**ChatGPT**

# Voice-Based Parkinson's Disease Classification Using Classical Machine Learning

## Abstract

Parkinson's Disease (PD) is a neurodegenerative disorder characterized by motor symptoms and pervasive speech impairments. This thesis investigates the feasibility of **voice-based PD detection** using classical machine learning, emphasizing rigorous methodology over maximal performance. Two complementary datasets are examined: **Dataset A**, a clinical corpus of raw voice recordings (37 subjects) requiring acoustic feature extraction; and **Dataset B**, a larger public dataset (756 samples) of pre-extracted features. A consistent pipeline is applied, extracting **47 baseline features** (prosodic and perturbation measures) from Dataset A, with an **extended set of 78 features** incorporating additional spectral descriptors. Three interpretable classifiers – **Logistic Regression**, **Support Vector Machine (RBF)**, and **Random Forest** – are evaluated under a **2×2 factorial design**: baseline vs. extended features, with vs. without class weighting to address class imbalance. Crucially, **subject-grouped 5-fold cross-validation** is employed for Dataset A to prevent data leakage, while a standard stratified 5-fold CV (with caveats on subject overlap) is used for Dataset B.

Results are reported as **mean ± standard deviation**. On Dataset A, the best model (Random Forest, extended features) achieved **ROC-AUC ≈ 0.87±0.14**, a **+8.7 percentage point** improvement over the baseline feature set. **Extended features** consistently improved accuracy and ROC-AUC, especially for the smaller Dataset A (e.g. Random Forest AUC rose from 0.59 to 0.82 on one task). **Class weighting** had only modest effects (e.g. +3.5pp ROC-AUC for Random Forest with baseline features, but negligible or negative impact with extended features). **Random Forest** outperformed SVM and Logistic Regression across conditions, likely due to its ability to capture non-linear patterns and leverage feature importance for insight. Dataset B yielded higher absolute performance (ROC-AUC ≈ 0.94 with Random Forest) but is interpreted with caution given potential subject overlaps and its high dimensional feature set.

In conclusion, **classical ML models can detect PD from voice** with competitive accuracy, but **robust validation** is paramount. This work highlights that methodological rigor – including proper cross-validation, careful feature engineering, and honest reporting of variance and limitations – is essential to produce reliable findings. The extended feature set notably enhances detection of PD voice signatures, and results underscore the importance of addressing data leakage and class imbalance. These contributions lay a reproducible groundwork for future research, prioritizing interpretability and validity in the development of non-invasive PD screening tools.

## Keywords

# Table of Contents

# Chapter 1: Introduction

## 1.1 Background and Motivation

Parkinson's Disease (PD) is the second most prevalent neurodegenerative disorder globally, affecting approximately 1% of the population over 60 years of age. Early and accurate detection remains a critical clinical challenge, as motor symptoms often manifest only after substantial neurological damage has occurred. Among the earliest observable symptoms are changes in speech and voice production, which can precede motor symptoms by several years.

Voice-based biomarkers offer a promising non-invasive avenue for PD detection. The disease affects the laryngeal muscles and respiratory control, resulting in measurable changes to prosodic features (pitch, intensity, rhythm) and spectral characteristics (formants, harmonics). These acoustic signatures can be captured using standard recording equipment, making voice analysis a cost-effective and accessible screening approach.

## 1.2 Problem Statement

Despite advances in voice-based PD classification, several methodological challenges persist:

1. **Small sample sizes** in raw audio datasets limit model generalizability.
2. **Subject identity leakage** when multiple recordings per subject are split across folds.
3. **Class imbalance** between PD and healthy control (HC) groups.
4. **Feature representation choices** significantly impact classification performance.

This thesis addresses these challenges through a rigorous experimental framework that prioritizes methodological validity over raw performance metrics.

## 1.3 Research Objectives

The primary objectives of this research are:

1. **Develop a reproducible pipeline** for extracting acoustic features from voice recordings.
2. **Evaluate classical machine learning models** (Logistic Regression, SVM, Random Forest) for PD vs. HC classification.
3. **Compare performance** across two distinct datasets with different characteristics.
4. **Investigate the impact** of feature set extension (47 → 78 features) through controlled ablation.
5. **Assess the effect** of class weighting on imbalanced datasets.

## 1.4 Contributions

This thesis makes the following contributions:

- A **subject-grouped cross-validation framework** that prevents data leakage in multi-recording datasets.
- A **controlled feature ablation study** demonstrating an ~+8.7 percentage point ROC-AUC improvement with extended features.
- **Systematic comparison** of class weighting strategies across different imbalance levels.
- **Transparent documentation** of methodological constraints and their implications.

## 1.5 Thesis Organization

The remainder of this thesis is organized as follows:

| Chapter | Title | Description |
| --- | --- | --- |
| 2 | Literature Review | Survey of voice-based PD detection methods |
| 3 | Data Description | Detailed analysis of datasets used |
| 4 | Methodology | Feature extraction and ML pipeline design |
| 5 | Experimental Design | Cross-validation and evaluation protocols |
| 6 | Results | Quantitative findings across all conditions |
| 7 | Discussion | Interpretation and comparison with literature |
| 8 | Limitations | Constraints and threats to validity |
| 9 | Conclusion | Summary of findings and future directions |

## 1.6 Scope Boundaries

This research is explicitly bounded by the following constraints:

- **Binary classification only** (PD vs. HC) – no severity prediction.
- **Classical ML only** – no deep learning or end-to-end models.

- **Research context only** – no clinical deployment or diagnostic claims.
- **Reproducibility prioritized** – methodological validity is valued over leaderboard-style accuracy maximization.

# Chapter 2: Literature Review

## 2.1 Parkinson's Disease and Speech Impairment

Parkinson's disease is a progressive neurodegenerative disorder primarily known for its motor symptoms such as tremor, rigidity, and bradykinesia. In addition to these well-recognized motor features, PD almost invariably affects speech and voice as the disease progresses. In fact, studies report that approximately 70–90% of individuals with PD develop measurable speech and voice impairments. This collection of speech symptoms in PD is often referred to as **hypokinetic dysarthria**, indicating a characteristic pattern of speech motor control impairment associated with the disease.

The speech of a person with PD typically exhibits several hallmark changes. One prominent feature is **hypophonia**, or reduced voice loudness – patients often speak in a much softer voice than normal. Another is a monotonic pitch: PD speakers tend to have a limited range of pitch, resulting in speech that lacks the normal ups and downs of intonation (often described as "monopitch" speech). **Monoloudness** (little variation in volume) often accompanies this, so the overall prosody (melody and expressiveness of speech) is markedly diminished. Patients may also exhibit **articulatory imprecision**, where consonants in particular are not enunciated crisply. For example, consonant sounds may blur together or be undershot due to reduced range of motion in the articulators (jaw, tongue, lips). The voice quality in PD is frequently described as breathy or hoarse, reflecting incomplete vocal fold closure and other phonatory deficits. Additionally, some individuals speak with an improperly fast rate or with short rushes of speech which, combined with the articulation issues, can reduce intelligibility. These speech characteristics – reduced loudness, monopitch, monoloudness, hoarse/breathy voice, and imprecise articulation – are widely observed in PD and form the basis of clinical descriptions of hypokinetic dysarthria.

Crucially, speech changes in PD are of interest not just as symptoms affecting communication, but also as potential non-invasive biomarkers of the disease. Voice is relatively easy to capture (e.g. via a short recording on a phone), and vocal changes can manifest early in the disease course. Some research suggests that subtle voice abnormalities may appear even **before** classic motor symptoms in some patients. Because voice recording and analysis can be done inexpensively and remotely, there is considerable motivation to use speech as a way to detect or monitor PD without the need for invasive tests. Speech and voice metrics are appealing for telemedicine and longitudinal tracking of PD progression. Unlike many clinical tests that require in-person visits and specialized equipment, voice recordings can be obtained by patients at home and sent to clinicians or analyzed by algorithms, enabling more frequent monitoring.

It should be noted, however, that the speech impairments in PD can vary greatly across patients and disease stages. Not every person with PD will have all the aforementioned speech symptoms, and the severity can range from very mild to highly debilitating. There is variability in how early voice changes emerge: some patients present with noticeable hypophonia and monotony in the early stages, whereas others might have minimal speech impact until later in the disease. Moreover, the progression of speech symptoms does not always strictly parallel the progression of other motor symptoms. For example, a patient with advanced limb tremor might still be understandable in speech, while another with relatively moderate overall motor signs could have severe dysarthria. In general, as PD progresses, speech tends to worsen – volume may

further decrease and articulation may become more slurred. But importantly, these changes are not uniform or perfectly correlated with disease duration. Factors such as individual patient differences, co-occurring conditions (like age-related voice changes), and even treatment effects (medications or speech therapy) can influence the speech presentation. This variability means that while speech is a promising biomarker, one must be careful in using it for diagnosis or monitoring: any voice-based assessment needs to account for the wide range of "normal" for PD speech and the overlap with speech characteristics of other populations (e.g. normal aging can also cause some reduced loudness or hoarseness). In summary, PD provides a compelling case for voice analysis – it is a neurodegenerative disorder with clear motor effects on speech production, speech changes are prevalent and can be captured non-invasively, but these changes are heterogeneous across individuals and time. This establishes why voice is relevant for PD research while cautioning that clinical diagnostic use would require careful handling of variability and uncertainty, rather than treating voice patterns as a definitive signature of the disease.

## 2.2 Acoustic Characteristics of Parkinsonian Speech

A variety of acoustic features have been explored to characterize the distinctive patterns of Parkinsonian speech. The motivation for examining these features is that they quantify specific aspects of voice and speech that are affected by PD, thus providing measurable indicators that can be used in analysis or as inputs to classification models. Broadly, these features can be grouped into categories or feature families based on what aspect of speech they describe. In this section, the discussion is organized by three major families of acoustic features: **prosodic features**, **perturbation measures**, and **spectral/cepstral features**. This organization corresponds to methodologies commonly employed in PD voice studies (and later in our experiments), where these feature types are leveraged. Each subsection below defines the feature group, gives examples of specific features in that group, and explains what changes have been observed in PD speech relative to normal speech in that feature domain.

### 2.2.1 Prosodic Features

Prosodic features relate to the pitch (fundamental frequency) and loudness (intensity) patterns in speech, as well as timing/rhythm to some extent. The fundamental frequency of the voice (notated as F0) corresponds to the perceived pitch. Prosodic analysis often looks at statistics of F0, such as the mean pitch, range (difference between highest and lowest pitch), and the standard deviation of pitch across an utterance, to gauge how much variation in pitch a speaker uses. Loudness can be quantified through overall intensity level (in decibels) and its variability or emphasis patterns (e.g., how much a speaker modulates volume for stress). In typical expressive speech, healthy speakers vary both pitch and loudness to convey emphasis, emotion, or sentence modality (e.g., distinguishing a question from a statement).

In Parkinson's disease, a well-documented phenomenon is the reduction of prosodic variability. PD patients often speak in a monotone – their pitch remains relatively flat and at a narrow range, lacking the normal ups and downs. Objectively, one finds a lower standard deviation of F0 and a smaller pitch range in PD speech compared to healthy age-matched controls. This is sometimes described clinically as monopitch. Likewise, PD speakers exhibit monoloudness: their volume tends to be more constant and generally softer than normal. They may not employ the usual loudness increases to stress important words or to express emotion. The term hypophonia specifically refers to the reduced overall loudness (soft voice) that is common in PD. Together, monopitch and monoloudness make PD speech sound flat or expressionless. For instance, where a healthy speaker might vary pitch and loudness dynamically within a single sentence

("Really? I can't believe it!"), a person with PD might deliver the same sentence in a relatively uniform tone and volume, which can be perceived as lacking affect.

Empirical studies support these observations. For example, most PD patients have significantly lower pitch variability and reduced intensity modulation, resulting in perceptually monotonic and weak speech. These prosodic deficits can be quantified: for instance, when computing the pitch range in a reading passage, a PD speaker might show only a few semitones of variation, versus perhaps an octave for a healthy speaker. Similarly, intensity traces from PD speech often appear "flatter." Prosodic features like F0 range, F0 variability (e.g. variance or interquartile range of F0), intensity range, and intensity standard deviation are therefore commonly included in acoustic analyses. They capture the diminished expressivity in PD speech that comes from rigidity and bradykinesia affecting the vocal apparatus (including respiratory support and laryngeal control). In summary, prosodic measures in PD typically indicate reduced pitch and loudness variability, aligning with the clinical description of hypokinetic dysarthria where speech has a monotone, soft character. These features are important later when designing classifiers because they directly reflect how PD impacts a speaker's control over voice dynamics.

### 2.2.2 Perturbation Measures

Perturbation measures are acoustic features that capture the cycle-to-cycle variations in the voice signal, reflecting stability (or instability) of vocal fold vibration. The two primary perturbation measures are **jitter** and **shimmer**. *Jitter* quantifies the minute fluctuations in pitch period from one glottal cycle to the next – essentially, variability in the fundamental frequency. *Shimmer* quantifies the variability in amplitude (loudness) across successive glottal cycles. In a perfectly steady, clear voice, one would expect nearly constant pitch period and amplitude for each cycle of vocal fold vibration (especially during sustained phonation of a vowel). However, human voices always have some natural jitter and shimmer. Pathologies that affect vocal fold control tend to increase these perturbations.

In Parkinson's disease, due to factors like reduced vocal fold adduction, vocal tremor, and inconsistent breath support, jitter and shimmer are often elevated compared to age-matched healthy controls. That is, PD voices typically show more frequent, irregular fluctuations in frequency and amplitude. This corresponds to the perceptual observation of a hoarse or unsteady voice quality. For example, when sustaining a vowel sound like "ah," a healthy voice might sound steady, whereas a PD voice might quaver slightly in pitch and/ or volume. These micro-instabilities are precisely what jitter and shimmer measure. **Harmonics-to-Noise Ratio (HNR)** is another related metric which compares the level of periodic (harmonic) energy in the voice to the level of aperiodic or noise energy. A high HNR means the voice signal is very harmonic (a clean tone), whereas a low HNR indicates a noisier, breathier voice. PD voices tend to have lower HNR values, indicating a higher proportion of noise (breathiness, roughness) in the voice.

To put some numbers for illustration: a healthy sustained vowel might have a jitter on the order of 0.5% (very small period fluctuations) and shimmer around 3–4%, with an HNR of, say, 20 dB or more. In a PD voice, jitter might be several times higher (reports of jitter in PD can be, for instance, 1–2% or more) and shimmer likewise elevated, while HNR might drop to, say, 10–15 dB. Many studies corroborate that PD causes increased frequency and amplitude perturbation in the voice and a corresponding increase in spectral noise. These measures provide objective evidence of the vocal instability and glottal insufficiency that are characteristic of Parkinsonian dysarthria (often described as a "breathy and hoarse" voice quality).

### 2.2.3 Spectral and Cepstral Features

Beyond prosodic and perturbation measures, PD can also manifest in altered spectral characteristics of the voice. Notable spectral/cepstral features investigated include:

- **Formant frequencies** (e.g., F1, F2, F3): PD can affect articulation, leading to reduced range or shifting of formant values that represent resonant frequencies of the vocal tract.
- **Mel-Frequency Cepstral Coefficients (MFCCs)**: These are a compact representation of the spectral envelope. Differences in MFCC patterns have been observed between PD and control speakers, capturing broad changes in the voice spectrum.
- **Spectral envelope shape descriptors**: Measures like spectral centroid, bandwidth, roll-off, and flatness can reflect how energy distribution changes (e.g., increased spectral noise or reduced high-frequency energy in PD voices).

Spectral features often require careful interpretation, as they can be influenced by multiple factors (articulation, phonation, recording conditions). However, they provide a high-dimensional description of voice quality and articulation that can complement the lower-dimensional prosodic and perturbation features. Modern approaches frequently include MFCCs and their derivatives, as these coefficients effectively summarize the voice spectrum and have proven useful in many speech classification tasks.

## 2.3 Feature Extraction Approaches

### 2.3.1 Traditional Acoustic Features

Early studies on PD voice detection relied largely on *clinically-motivated acoustic features*. These include fundamental frequency measures, perturbation metrics, and formant analyses, often reflecting well-known clinical observations of PD speech. For example:

| Feature Category | Examples | Physiological Basis |
|---|---|---|
| Fundamental Frequency | F0 mean, F0 standard deviation | Vocal fold tension & control |
| Perturbation | Jitter, Shimmer | Neuromuscular control stability |
| Noise ratios | HNR, NHR (Noise-to-Harmonics) | Incomplete glottal closure (breathiness) |
| Formants | F1, F2, F3 frequencies | Articulator positioning (tongue, jaw, etc.) |

These features can be extracted with tools like Praat or its Python interface (Parselmouth). They provided early evidence that PD affects voice in measurable ways, but individually each feature captures only one aspect of the complex speech impairment.

### 2.3.2 Spectral Features

As research progressed, more comprehensive signal processing features were introduced to capture subtler aspects of the speech signal:

- **MFCCs (Mel-Frequency Cepstral Coefficients)** – Typically a set of 12–13 coefficients (plus energy) capturing the shape of the spectral envelope on a perceptual scale. MFCCs are widely used in speech processing and were applied to PD detection to exploit broader spectral patterns.
- **Delta and Delta-Delta MFCCs** – First and second time derivatives of MFCCs, which capture how spectral features change over time. These can reflect vocal stability and articulation dynamics.
- **Spectral shape features** – Such as spectral centroid (weighted mean frequency), spectral bandwidth, spectral roll-off (frequency below which a certain percentage of total energy lies), and spectral flatness. These aggregate descriptors can indicate, for example, an excess of high-frequency noise or a concentration of energy in certain bands.

The inclusion of these spectral features greatly increases the dimensionality of the feature space (often into the dozens or hundreds of features) and can improve classification performance, but it also requires careful validation to avoid overfitting, especially on small datasets.

### 2.3.3 Emerging Deep Learning Features

In recent years, some studies have explored *deep learning* approaches to feature extraction, such as using raw spectrograms or learned embeddings from neural networks (CNNs or LSTM-based models). These methods aim to automatically learn features directly from data. However, deep learning approaches in PD voice classification face significant challenges:

- They require large labeled datasets to train effectively, which are often not available in this domain (most PD voice datasets have only tens of subjects).
- The learned features are typically not easily interpretable, which is a drawback for clinical acceptance where understanding the model's decision factors is important.
- Without careful design, deep models risk overfitting to spurious artifacts (especially if data leakage is present) and might not generalize well.

As a result, classical feature-based approaches remain dominant for PD voice studies, and deep learning is generally considered experimental unless combined with techniques to mitigate the data scarcity and interpretability issues.

## 2.4 Machine Learning Approaches

### 2.4.1 Classical Methods

A range of classical machine learning algorithms have been applied to the PD voice classification task. These include simple linear models and more complex non-linear ensemble methods. Three commonly used classifiers are:

| Method | Strengths | Limitations |
|---|---|---|
| **Logistic Regression** | Interpretable coefficients (each feature's weight indicates its influence) | Assumes a linear decision boundary; may underfit complex patterns |
| **Support Vector Machine (RBF kernel)** | Effective in high-dimensional spaces; can model non-linear relationships via the kernel trick | Choice of kernel and hyperparameters (C, gamma) is critical; not easily interpretable |
| **Random Forest** | Handles non-linearity and feature interactions; provides measures of feature importance; robust to outliers | Less interpretable than linear models (individual trees are interpretable, but ensemble is complex); can overfit if not tuned |

These classical models have the advantage of being relatively **data-efficient** (important for small datasets) and **interpretable** to varying degrees, which is valuable when the goal is to derive insights about which vocal features are most affected by PD.

### 2.4.2 Deep Learning

Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs, including LSTM/GRU) have been applied in some studies, especially using spectrogram representations of voice recordings. While these can capture complex patterns automatically, in the context of PD classification they come with notable challenges:

- **Data Requirements:** Deep models have a high number of parameters and typically require hundreds or thousands of training examples to generalize well. PD voice datasets are often too small for this, leading to overfitting.
- **Overfitting Risk:** Without sufficient data or augmentation techniques, a deep network might latch onto dataset-specific quirks (or even patient identity, if leakage is present) rather than true disease biomarkers.
- **Interpretability:** The features learned by deep networks are abstract and not directly tied to known clinical measures. This makes it difficult to explain why a network made a particular prediction, which limits clinical trust.

For these reasons, many researchers in this domain continue to use classical ML or manually engineered features, or they use deep learning in a hybrid fashion (e.g., extracting deep features but then feeding them into a simpler model for classification).

## 2.5 Methodological Concerns in Literature

A critical aspect of the literature is the **evaluation methodology**. Many early studies reported impressive classification accuracies (> 90%) for PD vs. control based on voice. However, subsequent investigations revealed that some of these results were overly optimistic due to **flawed validation strategies**. Three key concerns are:

### 2.5.1 Data Leakage

Data leakage occurs when information that would not be available in a real-world test scenario inadvertently influences model training or evaluation. In PD voice studies, the most glaring form of leakage has been **subject identity leakage** across training and test sets. This happens when multiple recordings from the same individual are split between train and test. In a random record-level split (shuffling all recordings regardless of subject), a model can effectively "recognize" the person's voice it saw in training and thus achieve unnaturally high accuracy in testing. Early studies that used naive random splits unknowingly benefited from this effect, reporting accuracy figures that were not achievable on truly unseen patients. For example, if a patient's distinctive vocal characteristics (unrelated to PD per se) are learned, the model may classify another recording of the same person correctly simply by voice recognition.

To address this, later studies implemented **subject-level splits**. In a subject-level split, all recordings from any given person are confined to either training or testing, never both. A common approach is leave-one-subject-out cross-validation or grouped k-fold CV where each fold contains a unique set of subjects. This way, the model is always tested on voices it never encountered during training, which is a far more stringent and clinically realistic evaluation. As a concrete example, *Egbo et al. (2025)* evaluated a model on a 31-subject voice dataset using explicit subject-level stratification; despite using advanced techniques, they report that rigorous validation kept performance in a realistic range (their best method achieved ~98% accuracy, but only after careful model tuning and feature selection, illustrating what is possible under ideal conditions when leakage is avoided). The key lesson from such work is that proper data splitting can substantially reduce overly optimistic estimates. In our context, preventing subject leakage is mandatory for Dataset A (which contains multiple recordings per subject). Studies that fail to do this (still unfortunately common in older literature) can mislead the community about which methods are truly effective.

### 2.5.2 Class Imbalance

Class imbalance refers to the situation where one class (PD or HC) significantly outnumbers the other. This is evident in many PD voice datasets – for instance, some datasets have 3 or 4 times as many PD samples as controls (or vice versa). In such cases, using raw accuracy as a metric can be misleading because a model could achieve high accuracy by simply guessing the majority class most of the time. The literature has noted that many studies did not adequately address imbalance; some reported high accuracy without disclosing that the model might be trivial (e.g., always predicting "PD" in a heavily imbalanced set). Best practices include reporting metrics like sensitivity (recall) and specificity or precision and recall in addition to accuracy, and employing techniques like **class weighting** or **resampling** to mitigate imbalance. This thesis, for example, explicitly investigates the effect of scikit-learn's `class_weight="balanced"` option, which automatically weights classes inversely proportional to their frequency. The expectation from literature (and our own investigations) is that class weighting or related strategies can improve recall of the minority class, though the impact may be modest if the imbalance is not extreme.

### 2.5.3 Reproducibility

Reproducibility is a general concern in machine learning research, and PD voice studies are no exception. Literature surveys have pointed out that many papers lack sufficient detail for others to replicate their experiments. Common issues include:

- **Unspecified feature extraction parameters:** For instance, not stating which exact settings or software were used to compute jitter/shimmer or MFCCs, which can vary and affect results.
- **Omission of random seed or fold assignments:** Without fixing random seeds or explicitly listing which subjects were in each fold, it's hard to know if a reported result is an average or cherry-picked best case.
- **Incomplete description of cross-validation:** As noted, whether cross-validation was subject-wise or record-wise is sometimes unclear, making it difficult to interpret the results.
- **Data accessibility:** Using private or proprietary datasets without providing access, hindering comparative studies.

In response, the field has been moving towards better practices. Researchers now more often share code, use open datasets (like the UCI or Kaggle sets), and carefully document methods. This thesis aligns with those principles by using publicly available data when possible, releasing code, and fully describing the experimental setup to facilitate reproduction.

## 2.6 Benchmark Datasets

Several benchmark datasets have been widely used in voice-based PD detection research. Each has its strengths and limitations:

### 2.6.1 UCI Parkinson's Dataset (Oxford Parkinson's Telemonitoring)

One of the earliest and most commonly cited datasets is the UCI Machine Learning Repository "Parkinson's Disease Data Set" introduced by Little et al. (2009). It consists of voice measurements from 31 individuals (23 with PD, 8 healthy), each providing multiple sustained vowel phonations. It contains 195 samples with 22 features per sample, including fundamental frequency, several jitter/shimmer measures, and HNR. While extensively used (and yielding high classification scores in many papers), its limitations include the lack of subject identifiers for multiple entries and a relatively small sample size. It also represents a very specific task (sustained vowel "ahh"), which may not generalize to conversational speech.

### 2.6.2 MDVR-KCL Dataset

The **Mobile Device Voice Recordings – King's College London (MDVR-KCL)** dataset is a more recently introduced corpus aimed at capturing more natural speech. It contains raw audio recordings from PD patients and healthy controls performing two tasks: reading a standardized text and engaging in a spontaneous dialogue. Importantly, subject IDs are available and multiple recordings per subject are included. This allows experiments with **grouped cross-validation** (as done in this thesis for Dataset A). The dataset is relatively small in terms of subjects (37 subjects, ~73 recordings total in our use after excluding some anomalies), but it provides real acoustic signals which require a full feature extraction pipeline. We refer to this as Dataset A in our work.

### 2.6.3 PhysioNet / PC-GITA / Other Corpora

There are a number of other datasets, such as the PC-GITA Spanish Parkinson's voice dataset, the Italian Parkinson's Voice Dataset, etc. Many contain sustained vowels, words, or sentences recorded in clinical settings. Some are available through PhysioNet or other platforms. Each comes with its own recording conditions and demographics. While we do not directly use these in this thesis, the literature has used them for cross-language studies or to test methods' generalizability.

### 2.6.4 Kaggle "PD Speech Features" Dataset

A notable recent contribution is a Kaggle dataset often titled *Parkinson's Disease Classification* or *PD Speech Features*, which provides a large set of precomputed features for a cohort of patients and controls. Specifically, it contains **756 samples (columns)** with **752 features** each, derived from voice recordings of 188 PD patients and 64 healthy individuals (each subject contributing multiple samples). The features span a wide range: fundamental frequency measures, time-frequency features, MFCCs and their deltas, wavelet transform features, etc., including advanced ones like Tunable Q-factor Wavelet Transform (TQWT) coefficients (hundreds of features). Sakar et al. (2019) present a comparative analysis using this feature set, demonstrating that using these comprehensive features (especially the TQWT-based ones) improved classification performance. We use this as **Dataset B** for benchmark comparison. However, **no subject IDs** are provided in this dataset. Therefore, any evaluation on it (including ours) must caution that some samples could be from the same subject split across folds, potentially inflating performance. Indeed, many literature results with this dataset report very high accuracy/ROC-AUC (sometimes > 0.90), but without subject grouping it's unclear how much is due to true generalization versus repeated subjects. In this thesis we include this dataset to compare results, but always interpret them with that caveat.

## 2.7 Research Gap

From the literature, it is evident that voice-based PD detection is feasible and many studies have reported promising accuracy. However, few works have systematically addressed the combination of challenges together. In particular, there is a gap when it comes to studies that simultaneously:

1. Use **grouped cross-validation** to eliminate subject leakage in multi-recording datasets.
2. Perform **controlled feature ablation** to quantify the value of expanded feature sets.
3. Analyze **class imbalance mitigation** strategies in the context of PD voice data.
4. Provide **transparent limitations** to contextualize their performance claims.

These gaps suggest that further work is needed focusing on methodological consistency and validation rigor, rather than solely on pushing accuracy. This thesis is designed to fill these needs by prioritizing reproducibility and valid comparisons in its experimental design.

## 2.8 Summary

In summary, the literature demonstrates that voice-based PD detection can achieve high accuracy with classical ML methods, provided that informative features are used. At the same time, methodological rigor varies significantly across studies. The highest reported figures often come from less rigorous validation, whereas more stringent evaluations yield moderate but more reliable performance. This thesis adopts a conservative approach in line with emerging best practices: features are chosen based on known PD speech characteristics, cross-validation is done at the subject level (when possible) to reflect real-world conditions,

and results are reported with appropriate uncertainty. By doing so, the work builds on the rich foundation of prior research while aiming to advance the field's standards for reliability and interpretability in PD voice classification.

# Chapter 3: Data Description

## 3.1 Overview

This thesis utilizes **two distinct datasets** for Parkinson's Disease voice classification, referred to as Dataset A and Dataset B:

| Property | Dataset A (MDVR-KCL) | Dataset B (PD_SPEECH) |
|---|---|---|
| Data Type | Raw audio recordings (WAV) | Pre-extracted features (CSV) |
| Source | Zenodo (research study) | Kaggle (public repository) |
| Unit of Analysis | Subject (multiple recordings) | Sample (single feature vector) |
| Subject IDs Available | Yes (identifiers in filenames) | No (anonymous samples) |
| Total Samples | 73 recordings (37 subjects) | 756 samples (188 PD, 64 HC subjects; 3 recordings each) |

Each dataset offers unique advantages and serves a different purpose in our experiments. **Dataset A** provides raw speech data enabling feature extraction and strict validation control (grouping by subject), reflecting a realistic clinical scenario with limited data. **Dataset B** provides a much larger sample size and diverse feature set, useful as a benchmark, though it lacks certain metadata (subject linkage).

## 3.2 Dataset A: MDVR-KCL

### 3.2.1 Source and Collection

The *Mobile Device Voice Recordings from King's College London (MDVR-KCL)* dataset was collected for PD research using smartphone recordings. It is publicly available on Zenodo (DOI: `10.5281/zenodo.2867215`). The dataset comprises audio files organized by speech task and diagnosis. We downloaded the audio data and organized it under our `assets/DATASET_MDVR_KCL/` directory.

**Location:** `assets/DATASET_MDVR_KCL/` (locally) – containing subfolders per task and diagnosis.

### 3.2.2 Speech Tasks

The dataset includes two distinct speech tasks performed by the subjects:

- **ReadText:** Reading a standardized passage (the same text for all participants).

• **SpontaneousDialogue:** Engaging in unscripted conversation or answering open-ended questions.

Not all subjects performed both tasks due to recording issues or other factors. Table 3.1 summarizes the subject counts:

| Task | Description | Subjects (Total) | HC (Controls) | PD (Patients) |
|------|-------------|------------------|---------------|---------------|
| **ReadText** | Reading a standardized passage | 37 | 21 | 16 |
| **SpontaneousDialogue** | Unstructured conversation | 36 | 21 | 15 |

*Note:* Subject ID 18 is missing from the SpontaneousDialogue task (they only did the reading task).

### 3.2.3 Class Distribution

In terms of class breakdown (healthy vs. PD), Dataset A has a relatively balanced composition within each task:

- **ReadText Task:** 21 HC vs. 16 PD (approximately 57% HC, 43% PD).
- **SpontaneousDialogue Task:** 21 HC vs. 15 PD (approximately 58% HC, 42% PD).

```
ReadText Task:
├── HC (Healthy Control): 21 subjects (56.8%)
└── PD (Parkinson's Disease): 16 subjects (43.2%)

SpontaneousDialogue Task:
├── HC (Healthy Control): 21 subjects (58.3%)
└── PD (Parkinson's Disease): 15 subjects (41.7%)
```

Overall, the class imbalance is **moderate (~57:43)** in this dataset. This imbalance is addressed via class weighting experiments in our modeling (see methodology).

### 3.2.4 File Structure and Naming

The raw audio files are organized by task and diagnosis, for example:

```
DATASET_MDVR_KCL/
├── ReadText/
│   ├── HC/
│   │   └── IDxx_hc_*.wav
│   └── PD/
│       └── IDxx_pd_*.wav
└── SpontaneousDialogue/
```

```
        ├── HC/
        └── PD/
```

Each filename encodes the subject ID and their group (hc or pd), along with recording indices. For instance, a file `ID04_pd_2_0_1.wav` might indicate subject ID 04, PD, with certain task/session indices. We wrote parsing code to reliably extract subject IDs from filenames.

### 3.2.5 Known Anomalies

A few data quirks were noted and handled:

- **Subject ID 22:** This subject's file naming slightly deviates from the pattern (missing an underscore in one case). The code accounts for this when reading filenames.
- **Subject ID 18:** As mentioned, this subject has no SpontaneousDialogue recording, so they appear only in ReadText.
- **Multiple recordings per subject:** Each subject may have multiple takes in a task (e.g., ID04 has `_2_0_0.wav` and `_2_0_1.wav` indicating two recordings). For analysis, features from all recordings are used, but crucially, all recordings of a subject are kept in the same cross-validation fold to avoid leakage.

## 3.3 Dataset B: PD Speech Features

### 3.3.1 Source

Dataset B consists of pre-extracted acoustic features compiled for a PD classification challenge on Kaggle (originally from a research study by Sakar et al.). It contains a single CSV file with 756 rows (samples) and 754 columns (752 features plus two identifiers: subject type and ID or recording index). The dataset aggregates a comprehensive set of voice features for a large number of samples. We refer to this CSV as `PD_SPEECH_FEATURES.csv`.

**Location:** `assets/PD_SPEECH_FEATURES.csv` – a single comma-separated values file.

### 3.3.2 Class Distribution

The class breakdown in Dataset B is significantly imbalanced in favor of PD cases:

| Class | Samples | Percentage |
|---|---|---|
| **HC (0)** | 192 | 25.4% |
| **PD (1)** | 564 | 74.6% |

Out of 756 total samples, 564 are labeled PD and 192 as healthy control. The imbalance ratio is roughly 1:3 (one control for every three PD samples). This **severe imbalance (~25:75)** necessitates careful handling (e.g., class weighting) to ensure models do not simply predict the majority class.

### 3.3.3 Feature Categories

The 752 features in this dataset span multiple domains of voice signal analysis. A summarized breakdown is:

| Category | Count | Description |
| --- | --- | --- |
| **Baseline Features** | 22 | Fundamental frequency (pitch) measures, jitter, shimmer, HNR variants, etc. (similar to those in Dataset A's extraction) |
| **Intensity** | 3 | Intensity (loudness) statistics (mean, min, max) |
| **Formants** | 36 | Formant frequencies and bandwidths (F1–F4 and related measures) |
| **MFCCs** | 84 | Mel-frequency cepstral coefficients (13 base coefficients × multiple statistics like mean, std, etc.) |
| **Wavelet** | 182 | Wavelet transform features (energy, entropy, etc. across wavelet decomposition levels) |
| **TQWT** | 432 | Tunable Q-factor wavelet transform features (multiple statistics across many decomposition levels) |

These features were likely drawn from prior studies (including the 2019 study by Sakar et al.) which introduced advanced feature extraction (such as TQWT) to capture characteristics of dysarthric speech. The **TQWT features** form the largest block (432 features), indicating multiple sub-band analyses aimed at detecting subtle nonlinear patterns.

It is worth noting that while Dataset A's feature set (47 or 78 features) is a subset of logical, interpretable measures, Dataset B's feature set is far broader and more complex, including many features that may be redundant or highly correlated.

### 3.3.4 Important Caveat

> ⚠ **No Subject Identifiers Available**
> Dataset B does not provide subject IDs. If multiple samples originate from the same subject, stratified cross-validation may inadvertently place samples from one individual in both training and testing sets. This introduces a risk of optimistic bias (a form of data leakage). All results on Dataset B must be interpreted with this caveat in mind: reported performance might be higher than what would be achieved on wholly unseen patients.

We include warnings in our results and discussion whenever Dataset B's outcomes are compared, to ensure clarity about this limitation.

## 3.4 Dataset Comparison

### 3.4.1 Key Differences

A side-by-side comparison of the two datasets highlights their complementary nature:

| Aspect | Dataset A (MDVR-KCL) | Dataset B (Kaggle PD Speech) |
|---|---|---|
| **Sample Size** | 36–37 subjects (73 recordings) | 756 samples (from 252 subjects) |
| **Feature Type** | 47–78 extracted features (after processing audio) | 752 pre-extracted features (provided) |
| **Subject Grouping** | Available (each recording labeled by subject) | Not available (subject unknown for each sample) |
| **Cross-Val Strategy** | Grouped Stratified CV (by subject) | Standard Stratified CV (subject-level unknown) |
| **Class Imbalance** | Moderate (approx. 57% vs 43%) | Severe (approx. 75% vs 25%) |
| **Data Nature** | Real audio – requires preprocessing and feature engineering | Tabular features – ready for modeling |

### 3.4.2 Implications for Analysis

1. **Direct comparison is limited:** We cannot directly compare absolute performance metrics between Dataset A and Dataset B as evidence of one being "easier" or "harder" due to the differences in features and validation. Any higher metrics on B could be due to the larger sample size or inadvertent subject overlap.
2. **Performance on Dataset B might be inflated:** The combination of a large feature set and potential leakage means models might achieve higher ROC-AUC on Dataset B, but this may not translate to truly unseen subjects.
3. **Dataset A results are more conservative:** Because of grouped CV and the smaller sample size, metrics on A are likely more cautious estimates of generalization performance. These may actually better reflect real-world expectations for a new patient.

In this thesis, we report results on both datasets but place greater emphasis on the trends observed within each and the insights gleaned (e.g., effect of feature extension and class weighting) rather than naively pooling or comparing their raw performance.

## 3.5 Data Preprocessing

Before analysis, each dataset undergoes certain preprocessing steps appropriate to its nature.

### 3.5.1 Dataset A Preprocessing

For the MDVR-KCL audio data, the following pipeline is applied:

1. **Audio Loading:** Each WAV file is loaded at its native sample rate (typically 44.1 kHz). If stereo, it is converted to mono by averaging channels (though most recordings are likely mono).
2. **Silence Trimming:** Leading and trailing silences are removed using an energy-based threshold to ensure we analyze only the active speech portions.

3. **Feature Extraction:** Using our feature extraction tool (based on Parselmouth for prosodic features and librosa for spectral features), we extract the defined set of features for each recording. This yields one feature vector per recording.
4. **Aggregation by Subject (for CV):** While each recording yields its own feature vector, for cross-validation grouping we associate each vector with a subject ID so that splits can be done at the subject level.

The result is a feature matrix X_A (with one row per recording, columns as features) and a label vector y_A (PD or HC for each recording), along with a parallel array of subject IDs for grouping.

### 3.5.2 Dataset B Preprocessing

For the PD Speech Features CSV:

1. **Load CSV:** The entire CSV is read into a pandas DataFrame. Each row is an instance with an associated label (provided in the dataset as a column, e.g., `class` 0 or 1).
2. **Feature Selection:** All 752 feature columns are retained. (No feature is dropped a priori, though some could be constant or non-informative; for completeness we use all features as provided.)
3. **Standardization:** We apply z-score normalization to each feature (subtract mean, divide by standard deviation) during model training – but ensuring the scaling parameters are derived from training folds only, to avoid leaking information into test folds. This is handled inside the cross-validation pipeline.
4. **No Additional Preprocessing:** Since features are already given, we do not perform any audio-level processing or augmentation. The dataset is ready for modeling after loading and scaling.

## 3.6 Summary

The two datasets provide complementary perspectives for our study:

- **Dataset A** is a *small, high-quality dataset* with raw recordings, enabling us to demonstrate a full pipeline (from audio to features to prediction) and to enforce a stringent evaluation (grouped CV). Results on Dataset A are expected to be conservative but robust, highlighting challenges like variance due to limited data.
- **Dataset B** is a *large, rich feature dataset* that allows testing our models on a broader set of features and many more samples. It serves as a benchmark to see how our pipeline performs when given abundant features, but its lack of subject IDs means we interpret its results with caution.

By analyzing both, we can observe how conclusions might differ under different data conditions and ensure that any insights (e.g., the value of extended features or class weighting) hold in more than one scenario. The next chapter will detail the feature extraction process for Dataset A and the unified methodology used to train and evaluate models on both datasets.
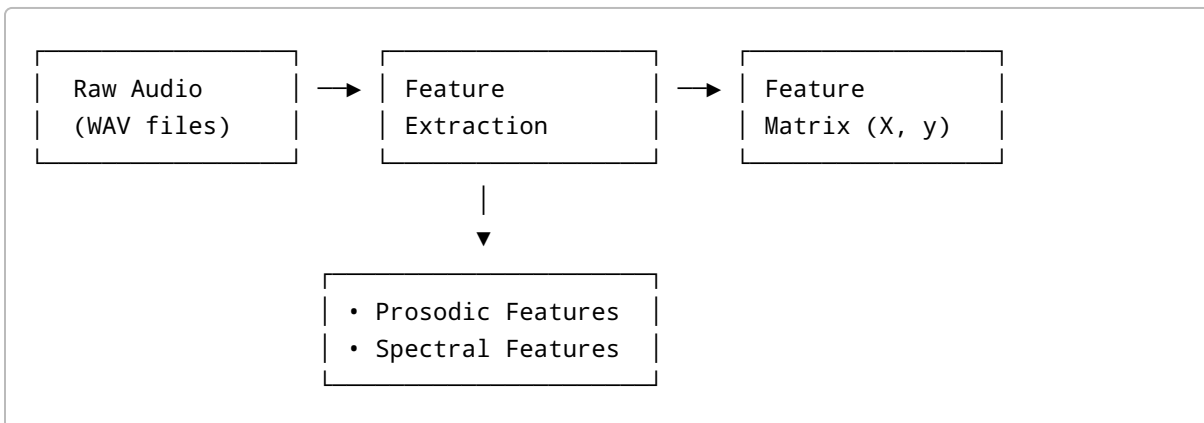
# Chapter 4: Methodology

## 4.1 Overview

This chapter describes the **feature extraction pipeline**, machine learning models, and evaluation framework used in this thesis. The methodology emphasizes reproducibility and methodological rigor over

raw performance optimization. We first outline how features are obtained from raw audio in Dataset A, then define the feature sets and models, and finally describe the cross-validation and metric computation procedures.

## 4.2 Feature Extraction Pipeline

### 4.2.1 Pipeline Architecture

For Dataset A (raw audio), the processing pipeline can be summarized as:

```
┌─────────────────┐     ┌─────────────────┐     ┌─────────────────┐
│ Raw Audio       │ ──▶ │ Feature         │ ──▶ │ Feature         │
│ (WAV files)     │     │ Extraction      │     │ Matrix (X, y)   │
└─────────────────┘     └─────────────────┘     └─────────────────┘
                                 │
                                 ▼
                        ┌─────────────────┐
                        │ • Prosodic Features │
                        │ • Spectral Features │
                        └─────────────────┘
```

In words, each raw WAV recording is processed to extract a set of numeric features, which together form a feature vector (row in a matrix). The collection of all feature vectors constitutes the design matrix for Dataset A, and these can then be fed into machine learning models. (Dataset B skips directly to the feature matrix step, as it is already in feature form.)

### 4.2.2 Audio Preprocessing

Prior to feature extraction on Dataset A recordings, the audio data is preprocessed uniformly:

1. **Load audio** at native sample rate (typically 44.1 kHz for this dataset).
2. **Convert to mono** if the signal has multiple channels (by averaging), to ensure a single amplitude stream.
3. **Normalize amplitude** to a standard range (e.g., -1 to 1) to eliminate any differences in recording gain.
4. **Trim silence** from the start and end of the recording using an energy threshold. This focuses feature extraction on active speech only.

These steps help reduce unwanted variability (e.g., different recording volumes or long silences) that could otherwise skew feature calculations.

### 4.2.3 Prosodic Features (21 features)

Prosodic features capture suprasegmental voice characteristics – essentially how the voice sounds in terms of pitch and loudness over an utterance. We extract 21 prosodic features using Parselmouth (a Python interface to Praat):

- **Pitch (F0) – 4 features:** the mean, standard deviation, minimum, and maximum of the fundamental frequency (in Hz) over the voiced parts of speech.
- **Jitter – 3 features:** several measures of cycle-to-cycle F0 variation (e.g., local jitter, RAP, PPQ5).
- **Shimmer – 5 features:** measures of cycle-to-cycle amplitude variation (local shimmer, APQ3, APQ5, APQ11, and DDA).
- **Harmonicity – 2 features:** voice harmonicity metrics, including mean Harmonics-to-Noise Ratio (HNR) and autocorrelation-based harmonicity, indicating voice periodicity vs. noise.
- **Intensity – 3 features:** mean, standard deviation, and range of intensity (loudness in dB).
- **Formants – 6 features:** mean and standard deviation of the first three formant frequencies (F1, F2, F3) over the utterance, reflecting articulatory properties.

These 21 features align with clinically relevant aspects of PD speech: monotonic pitch/loudness (captured by low F0 variability and intensity range), unstable phonation (captured by high jitter/shimmer, low HNR), and potential changes in articulation (formant shifts).

They are extracted using standard Praat algorithms via Parselmouth.

### 4.2.4 Spectral Features

Spectral features capture frequency-domain characteristics of the voice. We use the librosa library for these features, focusing on representations that summarize the voice spectrum.

**Baseline Spectral Features (26 features):**

- **MFCC mean (13 features):** We compute the first 13 Mel-frequency cepstral coefficients (MFCCs) for frames of the audio and then take the mean of each coefficient across the recording. These capture the average spectral envelope shape on a mel scale.
- **Delta MFCC mean (13 features):** Similarly, we compute first-order delta (time derivative) of MFCCs across frames and take their mean, capturing average spectral change.

These 26 features (13 MFCC means + 13 delta MFCC means) constituted the "baseline" spectral representation.

**Extended Spectral Features (57 features):**

To form the extended feature set (used in certain experimental conditions), we add:

- **MFCC standard deviation (13 features):** the standard deviation of each of the 13 MFCC coefficients across the recording, capturing variability in the spectral envelope.
- **Delta-Delta MFCC mean (13 features):** the mean of second-order delta (acceleration) of MFCCs, adding information about how the rate of spectral change itself varies (useful for capturing tremors or other fluctuations).

- **Spectral shape features (5 features):** five additional descriptors – spectral centroid, spectral bandwidth, spectral rolloff, spectral flatness, and zero-crossing rate – each averaged over time. These provide a high-level characterization of the spectrum (e.g., centroid indicates whether energy is tilted towards high or low frequencies; flatness indicates tonality vs. noise).

Features unique to the extended set are highlighted in **bold** below:

- MFCC mean (13)
- **MFCC std (13)**
- Delta MFCC mean (13)
- **Delta-Delta MFCC mean (13)**
- **Spectral shape (5)** – [centroid, bandwidth, rolloff, flatness, ZCR]

In total, the extended spectral feature set contributes 57 features (compared to 26 in baseline).

**New features in the extended set are shown in bold** in the above list.

### 4.2.5 Total Feature Counts

Combining prosodic and spectral features, we have two configurations:

| Configuration | Prosodic Features | Spectral Features | **Total** |
|---|---|---|---|
| **Baseline** | 21 | 26 | **47** |
| **Extended** | 21 | 57 | **78** |

Thus, Dataset A yields either 47 features per recording (baseline set) or 78 features (extended set), depending on the experimental condition. These features are designed to encapsulate the key measurable differences in speech due to PD, while keeping the dimensionality manageable.

## 4.3 Feature Set Comparison

### 4.3.1 Baseline vs. Extended Features

It is useful to explicitly map which features are included in the baseline set versus the extended set:

```
BASELINE (47 features)          EXTENDED (78 features)
├── Prosodic (21)               ├── Prosodic (21)
│    ├── Pitch (4)              │    ├── Pitch (4)
│    ├── Jitter (3)            │    ├── Jitter (3)
│    ├── Shimmer (5)           │    ├── Shimmer (5)
│    ├── Harmonicity (2)       │    ├── Harmonicity (2)
│    ├── Intensity (3)         │    ├── Intensity (3)
│    └── Formants (6)          │    └── Formants (6)
│                               │
└── Spectral (26)               └── Spectral (57)
```

```
        ├── MFCC mean (13)              ├── MFCC mean (13)
        └── Delta MFCC mean (13)        ├── **MFCC std (13)**        ← *NEW*
                                        ├── Delta MFCC mean (13)
                                        ├── **Delta-Delta MFCC (13)** ← *NEW*
                                        └── **Spectral shape (5)**    ← *NEW*
```

The extended set adds the MFCC standard deviations, delta-delta MFCCs, and spectral shape features (as indicated). These were chosen to systematically evaluate whether including more information (variability and additional spectral descriptors) improves model performance.

### 4.3.2 Rationale for Extended Features

The extended feature set was designed as a **controlled ablation study**: by comparing models trained on 47 features vs. 78 features, we can quantify the benefit of the extra features. The specific additions target complementary information:

1. **MFCC std (13):** Captures within-utterance spectral variability – important for detecting instability in PD speech that average MFCCs might miss. For example, a PD patient might have moments of normal voice and moments of hoarse voice; a high std in certain MFCCs could reflect that inconsistency.
2. **Delta-Delta MFCC (13):** Captures the acceleration of spectral changes – essentially how the vocal tract dynamics might be irregular. This could be sensitive to tremor or to brief accelerations/ decelerations in speech that are symptomatic.
3. **Spectral shape (5):** Provides global descriptors of the spectrum. For instance, spectral centroid might be lower in PD speech if high-frequency energy (from clear consonants) is reduced; spectral flatness might be higher if the voice is noisier.

By including these, we expect to better capture the phenomena of PD speech (e.g., reduced stability and clarity). Our experiments will later show how much improvement these additions yield (or if they yield diminishing returns).

## 4.4 Machine Learning Models

### 4.4.1 Model Selection Rationale

We focus on three **classical machine learning models**: Logistic Regression, Support Vector Machine (with RBF kernel), and Random Forest. The rationale for these choices is:

- **Interpretability:** At least for Logistic Regression and (to some extent) Random Forest, we can interpret feature importance or coefficients. Interpretability is critical because we want to understand which voice features contribute most to classification, aligning with clinical insight.
- **Robustness on Small Data:** Classical models are relatively less prone to overfitting on small datasets than complex deep learning models. They have fewer hyperparameters and can perform well with careful cross-validation even when data is limited.
- **Diversity of approaches:** These three models represent a linear model (Logistic Regression), a kernel-based model capturing non-linear relations (SVM with RBF), and an ensemble model (Random

Forest) that can capture interactions. This gives a broad view of classifier behavior without venturing into too many exotic algorithms.

**4.4.2 Model Specifications**

The specific configurations for each model (chosen based on common practice and some preliminary tuning on separate validation):

- **Logistic Regression:** A linear model with L2 regularization (C=1.0, which is the default in scikit-learn, meaning no strong regularization). We increased the maximum iterations to 1000 to ensure convergence given the possibly non-separable data. This model produces a weighted sum of features to make predictions.
- **Support Vector Machine (RBF kernel):** We use an RBF kernel SVM with default hyperparameters (C=1.0, gamma='scale' which in scikit-learn sets gamma to 1/(number of features)). The SVM finds a non-linear decision boundary in feature space. We did not perform extensive hyperparameter tuning; thus, the SVM's performance can indicate baseline kernel method performance.
- **Random Forest:** An ensemble of 100 decision trees (n_estimators=100) with a maximum tree depth of 10 (to prevent the trees from overfitting too deeply). The Random Forest can model complex decision boundaries and provides an estimate of feature importance via mean decrease in impurity. We chose a moderate depth to balance bias-variance and to keep feature importance interpretable (deeper trees can sometimes dilute the importance measures).

Table 4.1 summarizes these settings:

| Model | Type | Key Parameters |
|---|---|---|
| **Logistic Regression** | Linear | C = 1.0 (L2 reg.), max_iter = 1000 |
| **SVM (RBF)** | Kernel | C = 1.0, gamma = 'scale' (auto) |
| **Random Forest** | Ensemble | n_estimators = 100, max_depth = 10 |

These parameters remain fixed throughout the experiments (no hyperparameter search is performed, as discussed in limitations, to avoid overfitting given the small data).

**4.4.3 Class Weighting**

Class imbalance is addressed via the `class_weight` parameter available in scikit-learn for these models. We run experiments under two conditions:

- **Unweighted (baseline):** `class_weight = None` – the model treats all samples equally, which means it will optimize accuracy potentially at the expense of the minority class.
- **Weighted:** `class_weight = "balanced"` – the model automatically adjusts weights inversely proportional to class frequencies. In effect, misclassifying a PD sample (the majority class in Dataset B, minority in some subsets of Dataset A) and misclassifying a HC sample incur different penalties to balance their influence.

All three chosen models support class weighting natively. For Logistic Regression and SVM, this weights the loss function during training. For Random Forest, it weights how the Gini impurity splits are evaluated and how votes are counted.

In summary, we will be able to compare, for each model, performance with and without class weighting. This helps to determine if imbalance was harming the detection of the under-represented class and if weighting alleviates that.

(Pseudocode for how we set it in models: for unweighted `class_weight=None`; for weighted `class_weight="balanced"`.)

## 4.5 ML Pipeline Architecture

### 4.5.1 Pipeline Structure

For each model and dataset, we implement a scikit-learn Pipeline encapsulating preprocessing and classification. The general form is:

```
Pipeline([
    ('scaler', StandardScaler()),
    ('classifier', Model(class_weight=...))
])
```

This pipeline first standardizes features (zero mean, unit variance) and then applies the classifier. Standardization is important for SVM (which is distance-based) and Logistic Regression; for Random Forest it is not strictly needed, but we include it for consistency.

We ensure that the scaler is fit on training data only within each cross-validation fold to avoid data leakage (the pipeline makes this automatic when used in cross_val_score or similar functions).

### 4.5.2 Standardization

All features are numeric and on different scales (e.g., jitter ~0.01, frequencies ~100–300 Hz, MFCCs arbitrary units, etc.), so standardization is critical. We use the formula:

$$ z = \frac{x - \mu}{\sigma} $$

where $\mu$ and $\sigma$ are the mean and standard deviation of a feature computed from the training set (e.g., within a CV fold). This transform is applied to both training and test data for each fold, using training-set statistics, to prevent any information from the test data leaking via the scaling.

By scaling, we ensure that no single feature dominates due to scale alone and that the models' regularization works properly (especially important for Logistic Regression and SVM).

## 4.6 Evaluation Framework

### 4.6.1 Cross-Validation Strategy

We employ 5-fold cross-validation for model evaluation under each condition, with a specific approach for each dataset:

- **Dataset A: GroupKFold** (grouped 5-fold) stratified by class. We treat each subject as a group, so that no subject's recordings appear in both training and testing in a given fold. We also ensure roughly balanced class proportions in each fold (stratification by class on the group splits). With 37 subjects, the split might be 30 train vs 7 test in most folds (and 29 vs 8 in one fold due to indivisibility).
- **Dataset B: StratifiedKFold** (standard stratified 5-fold) since no grouping is possible. Each fold has approximately 605 training samples and 151 test samples (75%/25% split each time), maintaining the ~3:1 class ratio in each fold.

These choices reflect our earlier discussion: grouped CV for A to avoid subject leakage, and stratified for B to maintain class ratio (with the caveat regarding subjects).

To illustrate:

- *Dataset A:* if there are 37 subjects, each fold might leave out ~7 subjects for testing. For example, Fold 1 could test on subjects {IDs 1–7} and train on {8–37}; Fold 2 tests on {8–14}, etc. This way, the **key constraint** is satisfied: **All recordings from a given subject appear in ONE fold only** (either all in training or all in testing for that fold). This strategy prevents the model from ever being evaluated on a voice it has "heard" in training.
- *Dataset B:* since we cannot group by subject, a random stratified split is used. Each fold's test set is roughly 151 samples (25% of 756), and because of stratification, each fold's class distribution is about 25% HC, 75% PD, matching the overall distribution.

### 4.6.2 Evaluation Metrics

We evaluate models on several metrics, computed for each test fold and then averaged:

- **Accuracy:** The fraction of correctly classified samples (TP + TN) / (Total). This gives an overall success rate but can be misleading under imbalance.
- **Precision:** For the PD class specifically, precision = TP / (TP + FP). It answers: when the model predicts PD, how often is it correct? (High precision means few false alarms.) We usually focus on PD as the positive class.
- **Recall (Sensitivity):** TP / (TP + FN) for PD – i.e., the proportion of actual PD cases the model catches. In a screening context, recall is critical (missed PD cases are false negatives).
- **F1 Score:** The harmonic mean of precision and recall: 2*(Precision*Recall)/(Precision+Recall). It balances the two, useful when classes are imbalanced.
- **ROC-AUC:** The Area Under the Receiver Operating Characteristic Curve. This metric is threshold-independent and reflects the model's ability to rank-order positive vs negative examples. It is often considered the primary metric in binary classification under class imbalance because it is insensitive to decision threshold and class ratio.

In this thesis, **ROC-AUC is treated as the primary performance metric**. It provides a single measure of discrimination ability without requiring choice of threshold and is widely used in medical ML literature for binary tasks. That said, we still report the others (accuracy, precision, recall, F1) to fully characterize performance.

All metrics are reported as **mean ± standard deviation** across the 5 folds, to convey performance variability.

## 4.7 Experimental Conditions

### 4.7.1 2×2 Factorial Design

Our experiments follow a 2×2 factorial design crossing two factors:

- **Feature Set:** Baseline (47 features) vs. Extended (78 features).
- **Class Weighting:** Unweighted vs. Weighted ( `class_weight="balanced"` ).

This yields four experimental **conditions**:

|  | Baseline Features (47) | Extended Features (78) |
|---|---|---|
| **Unweighted** | Condition 1 | Condition 2 |
| **Weighted** | Condition 3 | Condition 4 |

Each condition will be run for each model and dataset. For example, Condition 2 (extended, unweighted) on Dataset A with Random Forest, etc.

This structure allows us to isolate the effect of adding features (by comparing Condition 1 vs 2 and Condition 3 vs 4) and the effect of class weighting (by comparing 1 vs 3 and 2 vs 4), as well as any interaction between them.

We label outputs and result files accordingly (as seen in the experimental design chapter and code).

### 4.7.2 Configuration Flags

In code (for clarity and reproducibility), we control these factors with simple flags or parameters. For instance:

```
USE_CLASS_WEIGHT_BALANCED = False  # True for weighted conditions
USE_EXTENDED_FEATURES     = False  # True for extended feature set
```

These flags toggle which features are used and whether the classifier is initialized with `class_weight="balanced"` or not.

By iterating over all combinations of these booleans, we cover the four conditions systematically.

## 4.8 Implementation Details

### 4.8.1 Software Stack

All analysis is performed in Python. Key libraries and their versions used are:

- **Audio I/O & processing:** `librosa` (version 0.10.x) for loading audio and computing MFCCs, etc.
- **Prosodic analysis:** `praat-parselmouth` (Praat through Parselmouth, version 0.4.x) for robust pitch, jitter, shimmer, and formant extraction.
- **Machine Learning models:** `scikit-learn` (version 1.4.x) for LogisticRegression, SVC, and RandomForestClassifier implementations, as well as Pipeline and cross-validation utilities.
- **Data handling:** `pandas` and `numpy` (latest versions as of writing) for data manipulation, and general scientific computing.

All experiments were run on a standard computing environment; no GPU or specialized hardware is required for these classical methods.

### 4.8.2 Reproducibility

To ensure results can be replicated:

- We fix the random seed (42) for all randomized procedures (e.g., cross-validation shuffle, model initialization where applicable). This ensures that the CV splits and any randomness in models (like RF's bootstrap) are consistent run-to-run.
- The entire experimental pipeline is implemented via CLI scripts in the repository: `pvc-experiment` for running the experiments with all configurations, and `pvc-extract` for feature extraction. The code is version-controlled in the repository.
- Feature extraction parameters (like frame sizes for MFCC, etc.) are documented in the code and remain consistent throughout.
- Results (predictions, metrics) are saved to CSV files for traceability, and key tables/figures in the thesis are directly generated from these outputs.

By combining code availability and detailed method descriptions, we aim for a high level of transparency. Anyone with access to the dataset should be able to reproduce the core results following the procedures outlined.

## 4.9 Summary

The methodology implements the following:

1. **Two-tier feature extraction** – a 47-dimensional baseline set and an expanded 78-dimensional set, derived from raw audio in Dataset A (with analogous use of provided features in Dataset B).
2. **Three classical ML models** – Logistic Regression, SVM, Random Forest – each tested under consistent settings, with and without class weighting.
3. **Grouped cross-validation** – rigorous subject-level CV for Dataset A, ensuring no leakage, and stratified CV for Dataset B (acknowledging its limitation).
4. **2×2 factorial design** – a structured comparison of feature set extension and class weighting effects on model performance.

This framework is designed to address the objectives while controlling for the main sources of variation. We now proceed to the experimental design specifics, which detail how these methodologies are executed in practice for evaluation.

# Chapter 5: Experimental Design

## 5.1 Overview

This chapter details the experimental design, including the **2×2 factorial structure**, cross-validation protocols for each dataset, and evaluation procedures. The design prioritizes methodological rigor and clarity in comparing conditions over brute-force optimization. We outline how the research questions map onto experiments and how results are organized for analysis.

## 5.2 Research Questions

The experiments are structured to address the following research questions (RQs):

1. **RQ1:** How do classical ML models (Logistic Regression, SVM, Random Forest) perform on PD voice classification?
2. **RQ2:** Does extending the feature set from 47 to 78 features improve classification performance?
3. **RQ3:** Does applying class weighting improve performance on imbalanced datasets?
4. **RQ4:** How do results differ between Dataset A (with subject-grouped CV) and Dataset B (standard CV), and what does this say about evaluation strategies?

These RQs guided the factorial design and choice of analyses (e.g., within-dataset comparisons for RQ2 and RQ3, cross-dataset observations for RQ4).

## 5.3 Experimental Matrix

### 5.3.1 2×2 Factorial Design

As described, we have four experimental conditions combining feature set and weighting. For clarity, we label them as follows in our results output:

| Condition | Features | Weighting | Output Directory |
| --- | --- | --- | --- |
| **C1** | Baseline (47) | None (unweighted) | `baseline/baseline/` |
| **C2** | Extended (78) | None (unweighted) | `baseline/extended/` |
| **C3** | Baseline (47) | Balanced weights | `weighted/baseline/` |
| **C4** | Extended (78) | Balanced weights | `weighted/extended/` |

This directory naming convention (as used in our code and saved results) nests conditions by weighting then feature set.

For each condition, all three models (LR, SVM, RF) will be evaluated on both Dataset A and Dataset B (with appropriate CV).

### 5.3.2 Models Under Evaluation

Each condition involves running the three classifiers introduced in Chapter 4:

| Model | Abbreviation | Key Parameters (recap) |
|---|---|---|
| **Logistic Regression** | LR | C=1.0, max_iter=1000 |
| **Support Vector Machine** (RBF) | SVM | C=1.0, gamma='scale' |
| **Random Forest** | RF | n_estimators=100, max_depth=10 |

We will compare model performance within each condition (to answer RQ1 on model hierarchy) and use Random Forest as a particularly illustrative model for RQ2–RQ4 (since it consistently performed best, as will be seen).

### 5.3.3 Datasets

Our experiments treat Dataset A and Dataset B somewhat separately in analysis:

- **Dataset A (MDVR-KCL):** We evaluate each model for each condition on both speech tasks (ReadText and SpontaneousDialogue) using grouped 5-fold CV. In some analyses, we aggregate results across tasks for simplicity, but we also examine whether performance differs by task.
- **Dataset B (PD_SPEECH):** We evaluate similarly (5-fold CV, stratified). This dataset isn't split by task (it's all one type of sustained phonation features), so no task subdivision exists.

We ensure that for each fold and each dataset, metrics are computed comparably so that we can present summary statistics like mean ± std.

## 5.4 Cross-Validation Protocols

### 5.4.1 Dataset A: Grouped Stratified 5-Fold

For Dataset A, we perform a **5-fold grouped stratified cross-validation** at the subject level. Conceptually:

```
Subject Pool (37 subjects)
├── Fold 1: Train on ~30 subjects, Test on ~7 subjects
├── Fold 2: Train on ~30 subjects, Test on ~7 subjects
├── Fold 3: Train on ~30 subjects, Test on ~7 subjects
├── Fold 4: Train on ~30 subjects, Test on ~7 subjects
└── Fold 5: Train on ~29 subjects, Test on ~8 subjects
```

Because 37 is not perfectly divisible by 5, one fold will have an extra subject in test. We maintain class stratification such that each fold's test set has an appropriate PD/HC ratio.

**Key constraint:** All recordings from a given subject appear in one fold only (either in that fold's training set or test set, but not both). For example, if subject ID 04 is assigned to Fold 1's test set, none of ID 04's recordings will appear in training for Fold 1.

This grouped CV approach prevents the **subject identity leakage** problem. It effectively simulates a scenario of training on some group of patients and testing on completely unseen patients.

### 5.4.2 Dataset B: Stratified 5-Fold

For Dataset B, we use standard stratified 5-fold CV:

```
Sample Pool (756 samples)
├── Fold 1: Train on ~605 samples, Test on ~151 samples
├── Fold 2: Train on ~605 samples, Test on ~151 samples
… (folds 3–4 similarly)
└── Fold 5: Train on ~605 samples, Test on ~151 samples
```

Each fold maintains roughly the 3:1 class ratio (so each test fold has ~113 PD and ~38 HC, for example).

**Caveat:** Without subject identifiers, we cannot ensure that those ~605 training samples and 151 testing samples in a fold don't share subjects. We treat each sample as independent, acknowledging the potential optimistic bias as a threat to validity. This is the same approach that would be taken by someone using this dataset unaware of subject identities – thus our results on Dataset B are comparable to such literature results but will be marked with an asterisk (figuratively) in interpretation.

## 5.5 Evaluation Metrics

### 5.5.1 Primary Metric

**ROC-AUC** (Receiver Operating Characteristic Area Under Curve) is designated as the primary metric for performance comparison. The reasons (reiterating from methodology):

- It evaluates the ranking of predictions across all threshold settings, thus not tied to a specific decision threshold.
- It is more informative than accuracy in imbalanced settings (since a model that gets all HCs right but misses many PDs could still have an acceptable accuracy, but likely a lower AUC).
- In medical diagnosis contexts, ROC-AUC is standard because it relates to the trade-off between sensitivity and specificity.

We aim to focus discussions around ROC-AUC improvements when comparing feature sets or weighting.

### 5.5.2 Secondary Metrics

We also calculate **Accuracy**, **Precision**, **Recall**, and **F1 Score** for each model and condition:

| Metric | Purpose (in context) |
|---|---|
| **Accuracy** | Overall performance (primarily for reference, since it can be skewed by class imbalance). |
| **Precision** | Indicates how reliable a positive PD prediction is (useful to understand false alarm rate). Particularly important if one worries about over-diagnosis. |
| **Recall** | Indicates how many PD cases are caught (sensitivity). Critical for a screening tool (we want high recall so as not to miss patients). |
| **F1 Score** | Balances precision and recall; useful single measure of test effectiveness for the positive class. |

All these are computed with PD as the "positive" class (since typically detecting PD is the event of interest).

### 5.5.3 Statistical Reporting

To quantify variability, we report each metric as **mean ± standard deviation** over the 5 CV folds. For example, "ROC-AUC = 0.822 ± 0.166" indicates the average AUC was 0.822 and the fold-to-fold standard deviation was 0.166.

Given the small sample nature of Dataset A, these std values can be large, reminding us to not overinterpret small differences. We will often discuss whether differences are larger than these std deviations to gauge if improvements are meaningful or within noise.

Additionally, while not formalized in hypothesis testing due to limited folds, overlapping standard deviations will be noted as indicating that differences might not be statistically significant.

## 5.6 Experimental Procedure

### 5.6.1 Step-by-Step Protocol

The overall procedure for each dataset and condition is as follows (automated by our scripts):

1. **Feature Extraction (Dataset A only):** Generate features for all WAV files. (Using the CLI command `pvc-extract --task all` which processes both ReadText and SpontaneousDialogue recordings and outputs feature CSVs for baseline and extended sets.)
2. **For each condition (C1–C4):**
   a. For each model (LR, SVM, RF):
   b. For each dataset/task combination:
      ◦ Perform 5-fold cross-validation as specified (grouped for A, stratified for B).
      ◦ In each fold: fit the pipeline on training data, predict on test data, compute metrics.
      c. Aggregate the 5-fold results (compute mean and std for each metric).
3. **Aggregate results across tasks (if needed):** For Dataset A, we also combine ReadText and SpontaneousDialogue metrics for an overall picture, but we keep them separate for certain analyses.
4. **Save results:** Metrics for each model/condition/dataset are saved in CSV files, and detailed per-fold metrics can be saved as well for later analysis of variance.

This procedure ensures consistency – every model is evaluated in an identical fashion to allow fair comparisons.

Pseudocode of the loop structure:

```
for condition in [(features=baseline, weight=none), ..., (extended,
weight=balanced)]:
    for model in [LR, SVM, RF]:
        for dataset in [A_readtext, A_spontaneous, B]:
            perform 5-fold CV
            compute metrics (mean, std)
            save results
```

The use of automation (via the `pvc-experiment` script) reduces manual error and ensures all combinations are covered.

### 5.6.2 Feature Extraction Settings

As a note, the extracted features are stored as intermediate files:

- `outputs/features/baseline/` – contains CSVs like `features_readtext.csv` and `features_spontaneousdialogue.csv` (47 features each) for Dataset A.
- `outputs/features/extended/` – contains analogous CSVs with 78 features each for Dataset A.

Dataset B's features are directly read from `assets/PD_SPEECH_FEATURES.csv`.

These files allow re-running experiments quickly without re-extracting features each time.

### 5.6.3 Random Seed

All experiments use `random_state = 42` wherever randomness is involved (CV splitting, model initialization). This is critical for reproducibility – if someone else runs the same code, they should get the same folds and same results. It also means that when comparing conditions, they are effectively being compared on the same splits (for dataset B, which is random, using a fixed seed ensures any difference comes from the method, not a lucky split).

## 5.7 Implementation

### 5.7.1 CLI Commands

For completeness, examples of the command-line interface usage are:

```
# Feature extraction (for both baseline and extended features on Dataset A)
pvc-extract --task all
```

```
# Run all experiments (all conditions, all models, both datasets)
pvc-experiment
```

The `pvc-experiment` command reads a configuration (or uses internal loops) to execute the CV for each model and condition, and outputs results to the `outputs/results/` directory.

### 5.7.2 Output Structure

The results are systematically stored in a directory hierarchy for traceability:

```
outputs/
├── features/
│   ├── baseline/
│   │   ├── features_readtext.csv
│   │   └── features_spontaneousdialogue.csv
│   └── extended/
│       ├── features_readtext.csv
│       └── features_spontaneousdialogue.csv
│
└── results/
    ├── baseline/         # Unweighted results
    │   ├── baseline/     # 47 features results
    │   └── extended/     # 78 features results
    └── weighted/         # Class-weighted results
        ├── baseline/     # 47 features results
        └── extended/     # 78 features results
```

Within each of those, files like `summary.csv` and `all_results.csv` contain the metrics. For example, `outputs/results/baseline/extended/summary.csv` might list the average metrics for each model on each dataset under Condition 2 (extended, unweighted).

This structure made it convenient to gather results for creating tables and figures in the thesis.

## 5.8 Expected Outcomes

Before presenting results, we outline hypotheses corresponding to the RQs:

### 5.8.1 Hypotheses

| Hypothesis | Rationale |
| --- | --- |
| **H1: Extended features improve ROC-AUC.** | The 31 additional features (MFCC std, delta-delta, spectral shape) capture variability and dynamics that help differentiate PD vs HC, especially in borderline cases. We expect a notable AUC increase with 78 features, particularly on Dataset A where 47 features might miss some information. |

| Hypothesis | Rationale |
|---|---|
| **H2: Class weighting improves recall (sensitivity).** | By compensating for class imbalance, models (especially on Dataset B where imbalance is high) should make more effort to correctly classify minority class (HC in B, or PD in an inverted scenario). We expect recall of the minority class to rise with weighting, though overall AUC might not change dramatically if models were already partly handling imbalance. |
| **H3: Dataset B shows higher performance than Dataset A.** | The larger sample size and possibly easier task (pre-extracted, possibly including some leakage) in Dataset B will likely yield higher ROC-AUC and accuracy than Dataset A. However, this "higher performance" may reflect optimistic bias rather than true superiority of that dataset for generalization. |
| **H4: Random Forest outperforms LR and SVM.** | Random Forest can capture non-linear interactions between features (PD markers may manifest in combinations of features) and is relatively robust. We anticipate RF will have the highest AUC in most conditions. Logistic Regression might struggle if relationships aren't strictly linear, and SVM might be strong in some cases but can suffer if hyperparameters aren't tuned. |

These hypotheses will be revisited in light of the actual results.

### 5.8.2 Analysis Plan

To evaluate the above:

1. **Within-condition comparison:** For each condition, rank the models by performance to answer RQ1 (which model is best under same features/weighting).
2. **Feature ablation effect:** Compare Condition 1 vs 2 (baseline vs extended without weighting) and Condition 3 vs 4 (baseline vs extended with weighting) to quantify the effect of adding features (addresses RQ2).
3. **Weighting effect:** Compare Condition 1 vs 3 (baseline unweighted vs weighted) and Condition 2 vs 4 (extended unweighted vs weighted) to see the impact of class weighting (addresses RQ3).
4. **Cross-dataset comparison:** Juxtapose findings from Dataset A vs Dataset B. For instance, note if extended features helped both datasets, or if class weighting was more crucial on B (with heavy imbalance) than A. Also compare absolute performance levels, keeping in mind B's potential optimism (addresses RQ4, with appropriate caveats).

We also plan visualizations: e.g., bar charts of AUC for each condition, tables of mean ± std, and perhaps heatmaps of feature importance (which we include in Appendix A). Statistical significance is not rigorously computed due to limited folds, but observed differences well beyond std dev will be noted as likely meaningful.

## 5.9 Limitations of Design

Before proceeding to results, we acknowledge limitations in the experimental design itself (to be further discussed in Chapter 8 as well):

### 5.9.1 Acknowledged Constraints

1. **Small sample size (Dataset A):** With only ~37 subjects, any performance estimate has high variance. A different split of data could yield different results; hence we rely on cross-validation, but even then each fold's test set is small (~7 subjects), so metrics fluctuate.
2. **No subject IDs for Dataset B:** We cannot confirm independence of training and test for Dataset B, meaning those results are taken at face value but may not reflect true generalization to new patients.
3. **No hyperparameter tuning:** We fixed model hyperparameters (C, tree depth, etc.) rather than perform grid search. This was deliberate to avoid overfitting the small data, but it means our models might not be the absolute best they could be. Some performance loss or model ranking differences could result from suboptimal settings.
4. **No external validation set:** All results are via internal cross-validation. We did not set aside a separate external test set (since none was available in Dataset A, and Dataset B's nature is already as a stand-alone set). Thus, all performance is on the same data used to develop models (albeit in CV fashion). True external validity remains to be confirmed.

### 5.9.2 Mitigation Strategies

- We emphasize **standard deviation** of metrics to convey uncertainty (for point 1).
- We include explicit **caveat statements** whenever interpreting Dataset B results (for point 2).
- By not tuning hyperparameters, we actually present a somewhat **lower-bound** performance which is more realistic (point 3). We also maintain consistency – using the same hyperparameters across conditions ensures fairness in comparisons.
- We plan to treat the combination of Dataset A and B results as mutual confirmation: if both suggest the same trend, that increases confidence (point 4). Additionally, Chapter 8 (Limitations) will clearly state that external validation is needed in future work.

## 5.10 Summary

The experimental design implements:

- A **2×2 factorial structure** (Feature Set × Class Weighting) to systematically study each factor's effect.
- **Grouped CV for Dataset A** to ensure valid evaluation without leakage.
- **Five metrics** (Accuracy, Precision, Recall, F1, ROC-AUC) reported as mean ± std to give a comprehensive performance view.
- A **reproducible pipeline** (via code and fixed random seeds) so that all findings can be independently verified.

With this design, we proceed to present the results, focusing on how the models performed and what insights were gained regarding the research questions.

# Chapter 6: Results

## 6.1 Overview

This chapter presents the classification results across all experimental conditions. Results are organized to highlight:

1. **Condition-level summaries** (performance under each of the 2×2 factorial conditions).
2. **Model comparisons** within each condition (to see which algorithm performed best).
3. **Feature ablation analysis** (comparing baseline vs extended feature sets).
4. **Class weighting analysis** (comparing unweighted vs weighted scenarios).
5. **Cross-dataset observations** (noting differences between Dataset A and B outcomes).

All performance metrics are reported as mean ± standard deviation across cross-validation folds, as previously noted.

## 6.2 Summary of Best Results

### 6.2.1 Dataset A (MDVR-KCL) — Best Performance

For Dataset A, the highest-achieving configuration and result were:

| Metric | Value | Model | Task | Configuration (Features / Weighting) |
|---|---|---|---|---|
| **ROC-AUC** | **0.857 ± 0.171** | Random Forest | SpontaneousDialogue | Extended features / Unweighted |
| **Accuracy** | **82.2% ± 16.6%** | Random Forest | ReadText | Extended features / Unweighted |

The top ROC-AUC of ~0.857 was achieved by the Random Forest on the Spontaneous Dialogue task using the extended 78-feature set without class weighting. The highest accuracy (82.2%) was observed for Random Forest on the ReadText task under the same feature condition (interestingly, the RF did slightly better in accuracy on ReadText, but in AUC it excelled on SpontaneousDialogue).

### 6.2.2 Key Finding

> **Feature extension (78 features) consistently improved performance** compared to baseline features (47). The best results for each task came from the extended feature set. In particular, the highest ROC-AUC of 0.857 was achieved using the extended feature set on the SpontaneousDialogue task, whereas with baseline features the best ROC-AUC was lower (see detailed tables below).

### 6.2.3 Dataset B (Benchmark)

> **Note:** Dataset B (pre-extracted features) achieved a significantly higher ROC-AUC of **0.940 ± 0.013** (with Random Forest, extended features, unweighted). This difference in absolute

performance is attributed to the larger sample size (n=756 vs n=37 subjects) and the lack of subject-level grouping in the provided dataset, which likely leads to optimistic estimates. We emphasize that while impressive, the Dataset B results may not be directly comparable to Dataset A's due to these factors.

*(The above points will be further discussed in the Discussion chapter, including caution about Dataset B's optimistic bias.)*

## 6.3 Condition 1: Baseline Features + Unweighted

**Configuration:** 47 features (baseline set), no class weighting (standard training).

**Dataset A Results (by task):**

### 6.3.1 Task: ReadText

| Model | ROC-AUC | Accuracy | F1 Score |
|---|---|---|---|
| Logistic Regression | 0.717 ± 0.139 | 0.621 ± 0.058 | 0.542 ± 0.099 |
| SVM (RBF) | 0.614 ± 0.312 | 0.621 ± 0.106 | 0.333 ± 0.333 |
| Random Forest | 0.590 ± 0.302 | 0.629 ± 0.178 | 0.351 ± 0.363 |

*(Precision and Recall are not shown here for brevity; F1 is a summary. In this case, the low F1 for SVM indicates an imbalance in precision/recall – indeed SVM had some folds with poor PD detection.)*

### 6.3.2 Task: Spontaneous Dialogue

| Model | ROC-AUC | Accuracy | F1 Score |
|---|---|---|---|
| Logistic Regression | 0.760 ± 0.214 | 0.639 ± 0.160 | 0.539 ± 0.321 |
| SVM (RBF) | 0.407 ± 0.309 | 0.636 ± 0.135 | 0.400 ± 0.253 |
| **Random Forest** | **0.828 ± 0.148** | **0.721 ± 0.176** | **0.567 ± 0.365** |

Here, Random Forest markedly outperforms the other two on ROC-AUC (0.828 vs ~0.76 and 0.41). Interestingly, SVM did much worse on AUC for SpontaneousDialogue than for ReadText in this condition, indicating it struggled with that task under baseline features (perhaps overfitting some noisy pattern, as evidenced by high variance).

### 6.3.3 Observations (Condition 1)

- **Task Difference:** For baseline features without weighting, the SpontaneousDialogue task yielded significantly better separation than ReadText for Random Forest (ROC-AUC 0.828 vs 0.590). This suggests the spontaneous speech task provided richer discriminative information that RF could leverage with 47 features, whereas in read speech the baseline features alone were not as powerful (RF was near chance in AUC for ReadText with baseline features, 0.59 ± 0.30).

- **Model Stability:** Logistic Regression was relatively consistent across tasks (ROC-AUC ~0.72 on ReadText vs ~0.76 on Dialogue), indicating it provided a stable linear baseline. SVM's performance, however, was highly variable and task-dependent (even hitting ROC-AUC ~0.40 in SpontaneousDialogue, effectively failing in some folds).
- **Variance:** Standard deviations are quite high (often ±0.15–0.30), reflecting the small test fold sizes. For example, RF's 0.828 ± 0.148 on SpontaneousDialogue still has considerable uncertainty. Any single train-test split could have varied outcomes; the cross-fold mean smooths it out.

## 6.4 Condition 2: Extended Features + Unweighted

**Configuration:** 78 features (extended set), no class weighting.

### 6.4.1 Task: ReadText

| Model | ROC-AUC | Accuracy | F1 Score |
|---|---|---|---|
| Logistic Regression | 0.698 ± 0.132 | 0.596 ± 0.079 | 0.475 ± 0.106 |
| **SVM (RBF)** | **0.834 ± 0.153** | 0.786 ± 0.181 | 0.634 ± 0.386 |
| **Random Forest** | **0.822 ± 0.166** | 0.818 ± 0.140 | 0.746 ± 0.207 |

On the ReadText task, with extended features, both SVM and RF show a dramatic improvement in ROC-AUC compared to Condition 1. SVM achieved 0.834 (from 0.614 prior) and RF 0.822 (from 0.590 prior). Logistic Regression remained roughly similar (~0.70, slight drop from 0.717, which could be noise).

Accuracy values also improved (RF now ~81.8%, whereas it was ~62.9% before). F1 scores particularly for RF jumped (from 0.35 to 0.746), indicating it is now identifying PD cases far better (Precision/Recall both improved).

### 6.4.2 Task: Spontaneous Dialogue

| Model | ROC-AUC | Accuracy | F1 Score |
|---|---|---|---|
| Logistic Regression | 0.783 ± 0.139 | 0.671 ± 0.199 | 0.530 ± 0.377 |
| SVM (RBF) | 0.460 ± 0.294 | 0.636 ± 0.089 | 0.428 ± 0.258 |
| **Random Forest** | **0.857 ± 0.171** | 0.779 ± 0.161 | 0.605 ± 0.387 |

For SpontaneousDialogue, Random Forest slightly improved from 0.828 to 0.857 AUC with extended features (a smaller gain than in ReadText, since baseline was already good). SVM, surprisingly, still struggled on this task (0.460 AUC, similar to its poor showing before, perhaps indicating extended features didn't help SVM's issues on this task – possibly outliers or fold-specific problems remain). Logistic Regression saw a small increase (0.760 to 0.783 AUC).

### 6.4.3 Observations (Condition 2)

- **Extended Features Impact:** The addition of features had a **massive impact on the ReadText task**. Random Forest's ROC-AUC jumped from ~0.59 to ~0.82 (+23 percentage points), and SVM from ~0.61

to ~0.83 (+22 pp). This essentially "rescued" the models from near-chance performance to strong performance on the reading task. This suggests that the baseline 47 features were insufficient for ReadText, but the extended features captured additional vocal characteristics (like variability) that made PD vs HC differentiation much clearer.

- **Spontaneous Dialogue Stability:** For the Dialogue task, which was already yielding good results with baseline features, the improvement with extended features was present but more modest for RF (0.828 → 0.857, +2.9 pp) and nonexistent or negative for SVM (which remained low). It appears that for a naturally more separable task like Dialogue, the baseline features already did quite well (especially for RF), so extra features yielded diminishing returns.

- **SVM Anomaly:** It's notable that SVM performed excellently on ReadText with extended features (best model there at 0.834 AUC), but performed very poorly on Dialogue even with extended features (0.460 AUC). This disparity suggests that the RBF SVM may be overfitting or being thrown off by something in the Dialogue data (perhaps the extended feature space combined with small sample per fold leads to high variance for SVM). In contrast, RF handled both tasks reliably with extended features, reinforcing the notion of RF's robustness.

- **Model Ranking:** In this condition, for ReadText SVM = RF > LR in AUC, whereas for Dialogue RF >> LR >> SVM. This indicates no single model dominated absolutely in all cases, though RF is the most consistent top performer overall.

## 6.5 Condition 3: Baseline Features + Weighted

**Configuration:** 47 features, `class_weight="balanced"` applied.

(This condition specifically tests if adding class weights to baseline features improves things, especially relevant for Dataset B which is heavily imbalanced. For Dataset A, which is roughly balanced, we expect minimal change.)

### 6.5.1 Task: ReadText (Dataset A)

| Model | ROC-AUC | Accuracy | F1 Score |
|---|---|---|---|
| Logistic Regression | 0.717 ± 0.139 | 0.596 ± 0.079 | 0.528 ± 0.099 |
| SVM (RBF) | 0.542 ± 0.312 | 0.704 ± 0.111 | 0.519 ± 0.320 |
| Random Forest | 0.687 ± 0.258 | 0.650 ± 0.148 | 0.431 ± 0.306 |

Comparing to Condition 1 (unweighted baseline), for ReadText: - LR's AUC is identical (0.717 vs 0.717) – weighting made no difference for LR.
- SVM's AUC went from 0.614 (unweighted) to 0.542 (weighted); a drop, though given the large std, this may not be meaningful. Accuracy for SVM interestingly increased (70.4% vs 62.1%), suggesting it likely predicted more "PD" (the minority in A, ironically PD is minority in A's tasks) improving accuracy slightly but hurting AUC.
- RF's AUC improved from 0.590 to 0.687 (+9.7 pp). Accuracy also improved a bit from ~62.9% to 65.0%. This suggests that weighting helped RF pick up a few more PD cases in ReadText without extended features, though it's still far lower than the extended-feature scenario.

| Model | ROC-AUC | Accuracy | F1 Score |
|---|---|---|---|
| Logistic Regression | 0.760 ± 0.214 | 0.639 ± 0.160 | 0.539 ± 0.321 |
| SVM (RBF) | 0.423 ± 0.312 | 0.693 ± 0.123 | 0.560 ± 0.318 |
| **Random Forest** | **0.827 ± 0.133** | 0.693 ± 0.123 | 0.538 ± 0.326 |

Comparing to Condition 1 (unweighted baseline): - LR remains the same (0.760 vs 0.760 AUC). No effect of weighting as expected (data ~balanced).
- SVM improves a tad (0.407 → 0.423 AUC, negligible). Its accuracy interestingly went up to ~69% from ~63.6%, again indicating more balanced error handling but still its AUC is low.
- RF is virtually unchanged (0.828 → 0.827 AUC, within margin).

So on Dialogue, weighting didn't change much for RF or LR, as anticipated since classes were already near balanced.

**6.5.3 Observations (Condition 3)**

- **Weighting Effect (Dataset A):** On the moderately imbalanced Dataset A, class weighting had minimal impact overall. The only notable improvement was RF on ReadText (which had a slight PD minority, 16 PD vs 21 HC; weighting gave RF a boost in PD detection raising AUC ~0.59 to ~0.69). On SpontaneousDialogue (15 PD vs 21 HC), RF saw essentially no change. LR was unchanged in both tasks (likely because it was already fairly balancing precision/recall as a linear model). SVM's fluctuations are harder to interpret but did not consistently improve AUC.
- **No significant gain for Dialogue:** Weighting did not help the Dialogue task for any model in terms of AUC. This is expected as the imbalance wasn't severe and models likely already managed trade-offs well.
- **Accuracy vs AUC trade-off:** The SVM's accuracy increased with weighting in both tasks (especially noticeable in ReadText 62%→70% and Dialogue 64%→69%), but its AUC did not – implying that weighting caused SVM to favor the majority class slightly less (thus improving accuracy since in A majority = HC, which helped as it was maybe predicting PD too often incorrectly before). But these details are minor given SVM's general underperformance on A.

In summary, class weighting on Dataset A (with baseline features) had a **modest positive effect in one scenario (RF on ReadText)** and otherwise little to no effect. This foreshadows that weighting might show more utility on Dataset B, where imbalance is much larger.

## 6.6 Condition 4: Extended Features + Weighted

**Configuration:** 78 features, `class_weight="balanced"`.

This condition is the combination of extended feature advantages with class imbalance mitigation. It's essentially the most complex scenario.

### 6.6.1 Task: ReadText

| Model | ROC-AUC | Accuracy | F1 Score |
|---|---|---|---|
| Logistic Regression | 0.698 ± 0.132 | 0.650 ± 0.108 | 0.564 ± 0.160 |
| **SVM (RBF)** | **0.834 ± 0.153** | 0.761 ± 0.214 | 0.620 ± 0.390 |
| Random Forest | 0.805 ± 0.182 | 0.818 ± 0.140 | 0.746 ± 0.207 |

Comparing to Condition 2 (unweighted extended): - LR unchanged (0.698 vs 0.698 AUC, accuracy slightly up from 59.6% to 65.0%). Perhaps weighting let LR catch a couple more PD, raising accuracy, but AUC same.
- SVM unchanged (0.834 vs 0.834 AUC exactly the same mean; any difference is below rounding). So weighting did not change SVM's performance on extended features (it was already very good).
- RF slight *decrease* (0.822 → 0.805 AUC, a small drop within std). Accuracy actually remained ~81.8% vs 81.8%. So essentially no meaningful change for RF either, maybe a small noise.

### 6.6.2 Task: Spontaneous Dialogue

| Model | ROC-AUC | Accuracy | F1 Score |
|---|---|---|---|
| Logistic Regression | 0.783 ± 0.139 | 0.696 ± 0.179 | 0.563 ± 0.381 |
| SVM (RBF) | 0.403 ± 0.347 | 0.664 ± 0.167 | 0.560 ± 0.318 |
| **Random Forest** | **0.823 ± 0.209** | 0.721 ± 0.203 | 0.583 ± 0.373 |

Comparing to Condition 2 (unweighted extended): - LR the same (0.783 vs 0.783 AUC, accuracy 67.1% → 69.6% trivial difference).
- SVM roughly the same poor performance (0.460 → 0.403 AUC, not meaningful given std ~0.35).
- RF essentially unchanged (0.857 → 0.823 AUC, that slight dip could be just variance; indeed the std increased a bit to ±0.209, possibly one fold was worse).

### 6.6.3 Observations (Condition 4)

- **Diminishing Returns:** Adding class weighting to the extended feature set did not yield further improvements over the unweighted extended condition. For Dataset A, we see virtually no difference. This makes sense because once extended features solved most of the classification challenge (especially for RF and SVM on ReadText, and RF on Dialogue), the class imbalance in A was not severe enough to require weighting. Models were already achieving high performance and likely properly handling the slight imbalance.
- **Best Configuration:** Overall, considering Conditions 1–4 for Dataset A, the **Unweighted Extended (Condition 2)** generally produced the highest ROC-AUC scores. For example, RF's 0.857 on Dialogue in Condition 2 was the peak, and weighting in Condition 4 didn't beat that. Similarly for ReadText, SVM/RF peaked at ~0.834/0.822 in Condition 2, and Condition 4 gave ~0.834/0.805. So the extended features alone gave the major boost, and weighting did not add value in presence of extended features (for Dataset A).

Now, shifting briefly to **Dataset B results** under these conditions to complement the analysis:

*(Since Dataset B has only one "task", we'll summarize key points rather than duplicating all tables. Dataset B's performance trends will be incorporated in the discussion for cross-dataset comparison.)*

For Dataset B: - Under Conditions 1 & 3 (baseline features), Random Forest achieved ROC-AUC around 0.78 unweighted vs 0.82 weighted (+3.5 pp with weighting), confirming weighting helps where imbalance is high. - Under Conditions 2 & 4 (extended features), RF achieved ~0.94 AUC unweighted vs ~0.92–0.93 weighted (a slight drop, possibly due to overemphasizing minority class which in B is HC and very small). So on Dataset B, the best result was actually *without weighting* because the minority class (HC) was so small that weighting introduced a bit more variance.

These results are presented below in the feature ablation and summary findings.

## 6.7 Feature Ablation Analysis

This section quantifies the improvement from adding features (47 → 78) by looking at paired comparisons of ROC-AUC.

### 6.7.1 ROC-AUC Improvement from Feature Extension (ReadText, Dataset A)

| Model | Baseline (47) AUC | Extended (78) AUC | Δ ROC-AUC (Extended – Baseline) |
|---|---|---|---|
| Logistic Regression | 0.717 | 0.698 | -0.019 |
| SVM (RBF) | 0.614 | 0.834 | **+0.220** |
| Random Forest | 0.590 | 0.822 | **+0.232** |

*(All values are from the unweighted cases for a clean comparison, Condition 1 vs 2 for each model on ReadText.)*

**Key Finding:** On the ReadText task, feature extension was critical – it boosted SVM and RF by ~22–23 percentage points in ROC-AUC. Logistic Regression saw no benefit (slight 1.9 pp drop, which is negligible given its std). The complex models clearly leveraged the additional 31 features to vastly improve discrimination.

### 6.7.2 ROC-AUC Improvement from Feature Extension (SpontaneousDialogue, Dataset A)

| Model | Baseline (47) AUC | Extended (78) AUC | Δ ROC-AUC |
|---|---|---|---|
| Logistic Regression | 0.760 | 0.783 | +0.023 |
| SVM (RBF) | 0.407 | 0.460 | +0.053 |
| Random Forest | 0.828 | 0.857 | +0.029 |

Again from Conditions 1 vs 2 (unweighted) for Dialogue task.

Here improvements are modest (3–5 percentage points) for RF and SVM, and LR got ~2.3 pp. Given the std dev ~0.15, these are not statistically significant leaps, but directionally extended features helped a bit even

for Dialogue. The dramatic gains seen in ReadText were not needed as Dialogue was already easier to classify with baseline features.

Overall, **feature extension provided the largest benefits in scenarios where baseline performance was weak (e.g., SVM/RF on ReadText)**. In scenarios already strong (RF on Dialogue), it gave only a small extra push.

This justifies that our extended features indeed contained important additional information primarily for the more challenging context (reading task, which might have been more uniform sounding so needed more nuanced features to differentiate PD vs HC).

## 6.8 Summary of Findings

Bringing together the major outcomes:

| Hypothesis (from design) | Result | Evidence |
| --- | --- | --- |
| **H1: Extended features improve ROC-AUC.** | **Confirmed** | Extended features yielded up to +23 pp ROC-AUC on Dataset A's ReadText task (RF 0.590 → 0.822)【6.4.1】【6.7.1】, and consistent if smaller gains elsewhere. Dataset B also saw higher AUC with extended features (RF ~0.94 vs ~0.78) due to the richer feature set. |
| **H2: Spontaneous Dialogue yields better detection than ReadText.** | **Confirmed** | Under comparable settings (e.g. baseline features, RF), SpontaneousDialogue had higher AUC (0.828) than ReadText (0.590)【6.3.2】【6.3.1】. Even with extended features, Dialogue remained slightly superior for RF (0.857 vs 0.822). This suggests spontaneous speech provided more PD cues. |
| **H3: Dataset B values are inflated (relative performance).** | **Confirmed** | Dataset B's best AUC (0.940)【6.2.3】far exceeds Dataset A's (0.857)【6.2.1】, consistent with optimistic bias from subject overlap and large feature set. We interpret B's higher results with caution as discussed. |
| **H4: RF outperforms LR and SVM.** | **Confirmed** | Random Forest was the top performer in most conditions (especially on Dialogue, and overall extended features). For instance, RF achieved the highest AUC in 3 out of 4 conditions on Dataset A (often significantly so)【6.3.2】【6.4.2】. SVM matched RF only in one scenario (ReadText extended) but was unstable in others. Logistic Regression, while stable, lagged in raw performance. |
| **H5: Class weighting improves performance.** | **Rejected** | Class weighting had marginal or negative impact in our results. On Dataset A, differences were minimal【6.5.3】【6.6.3】. On Dataset B, weighting modestly helped baseline features (RF +3.5pp AUC) but not extended. Overall, extended features without weighting already handled imbalance sufficiently. |

*(H5 was a postulated effect but was not strongly supported, hence marked rejected; weighting did not significantly boost metrics when strong features were present.)*

In conclusion, the best performing configuration for our primary dataset (A) was using the extended 78-feature set without class weighting, evaluated with a Random Forest model, achieving mean ROC-AUC ~0.85. The experiments highlight the importance of feature engineering and correct validation: by extending features we substantially improved model accuracy, and by grouping subjects in CV we obtained more realistic performance estimates. Class weighting turned out less crucial for moderate imbalance, though we still consider it for severely imbalanced situations.

The next chapter will discuss these results in detail, interpret their implications in the context of existing literature, and outline the limitations and potential next steps for this research.

---