

ParkinSense: A novel approach to non-intrusive & incipient Idiopathic Parkinsonism Diagnosis, Severity Profiling, and Telemonitoring via Ensemble Learning & Multimodal Data Fusion on Webcam-Derived Digital Biomarkers

For double blind review purpose

Please do not write your names and affiliations.
During the initial submission.

Abstract—Despite being the fastest growing neurodegenerative disease in the world, with over 10 million patients worldwide, there is no definitive diagnosis method for Parkinson's Disease currently. Current diagnosis technologies misdiagnose one in three patients, and rely on inaccurate technologies or are subjective to the administering clinicians. Furthermore they are inaccessible to billions due to immobility, geographic barriers, or associated costs. With early and accurate diagnosis being vital to effective treatment, a tremendous issue emerges. Parkinsense addresses these issues, serving as a web application that contactlessly analyzes three cardinal symptoms of Parkinson's Disease over a standard webcam and microphone: Hypomimia, Dysarthria, and Bradykinesia. With Ensemble Learning coupled with various Classifiers, such as Random Forest, and Support Vector Machines, Parkinsense was able to achieve accuracy rates of 99.72% when tested on synthetic patients. By examining, and combining multiple modalities, Parkinson's boasts accuracy rates higher than many advanced unimodal technologies, as patients often do not display all the symptoms of Parkinson's Disease substantially, leading to misdiagnosis when relying on only one symptom for diagnosis. Parkinsense strongly suggests that entirely contactless diagnosis of PD through digital biomarkers, can be effective, and web applications can be utilized for rapid, automated, and remote diagnosis, which can be crucial for those affected by barriers that prevent access to diagnosis.

Keywords—Parkinson's Disease; Diagnosis; Multimodal Models; Severity Analysis;

I. INTRODUCTION

Parkinson's Disease (PD) is the second most prevalent neurodegenerative disease in the world [X]. With Parkinson's Disease being the fastest growing neurodegenerative disease globally [X], the number of cases has surged to over 10 million in the last few decades, and this figure is expected to double by 2030 [X].

Parkinson's Disease results from the death of dopaminergic neurons in the *substantia nigra*. This neurodegeneration affects largely those over the age of 40, and leads progressive disability. Due to this progressive nature of PD, early and accurate diagnosis is vital for effective treatment, to stagnate or prevent further neuron loss. This neuron loss affects both motor and non-motor skills[X]. Some of the hallmarks of the disorder include *bradykinesia* (slowed movement), speech irregularities, tremors, visual impairment, postural instability, rigid muscles, and more[X]. Traditional and industry-standard diagnostic approaches typically take the form of extensive questionnaires that involve the manual examination and measurement of certain behaviors. However these are highly subjective to the administering clinician, relying on the detection of characteristics and behavior that may be too minute to be noticed by one's eyes [X]. While a few technologies, such as MRI scans and DaTScan can aid in diagnosis, they can offer little information beyond differentiating Idiopathic Parkinsonism, or Parkinson's Disease, from atypical forms of Parkinsonism or other diseases with similar symptoms. As such, they can only be used to exclude other causes, as opposed to directly diagnosing an individual for PD. Further exacerbating the issue, not only do many neurological diseases, such as Essential Tremor, share symptoms with PD, but many parkinsonism symptoms arise naturally in PD-susceptible, but healthy, individuals, as a result of old age, complicating PD diagnosis. Resultantly, nearly 1 in 3 patients are misdiagnosed for PD at least once, severely inhibiting the effectiveness of potential treatments, and enabling the disease time to progress. Furthermore such a diagnosis typically is conducted in a hospital or institution with trained clinicians [X], a single diagnosis can cost upwards of \$5,000, and results can take weeks to arrive [X]. Coupled together these





The linked image cannot be displayed. The file may have been moved, renamed, or deleted. Verify that the link points to the correct file and location.

Figure 1: Summary of sources, demographic information, content information, and citations. N/A represents categories where no information was available, and information could not be compute

difficulties essentially render PD diagnosis inaccessible to billions, specifically those mainly in developing regions, or those without access to affordable healthcare [X].

The goal of this study was to create an accurate, rapid, and accessible method of PD diagnosis, and severity analysis using machine learning for multiple symptoms. This method is to be accessible via a web platform. The study contributions are multifold:

- I. Multimodal Diagnosis PD can lead to accuracy gains when compared to unimodal diagnosis tools.
- II. Aggregating multiple state-of-the-art classifiers simply via late fusion with majority voting, without necessitating fusion dependent on specific classifier attributes or values, is successful in improving model performance.
- III. PD diagnosis to be accurately conducted entirely non-intrusively and contactlessly.
- IV. PD diagnosis software can be deployed in an engaging and accessible manner worldwide
- V. EBR is promising modality for PD diagnosis.
- VI. PD diagnosis can be done rapidly and in an automated manner with server-side feature extraction.
- VII. PD severity analysis can be conducted contactlessly and accurately, using the same feature set as the diagnosis models.

II. RELATED WORKS

With technology significantly emerging into the healthcare industry over the past few decades, a great deal of literature has been published regarding PD diagnosis with the assistance of computational resources. Multiple approaches have been tried, including wearable technology, MRIs, machine-learning-assisted analysis of questionnaires, and plenty more. However a proven and accurate method of diagnosis has been elusive. A promising portion of this healthcare-technology convergence is Machine Learning. Machine Learning has been seen to be effective in signaling for PD in an individual, by analyzing data such as vocal irregularities, MRI scans, sebum RNA, and much more. For the purpose of this study, a focus was placed on machine learning methods for PD diagnosis that can be conducted contactlessly (e.g. without wearable technology, or MRI scans).

While numerous symptoms are present in PD patients [X], three were chosen for this study, stemming from their ability to easily and accurately be assessed noninvasively, with minimal technology requirements, ensuring accessibility while retaining reliability.

2.1 Dysarthria

Dysarthria categorizes the vocal irregularities deriving from the hypokinetic nature of parkinsonism [X]. It can cause deficits and abnormalities in intonation, tone, volume, fluidity, and more. Often, it leads to monotonicity in pitch and volume, reduced stress, and a breathy and hoarse voice. Over 90% of PD subjects develop this symptom. *Dysarthria* typically is noticeable in the moderate stages of PD, a few years after onset.

The vast majority of approaches to PD diagnosis via machine learning, attempt to identify vocal irregularities, from dozens of features, to differentiate between PD patients and HC [X, X, X]. Sustained vowel phonation, of /a/, for example, have proven to be effective in differentiating PD patients from HC [X, X]. However, many past approaches are only tested on a small number of instances, which may make them unrepresentative of their performance on a large number of people.

Additionally, as not all PD patients exhibit *Dysarthria* [X], and since it is more recognizable in the later stages of PD [X], the analysis of this symptom alone is not adequate for a PD diagnosis platform that aims to be widespread.

2.2 Hypomimia

Also referred to as ‘Facial Masking’, *Hypomimia* is the reduction in pronunciation, and rapidness of facial expressions [X]. This symptom arises from the rigidity many PD patients encounter in their muscles. This rigidity limits the ability for PD patients to express proper emotions through facial expressions, and can give the impression to others that you are wearing a mask. *Hypomimia* is one of the earliest symptoms of PD, and over 70% of PD patients develop *Hypomimia*.

The human face can be partitioned into groups of Action Units (AU), that represent groups of facial muscles [X]. When making certain facial expressions, such as disgust or surprise, numerous AU’s are activated, and the variance of the magnitudes of this activation, reveal quantify the amount of muscle movement, which in turn, can divulge information about the presence of PD. *Hypomimia* has begun to be studied in much more detail, especially in regard to Machine Learning, in recent years [X, X]. A few studies have discovered that *Hypomimia* is a promising method for PD diagnosis [X], since the variance of AU Activation can be picked up by Machine Learning algorithms, but may be too minute for clinicians to measure.

2.3 Bradykinesia/akinesia

Bradykinesia describes the slowness-of-movement characteristic of many PD patients [X]. It involves reduced amplitude and speed of repetitive and voluntary actions. Eye Blink Rate (EBR) is a widely accepted product of *Bradykinesia*. Between 50- 80 percent of PD patients have significantly reduced EBR as compared to Healthy Controls

(HC). Additionally, EBR is one of the earliest signs of PD, making it a vital factor for incipient diagnosis.

Literature on *Bradykinesia*-induced reductions in EBR, for PD diagnosis is quite sparse relative to the amount

of literature on *Dysarthria* and *Hypomimia* [X, X]. However, numerous studies have shown statistically significant

Figure 1: Summary of sources, demographic information, content information, and citations. N/A represents categories where no information was available, and information could not be compute

differences in the EBR of PD patients, and HC, providing a promising space to exploit for PD diagnosis. During a reading task, it was found that PD individuals averaged only 2.4 blinks/minute, while HC averaged 7.8 blinks/per minute [X].

2.4 Multimodal Diagnosis

While many classification algorithms with Machine Learning, can perform quite well on PD diagnosis, they often rely on only one symptom. While this can be effective, the onset and severity of many PD symptoms vary greatly. Some symptoms may never manifest in a PD patient, while others typically only manifest in the middle to late stages of PD. As such, it is unreliable to use only unimodal classifiers. Furthermore, past multimodal studies, for not only PD diagnosis [X], but the diagnosis of other neurodegenerative diseases, have shown that multimodal models can offer significant performance improvements over unimodal classifiers [X, X,X].

2.5 Severity profiling

The Hoehn and Yahr Scale (H&Y Scale) is a 5-point scale that rates a PD patient's level of disability and disease progression, from the evaluation of dozens of tasks by a clinician. A few have managed to estimate disease severity from symptoms with Machine Learning, although the Root Mean Square Error (RMSE), is usually close to one full stage [X], and they usually rely on symptoms that are difficult or impractical to measure contactlessly and automatically.

III. DATA & METHODOLOGIES

Since we theorized that *Dysarthria*, *Hypomimia*, and *Bradykinesia* features could reliably be collected and analyzed via a web platform, these symptoms of PD were focused on. ParkinSense, a web application, was developed that enabled individuals to complete a 6 minute study. After the study was complete, data processing, feature extraction, model diagnosis, and H&Y scale estimation took under 5 minutes.

3.1 Data Sources

Data was gathered from a variety of sources. All data was gathered or collected from sources that obtained data in scientifically rigorous environments and methods. Speech Recording data with extracted features were gathered from numerous datasets in UC Irvine Database. Data was collected from PDs and HC while they phonated /a/, /o/, and /u/. Datasets A-E of Figure 1 outlines the size,

demographics, and information within each dataset used for this modality. The vowel phonantons of /a/ were aggregated into a single large dataset, with the features represented being the intersection of the features in Datasets A-E. A

similar procedure was done for the vowel phonation of /o/ from datasets B and C. Data was also collected from PD and HC while they spoke words and short sentences, and is found in dataset B.

Facial Expression Data can be seen in Dataset F of Figure 1, and contains information regarding the variance of the magnitude of AUs involved with Smile, Surprise, and Disgust facial expressions, when the AU was activated.

EBR data is represented in Dataset G of Figure, and provides information about the EBR of PD patients and HC during a reading exercise, and also contains information about the H&Y Scale of diagnosed individuals, and whether or not they wore glasses during the exercise. The EBR dataset was augmented via the Hastie Algorithm to match the size of the other datasets [x]

3.2 Website Diagnosis Framework

A diagram of the website, which is to be described in sections 4.2-4.4, can be found in Figure 2. A prototype example of this website can be found in Appendix A.

When users navigate to the web application, and begin the diagnosis, they are first asked to be in a seated position, with their laptop or computer one foot away from them, such that their entire head, and part of their torso are in the frame, not dissimilar from a passport photo. They are requested to be well-hydrated and not tired. Finally, they should be in a quiet location, with their face being well lit.

Before they begin the diagnosis, they are prompted with a screen that allows them to adjust their camera and microphone volume, so they are able to be in the optimal position relative to the above conditions, when they begin the diagnosis.

3.2.1 Vocal Analysis

When users begin the diagnosis through the web application, they are first asked to vocalize sustained vocal phonations, for as long as they can, or until 10 seconds elapse, whichever is shorter. It was emphasized that they should maintain constant volume and pitch. They are requested to complete three such phonations: , , and . These were chosen, as they have had success in the past with accurate results [X, X].

3.2.2 Facial Expressions

Then the user is asked to make 3 pronounced facial expressions: joy, surprise, and disgust. These were largely chosen since in [X], training models on these expressions proved to be effective when differentiating between PD patients and HC. To improve accuracy users were asked to make each face three times in succession, alternating with a neutral face.

3.2.3 Short Passage

As the next section of the diagnosis, users were asked to read as much as they could of a selection of 5 passages for 60 seconds. They were asked to read in their regular reading speed, volume, and voice. Text spanned 90% of the viewport width, centered in the screen, at a font-size of 15 pixels. The specific reading passages were chosen for having a diverse set of words, ranging from short to long, monosyllable to polysyllabic, containing most sounds in the English Language, and suspected phrases that might divulge important information for the machine learning models to pick up on. The specific passages can be found in Appendix A.

Paragraphs were ordered in order of expected relative importance, since users may not have time to read all five paragraphs in the allotted time with their regular reading speed.

From this section, data for 2 of the modalities were collected. Specifically, data was gathered for vocal analysis, to be split up into small sentences, phrases, and individual words. Furthermore video data was collected for EBR analysis.

3.2.4 Questionnaire

The final section of the diagnosis involved a short questionnaire, asking a few questions including gender, age, and whether the individual wore glasses during section 3.2.3.

3.3 Data Pre-processing & Feature Extraction

3.3.1 Vocal Data

Data processing and Feature Extraction was done entirely through an automated workflow, to mimic the requirements for automated data processing when used as a web application for rapid diagnosis.

First vocal data from section 3.2.1 was separated from webcam data, and webcam data discarded. Additionally, the recording was trimmed, such that the first 0.5 seconds of the recording was discarded, since it was noticed in [X] that the onset of speech results in a nonuniform levels of pitch and volume. Additionally, the last 2 seconds were thrown out, since during this timeframe, the user typically runs out of breath, resulting in diminishing volume and degrading pitch that does not represent the user's true vocalization. Features were extracted from vocal data that represents the intersection of Features in datasets A-E. These features were extracted using openSMILE 3.0

3.3.2 Facial Data

Each participant submits 3 videos, each containing 6 facial expressions (3 of the requested facial expression, and 3

neutral faces), for a total of 18 facial expressions (9 total facial expressions for smile, surprise, and disgust).

For facial data, each of the 18 facial expressions was run through DeepFace [X] to ensure the correct facial expression was made. If the incorrect facial expression was detected, the recording was discarded. Then with a server-hosted instance of OpenFace 2.0 [X], the magnitude of the AU's associated with each of the 3 requested facial expressions were calculated when the expression was being made, and the variance of these AU magnitudes was calculated, representing a metric of the amount of facial movement. Recordings were split at the midpoint troughs of AU variance, indicating a neutral face.

3.3.3 Reading Passage Data

For data derived from the reading passage, vocal and webcam data were split from each other. Vocal data was split up into individual words using the Pydub module [X], which was temporally paired with the words in the passage, to create a mapped transcript. RT-Bene was used to calculate the EBR, and is effective both with and without lenses.

3.4 Website Deployment

The diagnosis website was developed in React.js and deployed on a personal domain with a signed certificate. It can be accessed in Appendix A. Data is collected client-side, and processed server-side. The unprocessed video and audio data was stored server-side, then discarded after being processed to protect user privacy. Server-side was written in (PHP), connected to the front-end through (AJAX) requests. The trained models were hosted in Amazon Web Services (AWS), and accessed by the server through an AWS Gateway API.

3.5 Model Framework

A diagram of Parkinson's Diagnosis Model Framework can be seen in Figure 3. ParkinSense's diagnosis framework revolves around the use of multiple classifiers to provide a binary classification of PD presence in a subject

3.5.1 Dysarthria

Given the ample literature on Vocal-based methods of PD diagnosis, numerous classifiers and models have been well-tested and documented. In this step, the featureset was reduced to 10 using Genetic Algorithm-based feature set selection for all datasets. A Support Vector Machine (SVM) and Random Forest were used via Ensemble Learning with a Logistic Regression combination classifier to provide a single binary classification for . For , and , an Adaboost Classifier, and Gradient Boosted Classifier were used respectively, as they provided the best performance out of the tested classifiers. Finally, a random forest was used for both words and short-sentence classification. A final binary classification for Dysarthria-based methods was generated through hard majority-voting decision fusion, from the 5 individual classifiers. Weighting individual classification during the decision fusion stage did not seem to out-perform equal-weighting in this study.

3.5.2) Hypomimia

Since three sets of 3 facial expressions were recorded for each user. A Support Vector Machine with a Radial Basis Function Kernel was used on each set, and combined with a majority voting decision fusion to generate a single binary classification for this modality.

3.5.3) Bradykinesia

For EBR, a Supervised Multi-layer Perceptron was used to predict a binary classification

3.5.4) Modality Decision Fusion

To combine the decisions from each modality, Soft Voting was used, with weights of 30, 35, and 35 were used for Dysarthria, Hypomimia, and Bradykinesia respectively, to account for the fact that Dysarthria typically manifested later on than the other symptoms.

3.5.5) Severity Profiling

We predicted two scores: H&Y Scale score, which classifies the severity and progression of PD into a five point scale, and UPDRS total score, which is a 199 point scale for disease severity and progression. Both of these scores are conventionally calculated through lengthy questionnaires conducted by a clinician. H&Y Score was predicted via Stochastic Gradient Descent Regression using EBR data. UPDRS scores were predicted from datasets B and D individually, using Gradient Boosted Tree Regressors, and used a voting regressor, with 80% weight assigned to the dataset B, due the magnitude of data in that set.

3.6) Evaluation

Due to the difficulty in finding data, especially on all modalities simultaneously for overall model evaluation, high-fidelity synthetic patients were generated, with specific modality instances being paired by gender, age, and disease severity. As reported by [X] synthetic patients can provide near identical results when compared to live patients. K-fold-cross validation was used for all model evaluations, coupled with metrics discussed in selection 4.

IV. RESULTS & DISCUSSION

Accuracy (), precision (), recall(), F-1 score (), sensitivity, (), specificity (), and negative predictive value () scores were used as the metrics for this task. Formulas for these equations can be found below:

$$a = \frac{TN + TP}{TP + TN + FP + FN}$$

$$p = \frac{TP}{TP + FP}$$

$$r = \frac{TP}{TP + FN}$$

$$s = \frac{TP}{TP + FN}$$

$$S = \frac{TN}{TN + FP}$$

$$f = \frac{2 * TP}{2 * TP + FP + FN}$$

$$n = \frac{TN}{TN + FN}$$

Where (TP) defines the number of True Positives, (FN) defines the number of False Negatives, (TN) defines the number of True Negatives, and (FP) defines the number of False Positives.

The results for each modality, and specific internal classifiers, and the overall model, can be seen in table 2. In large, all modalities had accuracy rates above 85%, with vocal analysis being the most accurate, and EBR being the least accurate.

Given the accuracy rates of each modality algorithm, the lower bound on accuracy is 0.995. However, ParkinSense was able to obtain an overall accuracy rate of 0.9972, confirming it as an accurate method of diagnosis.

The high value of is important, as it is vital to reduce false-positives as much as possible in the medical diagnosis industry.

5.1.1 Comparative Analysis

ParkinSense out-performs every unimodal diagnosis software in literature that was trained on large datasets. Parkin-Sense's individual and facial analysis classifiers perform similarly to the top models that utilize similar data in literature, but Parkin-Sense is the first to analyze and , as well as words, short sentences, or EBR. Furthermore, Parkin-Sense performs better than or equal to most multimodal PD diagnosis models, even those that analyze gene, MRI scan, or other invasive or non contactless forms of data.

5.1.2 Severity Profiling & Telemonitoring

Root-Mean-Square-Error (RMSE), R-Square, and Mean Absolute Error(MAE) were used as metrics for both H&Y Scale and UPDRS Scale regression predictions. The equations can be seen below(adding later). The metrics can be found in table 3 below.

Overall, ParkinSense's severity profiling model can be used relatively accurately to predict H&Y Scale and UPDRS Scale scores, allowing users to monitor disease progression, and medication effectiveness over time.

5.1.3. Website Analysis

All methods of data collection, feature extraction, and model querying via the web application were validated to be effective and successful. Data collection depended on the user, but usually took around 6 minutes. Pre-processing and Feature Extraction usually took around 3 minutes, and prediction via the model took around 2 minutes. As such a rapid diagnosis was achieved, as the entire process took around 11 minutes.

6.1 Error Analysis

Additionally, for dataset aggregation in section 4.1, certain features likely were not collected in similar manners for each dataset, which may have introduced some noise that made the model underperform.

Finally, other approaches to PD diagnosis, or other neurodegenerative diseases, and diseases with cardinal behavioral symptoms can benefit from a similar contactless diagnosis framework. By analyzing symptoms rapidly, accurately, and entirely contactlessly, many diagnosis tools

ParkinSense makes an encouraging process for the development of PD contactless diagnosis software via machine learning technologies and multiple modalities. ParkinSense automates data collection and feature extraction on the web application, and regarding the unimodality models, have performed similarly to many top models. Furthermore, we introduce the first models that involve PD diagnosis classification through EBR, and phonations, words and short sentences. Additionally, we developed UPDMS and H&Y Scale predictors that perform quite well. ParkinSense integrates all of this together, to develop a PD diagnosis and severity-proofing tool, that outperforms all unimodal and multimodal diagnosis tools in literature for PD diagnosis, and delivers competitively accurate levels of PD severity and progression metrics. All of this is accessible through a web application that is accurate and rapid. ParkinSense suggests that contactless PD diagnosis and severity tracking can be conducted accurately, which may prove invaluable for those in which PD diagnosis is inaccessible due to barriers, immobility, or associated costs. The unparalleled accuracy of ParkinSense demonstrates the effectiveness of multimodal machine learning algorithms, and may allow millions worldwide to be diagnosed correctly, ensuring effective medication can be distributed immediately, and preventing PD progression, as well as improving the quality of life of PD patients.

[illegible]

[illegible]

[1] blah blah blah blah blah blah blah blah blah blah blah blah blah blah
 blah blah blah blah blah blah blah blah blah blah blah blah blah blah
 [1] blah blah blah blah blah blah blah blah blah blah blah blah blah blah
 blah blah blah blah blah blah blah blah blah blah blah blah blah blah

Appendix A

Diagnosis steps and information can be found at
<https://measure.parkin-sense.pulkith.com>

TO DO: label modalites in graphic

TO DO: fix colors in graphic

TO DO: fix resolution

TO DO: add citations

TO DO: add latex in table headers

TO DO: add regression metric formulas

