# Parkinson's Disease Prediction

**Sandhya Budhavale**
*Usha Mittal Institute of Technology*
*SNDT Women's University*
Mumbai 400049, India

**Siddhi Mane**
*Usha Mittal Institute of Technology*
*SNDT Women's University*
Mumbai 400049, India.

**Rutuja Shinde**
*Usha Mittal Institute of Technology*
*SNDT Women's University*
Mumbai 400049, India

**Prof. Sujata Kullur**
*Usha Mittal Institute of Technology*
*SNDT Women's University*
Mumbai 400049, India

*Abstract*—In 2019, Parkinson's disease (PD) affected more than 8.5 million people worldwide. Men are usually more affected than women. Parkinson's disease is a neurodegenerative disease that primarily affects the motor system (slow movement,tremor,rigidity). The aim is to develop a machine-learning model for predicting Parkinson's disease based on speech and drawings. The model collects and analyses data from Parkinson's patients and healthy individuals to identify patterns and correlations between various features. The model is trained on data (voice and drawing) and tested to achieve excellent accuracy and performance in predicting Parkinson's disease. This discovery may lead to an earlier diagnosis and more effective treatment regimens for patients with Parkinson's disease.

*Index Terms*—*Parkinson's, Naive Bayes, KNN, Logistic Regression, Random Forest, XGBoost.*

## I. INTRODUCTION

Parkinson's disease (PD), or simply Parkinson's disease, is a long term neurodegenerative disease of the central nervous system that mainly affects the motor movements. Symptoms began gradually, with a slight tremor in one hand and a feeling of body stiffness that worsened over time.Parkinson's sickness has each motor & non-motor symptoms . It is mainly caused by the destruction of cells in the nervous system. A major pathological indication of Parkinson's disease is cell death in the basal ganglia of the brain (affecting up to 70percent of dopamine-secreting neurons in the substantia nigra pars compacta).PD begins in the medulla and olfactory bulbs before migration to the substantia nigra pars compacta and the rest of the midbrain/basal region of the forebrain. When these disease begins to affect the substantia nigra pars compacta, motor symptoms begin to appear. Certain nerve cells (neurons) in the brain gradually deteriorate or die.Many symptoms are caused by the loss of neurons in the brain that produce dopamine, a chemical messenger. Low dopamine levels causes abnormal brain activity.

signs and symptoms of PD includes:

- Temor- A tremor or rhythmic tremor usually starts in the limbs.

- Slow movement (bradykinesia)- Bradykinesia is the most disabling symptom.
- Impaired Posture and Balance: Posture may be bent. Postural instability leads to loss of balance, frequent falls, fractures, and reduced mobility.
- Loss of automatic movements: reduced capacity to carry out unconscious actions like blinking, grinning, or waving your arms while walking.
- Change in speech style: speaking quietly, quickly, slurred, or hesitantly before speaking monotonous rather than normal language patterns.
- Writing Changes: writing can be difficult and appear small.
- Cognitive: This includes mood, behavior, and thinking problems. problems with planning, abstract thinking, inhibition of inappropriate behavior, initiation of appropriate behavior, working memory, and attention control.

There is no known cure for Parkinson's disease. Medicines, surgery, and physical therapy can improve people's quality of life. Drugs most commonly used for treatment:**levodopa, COMT inhibitors, dopamine agonists, MAO-B inhibitors.**

**Similar cousins:** Corticobasal Degeneration- A (CSB) rare type of parkinsonism that affects people over the age of 40, usually between the 50s and 70s. Corticobasal degeneration is a disorder where cells in specific regions of the brain (cortex and basal ganglia) are harmed due to the accumulation of tau protein over time. In the normal brain, tau is degraded to avoid accumulation, but this does not occur in corticobasal degeneration.

Multiple system atrophy- It (MSA) is a degenerative brain disorder that affects physical functions such as blood pressure, heart rate, and bladder function and is associated with Parkinson's disease.

Progressive supranuclear palsy- It (PSP) is a rare neurodegenerative disease that is often misdiagnosed as Parkinson's

disease due to similar symptoms. PSP is caused by the destruction of brain cells in a few small but very important areas at the base of the brain. Many of the symptoms of palsy become more apparent when this area of the brain is affected by the disease. It is not yet known why brain cells degenerate.

## II. Literature Survey

Shetty, S., and Rao, Y.S. (2016) proposed an "SVM-based machine learning approach for identifying Parkinson's disease using gait analysis"[1][2]. The set of statistical feature vectors considered here from time series response data is reduced using a correlation matrix. These feature vectors are then analysed individually to extract the best seven feature vectors.These feature vectors are classified using SVM, kernel-based classifiers, and Gaussian radial basis functions. The results show that the seven features selected for SVM achieve a good overall accuracy of 83.33%, a good Parkinson's disease detection rate of 75%, and a low false-positive rate of 16.67%.

K. Manoj, I. Manikanta, CH. Deekshith, and K.V. Mukesh proposed a paper "Parkinson's disease prediction using machine learning techniques"[2]. In this paper, the author uses various machine learning techniques, such as KNN, Naive Bayes, and Logistic Regression, and describes how these algorithms can be used to predict Parkinson's disease based on user input and how they work. Based on these features, authors predict more accurate algorithms. The Naive Bayes algorithm provides the highest precision (81%).

Anila M and Dr G Pradeepini proposed the paper "Diagnosis of Parkinson's disease using Artificial Neural network"[3][2]. The main purpose is to perform disease detection through voice analysis in Parkinson's disease patients. To do this, various machine learning techniques such as ANN, Random Forest, KNN, SVM, and XGBoost are used to classify the best model, compute the error rate, and measure the performance of all models used. The main drawback of this paper is that it is limited to ANN only two hidden layers.And such a neural network with two hidden layers is sufficient and efficient for simple data sets. They only used one feature of the selection method to reduce the number of functions.

SR Sonu, Vivek Prakash, Ravi Ranjan, and K. Saritha proposed the paper "Prediction of Parkinson's disease using data mining"[4]. Early-onset vocal cord disorders Voice impairment affects approximately 90% of people with Parkinson's disease. Our aim is to use a data mining algorithm decision tree (CART) to predict whether an individual will develop Parkinson's disease using a dataset of patient voice recordings. The patient's voice is recorded and converted into audio attributes such as jitter, simmer, and frequency using the PRAAT script. Recorded audio is tested to predict whether you have Parkinson's disease and to determine your disease status.

Anitha R, Nandhini T, Sathish Raj S, and Nikitha V proposed a paper titled "Early Detection of Parkinson's Disease Using Machine Learning"[5].Decision trees and k-means clustering are combined in the suggested predictive analytics paradigm. Furthermore, to facilitate the usage of real-time machine perception, the OpenCV (Open Source Computer Vision Library), a library of programming functions primarily targeted at real-time image processing, was created. In this way,the achievements demonstrate early detection of disease, and with appropriate treatment and medicine, we can extend the lives of sick patients and lead them to peaceful lives.

T. J. Wroge, Y. Ozkanca, C. Demiroglu, D. Si, D.C.Atkins and R. H. Ghomi, proposed paper titled "Parkinson's Disease Diagnosis Using Machine Learning and Voice"[6]. Uses supervised classification algorithms such as deep neural networks to identify individuals who can accurately diagnose a disease. It examines the effectiveness of doing analysis. For that, the paper compares the effectiveness of various machine learning classifiers in diagnosing diseases against noisy high-dimensional data. 85% peak accuracy The accuracy provided by the machine learning model exceeds the average clinical diagnostic accuracy of laymen (73.8%) and the average accuracy of movement disorder specialists (79.6% without follow-up and 83.9% after follow-up) in a pathological examination as ground truth.

Sabyasachi Chakraborty, Satyabrata Aich, Jong-Seong-Sim, Eunyoung Han, Jinse Park, and Hee-Cheol Kim proposed a paper "Parkinson's Disease Detection from Spiral and Wave Drawings Using Convolutional Neural Networks: A Multistage Classifier Approach"[7]. Identifying relevant biomarkers associated with specific health problems and detecting them is of paramount importance for the development of clinical decision support systems. The system developed in this study uses two different convolutional neural networks (CNNs) for analysis. Drawing patterns from both spiral and wave sketches Additionally, prediction probabilities are trained with a metal classifier based on ensemble votes to provide weighted predictions from both spiral and wave sketches. The full model was trained on data from 55 patients. The mean f1 score was 93.94%, the mean recall was 94%, the mean precision was 93.5%, and the overall accuracy was 93.3%.

Richa Mathur, Vibhakar Pathak, and Devesh Bandil proposed a paper titled "Parkinson Disease Prediction Using Machine Learning Algorithm"[8]. Data mining is the process of selecting, extracting, and modelling unknown, hidden patterns from large data sets. Machine learning (ML) algorithms can be used for early disease detection to extend life in the elderly and improve Parkinson's lifestyles. In this post, various MLs are used to improve dataset performance and play a key role in the early and timely prediction of disease. After comparing these

algorithms, choose the one that is most effective in terms of accuracy. Our experimental results indicate that the accuracy obtained from the combined operation of the ANN algorithm and ANN is superior compared to other algorithms.

## III. PROPOSED SYSTEM

The primary goal of this model is to improve prediction accuracy, which will be advantageous for Parkinson's patients and reduce the disease's prevalence. Only if Parkinson's is identified at an early stage can the necessary steps be taken for treatment.So it's important to find Parkinson's disease at an early stage. In order to locate this machine learning model, The outcomes of speech and drawing are displayed based on the model's training and testing, which utilised the following algorithms: KNN, Naive Bayes, Logistic Regression, Random Forest, and XGBoost. The datasets used were downloaded from Kaggle. The results of this will allow the doctor to determine whether something is normal or abnormal and to give the appropriate medication for the affected stage.

## IV. METHODOLOGY

The flow of execution of model includes six main steps:

- A collection of datasets proving Parkinson's disease.
- Pre-process the collected raw data into an understandable format.
- Divide the dataset into training and test data, then use the training data.
- Data are analysed using machine learning algorithms such as KNN, Naive Bayes, logistic regression for Speech and Random Forest, XGBoost Algorithms for drawing.
- After training data with these algorithms, tests are run.
- Finally, using the maximum accuracy for displaying results.
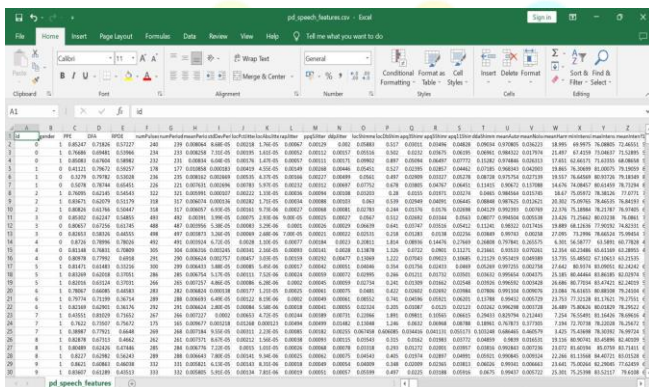
**SPEECH :**



Fig. 1. Speech dataset

The speech dataset that has collected from kaggle website. This dataset has around 756 patient's data and each row has

755 different voice features.Some are of no use when building a model. so need to remove unnecessary functions that are not responsible for generating output. Accuracy declines when more features are chosen. After identifying and dropping some features, only 11 main features are chosen.Those features are: **ID, Gender, PPE (Pitch Period Entropy), DFA (Detrended Fluctuation Analysis), RPDE (Recurrent Period Density Entropy) , NumPulses, NumPeriodPulses, MeanPeriod-Pulses, StdDevPeriodPulses, LoctPctJitter, Class**[9][10].
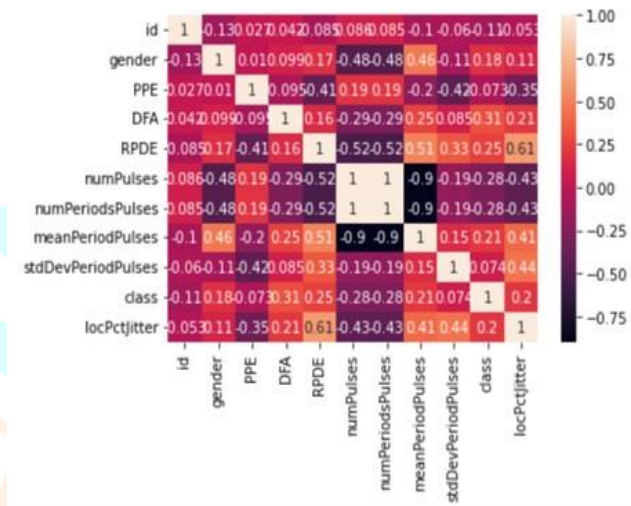


Fig. 2. Heatmap

After preprocessing the collected data, the best features are identified. These identified best features should be able to provide high efficiency. Testing and training can be done with the model now. The dataset splits into a train and test set in ratio of 7:3. KNN, Naive Bayes, and Logistic Regression algorithms were run on the training data set, and their performance in terms of accuracy was evaluated.The accuracy of these methods was **79.29% for Logistic regression, 81.06% for Naive Bayes, and 80.18% for KNN**. The Naive Bayes algorithm, which has the best accuracy, is employed to forecast the outcome.
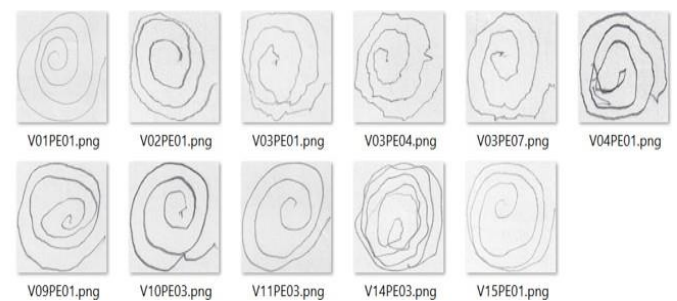
**DRAWING :**



Fig. 3. Spiral drawing dataset

The dataset itself consists of images, pre-split into a training set and a test set.

Train: 270 , Test: 60

Total the dataset consist of 330 images. The image has undergone pre-processing to enhance image quality, producing a picture that is superior to the original. Image acquisition's goal is to gather images that are less noisy than HD ones. Getting images with excellent clarity and less noise and distortion is the key benefit. The purpose of segmentation is to make the image easier to view or assess. The HOG descriptor is used to captures and quantifies local gradient changes in the input image. HOG can quantify how the direction of a spiral and wave changes. The resulting feature vector obtained is used to train the classifier.

The model is trained with a random forest and an XGBoost classifier. The one with greatest accuracy is considered for prediction. Accuracy achieved after model training

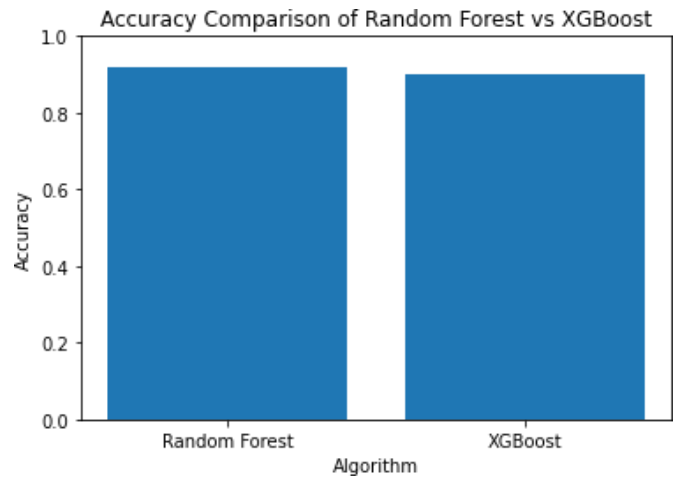**spiral: Random Forest: 93.33% and XGBoost: 91.67%**
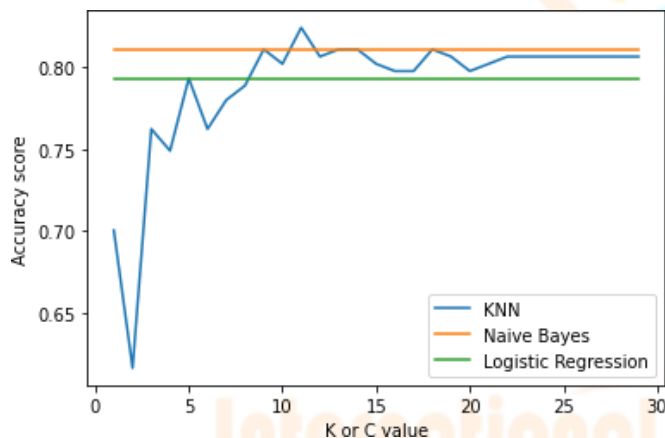
## V. RESULTS

**SPEECH :**



Fig. 4. Speech algorithm's accuracy

Naive Bayes has the highest accuracy, and this model is used on the front end. The model is loaded into a pickle file, and this file is opened in the frontend to compare the user input value with the corresponding model. Finally, it displays a text message indicating whether the patient has Parkinson's disease or not.

**DRAWING :**

Random forest has the highest accuracy for spirals and waves. As a result, it is applied to Parkinson's disease prediction.

## VI. CONCLUSIONS

Parkinson's disease is the second-deadliest neurodegenerative disease to date with no cure, so it is important to reduce the disease. The model is a combination of speech and drawing



Fig. 5. Drawing algorithm's accuracy

used for Parkinson disease prediction. For voice, the three algorithms employed are KNN, Naive Bayes, and Logistic Regression; for drawing, Random Forest, and XGBoost. The dataset is formed using these models, and we also compare the different models built using different methods to determine the most suitable one. The speech dataset consists of 700 features, of which only the best are selected, while the drawing dataset consists of 330 drawings of healthy and Parkinson's patients. From the above results, Nave Bayes stands out from the other two machine learning algorithms in speech, while Random Forest stands out in drawing, with 81% accuracy and 93.33% accuracy, respectively.In the future, the work can be extended by building hybrid models that can detect multiple diseases using accurate data sets that have common features between two diseases. The work can be extended to create a model that extracts more important features from all features in the dataset and improves accuracy.

## REFERENCES

[1] S. Shetty and Y. Rao, "Svm based machine learning approach to identify parkinson's disease using gait analysis," in *2016 International conference on inventive computation technologies (ICICT)*, vol. 2. IEEE, 2016, pp. 1–5.

[2] K. MANOJ, "Parkinson's disease prediction using machine learning techniques," 2021.

[3] M. Anila and G. Pradeepini, "Diagnosis of parkinson's disease using artificial neural network," *International Journal of Emerging Trends in Engineering Research*, vol. 8, no. 7, pp. 3700–3707, 2020.

[4] S. Sonu, V. Prakash, R. Ranjan, and K. Saritha, "Prediction of parkinson's disease using data mining," in *2017 international conference on energy, communication, data analytics and soft computing (ICECDS)*. IEEE, 2017, pp. 1082–1085.

[5] R. Anitha, T. Nandhini, S. Raj, and V. Nikitha, "Early detection of parkinson's disease using machine learning," *IEEE Access*, vol. 8, pp. 147 635–147 646, 2020.

[6] T. J. Wroge, Y. Özkanca, C. Demiroglu, D. Si, D. C. Atkins, and R. H. Ghomi, "Parkinson's disease diagnosis using machine learning and voice," in *2018 IEEE signal processing in medicine and biology symposium (SPMB)*. IEEE, 2018, pp. 1–7.

[7] S. Chakraborty, S. Aich, E. Han, J. Park, H.-C. Kim *et al.*, "Parkinson's disease detection from spiral and wave drawings using convolutional neural networks: a multistage classifier approach," in *2020 22nd Inter-*

*national Conference on Advanced Communication Technology (ICACT)*. IEEE, 2020, pp. 298–303.

[8]  R. Mathur, V. Pathak, and D. Bandil, "Parkinson disease prediction using machine learning algorithm," in *Emerging Trends in Expert Applications and Security: Proceedings of ICETEAS 2018*. Springer, 2019, pp. 357–363.

[9]  T. Sriram, M. V. Rao, G. Narayana, D. Kaladhar, and T. P. R. Vital, "Intelligent parkinson disease prediction using machine learning algorithms," *Int. J. Eng. Innov. Technol*, vol. 3, pp. 212–215, 2013.

[10] G. Abdurrahman and M. Sintawati, "Implementation of xgboost for classification of parkinson's disease," in *Journal of Physics: Conference Series*, vol. 1538, no. 1.  IOP Publishing, 2020, p. 012024.