# Survey on Parkinson Detection Using Machine Learning

Amit Ambadas Patil, Prof. Sudesh L Farpat,

[1,2] Department of Computer of Engineering Padmashri Dr. V.B. Kolte College of Engineering

**ABSTRACT**

*Diagnosis of Parkinson's  disease (PD) is very important for treatment. In this study we highligh the potential of machine learning algorithms in aiding the development of non-invasive diagnostic tools to aid diagnosis of PD. Analogously, we proposed a multi-response model based  on various clinical, motor function, and neurosurgical datasets. Machine learning methods, including support vector machine, random forest and neural networks, were used for PD and non-PD cases  classification. The sensitivity, specificity and accuracy of the model performance subsided traditional diagnostic methods. Cross-validation and hyperparameter tuning were needed for its  results. Results suggest that machine learning algorithms may have utility in the diagnosis of PD, which may allow earlier detection and improved outcomes for patients. Further work includes enlarging the database and optimizing the  algorithms for improved accuracy and applicability.*

*Index Terms—Parkinson's disease, machine learning, diagnosis, classification, support vector machine, random forest, neural network, feature selection, cross-validation, hyperparameter tuning, sensitivity, specificity, accuracy.*

## 1. INTRODUCTION

Parkinson 's disease Parkinson's disease (PD) is a neurodegenerative disorder, which is accompanied by slow motor function deterioration because of the loss of dopamine-generating neurons in the brain. Typical symptoms can include tremors, muscle stiffness, loss of balance, and problems with walking, coordination and speech. These symptoms generally develop gradually and become increasingly debilitating, leading to severe mobility and communication impairment. Other non-motor symptoms, such as behavioural changes, dementia and depression are also common. The presence of a combination of these core motor symptoms results in  a condition known as "Parkinsonism" or "Parkinsonian Syndrome." At present, PD is generally diagnosed  in late stages of the disease and after extensive dopamine loss has already occurred. This impaired detection creates a significant  clinical burden. Symptoms can be very different from one person to another and treatments decisions are complex because diagnosis is highly dependent on medical assessment of doctors. Because of this, the mental health aspects of PD are frequently under-recognized, leading to a series of secondary health implications. Standard clinical tools for PD diagnosis are  the following:

- MRI or CT (Magnetic Resonance Imaging or Computed Tomography): Traditional imaging techniques do not  notice early onset Parkinson's.
- Positron Emission Tomography (PET): Assesses the level of activity in the brain areas controlling movement. Single-
- Photon Emission Computed Tomography (SPECT):Shows chemical  activity in the brain (like a shortfall of dopamine).

Nevertheless, they cause frequent misdiagnoses, if applied by nonspecialists – misdiagnosis may  occur in 25 % of patients. Many patients sufferfrom the disease for  many years before having a cleandiagnosis and early interventions are ineffective. Machine learning (ML), a  subfield of artificial intelligence, has shown great potential to address these challenges. ML has the capability of letting machines learn from data and perform better in time, without being programmed explicitly to do a certain task. These models learn from training data to perform predictions or make decisions, and are increasingly  useful in areas like healthcare, speech recognition, computer vision and diagnostic medicine—where traditional rule-based algorithms might not be sufficient. By leveraging ML to medical datasets, such as speech, handwriting, and  imaging, researchers hope to elevate the accuracy and speed of PD diagnosis, which could allow for earlier detection, and more personalized treatment options. Parkinson's

disease (PD) advances as a neurodegenerative condition which generates several motor as well as non-motor manifestations that substantially decrease patient existence quality. Dividing detection and diagnosis at an early stage enables both effective disease management and treatment through improved patient outcomes and disease progression tracking abilities [1]. Strategies based on traditional diagnosis methods use clinical evaluation and subjective evaluation methods that commonly lead to prolonged diagnosis times and incorrect classifications. The recent developments in machine learning open new promising prospects to improve both accuracy and speed of Parkinson's disease detection. Machine learning models equipped with extensive clinical records and neurosurgical data and biomarkers detect complex patterns and minor transformations linked to PD [2]. The models rapidly process large datasets to identify individuals at risk who may display signs of Parkinson's disease before prominent symptoms appear [3].

## 2. LITREATURE SURVEY

**A. Synopsis of Machine Learning-Based Speech Prediction for Parkinson's Disease:** The investigation creates a Parkinson's disease predictive model through Support Vector Machines (SVM), Decision Trees, and Logistic Regression machine learning methods [4,5]. The methods today face challenges with their diagnostic tools because they lack effective noninvasive techniques for early detection and create spaces for further research [8]. A dataset consisting of speech-related features collects audio recordings between Parkinson's patients and individuals without Parkinson's disease. Researchers examine the audio recordings to find early Parkinson's disease indicators. The prediction models accept Mel-Frequency Cepstral Coefficients (MFCCs) as speech characteristic representations after feature extraction and pre-processing. Accuracy together with confusion matrix metrics determine model performance [6]. The models encounter obstacles because speech variation stems from environmental noise along with differences between individual speakers. The prediction models encounter performance limitations due to divergent speech sounds that arise from people-specific features or noisy environments.
Voice analysis coupled with machine learning demonstrates great potential for healthcare and neurological disorder diagnosis as it brings technological advancement to early disease detection and treatment of Parkinson's disease [16]. Voice analysis integrated with machine learning shows promise as a technological healthcare solution to enhance the early identification and treatment of neurological conditions including Parkinson's disease [16].

**B. Overview of Machine Learning-Based Parkinson's Detection Using Speech Features:** Researchers have evaluated the potential for diagnosing Parkinson's disease through machine learning algorithms which analyze vocal features. Support Vector Machines (SVM) and K-Nearest Neighbors (KNN) serve as the main algorithms while the research focuses on detecting Parkinson's disease at an early stage through non-invasive methods. The research underscores a requirement for new voice-based assessment instruments that can recognize PD symptoms in their initial stages. The research relies on sustained vowel sound "/a/" recordings found in its dataset to achieve its objective. The research dataset combines data from 188 Parkinson's patients (107 men and 81 women) with data from 64 healthy participants (23 men and 41 women) across the age range of 33 to 87 [7,8].
Research uses a machine learning feature matrix of extracted audio characteristics to measure 752 different elements between intensity metrics and formant frequencies and baseline vocal parameters. By applying the KNN algorithm the researchers obtained exceptional results since it detected Parkinson's disease perfectly in every case which confirmed speech features work well for diagnosis. Research examines how gender influences both diagnostic precision and disease development timelines. Researchers plan to investigate how gender along with PD disease progression rates influence prediction results in upcoming studies [9].

**C. Detection of Parkinson's Disease Using Voice Source Information**
A research study examines PD identification using traditional Support Vector Machines (SVM) and Multi-Layer Perceptron (MLP) and deep learning technique Convolutional Neural Networks (CNN) along with voice source analysis [10, 11]. Researchers focus on discovering robust speech features that effectively distinguish voice elements between PD patients and healthy people.
The research draws its data from the PC-GITA speech dataset which contains balanced audio recordings of 50 Parkinson's disease patients together with 50 healthy controls. Both raw speech and glottal flow waveforms are analyzed in two pipelines: Voice source features helped increase detection accuracy resulting in approximately 68% successful outcomes.
A binary classification process identifies regular speakers versus individuals with PD without additional context. Research demonstrates QCP achieves better glottal inverse filtering performance than IAIF because of its moderate accuracy level. Future research aims to enhance end-to-end systems while expanding the dataset in order to improve diagnostic performance. The upcoming research will concentrate on establishing bigger datasets and advancing end-

to-end systems to enhance diagnostic outcomes. Glottal source analysis shows potential for PD detection during early stages based on research in [14, 17].

**Table 1: literature survey from current years**

| Year | Gap Identified | Benefits | Demerits | Future Work |
|---|---|---|---|---|
| 2024 | Lack of timely, noninvasive detection methods | Early detection enables better outcomes | Variability in voice data affects results | Augment dataset, improve UI, explore transfer learning |
| 2024 | Need for robust handwriting analysis methods | High detection accuracy; new dataset strengthens research | Validated only on Arabic handwriting; needs adaptation | Explore additional parameters and learning techniques |
| 2023 | High-dimensional data risks overfitting | Achieved 100% accuracy in detection | Focused only on detection, not progression | Investigate effects of disease progression |
| 2023 | Limited research on vocal data vs. other biomarkers | Promotes telemedicine applications for PD classification | No single method perfect; analysis is time-consuming | Integrate vocal data with REM sleep data for better outcomes |
| 2022 | Need for effective detection in noisy environments | Suitable for real-world applications | Excludes accepts; only post-diagnosis participants used | Develop advanced models; improve noise robustness |
| 2022 | Need for better early-stage detection systems | Enhances patient outcomes through early identification | Small sample sizes limit generalizability | Expand datasets; explore new techniques |
| 2022 | Challenges with self-reported data accuracy | Introduces high-frequency voice-based screening tool | Misclassification due to label noise; device variability impacts results | Research on federated learning approaches |
| 2021 | Need for robust speech features for differentiation | Improved accuracy using voice source information | Modest accuracy (~68%); large datasets required | Explore alternative voice representations; transfer learning |
| 2021 | Limited use of advanced algorithms for early detection | High accuracy (100% for individuals) | Variability in voice data affects results | Explore neural networks for improved accuracy |

## D. Machine Learning Approaches for Parkinson's Disease Detection

A research project examines how Naive Bayes, XGBoost, and Decision Tree classifiers work to predict Parkinson's disease. The research addresses an existing performance deficit in previous models because their accuracy reached only 73.8% by examining new sophisticated algorithms for enhanced early detection results. The available dataset contains voice files with each recording having 24 numerical features describing vocal symptoms and a target field that marks patients as PD (value=1) or healthy (value=0). The tested model demonstrates 100% precision during single case tests while achieving 94.87% total success across every participant[12].

XGBoost showed the most effective predictive abilities among the tested algorithm models. The study mentions personal and environmental factors that create challenges in stabilizing speech data variations. Future research will examine multiple speech processing approaches alongside robust model development for improving detection performance in realistic challenging circumstances [5, 16].

## E. Smartphone-Based Phoneme Analysis in Real-World Settings

Researchers conducted an evaluation of phoneme recordings taken on smartphones in natural environments to detect Parkinson's disease [13]. Random Forest classifiers integrated with Support Vector Machines (SVM) form the research backbone to improve detection of PD for clinical practice together with telemedicine functions. The researchers obtained smartphone-based phoneme data from 72 participants including 36 individuals with Parkinson's disease and 36 healthy participants at clinical locations. The study recorded participants producing three phonemes (/a/, /o/ and /m/) for the analysis of speech production variation.

Using leave-one-out cross-validation, the study reports perfect classification performance: 100% accuracy, sensitivity, and specificity. Results from the study established significant phoneme production distinctions between patients with Parkinson's Disease and their age-equivalent healthy subjects despite sound interference. The study's limitations stem from using Melbourne-based participants and excluding minimal PD cases. Research improvements will direct toward the creation of personalized models combined with new biomarkers, long-term monitoring

systems while advancing noise reduction methods and implementing transregional tests. Phoneme analysis conducted on smartphones shows promise as a usable means for diagnosing Parkinson's disease remotely [14].

### F. Machine Learning-Based Early Parkinson's Disease Detection

The research examines machine learning techniques including Random Forest together with K-Nearest Neighbors (KNN) and Support Vector Machines (SVM) and Logistic Regression for identifying Parkinson's disease at its early stages. Voice data remains underutilized in Parkinson's disease diagnostics while MRI approaches serve as the primary focus among existing traditional diagnostic techniques[15].

The research consisted of 31 participants between 46 and 85 years old who included 23 PD-diagnosed participants. The participants submitted on average 195 phonation samples spanning 1 to 36 seconds in length. Three modeling strategies were evaluated: The research contained three separate experiments: a balanced dataset experiment followed by a complete dataset experiment and a final experiment with principal component analysis (PCA).

The Random Forest classifier emerged as the top model due to its 91.83% accuracy and 0.95 sensitivity rate. After performing PCA the KNN model produced accuracy results comparable to other models. The research notes that voice data demonstrates diagnostic potential but suggests that diagnostic accuracy could improve through the combination of voice data alongside REM sleep data.

Study results demonstrate mobile voice recording technology holds potential for establishing low-cost telemedicine solutions in remote PD diagnosis. While no model reached absolute accuracy it showed improved efficiency and reduced size compared to deep learning systems. The authors suggest Random Forest as the top choice among models for Parkinson's disease classification since it demonstrated the best performance accuracy [16].

### G. Overview of Parkinson's Disease Detection via Online Handwriting Analysis

A new diagnostic method for Parkinson's disease is studied which analyzes online handwriting by integrating beta-elliptical models with fuzzy perceptual detectors and Bidirectional Long Short-Term Memory (BLSTM) networks. During these tasks researchers documented pressure data and x and y coordinate information to support feature extraction and classification applications [17]. Handwriting samples from 30 PD patients along with 30 healthy control subjects were obtained through the usage of a Wacom digitizing tablet. Participants performed five tasks which consisted of writing Arabic words and drawing spiral shapes and ellipse designs. The collected data included pressure measurements alongside x and y coordinate positions for feature extraction during classification [4, 15]. The implemented methodology performs similarly to or even better than conventional diagnostic procedures. The results indicate that this diagnosis method competes effectively with existing approaches while displaying superiority in specific instances. A major drawback exists because the methods solely examine Arabic handwriting yet extensive validation is required to cover multiple languages and geographical regions. Additional research will utilize Hamilton–Jacobi–Bellman (HJB) equation-based learning methods to enhance classification by analyzing new variables. Overall, the findings lay the groundwork for expanding handwriting-based PD diagnostics and confirm the promise of this method for non-invasive, early-stage detection [8].

### H. A Review of Datasets Used for Parkinson's Disease Detection

Multiple research groups developed PD detection models using different datasets which generated distinct features and insight into the disease. The UCI Parkinson's dataset contains speech data from 188 patients with Parkinson's disease and 64 healthy individuals that is commonly used for research purposes. Physiological data from patients with Parkinson's disease and healthy controls shows sustained vowel phonations yielding acoustic features helping diagnose PD. The PC-GITA speech database stands as a crucial resource which captures voice data from 50 PD patients together with 50 healthy controls. The dataset contains speech recordings alongside glottal flow waveforms to provide support for traditional algorithms along with end-to-end deep learning methods.

Researchers are finding Smartphone-based phoneme datasets increasingly relevant for their work. Researchers collected phoneme samples from 72 participants including 36 people diagnosed with PD and 36 healthy controls while they recorded in real-life settings using iOS devices. Data sets demonstrate that telemedicine alongside non-clinical speech-based biomarker monitoring presents real-world possibilities. Recent research initiatives have explored handwriting data as a supplementary method for indicating Parkinson's disease detection. Good experimental data utilizing digitizing tablets acquire and monitor pressure alongside position coordinates from individuals conducting writing tasks including spiral drawings and word writing for motor symptom analysis. A notable large-scale initiative is the i-Prognosis project, which collected over 29,000 anonymized phone call recordings through a smartphone app. This dataset investigates voice as a continuous biomarker in uncontrolled, real-world settings, emphasizing scalability and ecological validity.

## 3. CHALLENGES

Medical systems that detect Parkinson's disease through voice features experience multiple significant problems which reduce their diagnostic precision and reliability and generalizability. Speech pattern variability represents a significant barrier because it occurs from various factors such as patient age and gender and their specific regional accent alongside potential health conditions. Different variations create barriers that prevent machine learning models from mapping knowledge across various populations. Insufficient availability along with dataset imbalance represents an essential hurdle within the field. The current datasets have insufficient sample sizes along with insufficient diversity which hinders the development of reliable models that perform effectively in real-life conditions. Data scarcity heightens the possibility that models will perform effectively on their training data yet struggle to apply learned knowledge to unidentified test data.

Intrinsic variations within Parkinson's patients' speech patterns become more challenging because they stem from time-dependent conditions together with medicine effects and emotional changes in patients' moods. Voice-based models experience deterioration in predictive stability when faced with these fluctuating alterations. The practice of selecting features creates complex problems during implementation. Voice data collection requires the establishment of proper ethical protocols and privacy safeguards. Experts recognize feature extraction of significant voice information as an essential active area of scientific inquiry. Assessing the ethical implications and ensuring privacy protection for voice data collection remains the final critical aspect. Improving the clinical usability and fairness alongside the system's robustness requires the solution of current challenges with voice-based Parkinson's disease detection systems. The clinical effectiveness and robustness of voice-based Parkinson's disease detection systems requires solutions to overcome these identified challenges.

## 4. CONCLUSION

This review highlights the significant progress made in the early, non-invasive diagnosis of Parkinson's disease (PD) through the application of machine learning techniques. From voice-based analysis using algorithms such as Support Vector Machines (SVM), Random Forest, and K-Nearest Neighbors, to handwriting-based detection utilizing Bidirectional Long Short-Term Memory (BLSTM) networks, these studies collectively demonstrate the transformative potential of integrating artificial intelligence into clinical diagnostics. The research underscores the strong diagnostic value of speech features, with several models achieving exceptionally high accuracy—some nearing or reaching 100% under specific conditions. Despite these advancements, key limitations persist, including limited dataset diversity, small sample sizes, and challenges related to speech variability influenced by demographic and environmental factors. Looking forward, expanding datasets to encompass broader populations, improving model generalization, and incorporating additional biomarkers such as REM sleep or motor data are crucial next steps. The growing use of smartphones and remote monitoring tools presents an exciting opportunity to enhance accessibility and efficiency through telemedicine, particularly for early-stage diagnosis and continuous tracking. In conclusion, the integration of machine learning into PD diagnosis offers a promising path toward more accurate, timely, and scalable healthcare solutions. Continued innovation and interdisciplinary collaboration will be essential in translating these methods into real-world clinical practice, ultimately improving patient outcomes and quality of life.

## REFERENCES

[1] Exemplar-based sparse representations for detection of parkinson's dis- ease from speech. IEEE, 2023.

[2] Parkinson's disease detection from voice and speech data using machine learning(2023 researchgate), 2023.

[3] Bilal Alatas, Shadi Moradi, Leili Tapak, and Saeid Afshar. Identification of novel noninvasive diagnostics biomarkers in the parkinson's diseases and improving the disease classification using support vector machine. BioMed Research International, 2022.

[4] L. Ali, C. Chakraborty, and Z. et al. He. A novel sample and feature dependent ensemble approach for parkinson's disease detection. Neural Comput Applic, 2022.

[5] F. Amato, I. Rechichi, L. Borz`ı, and G. Olmo. Sleep quality through vocal analysis: A telemedicine application. In 2022 IEEE International Conference on Pervasive Computing and Communications Workshops and other Affiliated Events (PerCom Workshops), pages 706–711, 2022.

[6] B. Anbalagan, S. Karnam Anantha, and R. Kalpana. Novel approach to prognosis parkinson's disease with wireless technology using resting tremors. Wireless Pers Commun, 2022.

[7] Adjunct Faculty at the University of Science and UAE Technology of Fujairah, Fujairah. Parkinson's disease detection using voice features and machine learning algorithms. In 2023 International Conference on Microelectronics (ICM). ©2023 IEEE, 2023. Visiting Scholar at the University of Sharjah, Sharjah, UAE from Sept. to Dec.2023.

[8] Muskan Bahrani, Meet Chhabria, Kaustubh Kharche, Sakshi Shinde, and Prashant Kanade. Voice-based parkinson's disease prediction using machine learning. IJNRD, 9(4):ISSN: 2456–4184, apr 2024.

[9] J. R. Barr, M. Sobel, and T. Thatcher. Upsampling, a comparative study with new ideas. In 2022 IEEE 16th International Conference on Semantic Computing (ICSC), pages 318–321, 2022.

[10] I. Gupta, V. Sharma, S. Kaur, and A. K. Singh. Pca-rf: An efficient parkinson's disease prediction model based on random forest classifica- tion. 2022.

[11] R. Islam, E. Abdel-Raheem, and M. Tarique. A study of using cough sounds and deep neural networks for the early detection of covid-19. Biomedical Engineering Advances, 3:100025, June 2022.

[12] R. Islam, E. Abdel-Raheem, and M. Tarique. Voiced features and artificial neural network to diagnose parkinson's disease patients. In Proceedings of the International Conference on Electrical and Com- puting Technologies and Applications, American University of Ras Al Khaimah, United Arab Emirates (UAE), November 2022.

[13] John Hopkins Medicine. Parkinson's symptoms.

[14] N. P. Narendra, Bjo¨rn Schuller, and Paavo Alku. The detection of parkinson's disease from speech using voice source information. IEEE Transactions on Biomedical Engineering, 70(11):3473–3483, 2023.

[15] N. D. Pah, M. Motin, and D. Kumar. Voice analysis for diagnosis and monitoring parkinson's disease. In Techniques for Assessment of Parkinsonism for Diagnosis and Rehabilitation, pages 119–133. Springer, Singapore, 2022.

[16] N. D. Pah, M. A. Motin, S. Raghav, and D. K. Kumar. Parkinson's disease detection using smartphone recorded phonemes in real world conditions. IEEE Access, 10:97600–97609, 2022.

[17] Nicolo´ G. Pozzi and Ioannis U. Isaias. Chapter 19 - Adaptive deep brain stimulation: Retuning Parkinson's disease, volume 184, pages 273–284. Elsevi