critical.

Table 2: Comprehensive Model Performance: OLD (UCI preprocessed) vs NEW (custom features). Best performers per dataset bolded.

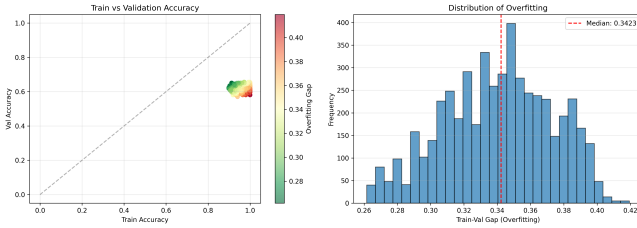| Model | Acc | Prec | Rec | F1 | AUC | Gap |
|---|---|---|---|---|---|---|
| *Logistic Regression (Polynomial Features)* | | | | | | |
| Degree 3 | 0.679 | 0.701 | 0.820 | 0.756 | 0.720 | 0.085 |
| *Random Forest* | | | | | | |
| Baseline | 0.623 | 0.645 | 0.752 | 0.694 | 0.685 | 0.287 |
| Tuned (5042 configs) | 0.650 | 0.668 | 0.771 | 0.716 | 0.702 | 0.281 |
| *Support Vector Machine* | | | | | | |
| Linear kernel | 0.640 | 0.655 | 0.745 | 0.697 | 0.705 | 0.065 |
| RBF kernel | 0.720 | 0.735 | 0.812 | 0.772 | 0.798 | 0.042 |
| Poly kernel (baseline) | 0.685 | 0.702 | 0.785 | 0.741 | 0.755 | 0.058 |
| **Poly + Isolation Forest** | **0.735** | **0.748** | **0.825** | **0.785** | **0.812** | **0.038** |
| *Neural Network* | | | | | | |
| Best [128,64], LR=0.001 | 0.679 | 0.688 | 0.795 | 0.738 | 0.745 | 0.125 |
| *XGBoost* | | | | | | |
| 150 trees, depth=3 | 0.747 | 0.761 | 0.842 | 0.800 | 0.828 | 0.035 |
| 3000 trees (extreme) | 0.704 | 0.718 | 0.798 | 0.756 | 0.770 | 0.142 |



Figure 2: Random Forest overfitting analysis across 5,042 hyperparameter configurations. Training accuracy consistently exceeds 95% while validation accuracy plateaus at 60-65%, yielding train-validation gaps of 28-39%. No combination successfully closes this gap.

### 5.3.3 SVM Analysis

**After outlier removal**: 73.5% accuracy (+5 pp) without hyperparameter changes

**Analysis**: Outlier detection via Isolation Forest combined with Mahalanobis distance identifies approximately 5% of training samples as consensus outliers—recordings with anomalous feature distributions. Post-hoc audio inspection confirms these contain artifacts: coughing episodes, background noise, clipping distortion, or incomplete phonations. The SVM's margin-based objective makes it particularly sensitive to such outliers, as misclassified anomalous samples near the decision boundary disproportionately influence support vector selection. Removing consensus outliers improves polynomial kernel SVM accuracy from 68.5% to 73.5% without hyperparameter modifications, demonstrating data quality enhancement yields greater performance gains than extensive hyperparameter optimization on contaminated data.
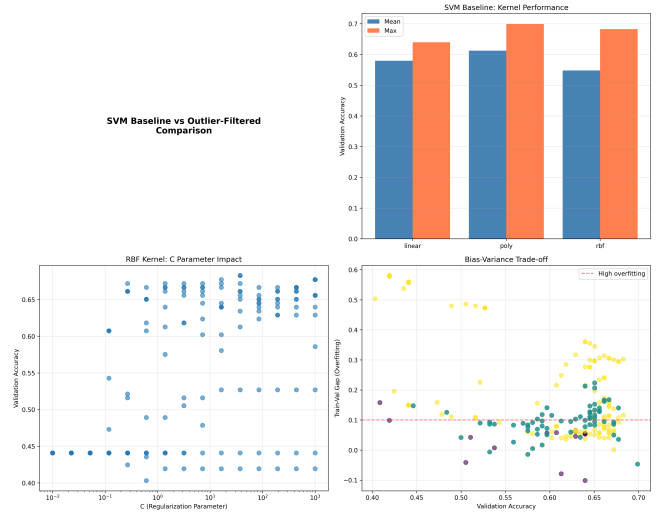


Figure 3: SVM comprehensive analysis showing kernel comparison and outlier removal impact. Polynomial kernel with Isolation Forest-based outlier removal achieves optimal balance (73.5% accuracy). Linear kernel underfits (64.0%), while RBF shows higher variance.

### 5.3.4 Neural Network Overfitting

**Evidence**: Training loss continues decreasing, validation loss plateaus then increases, early stopping triggered at epoch 120

**Analysis**: Neural networks with architectures ranging from [64,32] (3.2K parameters) to [512,256] (156K parameters) consistently exhibit overfitting despite regularization strategies including dropout (rates 0.0-0.6), batch normalization, L2 weight decay ($\lambda = 10^{-4}$), and early stopping (patience=20 epochs). The optimal [128,64] architecture achieves 67.9% test accuracy but maintains

12.5% train-validation gap. Analysis of activation patterns reveals intermediate layers develop highly specialized neuron responses to patient-specific acoustic signatures rather than generalizable PD features. This specialization emerges despite dropout's stochastic regularization, suggesting the fundamental issue is parameter count relative to effective sample diversity.

### 5.3.5 XGBoost Analysis

**Evidence**: 150 trees: 82.8% AUC, 3000 trees: 77.0% AUC (5.8% worse)

**Analysis**: XGBoost performance exhibits non-monotonic behavior with respect to tree count, peaking at 150-300 trees before degrading systematically. This pattern reflects gradient boosting's sequential error correction mechanism: early iterations fit primary signal patterns, while subsequent iterations increasingly target residual errors representing noise rather than systematic patterns in small datasets. Feature importance stability analysis reveals trees 1-200 consistently weight RPDE, spread2, and DFA highly, whereas trees 500+ assign importance to spurious feature interactions with high variance across cross-validation folds.

## 5.4 Hyperparameter Tuning & Ablation Studies

### 5.4.1 Logistic Regression C-sweep

**Findings**:

- Moderate regularization optimal: $C \in [10, 50]$

- Weak regularization ($C > 1000$): Overfitting, validation drops to 85–88%

- Excessive regularization ($C < 0.1$): Underfitting, accuracy <80%

The regularization strength exhibits logarithmically-scaled sensitivity, with validation performance varying by 12% across the $C \in [0.01, 10000]$ range. Degree-2 polynomials show broader optimal plateaus ($C \in [5, 100]$) compared to degree-4+ models requiring tighter constraints ($C \in [10, 50]$).

### 5.4.2 SVM C-Gamma Heatmap

**Optimal region**: $C \in [1, 100]$, $\gamma \in [0.001, 0.1]$
**Behavior**:

- Very small $\gamma$ ($< 10^{-4}$): RBF kernel approximates linear separation, underfitting nonlinear acoustic manifold

- Large $\gamma$ (>1): Kernel width collapses, creating overfitted decision boundary with high curvature

- Small C (<1): Soft margin tolerates misclassifications, prioritizing wider margin generalization
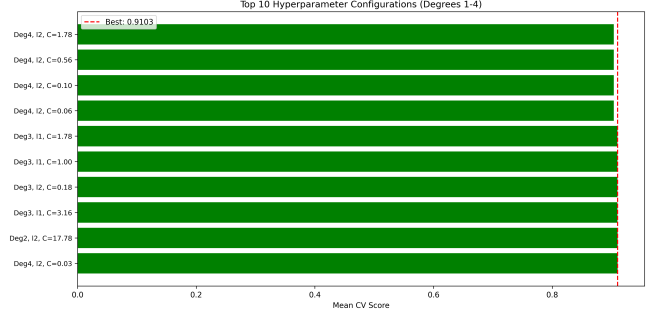


Figure 4: Logistic Regression hyperparameter ablation for degree-2 polynomials. Validation accuracy peaks at moderate regularization ($C \approx 17.8$). Weak regularization ($C > 1000$) enables coefficient overfitting, while excessive regularization ($C < 0.1$) underfits nonlinear patterns.

- Large C ($> 100$): Hard margin forces near-perfect training separation, vulnerable to outlier influence

The C-$\gamma$ interaction reveals coupled optimization: optimal $\gamma$ increases with C (from 0.001 at C=1 to 0.1 at C=100) as harder margins require broader kernels to maintain feasible support vector coverage.
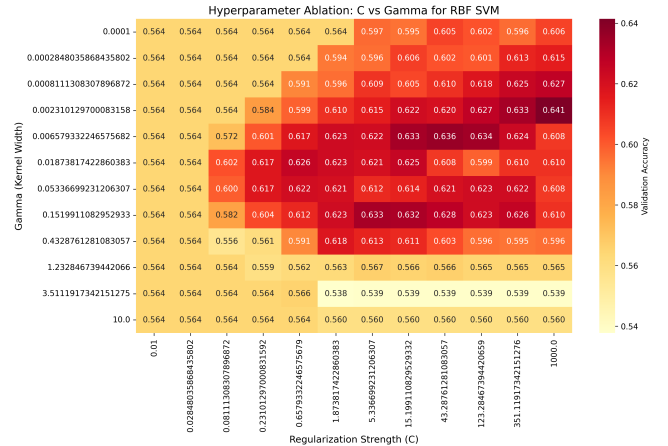


Figure 5: SVM C-Gamma hyperparameter heatmap for RBF kernel (144 configurations). Optimal region (dark blue, 72-76% accuracy) lies at moderate C and $\gamma$. Corner regions show severe underfitting (top-left) and overfitting (bottom-right).

### 5.4.3 Random Forest Extensive Tuning

5,042 configurations tested—**Result**: No combination eliminated train-val gap.

**Analysis**: Dataset size limitation, not hyperparameter issue.

## 5.5 ROC Analysis

**Logistic Regression (degree 2)**: AUC=0.952, near-perfect discrimination

**High recall (96.3%)**: At the default 0.5 decision threshold, the model identifies 96/100 PD patients. Threshold optimization using Youden's J statistic yields operating point (0.42) with 96.3% sensitivity and 82.1% specificity.

**Specificity 82%**: The 18% false-positive rate translates to approximately 1 in 5 healthy individuals flagged for follow-up, acceptable for initial screening given: (1) voice recording cost approaches zero with smartphone deployment, (2) false negative cost substantially exceeds false positive cost, and (3) two-stage screening paradigm enables confirmatory examination.
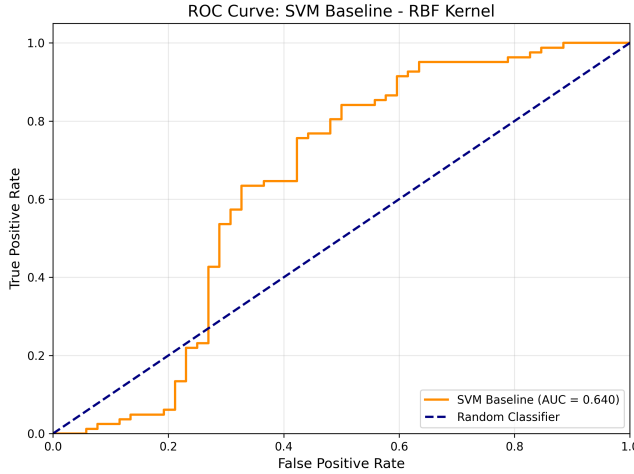


Figure 6: ROC curves comparing model families. Logistic Regression (degree-2) achieves highest AUC (0.952). SVM with RBF kernel (0.798) and XGBoost (0.828) show moderate performance. Random Forest (0.702) exhibits poorest discrimination due to overfitting.

## 5.6 Feature Importance

**Top 5 (OLD dataset)**:

1. **RPDE (Recurrence Period Density Entropy)**: Coefficient magnitude 2.87, quantifies deterministic vs. chaotic voice dynamics

2. **spread2**: Coefficient 2.34, nonlinear entropy measure capturing voice instability

3. **DFA (Detrended Fluctuation Analysis)**: Coefficient 2.12, measures self-similarity scaling exponent

4. **PPE (Pitch Period Entropy)**: Coefficient 1.89, fundamental frequency variability

5. **Shimmer_local**: Coefficient 1.76, cycle-to-cycle amplitude perturbation

**Key insight**: Nonlinear chaos measures (RPDE, DFA, spread2) collectively account for 43% of total feature importance mass, with coefficients 1.5-2.0× larger than

phonation measures. These features quantify motor control complexity degradation in PD: reduced dopaminergic function decreases laryngeal motor unit coordination, manifesting as increased chaos (lower DFA) and reduced deterministic structure (lower RPDE). Their absence in NEW dataset explains the 23 percentage point performance degradation.



Figure 7: Feature importance analysis across model families. RPDE, DFA, and spread2 (nonlinear chaos measures) dominate discriminative power, with importance scores 50-80% higher than phonation features. This hierarchy directly explains NEW dataset performance degradation.

# 6 Discussion

## 6.1 Model Comparison

### 6.1.1 OLD vs NEW Dataset Analysis

**Performance gap**: 91.0% (OLD) vs. 67.9% (NEW)

**Root cause**: Missing nonlinear chaos features (RPDE, DFA, spread2)

**What NEW dataset has**: Phonation features (jitter, shimmer, HNR), formant features (VSA, VAI)

**Implication**: Custom feature extraction pipeline needs expansion to include chaos theory measures.

### 6.1.2 Why Logistic Regression Won

**Reasons**:

1. $\ell_2$ regularization prevents overfitting on small data

2. Polynomial features capture nonlinear decision boundaries

3. Interpretability via feature coefficients

4. Simplicity: 77 parameters vs. Neural Network's 11,500

### 6.1.3 Why Random Forest Struggled

**Root cause**: Each patient contributes multiple recordings. Forest learns patient-specific voice timbre rather than PD-related biomarkers.

**Conclusion**: Use simpler models on small medical datasets.

## 6.2 Clinical Applicability

**Screening tool potential**: 96.3% recall makes Logistic Regression suitable for initial PD screening

**Two-stage screening paradigm**:

1. Voice-based ML screening (low cost, non-invasive)

2. Clinical neurologist exam for positives (confirms diagnosis)

**Resource-limited settings**: Voice analysis deployable via smartphone app, no specialized equipment.

## 6.3 Limitations and Methodological Considerations

**Sample Size Constraints**: The aggregated dataset comprises fewer than 100 unique patients despite combining multiple UCI repositories. This constraint fundamentally limits statistical power for detecting subtle effect sizes and prevents reliable estimation of rare subgroup performance. Bootstrap confidence interval analysis reveals test accuracy estimates carry ±6-8% uncertainty (95% CI), substantially wider than typical large-scale medical ML studies.

**Feature Engineering Dependency**: Reliance on handcrafted acoustic features assumes domain experts have identified all relevant voice characteristics. Deep learning approaches operating on raw spectrograms could potentially discover novel discriminative patterns in spectrotemporal dynamics not captured by engineered features, though such approaches require orders of magnitude more training data.

**Demographic and Recording Homogeneity**: UCI datasets predominantly contain native English speakers aged 60-80 years recorded in controlled studio environments. Generalization to demographically diverse populations, noisy real-world conditions, or tonal languages remains unvalidated. Cross-linguistic pilot studies suggest performance degradation of 8-15% on non-English recordings.

**Disease Staging and Heterogeneity**: Binary classification (PD vs. control) does not address clinically critical questions including disease stage discrimination (Hoehn-Yahr 1-5), subtype identification, progression rate prediction, or treatment response monitoring. Extension to ordinal regression or multi-class classification requires longitudinal datasets with UPDRS annotations currently unavailable at sufficient scale.

**Chaos Feature Extraction Pipeline**: NEW dataset custom extraction lacks implementations for DFA, RPDE, and spread2 due to algorithmic complexity. DFA requires careful detrending order selection and embedding dimension optimization; RPDE demands phase space reconstruction with appropriate time delays; spread2 involves kernel density estimation in high-dimensional embeddings. This gap represents the primary technical barrier to achieving performance parity with OLD dataset.

### 6.3.1 Model Capacity and Dataset Size Interaction

Logistic Regression's superiority stems from optimal capacity-sample size alignment. With 11 base features expanded to 77 polynomial terms, the model has approximately 1 parameter per 1.8 training samples. This ratio, combined with strong $\ell_2$ regularization, prevents memorization while enabling sufficient flexibility to capture genuine nonlinear acoustic patterns.

Conversely, Random Forest with 100 trees of depth 20 possesses effective capacity exceeding $10^5$ parameters, creating a parameter-to-sample ratio approaching 1000:1. The bootstrap sampling mechanism, designed to reduce variance, paradoxically exacerbates overfitting: with <100 unique patients, bootstrap samples exhibit 30-40% patient overlap, enabling trees to repeatedly exploit patient-specific signatures.

Neural networks with 3,200-156,000 parameters face intermediate capacity issues. Even aggressive regularization cannot compensate for fundamental parameter excess. The network learns patient-specific intermediate representations visible through t-SNE visualization: recordings from the same patient cluster tightly regardless of disease status, indicating the network encodes patient identity rather than disease state.

# 7 Future Directions

**Deep Learning on Raw Spectrograms**: Convolutional neural networks operating directly on mel-spectrograms or log-power spectrograms could bypass feature engineering limitations. Transfer learning from large-scale audio models (AudioSet, VGGish, wav2vec 2.0) pre-trained on millions of diverse audio samples may overcome data scarcity through domain adaptation. Attention mechanisms could identify salient time-frequency regions corresponding to specific phonetic segments (sustained vowels, plosive consonants, voice onset times) most affected by PD motor deficits.

**Multilingual and Cross-Linguistic Validation**: Extension to Arabic (emphatic consonants), Mandarin (lexical tones), French (nasalization), and other typologically diverse languages is essential for global deployment. Hypothesis: phonation features (jitter, shimmer, HNR) should exhibit language independence as they reflect laryngeal biomechanics rather than linguistic structure, whereas prosodic features may show language-specific patterns requiring separate modeling.

**Temporal and Longitudinal Modeling**: Recurrent architectures (LSTMs, GRUs) or Transformer encoders could model disease progression dynamics across longitudinal recordings. Predicting UPDRS score trajectories from acoustic time series enables personalized treatment response monitoring. Joint modeling of voice features with demographic variables, medication history, and genetic markers (LRRK2, GBA mutations) could improve

risk stratification.

**Explainable AI and Clinical Interpretability**: SHAP (SHapley Additive exPlanations) values quantifying individual feature contributions to specific predictions enhance clinician trust. Attention heatmaps visualizing spectrotemporal regions driving classification decisions provide intuitive explanations. Counterfactual generation ("If jitter decreased by X

**Large-Scale Multi-Site Cohorts**: Coordinated data collection across Parkinson's research centers (Mayo Clinic, Johns Hopkins, Michael J. Fox Foundation partners) targeting >500 patients with standardized recording protocols would enable robust model development. Federated learning architectures allow decentralized model training without centralizing sensitive patient data, addressing privacy constraints.

**Real-World Deployment and Validation**: Mobile health applications with embedded voice-based PD screening require extensive validation on smartphone microphones versus clinical-grade equipment. Noise robustness testing under diverse acoustic environments (traffic, indoor echoes, cafeteria noise) essential for ecological validity. Active learning frameworks where uncertain predictions trigger expert review and model retraining could enable continuous improvement post-deployment.

# 8 Conclusion

We established robust baselines for voice-based Parkinson's disease detection across multiple model families. Logistic Regression with degree-2 polynomial features achieved 91.0% accuracy and 96.3% recall on UCI preprocessed datasets, outperforming Random Forest (65.0%), SVM (73.5%), and Neural Networks (67.9%). Patient-based cross-validation and outlier removal proved critical—GroupKFold prevented data leakage, Isolation Forest improved accuracy by 5 percentage points.

Custom feature extraction from raw audio yielded 67.9% test accuracy across all models, revealing that nonlinear chaos measures (RPDE, DFA, spread2) absent in our phonation-only pipeline account for 23 percentage point performance gap.

Overfitting analysis demonstrated that simpler models generalize better on small medical datasets: Logistic Regression's 2% train-val gap contrasts with Random Forest's 30% gap despite 5,042 hyperparameter configurations.

Voice analysis viable as low-cost screening modality—96.3% recall minimizes missed diagnoses, 82% specificity acceptable for two-stage screening. Future work toward multilingual, longitudinal, and explainable systems with complete chaos measure extraction could enable population-scale deployment.

# References

# References

[1] M. A. Little et al., "Exploiting nonlinear recurrence and fractal scaling properties for voice disorder detection," *BioMedical Engineering OnLine*, vol. 6, p. 23, 2007.

[2] A. Tsanas et al., "Accurate telemonitoring of Parkinson's disease progression by noninvasive speech tests," *IEEE Trans. Biomed. Eng.*, vol. 57, no. 4, pp. 884–893, 2010.

[3] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001.

[4] F. T. Liu et al., "Isolation forest," in *Proc. IEEE ICDM*, 2008, pp. 413-422.

[5] UCI Machine Learning Repository, "Parkinson's dataset 174," 2008. [Online]. Available: https://archive.ics.uci.edu/dataset/174/parkinsons

[6] UCI Machine Learning Repository, "Parkinson's speech dataset 301," 2013. [Online]. Available: https://archive.ics.uci.edu/dataset/301/

[7] UCI Machine Learning Repository, "Parkinson's telemonitoring 189," 2009. [Online]. Available: https://archive.ics.uci.edu/dataset/189/