A
Mini Project Report
On

# "Next-Generation Parkinson's Detection: Leveraging Machine Learning Innovations"

Submitted to JNTU Hyderabad

In Partial Fulfilment of the requirements for the Award of Degree of

**BACHELOR OF TECHNOLOGY**

**IN**

**INFORMATION TECHNOLOGY**

Submitted By

**V. BHANU PRAKASH**                                   **228R5A1211**

Under the Esteemed guidance of

**MR. K. ANIL KUMAR**

Assistant Professor, Department of IT

**Department of Information Technology**



# CMR ENGINEERING COLLEGE
# (UGC AUTONOMOUS)

(Accredited by NAAC & NBA, Approved by AICTE NEW DELHI, Affiliated to JNTU, Hyderabad)

(Kandlakoya, Medchal Road, R.R. Dist, Hyderabad-501 401)

**(2024-2025)**

# CMR ENGINEERING COLLEGE
# (UGC AUTONOMOUS)

(Accredited by NAAC & NBA, Approved by AICTE NEW DELHI, Affiliated to JNTU, Hyderabad)

(Kandlakoya, Medchal Road, R.R. Dist, Hyderabad-501 401)

# Department of Information Technology



# CERTIFICATE

This is to certify that the project entitled **"Next-Generation Parkinson's Detection: Leveraging Machine Learning Innovations"** is a bonafide work carried out by

**V. BHANU PRAKASH**                                        **228R5A1211**

In Partial fulfilment of the requirement for the award of the degree of **BACHELOR OF TECHNOLOGY** in **INFORMATION TECHNOLOGY** from CMR Engineering College, affiliated to JNTU, Hyderabad under our guidance and supervision.

The results presented in this project have been verified and are satisfactory. The results embodied in this project have not been submitted to any other university for the award of any other degree or diploma.

Internal Guide                                                                     Head of the Department

**Mr. K. ANIL KUMAR**                                                **Dr. MADHAVI PINGILI**

Assistant Professor                                                                 Professor & HOD

Department of IT                                                                       Department of IT

CMREC, Hyderabad                                                               CMREC, Hyderabad

# DECLARATION

This is to certify that the work reported in the present project entitled **"Next-Generation Parkinson's Detection: Leveraging Machine Learning Innovations"** is a record of bonafide work done by us in the Department of Information Technology, CMR Engineering College, JNTU Hyderabad. The reports are entirely based on the project work we did and are not copied from any other source. We submit our project for further development by any interested students who share similar interests to improve the project in the future.

The results embodied in this project report have not been submitted to any other University or Institute for the award of any degree or diploma to the best of our knowledge and belief.


**V. BHANU PRAKASH**     **228R5A1211**

# ACKNOWLEDGEMENT

We are extremely grateful to **Dr. A. Srinivasula Reddy,** Principal, and **Dr. Madhavi Pingili,** HOD, **Department of IT, CMR Engineering College** for their constant support.

We are extremely thankful to **Mr. K. ANIL KUMAR,** Assistant Professor, Internal Guide, Department of IT, for his constant guidance, encouragement, and moral support throughout the project.

We will be failing in duty if we do not acknowledge with grateful thanks to the authors of the references and other literature referred to in this project.

We express our thanks to all staff members and friends for all help and coordination extended in bringing out this project successfully in time.

Finally, we are very much thankful to our parents who guided us through every step.

**V. BHANU PRAKASH**          **228R5A1211**

# CONTENTS

# ABSTRACT

Parkinson's Disease (PD) is the second most common age-related neurological disorder that leads to a range of motor and cognitive symptoms. A PD diagnosis is difficult since its symptoms are quite similar to those of other disorders, such as normal aging and essential tremor. When people reach 50, visible symptoms such as difficulties walking and communicating emerge. Even though there is no cure for PD, certain medications can relieve some of the symptoms. Patients can maintain their lifestyles by controlling the complications caused by the disease. At this point, it is essential to detect this disease and prevent it from progressing. The diagnosis of the disease has been the subject of much research. In our project, we aim to detect PD using different types of Machine Learning (ML), and Deep Learning (DL) models such as Support Vector Machine (SVM), Random Forest (RF), Decision Tree (DT), K-Nearest Neighbor (KNN), and Multi-Layer Perceptron (MLP) to differentiate between healthy and PD patients by voice signal features. The dataset taken from the University of California at Irvine (UCI) machine learning repository consisted of 195 voice recordings of examinations carried out on 31 patients. Moreover, our models were trained using different techniques such as Synthetic Minority Over-sampling Technique (SMOTE), Feature Selection, and hyperparameter tuning (GridSearchCV) to enhance their performance. In the end, we found that MLP and SVM with a ratio of 70:30 train/test split using GridSearchCV with SMOTE gave the best results for our project. MLP performed with an overall accuracy of 98.31%, an overall recall of 98%, an overall precision of 100%, and an f1-score of 99%. In addition, SVM performed with an overall accuracy of 95%, an overall recall of 96%, an overall precision of 98%, and f1-score of 97%. The experimental results of this research imply that the proposed method can be used to reliably predict PD and can be easily incorporated into healthcare for diagnosis purposes.

**Keywords:** Parkinson's disease, machine learning, deep learning, diagnosis, differential diagnosis, GridSearchCV, SMOTE, feature selection.

# INTRODUCTION

## 1.1 Introduction

Parkinson's disease is a progressive neurodegenerative disorder that primarily affects movement, leading to symptoms such as tremors, stiffness, and impaired balance. Early diagnosis is crucial for managing the disease and improving the quality of life for those affected. Traditional diagnostic methods rely heavily on clinical observations and assessments, which can be subjective and often result in late-stage diagnosis. In recent years, the application of machine learning techniques in healthcare has opened new avenues for early and accurate disease detection. This mini-project focuses on utilising machine learning algorithms to detect Parkinson's disease from biomedical data. By analysing features such as vocal characteristics, handwriting dynamics, or other physiological signals, machine-learning models can identify patterns and anomalies that indicate the presence of Parkinson's disease. The goal of this project is to develop a machine learning model that can assist in the early detection of Parkinson's, offering a non-invasive, cost-effective, and efficient diagnostic tool. Through the exploration of various algorithms and data preprocessing techniques, this project aims to contribute to the growing body of research in the intersection of machine learning and healthcare, demonstrating the potential of AI-driven diagnostics in combating neurodegenerative diseases. Early diagnosis of Parkinson's disease is critical for effective treatment and management. However, the current diagnostic process is predominantly clinical, relying on the observation of motor symptoms and patient history. This method can be subjective, often leading to delayed diagnosis, as symptoms typically become noticeable only after significant neurodegeneration has occurred. Furthermore, the overlap of Parkinson's symptoms with other neurodegenerative disorders complicates the diagnosis, increasing the risk of misdiagnosis. In recent years, advancements in machine learning and data science have provided new opportunities to enhance the accuracy and timeliness of Parkinson's disease diagnosis. Machine learning, a subset of artificial intelligence, involves the development of algorithms that can learn from and make predictions based on data. By applying these techniques to biomedical data, such as voice recordings, handwriting samples, gait analysis, and neuroimaging, it is possible to identify subtle patterns and biomarkers that may not be easily detectable by human observation. This mini-project, titled "Parkinson's Detection Using Machine Learning," aims to explore the potential of machine learning models in diagnosing Parkinson's disease at an earlier stage, thereby improving patient outcomes. The project involves collecting and preprocessing relevant datasets, selecting appropriate features, and training various machine learning algorithms to detect Parkinson's disease with high accuracy. The focus will be on evaluating different models, such as Support Vector Machines (SVM), Random Forests, Neural Networks, and other classifiers, to determine the most effective approach for this application. One of the key challenges in this project is the handling of imbalanced datasets, as the prevalence of Parkinson's disease is relatively low compared to the general population. Techniques such as oversampling, undersampling, and the use of synthetic data generation (e.g., SMOTE) will be employed to address this issue. Additionally, feature selection and dimensionality reduction methods, such as Principal Component Analysis (PCA), will be utilised to enhance the model's performance and interpretability.

## 1.2 Project Objectives

The objective of this project is to develop a machine-learning model capable of accurately detecting Parkinson's disease at an early stage by analysing relevant biomedical data. By leveraging advanced algorithms and data preprocessing techniques, the project aims to create a non-invasive, efficient, and reliable diagnostic tool that can assist healthcare professionals in identifying Parkinson's disease sooner, thereby enabling timely intervention and better patient outcomes.

## 1.3 Purpose of the Project

The purpose of this project is to harness the power of machine learning to enhance the early detection of Parkinson's disease, providing a more accurate and objective diagnostic method. This approach aims to improve the timeliness of diagnosis, allowing for earlier treatment and potentially better management of the disease.

## 1.4 Existing System with Disadvantages

The existing system for diagnosing Parkinson's disease primarily relies on clinical assessments and the observation of motor symptoms by healthcare professionals. This method is often subjective and can lead to delayed diagnosis, as symptoms usually become evident only after significant neurological damage has occurred. Additionally, the overlap of Parkinson's symptoms with those of other neurodegenerative disorders increases the risk of misdiagnosis. The lack of objective, quantifiable biomarkers in traditional approaches further complicates early detection, potentially limiting treatment options and affecting patient outcomes.

**Disadvantages**
- Limited Early Detection
- Data Imbalance
- Feature Overlap
- Generalization Issues
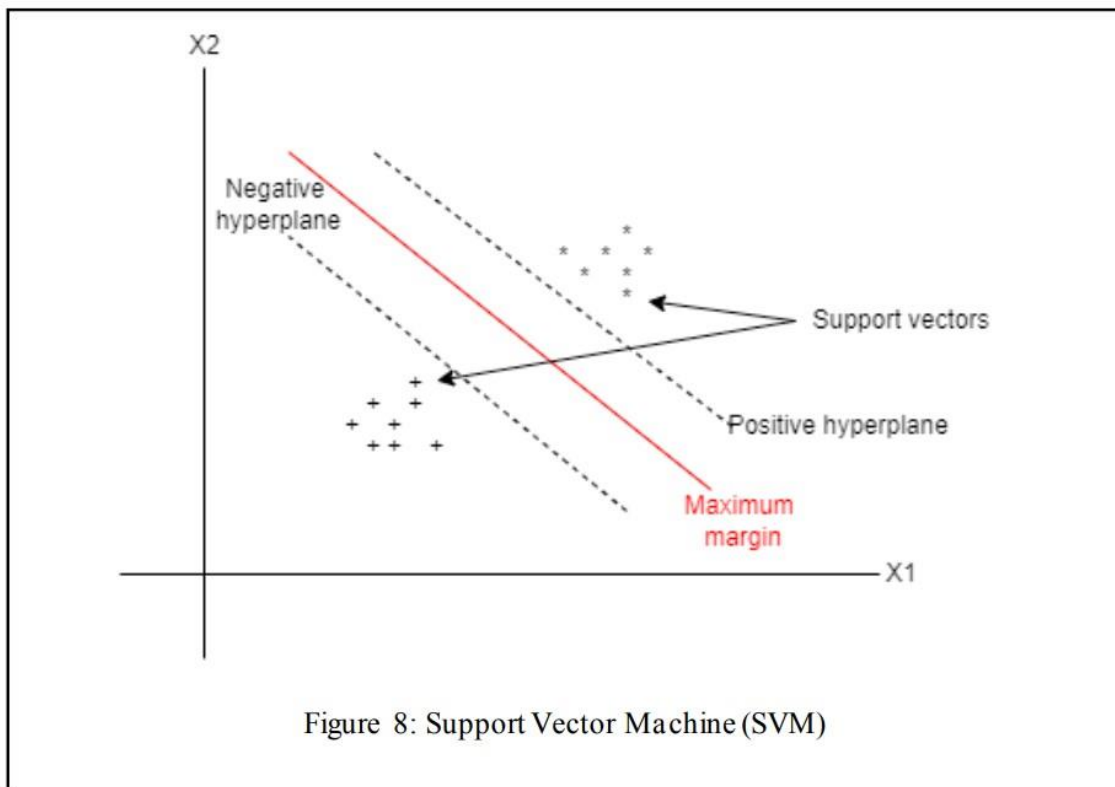- Limited Real-Time Monitoring
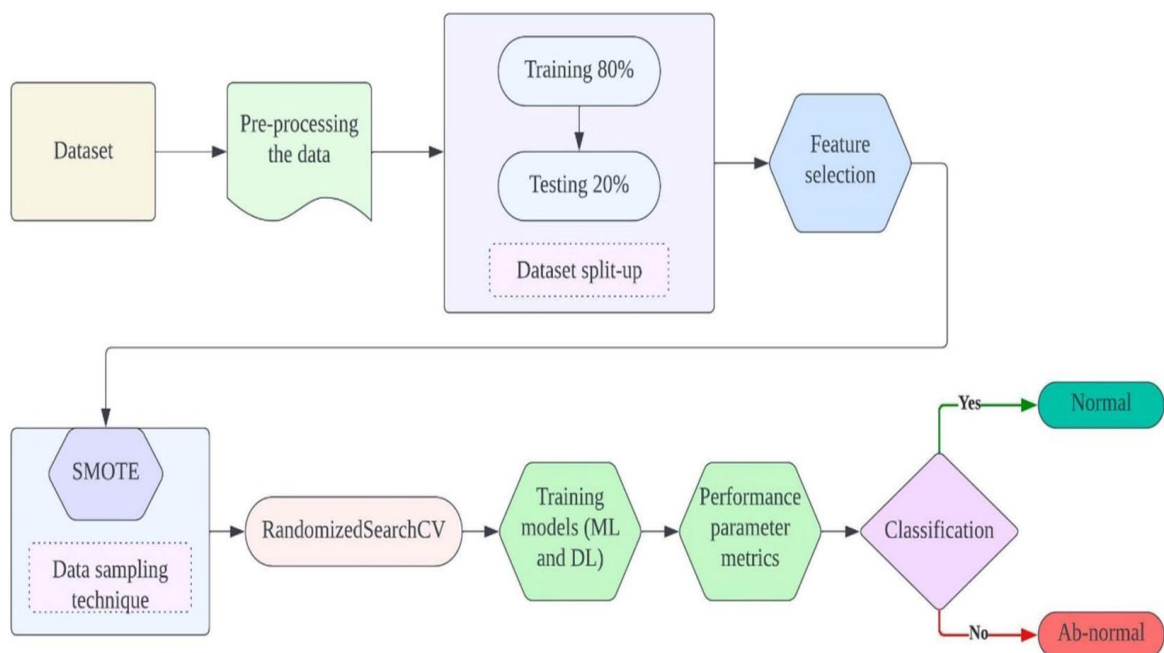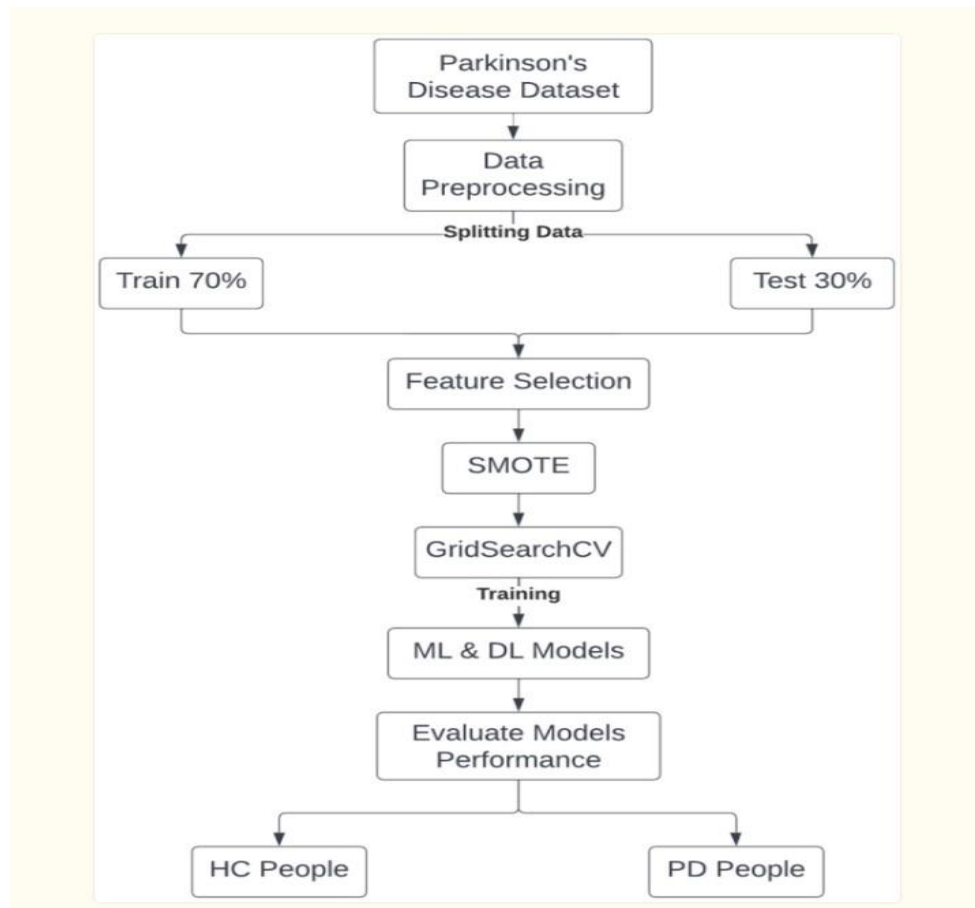
## 1.5 Proposed System with Features

The proposed method is designed to classify whether the patient has PD or not by using the Google Colab environment and Python language. The methodology of the proposed model is structured into six steps: data preprocessing, features selection, Synthetic Minority Over-sampling Technique (SMOTE), hyperparameter tuning (GridSearchCV), machine and deep learning classification models, and performance evaluation. It begins with data collection through multiple sources like wearable devices and clinical records, ensuring real-time acquisition and robust storage solutions. Data preprocessing involves cleaning, normalisation, and feature extraction to prepare for model training. Machine learning algorithms are selected and fine-tuned, with models evaluated for performance using

metrics such as accuracy and precision. For deployment, the system integrates real-time prediction capabilities with a user-friendly interface and a feedback loop for continuous improvement. Additionally, it includes performance monitoring, adherence to security and compliance standards, and thorough documentation to support ongoing maintenance and stakeholder engagement using **Support Vector Machine.**
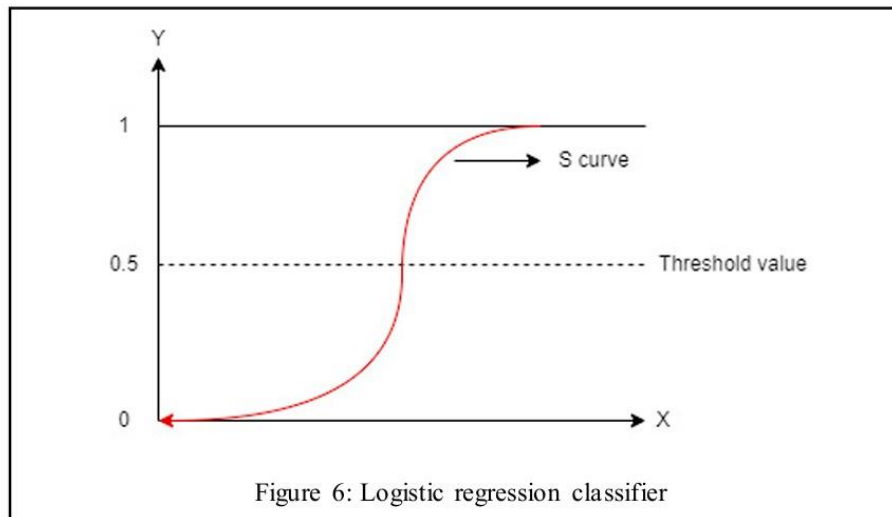
## Support vector Machine

- SVM is a supervised learning algorithm used for classification and regression tasks. It aims to find the optimal hyperplane that best separates the data into different classes in a high-dimensional space. The hyperplane is chosen to maximise the margin between the classes.
- SVM includes a regularisation parameter (C) that controls the trade-off between achieving a low error on the training data and minimising the margin of the hyperplane. A higher C value aims for a smaller margin with fewer classification errors, while a lower C value tries to achieve a larger margin even if it means more classification errors.
- Support vectors are the data points that lie closest to the decision boundary and are critical in defining the optimal hyperplane. They are the only data points used to define the hyperplane and margins, making them crucial for the model's performance and stability.



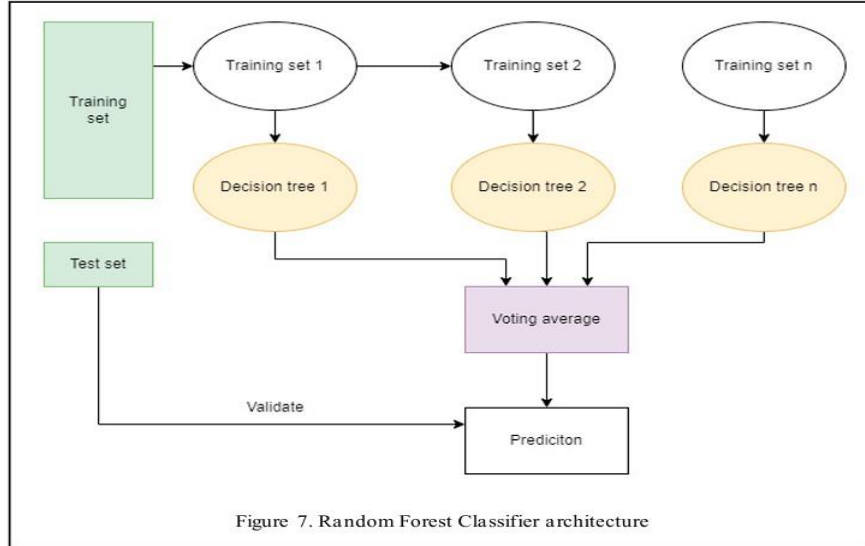Figure 8: Support Vector Machine (SVM)

## Logistic regression for classifier

Logistic regression is a prevalent supervised ML algorithm that predicts categorically dependent variables using a set of independent variables. It uses a curve fitting method to predict a probabilistic value in the range of 0 to 1, as the outcome of a categorical or discrete input. Compared to Linear regression where a line is fit to linearly predict one or more dependent variables, logistic regression predicts an S shaped logistic curve for values in range 0 to 1. This is ideal for audio data, as attributes affecting classification of PD are not linearly correlated, rather follow an exponential pattern.



Figure 6: Logistic regression classifier

## Random Forest classifier

Random Forest is a supervised machine learning algorithm that is applicable to classification and regression problems. This paper implements random forest classifiers to train a number of decision trees on subsets of the dataset and consider the average to increase the predictive accuracy of the results. It is a democratic model, where no single decision tree model is treated as superior, instead the majority vote of prediction from all models is considered to give an average prediction of the output. As no. of trees increases, the chances of overfitting decrease.

Figure 7. Random Forest Classifier architecture

## K Nearest Neighbour

K nearest neighbours (KNN) is a non-parametric, supervised machine learning algorithm that groups data into clusters based on underlying similarities. It works best for balanced audio data of 109 records due to small dataset size. Two clusters for PWP and healthy data are created in an efficient manner. KNN is a lazy learning algorithm, implying no presumptions of data are applied, ensuring novel patterns are learnt from training data.

## Model Evaluation

To identify the best model, we compare the results of 3 approaches and 9 models trained. For comparison, metrics chosen are ROC-AUC curve, confusion matrix, accuracy, precision, recall and F1 score. Formulae for these metrics are illustrated in equations 1-3, where TP stands for True Positives, FP for False Positives, TN for True Negatives and FN for False Negatives.
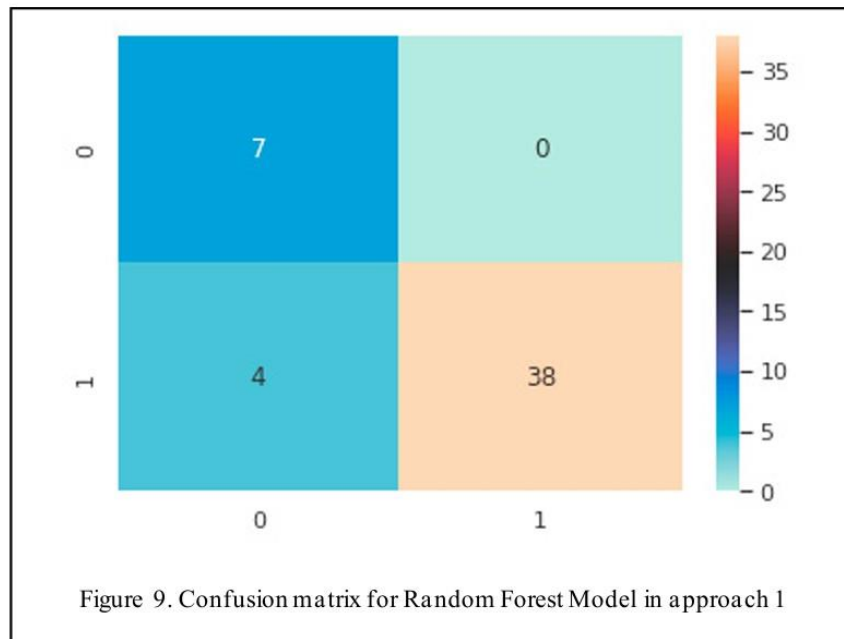
$$Precision = TP\overline{^{TP + FP}}$$

$$Recall = TP\overline{^{TP + FP}}$$

$$Accuracy = 2\,\overline{Precision^* Precision + Recall^* Recall}$$

Receiver Operating Characteristics or ROC curve is a probability curve and AUC measures area under this curve. It is a measure of separability or ability of the model to distinguish the classes. This metric measures the trade-off between clinical sensitivity and specificity for a set of tests, in question.

| Metric | Logistic Regression | Random Forest | SVM | KNN |
|---|---|---|---|---|
| Accuracy | 83.67% | 91.83% | 85.71% | 85.71% |
| Precision | 1.0 | 0.95 | 1.0 | 0.95 |
| Recall | 0.83 | 0.86 | 0.84 | 0.86 |
| ROC AUC curve | 0.636 | 0.701 | 0.682 | 0.701 |

Random Forest classifier is ideal for complete dataset, as it is an ensemble model. It evaluates the average of 100 decision trees, before the result is predicted. Every attribute is equally weighted during the classification process. The confusion matrix of this model has been illustrated in figure 9 below where, model classifies 7 true negatives (no PD), 4 false negatives, 38 true positives (PWP) and 0 false positives.



Figure 9. Confusion matrix for Random Forest Model in approach 1

Support Vector classifier with L1-support and linear kernel is ideal for PCA data, as it identifies ideal hyperplanes in less time and greater accuracy. Confusion matrix of this model has been depicted in figure 10 below, where SVM classifies 5 TN (no PD), 6 FN, 38 TP (PWP) and 0 FP.

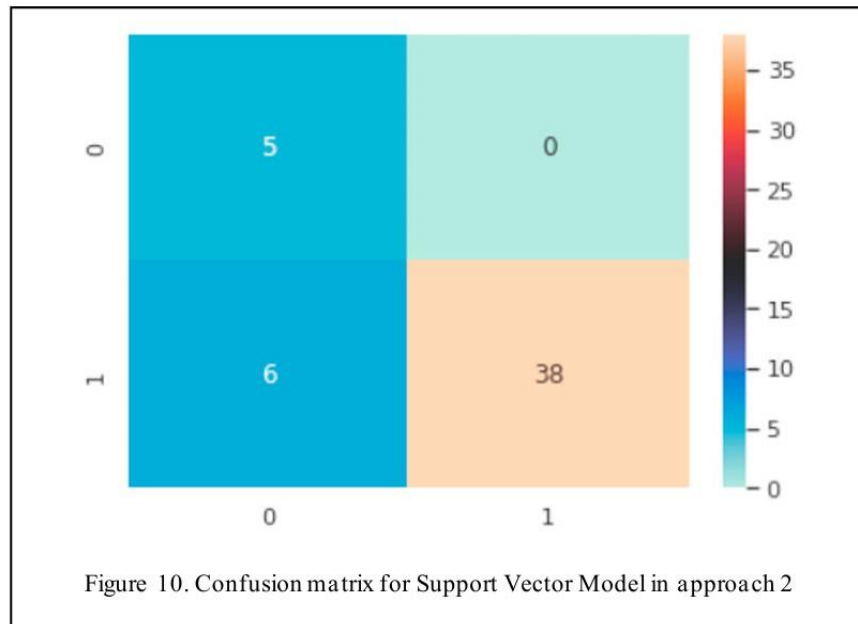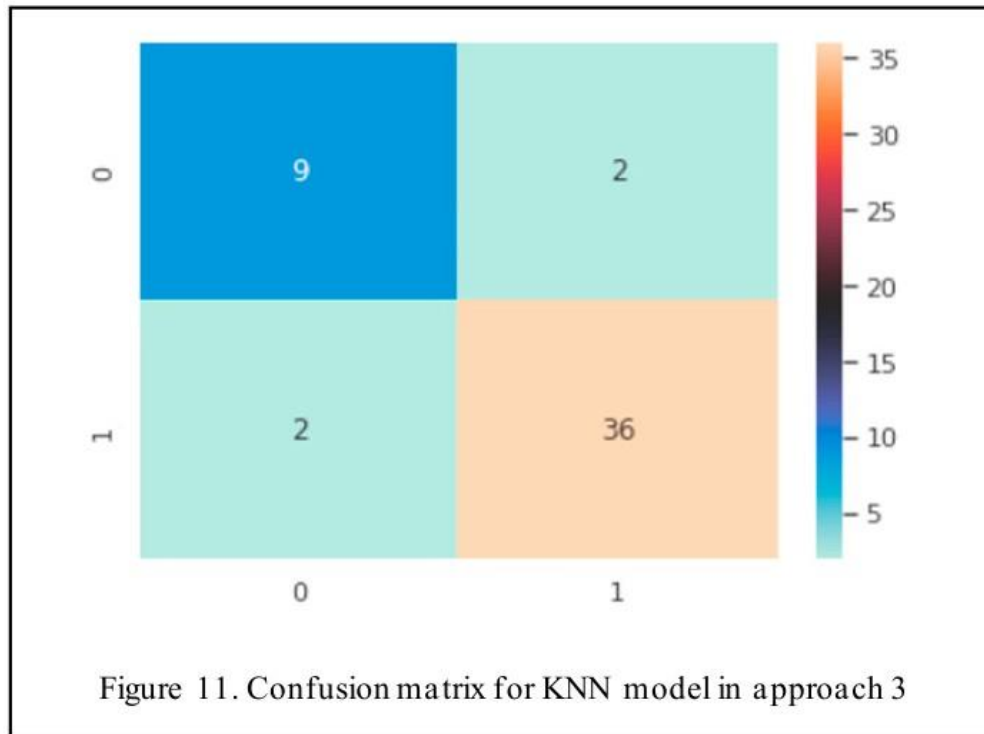Figure 10. Confusion matrix for Support Vector Model in approach 2

Table 4 below depicts the results for approach 3, using a balanced dataset. Models were trained on equal no. of records of normal and Parkinson's patients' data. Balancing ensures equal weightage is given to PWP and non Parkinson's patients. The results are as follows:

| Metric | Logistic Regression | Random Forest | SVM | KNN |
|---|---|---|---|---|
| Accuracy | 85.71% | 85.71% | 81.63% | 91.83% |
| Precision | 0.89 | 0.89 | 0.82 | 0.95 |
| Recall | 0.92 | 0.92 | 0.94 | 0.95 |
| ROC AUC curve | 0.811 | 0.811 | 0.817 | 0.883 |

K nearest neighbour's model performs best for balanced dataset, with highest precision and recall of 0.95. Due to equal distribution of data, identification of similarity in PWP and non-Parkinson's patients is faster. Classification results are illustrated in the confusion matrix in figure 11. KNN model classifies data into 9 TN (no PD), 2 FN, 36 TP (PWP) and 2 FP.

Figure 11. Confusion matrix for KNN model in approach 3

## Results and discussion

Parkinson's disease classification using vowel phonation data gives an 91.835% accuracy and 0.95 sensitivity for Random Forest classifier. Results of the Random Forest model are ideal, due to equal importance given to all 22 attributes in the MDVP dataset. This paper also highlights the results of the SVM model that gives an accuracy of 91.836% and sensitivity of 0.94, after PCA is applied to the dataset. Both SVM and Random Forest models perform well for outliers and are robust models. The models predict no false positives in the results. K nearest neighbour (KNN) model also performs well for balanced dataset, as classification into 2 categories without presumptions of data is favoured. Thus, we recommend the use of the Random Forest model to classify progress of the disease.

## 1.6 Input and Output Design

### Input Design

The system's input framework is meticulously structured to capture comprehensive patient data crucial for accurate detection of Parkinson's Disease. This includes demographic details such as age and gender, alongside detailed medical histories that encompass familial health patterns and prior diagnoses. Integral to the detection process are biometric inputs like speech recordings, which help identify vocal biomarkers through analysis of voice tremors and speech irregularities. Additionally, handwriting samples are collected to detect micrographia, a common symptom characterised by abnormally small handwriting.

9

Advanced sensor data capturing tremor patterns and gait abnormalities provide quantitative measures of motor function impairments. Coupled with standardised neurological assessment scores and lifestyle information covering diet, exercise, and sleep habits, the diverse dataset ensures a holistic evaluation of each patient's health status pertinent to Parkinson's Disease indicators.

## Objectives

- Input Design is the process of converting a user-oriented description of the input into a computer-based system. This design is important to avoid errors in the data input process and show the correct direction to the management for getting correct information from the computerised system.
- It is achieved by creating user-friendly screens for the data entry to handle large volumes of data. The goal of designing input is to make data entry easier and to be free from errors. The data entry screen is designed in such a way that all the data manipulations can be performed. It also provides record viewing facilities.
- When the data is entered it will check for its validity. Data can be entered with the help of screens. Appropriate messages are provided when needed so that the user will not be in maize instant. Thus, the objective of input design is to create an input layout that is easy to follow.
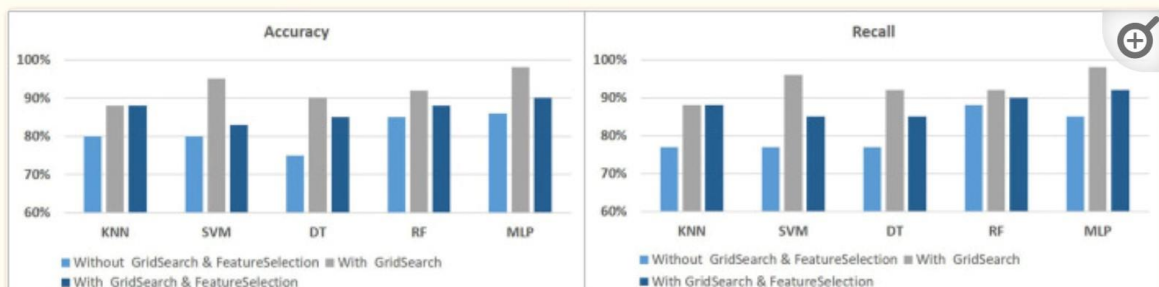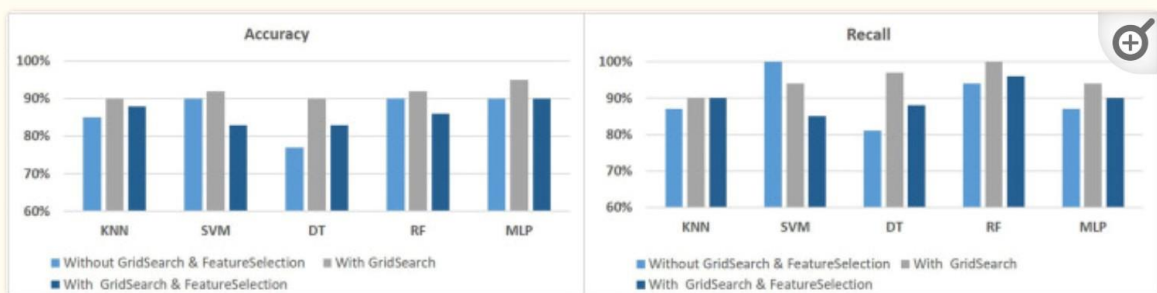
## Output Design

On the output side, the system provides a Parkinson's Disease risk score, which quantifies the likelihood of the disease's presence. This score is supplemented by a confidence level that reflects the model's certainty in its prediction. The output may also include a binary classification indicating whether the disease is present and, if so, the stage of progression (early, moderate, or advanced). Along with these predictive outputs, the system may offer recommendations for further diagnostic tests and possible treatment options. Visual aids, such as graphs and timelines, track the patient's symptom progression, enabling both patients and doctors to better understand the disease's course.

- **Diagnosis Probability:** A numerical score indicating the likelihood of Parkinson's Disease, ranging from 0 to 1 or as a percentage.
- **Classification:** Binary output (Yes/No) indicating whether the patient is likely to have Parkinson's Disease.
- **Recommendation System:** Suggestions for additional diagnostic tests if the model detects early signs of Parkinson's. General recommendations for treatment plans, such as medication, physical therapy, or lifestyle changes.
- **Visualisation:** Visual representations of key features like tremor patterns, gait analysis, and speech feature variations. A timeline or historical data visualisation showing the progression or stability of symptoms over time.

| Voice measure | Meaning |
| --- | --- |
| Name | ASCII name of subject and recording number (categorical variables). |
| MDVP:Fo(Hz) | Average vocal fundamental frequency (Numerical variables). |
| MDVP:Fhi(Hz) | Maximum vocal fundamental frequency (Numerical variables). |
| MDVP:Flo(Hz) | Minimum vocal fundamental frequency (Numerical variables). |
| MDVP:Jitter(%) | |
| MDVP:Jitter(Abs) | |
| MDVP: RAP | Several measures of variation in fundamental frequency (Numerical variables). |
| MDVP: PPQ | |
| Jitter:DDP | |
| MDVP:Shimmer | |
| MDVP:Shimmer(dB) | |
| Shimmer: APQ3 | Several measures of variation in amplitude (Numerical variables). |
| Shimmer: APQ5 | |
| MDVP: APQ | |
| Shimmer:DDA | |
| NHR | Measures of the ratio of noise to tonal components in |
| HNR | the voice (Numerical variables). |
| status | 0 for HC and 1 for PD (Numerical variables). |
| RPDE | Nonlinear dynamical complexity measures (Numerical variables). |
| D2 | |
| DFA | Signal fractal scaling exponent (Numerical variables). |
| spread1 | |
| spread2 | Nonlinear measures of fundamental frequency variation (Numerical variables). |
| PPE | |

Model's performance.

| | Without gridsearch and feature selection | With GridSearch | With GridSearch and feature selection |
|---|---|---|---|
| KNN | Accuracy = 80% | Accuracy = 88% | Accuracy = 88% |
| | Recall = 77% | Recall = 88% | Recall = 88% |
| | Precision = 97% | Precision = 98% | Precision = 98% |
| | F1-Score = 86% | F1-Score = 92% | F1-Score = 92% |
| SVM | Accuracy = 80% | Accuracy = 95% | Accuracy = 83% |
| | Recall = 77% | Recall = 96% | Recall = 85% |
| | Precision = 97% | Precision = 98% | Precision = 93% |
| | F1-Score = 86% | F1-Score = 97% | F1-Score = 89% |
| DT | Accuracy = 75% | Accuracy = 90% | Accuracy = 85% |
| | Recall = 77% | Recall = 92% | Recall = 85% |
| | Precision = 90% | Precision = 96% | Precision = 95% |
| | F1-Score = 83% | F1-Score = 94% | F1-Score = 90% |
| RF | Accuracy = 85% | Accuracy = 92% | Accuracy = 88% |
| | Recall = 88% | Recall = 92% | Recall = 90% |
| | Precision = 93% | Precision = 98% | Precision = 96% |
| | F1-Score = 90% | F1-Score = 95% | F1-Score = 92% |
| MLP | Accuracy = 86% | Accuracy = 98% | Accuracy = 90% |
| | Recall = 85% | Recall = 98% | Recall = 92% |
| | Precision = 98% | Precision = 100% | Precision = 96% |
| | F1-Score = 91% | F1-Score = 99% | F1-Score = 94% |

# 2. LITERATURE SURVEY

**Analysis of Tremors in Parkinson's Disease Using Accelerometer:** Niya Romy Markose, Priscilla Dinkar Moyya, and Mythili Asaithambi, Analysis of tremors in Parkinson's Disease using accelerometer, IEEE, July 2021. This prototype was designed to observe and quantify the tremor signal from Parkinson's disease patients. The prototype is based on Arduino Uno programming and interfacing, and the ADXL335 tri-axial accelerometer is used as a sensor. The resting tremor signal was acquired in the form of acceleration using the sensor accelerometer from the fingertip, wrist, and forearm of the patient. The Arduino processed the data which was transferred to MATLAB for further processing. The resting tremor was observed in terms of amplitude and spectral density. For the three parts considered, the amplitude values of acceleration ranged from 40 dB/Hz to 80 dB/Hz, and spectral density was observed and compared. Hence this basic prototype could be useful and developed further to assist Parkinson's Disease patients

**Building a Machine Learning Framework to Remotely Assess Parkinson's Disease Using Smartphones:** Oliver Y. Chen, Florian Lipsmeier, Huy Phan, and John Prince, IEEE, May 2020. Using this framework, we map the PD-specific architecture of behaviours using data obtained from both PD participants and healthy controls (HCs). Utilising these atlases of features, the framework shows promises to (a) discriminate PD participants from HCs, and (b) estimate the disease severity of individuals with PD.

**Deep Convolutional Neural Network for Parkinson's Disease Based on Handwriting Screening:** Mohamed Shaban, IEEE, Sep 2020. In this paper, the use of a fine-tuned VGG-19 for screening Parkinson's Disease (PD) based on a Kaggle handwriting dataset is investigated and experimented with. The dataset including 102 wave and 102 spiral handwriting patterns was pre-processed where images were resized and a data augmentation based on image rotation was adopted to minimise overfitting. The Convolutional Neural Network (CNN) model was then trained on the pre-processed dataset and validated using both 4-fold and 10-fold cross-validation techniques. The CNN model achieved an accuracy of 88%, 89%, and a sensitivity of 89%, and 87% on the wave and spiral patterns respectively when a 10-fold cross-validation was used.

13

**Detection of Parkinson's Disease through Smell Signatures:** <u>Shrinidhi Kulkarni, Neenu George Kalayil, Jinu James, Sneha Parsewar and Revati Shriram</u>**,** <u>IEEE, Sep 2020.</u> The intensity of these symptoms differs from person to person. Amongst these two types of symptoms, non-motor symptoms are identifiable at an early stage. Hence detection of these symptoms helps in recognizing whether a person has Parkinson's Disease at an early stage. Patients diagnosed with Parkinson's Disease give out a distinguishable musky smell.

**Effectiveness Analysis of Bio-electric Stimulation Therapy to Parkinson's Disease via Discrete Fourier Transform Approach:** <u>Yuxin Lin, Bingo Wing-Kuen Ling, Nuo Xu, Ringo Lam, and Charlotte Ho, IEEE, May 2021:</u> This paper studies the effectiveness of applying the bio-electric stimulation therapy to the patients suffering from the Parkinson's diseases via the resting tremor signals. Eight patients suffering from Parkinson's disease are invited to wear the wearable devices before and after applying the bio-electric stimulation therapy to obtain the resting tremor signals. Then, the resting tremor signals are classified into two groups, namely before and after applying the bio-electric stimulation therapy. To perform the classification, first, the signals are represented in the frequency domain via performing the fast Fourier transforms.

**Efficient Pre-diagnosis Approach for Parkinson's Disease with Machine Learning:** <u>Yuqi Qiu, IEEE, Nov 2020.</u> This research proposes deep learning methods to approach the diagnosis of Parkinson's Disease (PD) with more accurate results than those of existing methods. This automation approach also reduces the diagnosis expenses and makes it possible to be applied in developing countries. To find a more accurate and cheaper solution to the problem, I have employed and contrast three different experimental algorithms-K Means clustering, Support Vector Machine (SVM), and Convolutional Neural Network (CNN). The best performance, with an accuracy of 100%, can be achieved by the SVM model and convolutional neural network model.

**Feature Selection-Based Twin-Support Vector Machine for the Diagnosis of Parkinson's Disease:** <u>Surendra Bikram, Thapa, Surabhi Adhikari, Awishkar Ghimire and Anshuman Aditya, IEEE, June 2021.</u> With the growing ageing population, Parkinson's disease has become a serious problem for a huge fraction of people above 60. The disease severely affects the motor system and can lead to the death of the patients. There is no cure available for the disease. The symptoms in the motor system are seen very late which leads to difficulty in the management of the

disease. There is no cure for the disease which makes it more difficult when the disease is diagnosed later.

**Prediction of Parkinson's disease and severity of the disease using Machine Learning and Deep Learning algorithm:** Pooja Raundale, Chetan Thosar and Shardul Rane, IEEE, Aug 2021. People have trouble vocally, writing, strolling, or completing other simple tasks when dopamine-generating neurons in parts of the brain become impaired or expire. These symptoms worsen over time, increasing the severity of the condition in patients. We have suggested a methodology in this article for the prediction of Parkinson's disease severity using deep neural networks on UCI's Parkinson's Telemonitoring Vocal Data Set of patients. We have created a neural network to predict the severity of the disease and a machine learning model to detect the disorder. Classification of Parkinson's Disease is done by a Neural network, a Random Forest Classifier.

**Parkinson's Disease Detection from Spiral and Wave Drawings Using Convolutional Neural Networks A Multistage Classifier Approach:** Sabyasachi Chakraborty, Satyabrata Aich, Jong-Seong-Sim, Eunyoung Han, Jinse Park and Hee-Cheol Kim, IEEE, June 2020. In this paper, a system design is proposed for analysing Spiral drawing patterns and wave drawing patterns in patients suffering from Parkinson's disease and healthy subjects. The system developed in the study leverages two different convolutional neural networks (CNN), for analysing the drawing patterns of both spiral and wave sketches respectively. Further, the prediction probabilities are trained on a metal classifier based on ensemble voting to provide a weighted prediction from both the spiral and wave sketch. The complete model was trained on the data of 55 patients and has achieved an overall accuracy of 93.3%, average recall of 94 $^{\%}$, average precision of 93.5%, and average f1 score of 93.94%

**A Speech-Based System for Parkinson's Disease Analysis and Monitoring:** Daniel Palacios-Alonso, Guillermo Melendez-Morales, Agustin Lopez-Arribas, Carlos Lazaro-Carrascosa, Andres Gomez-Rodellar and Pedro Gomez-Vilda, MonParLoc, IEEE, Oct 2020. This study aims to describe an added-value solution considering the cooperation of both previously mentioned methods: speech analysis-based monitoring called within the project, Monitoring Parkinson using Locution (MonParLoc), and acoustical neurostimulation, called within project neuro-Acoustic-stimulation Parkinson (AcousticPar). The applications designed in both projects are embedded into a global solution denominated Teca-Park which consists of four main activities: speech evaluation, neurostimulation, motor symptom longitudinally, and questionnaires. This framework is conceived to be a powerful tool

for treating and monitoring longitudinally remotely and contact-free. MonParLoc was tested and validated in real scenarios involving patient associations. Validation results produced in these associations demonstrating the utility of this approach are given in the study, particularly about protocol vulnerability and robustness. This paper proposes a complete framework (a mobile app and a scorecard solution) including different services for Parkinson's clinical monitoring and patient management using speech, movement, and acoustic stimulation.

# 3. SOFTWARE REQUIREMENTS ANALYSIS

## 3.1 Problem Statement

This project aims to address the spread of misleading information online poses a serious challenge, compromising public trust and perpetuating false narratives. Conventional detection methods often fall short due to the rapid evolution and vast scale of digital content. Advanced, scalable solutions are required to accurately classify and manage misinformation in real time. Addressing this issue involves deploying sophisticated algorithms and large-scale data processing to enhance accuracy and effectiveness. Innovative approaches are crucial for effectively countering the widespread dissemination of misleading content.

## 3.2 Modules and their Functionalities

## Data preprocessing

Data Preprocessing is an essential phase in preparing data for machine learning models, ensuring that it is clean, relevant, and well-structured for effective training and evaluation. The Process begins with collecting the dataset and importing the necessary libraries for data manipulation. The raw data is then processed through several stages: cleaning, where redundant attributes and missing values are addressed to improve accuracy, splitting, which divides the dataset into training and testing datasets to evaluate model performance, and finally training and testing to assess its performance on unseen data. The final step involves evaluating the classifier using accuracy, precision, and recall metrics to ensure it meets performance standards and generalises well to new data.

- **Data Collection**: Gather raw data from various sources.
- **Import Libraries**: Load necessary libraries for data processing and machine learning.
- **Data Cleaning**: Remove redundant attributes and handle missing values.
- **Data Splitting**: Divide the dataset into training and testing datasets.
- **Model Training and Testing**: Train the classifier with the training data and evaluate it using the test data.
- **Model Evaluation**: Assess model performance with metrics such as accuracy and precision.

### 3.3 Functional Requirements

1. **Data Collection and Preprocessing:**
   - **Input Data Acquisition**: The system should be able to accept input data such as voice recordings, gait analysis, handwriting samples, or medical history.
   - **Data Cleaning**: The system must clean the data to remove noise, handle missing values, and normalise/standardise the input features.

2. **Model Training:**
   - **Model Selection**: The system should allow for the selection of various machine learning models (e.g., SVM, Random Forest, Neural Networks) for training.
   - **Training Process**: The system should train the chosen model using the processed data, optimising for accuracy and other relevant metrics.

3. **Model Evaluation:**
   - **Performance Metrics**: The system must calculate and display performance metrics such as accuracy, precision, recall, F1-score, and ROC-AUC curve.
   - **Confusion Matrix**: The system should generate a confusion matrix to visualise true positives, false positives, true negatives, and false negatives.

4. **Prediction:**
   - **Real-Time Prediction**: The system should be capable of making real-time predictions based on new input data.
   - **Batch Prediction**: The system should support batch processing to predict Parkinson's disease for multiple input samples at once.

5. **User Interface:**
   - **Input Interface**: A user-friendly interface for inputting patient data.
   - **Visualisation**: The system should provide visualisations of the input data, feature importance, and model predictions.
   - **Results Display**: The system should display the diagnosis result (e.g., likelihood of Parkinson's disease) along with confidence levels.

6. **Integration:**
   - **API for External Systems**: The system should expose an API for integration with other healthcare systems or data sources
   - **Database Connectivity**: The system should connect to a database for storing patient data and model outputs.

7.  **Security and Privacy:**
    ● **Data Encryption**: Ensure that all patient data is encrypted during transmission and storage.
    ● **Access Control**: Implement user authentication and authorization to control access to sensitive data and functionalities.
    ● **Anonymization**: The system should anonymize patient data to comply with healthcare regulations (e.g., HIPAA).
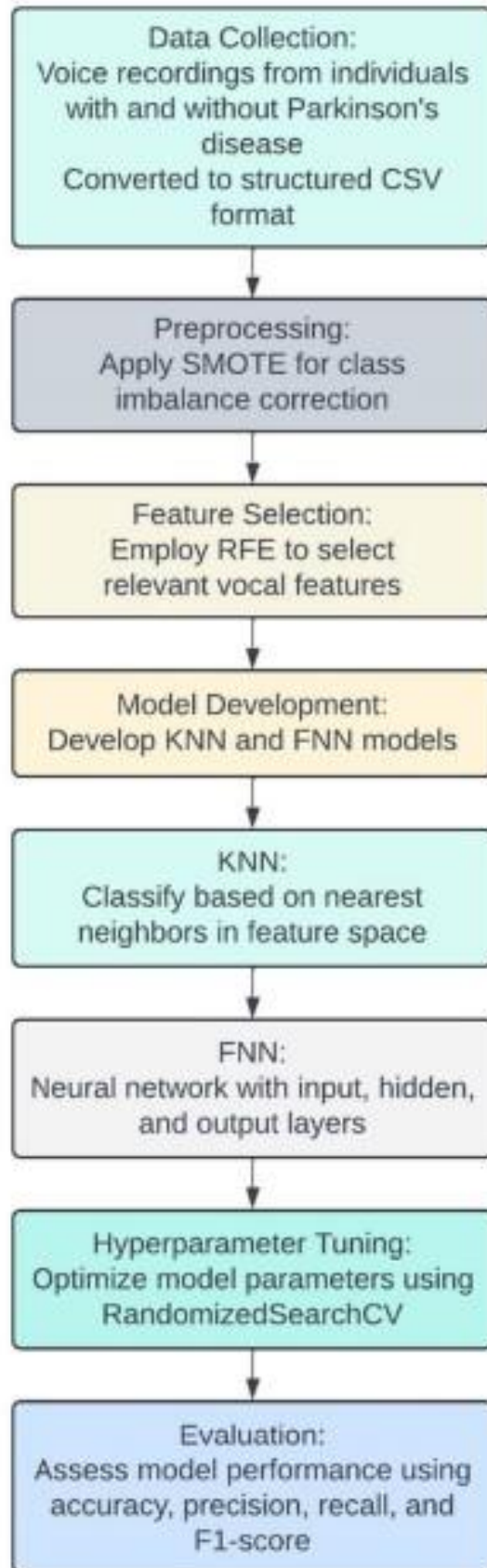
**Figure 2.** Operational flow diagram of proposed model.

## 3.4 Non-Functional Requirements

The Online Content Credibility Evaluation System must meet several key non-functional requirements to ensure it operates effectively and efficiently. It should deliver high performance, processing data quickly to provide timely results. Scalability is essential to accommodate increasing data and user demands without performance degradation. Robust security measures are required to protect user data through encryption and secure authentication. Additionally, the system must be reliable, ensuring consistent operation with minimal downtime and errors

This is done through 5 steps:

- **Performance:** Ensure rapid data processing for timely results.
- **Scalability:** Support increasing data volumes and user activity.
- **Security:** Implement robust encryption and secure authentication.
- **Usability:** Provide an intuitive and user-friendly interface.
- **Reliability:** Guarantee consistent operation with minimal downtime.

## 3.5 Feasibility Study

The Online Content Credibility Evaluation System confirms that the project is technically and operationally viable with manageable costs. Compliance with legal standards for data protection and privacy ensures the system can be implemented effectively.

- Technical Feasibility
- Operational Feasibility
- Economical Feasibility
- Legal Feasibility

## Technical Feasibility

Sufficient datasets (e.g., UCI Machine Learning Repository, Parkinson's Progression Markers Initiative) are available, containing patient data related to Parkinson's disease.

## Operational Feasibility

The system should integrate smoothly with existing healthcare infrastructure, such as Electronic Health Record (EHR) systems. This might involve API development and ensuring interoperability.

## Economical Feasibility

The system can provide economic benefits by improving early diagnosis, reducing the need for more expensive treatments later, and potentially reducing the burden on healthcare systems.

# 4. SOFTWARE AND HARDWARE REQUIREMENTS

## 4.1 Software Requirements

The functional requirements or the overall description documents include the product perspective and features, operating system and operating environment, graphics requirements, design constraints, and user documentation.

The appropriation of requirements and implementation constraints gives the general overview of the project regarding what the areas of strength and deficit are and how to tackle them.

Operating System        :        Windows 10

Coding Language        :        Python

Tool                :        Visual Studio Code

Server                :        Streamlit

## 4.2 Hardware Requirements

Minimum hardware requirements are very dependent on the particular software being developed by a given Enthought Python / Canopy / VS Code user. Applications that need to store large arrays/objects in memory will require more RAM, whereas applications that need to perform numerous calculations or tasks more quickly will require a faster processor.

System        :        Intel Core

Hard Disk    :        120 GB

Monitor        :        15" LED

Input devices :        Keyboard, Mouse

Ram            :         4 GB

# 5. SOFTWARE DESIGN
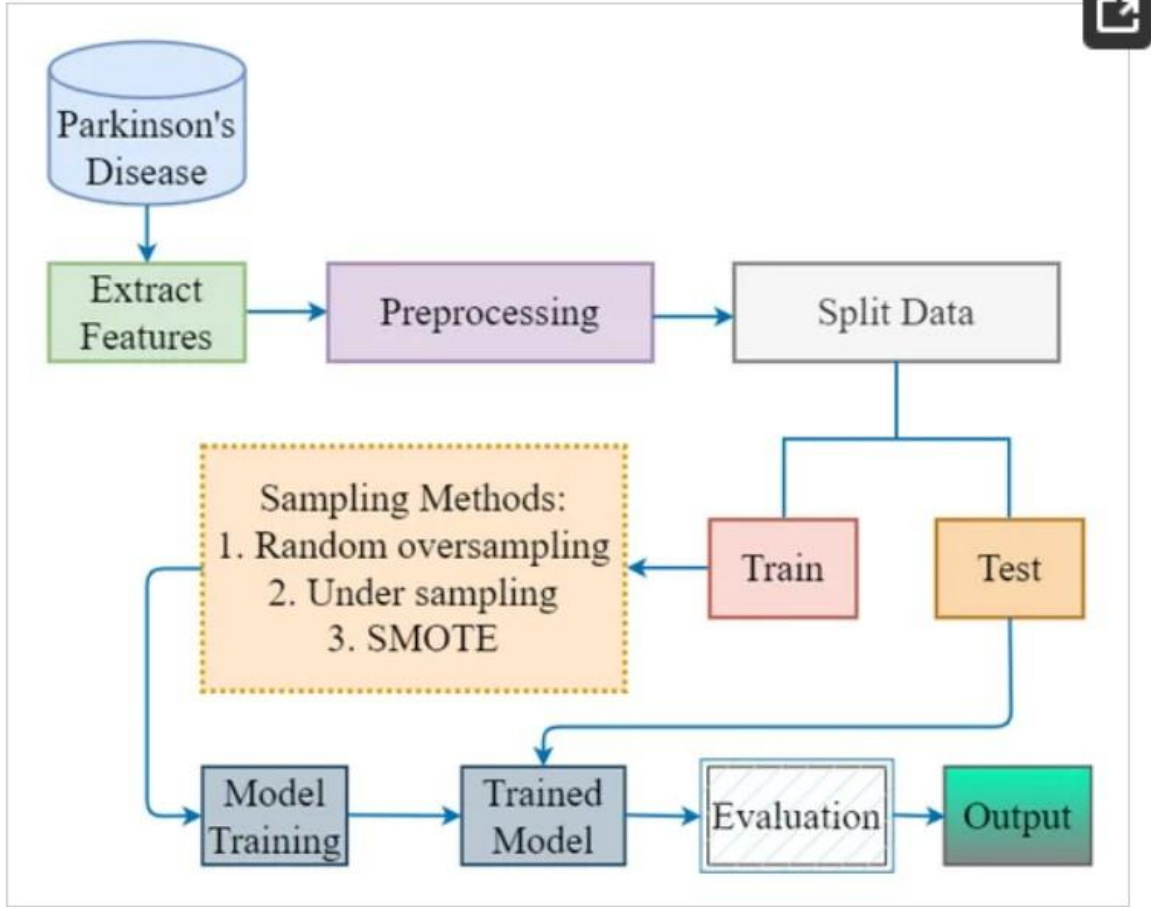
## 5.1 System Architecture



**Figure 1.** Work Flow of Proposed Methodology.

The system architecture for the "Parkinson's Disease Detection using Machine Learning" project is designed to efficiently process and analyse patient data to detect signs of Parkinson's Disease. The architecture begins with **Data Collection**, where wearable devices or clinical tools gather sensor data, such as motion or voice patterns, from patients. This data is then transmitted to a central database in the **Data Acquisition** phase.

In the **Data Preprocessing** stage, the collected data undergoes cleaning, feature extraction, and augmentation to ensure it is suitable for analysis. The preprocessed data is then used in the **Model Development** phase, where machine learning models are trained, validated, and selected based on their performance in detecting Parkinson 's-related features.

Once the model is trained, it is deployed in the **Model Deployment** stage, where it is integrated into an API layer and connected to a user interface for real-time predictions. The **Prediction & Monitoring** phase utilises the model to continuously monitor patient data.
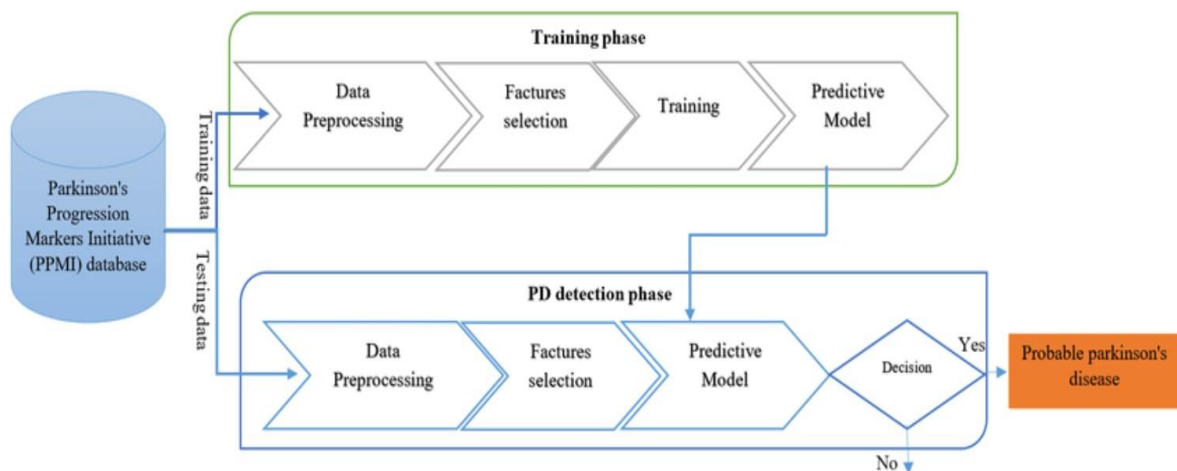
## 5.2 Data Flow Diagram

A data flow diagram (DFD) is a graphical representation of the flow of data through a system. It provides a visual overview of how data is input, processed, and output. In the context of machine learning projects, DFDs are particularly useful for understanding the entire workflow, from data acquisition to model deployment.

By breaking down the project into logical components and illustrating how data moves between them, DFDs help in:

- **Identifying bottlenecks:** Visualising the data flow can highlight areas where data processing might be slow or inefficient.
- **Understanding dependencies:** DFDs show how different components rely on each other, making it easier to identify potential risks or issues.
- **Communicating the process:** DFDs can be used to communicate the project's workflow to stakeholders, team members, or clients.
- **Ensuring data integrity:** By tracking the flow of data, DFDs can help maintain data quality and prevent errors.

For a machine learning project, a typical DFD might include the following components:

- **Data acquisition:** This involves collecting raw data from various sources.
- **Data preprocessing:** This stage involves cleaning, transforming, and preparing the data for modelling.
- **Feature engineering:** Here, relevant features are extracted or created to improve model performance.
- **Model training:** The training data is used to build and train the machine learning model.
- **Model evaluation:** The model's performance is assessed using appropriate metrics.
- **Model deployment:** The trained model is deployed into a production environment for real-world use

## 5.3 UML Diagrams

UML is a standard language for specifying, visualising, constructing, and documenting the artifacts of software systems. UML was created by the Object Management Group (OMG) and the UML 1.0 specification draft was proposed to the OMG in January 1997.

There are several types of UML diagrams, each of which serves a different purpose regardless of whether it is being designed before the implementation or after (as part of documentation).

UML is directly related to object-oriented analysis and design. After some standardisation, UML has become an OMG standard. The two broadest categories that encompass all other types are:

- Behavioural UML diagram and
- Structural UML diagram.

As the name suggests, some UML diagrams try to analyse and depict the structure of a system or process, whereas others describe the behaviour of the system, its actors, and its building components.

**Goals**: The Primary goals in the design of the UML are as follows:

- Provide users with a ready-to-use, expressive visual modelling Language to develop and exchange meaningful models.

- Provide extensibility and specialisation mechanisms to extend the core concepts.

- Be independent of particular programming languages and development processes.

- Provide a formal basis for understanding the modelling language.

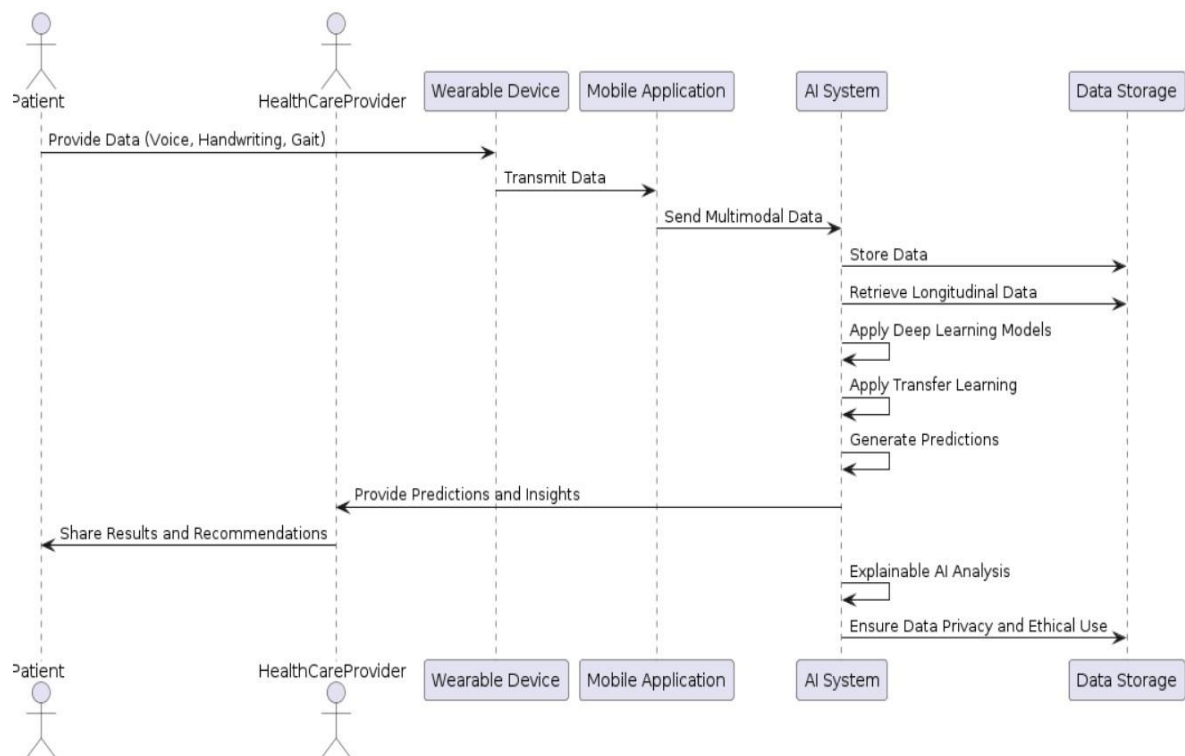- Encourage the growth of the tools market.

- Support higher-level development concepts such as collaborations, frameworks, patterns, and components.

- Integrate best practices.

**The different types are as follows:**

- Sequence diagram

- Use case Diagram

- Activity diagram

- Class diagram

- Collaboration diagram

## Sequence Diagram

A sequence diagram simply depicts interaction between objects in a sequential order i.e., the order in which these interactions take place. We can also use the terms event diagrams or event scenarios to refer to a sequence diagram. Sequence diagrams describe how and in what order the objects in a system function. These diagrams are widely used by businessmen and software developers to document and understand requirements for new and existing systems.

**List of Actions**

**User:**
The user needs to press any of the given three (i.e., prediction data, prediction
Skills) then he will get the output accordingly

**System:**
The system will give output as he enters according to the given data.

**Results:**

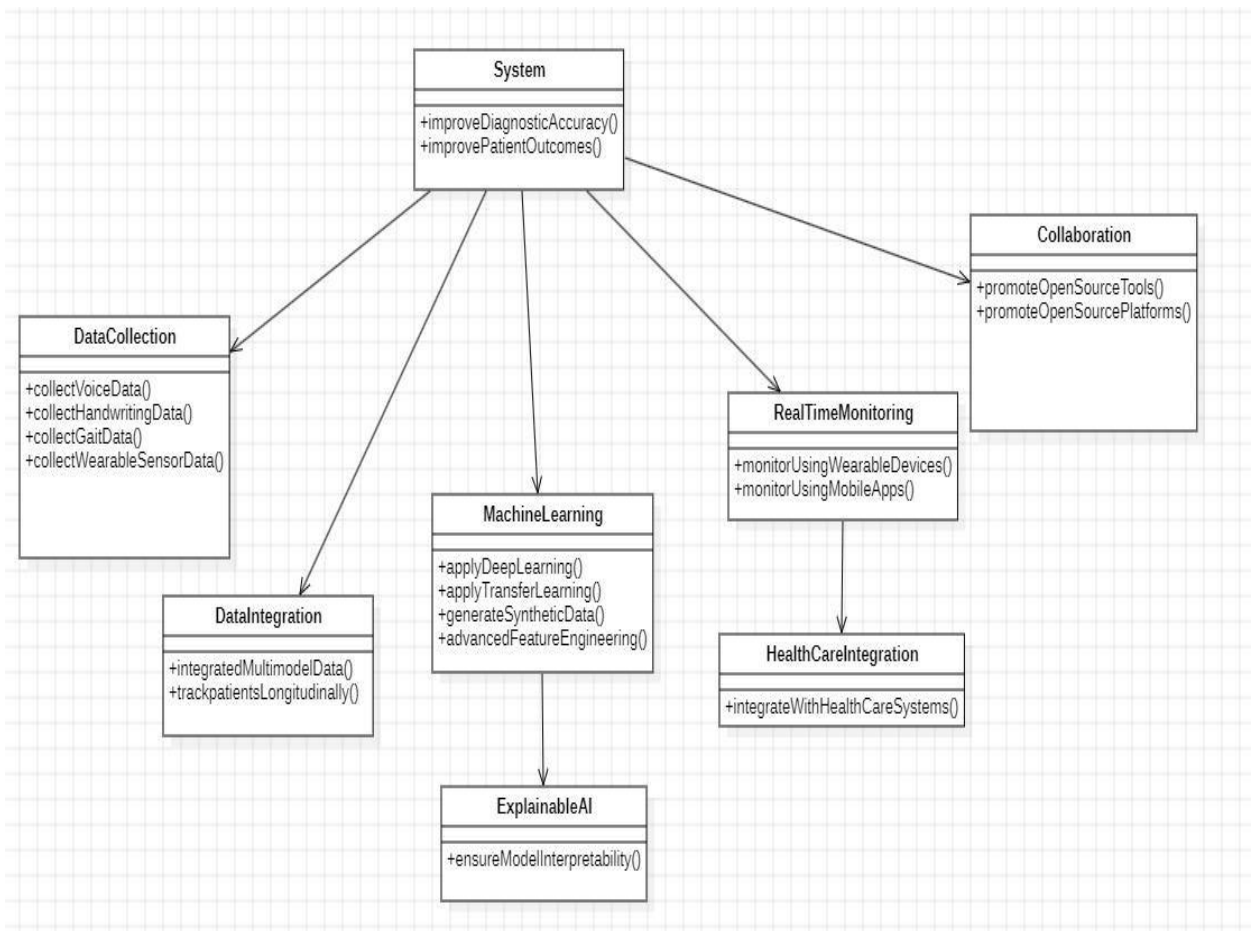As the user enters the data it will give whether the patient is positive or negative.

## Use Case Diagram

A use case diagram at its simplest is a representation of a user's interaction with
the system that shows the relationship between the user and the different use
cases in which the user is involved. A use case diagram is used to structure the
behaviour thing in a model. The use cases are represented by either circles or
ellipses.

## Class Diagram

In software engineering, a class diagram in the Unified Modelling Language (UML) is a type of static structure diagram that describes the structure of a system by showing the system's classes, attributes, operations (or methods), and the relationships among the classes. It explains which class contains information.
.



**Potential classes and their relationships in a Parkinson's disease detection system might include:**
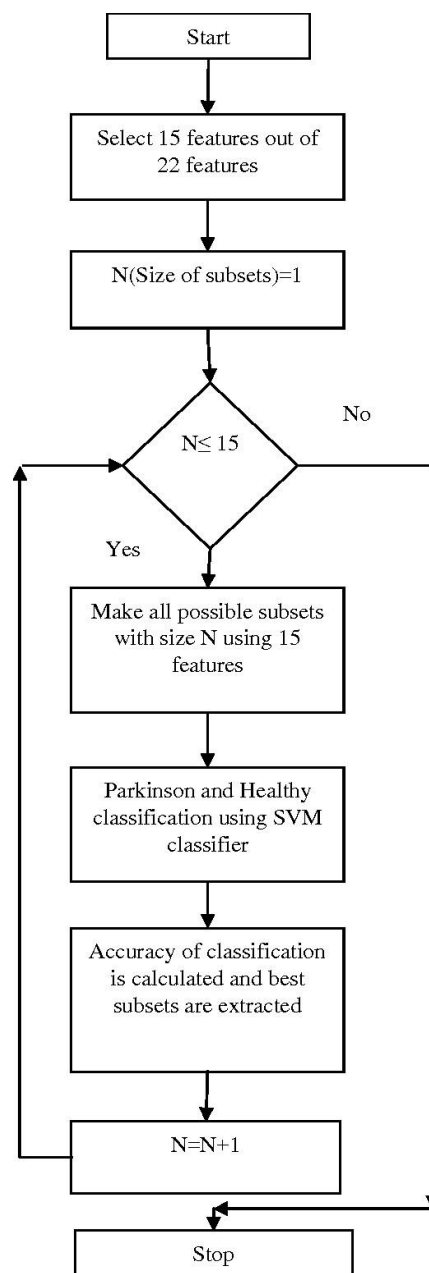
- **Patient:** Represents a patient with or without Parkinson's disease.
  - Attributes: Patient ID, name, age, gender, medical history, etc.
    - Operations: Get patient data, update patient information, etc.
- **DataCollector:** This represents a component responsible for collecting patient data, such as voice recordings, gait analysis, or medical records.
  - Operations: Collect voice data, collect gait data, collect medical history, etc.
- **DataPreprocessor:** Represents a component responsible for cleaning, normalising, and transforming the collected data.

○        Operations: Remove noise, handle missing values, extract features, etc.

## Activity Diagram

An activity diagram visually presents a series of actions or flow of control in a system similar to a flowchart or a data flow diagram. Activity diagrams are often used in business process modelling. They can also describe the steps in a use case diagram. Activities modelled can be sequential and concurrent.

Activity Diagrams are used to illustrate the flow of control in a system and refer to the steps involved in the execution of a use case. We can depict both sequential processing and concurrent processing of activities using an activity diagram that focuses on the condition of flow and the sequence in which it happens.

```
                    ┌──────────────┐
                    │    Start     │
                    └──────┬───────┘
                           │
                    ┌──────▼───────┐
                    │ Select 15 features out of │
                    │   22 features │
                    └──────┬───────┘
                           │
                    ┌──────▼───────┐
                    │ N(Size of subsets)=1 │
                    └──────┬───────┘
                           │
                        ◇ N≤ 15 ◇ ──── No
                        Yes │
                    ┌──────▼───────┐
                    │ Make all possible subsets │
                    │ with size N using 15 features │
                    └──────┬───────┘
                    ┌──────▼───────┐
                    │ Parkinson and Healthy classification using SVM classifier │
                    └──────┬───────┘
                    ┌──────▼───────┐
                    │ Accuracy of classification is calculated and best subsets are extracted │
                    └──────┬───────┘
                    ┌──────▼───────┐
                    │   N=N+1      │
                    └──────────────┘
                    ┌──────────────┐
                    │    Stop      │
                    └──────────────┘
```

# 6. CODING AND IMPLEMENTATION

## 6.1 Source Code

```python
import streamlit as st import numpy as np import

pandas as pd from sklearn import svm from

sklearn.metrics import accuracy_score from

sklearn.preprocessing import StandardScaler from

sklearn.model_selection import train_test_split


# Set up Streamlit session state for login

if 'logged_in' not in st.session_state:

    st.session_state['logged_in'] = False


# Define valid credentials (you can customize this or fetch from a secure source)

VALID_USERNAME = "admin"

VALID_PASSWORD = "password123"


#   Login   function

def login():

    st.title('Login Page')

    st.write("Please enter your username and password to access the Parkinson's Disease
Detection tool.")


    username = st.text_input("Username") password =

    st.text_input("Password", type="password") if

    st.button("Login"):
```

```python
    if username == VALID_USERNAME and password == VALID_PASSWORD:

        st.session_state['logged_in']    =    True    st.success("Login

        successful! You now have access to the app.")

    else:    st.error("Invalid    username    or

        password")


# If the user is not logged in, show the login page

if    not    st.session_state['logged_in']:

    login()


# If the user is logged in, show the Parkinson's Disease Detection tool

if st.session_state['logged_in']:

    st.title('Parkinson\'s    Disease    Detection')

    st.write("""

    This is a simple web app to predict whether a person has Parkinson's Disease based on various vocal measurements.

    """)


    # Load the data

    @st.cache def

    load_data():

        df  =  pd.read_csv(r"C:\Users\syeda\OneDrive\Desktop\Mini-Project\parkinsons.csv")

        return df


    df = load_data()
```

```python
# Show dataset if

st.checkbox('Show raw data'):

st.write(df)


# Preprocess the data

X = df.drop(columns=['name', 'status'], axis=1)

Y = df['status']


# Split the data

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=2)


# Standardize the data

ss = StandardScaler()

X_train = ss.fit_transform(X_train)

X_test = ss.transform(X_test)


# Train the model model =

svm.SVC(kernel='linear')

model.fit(X_train, Y_train)


# Model evaluation

X_train_pred = model.predict(X_train) train_data_acc =

accuracy_score(Y_train, X_train_pred)


X_test_pred = model.predict(X_test)
```

```
test_data_acc = accuracy_score(Y_test, X_test_pred)


st.write("### Model Accuracy") st.write(f"Training

Data Accuracy: {train_data_acc}") st.write(f"Testing

Data Accuracy: {test_data_acc}")


# User input for prediction st.write("##

Predict Parkinson's Disease")


def user_input_features():

    data = {} for col in

    X.columns:

        data[col] = st.number_input(f'Enter {col}', value=0.0)
    features = pd.DataFrame(data, index=[0]) return
    features


input_data = user_input_features()


#    Prediction    if

st.button('Predict'):

    input_data_std = ss.transform(input_data)

    prediction = model.predict(input_data_std)

    if prediction[0] == 0:

        st.write("Negative, No Parkinson's Found - Patient is HEALTHY")

    else:

        st.write("Positive, Parkinson's Found - Patient is UNHEALTHY")
```

## 6.2 Implementation

### 6.2.1 Python

Python is an interpreted high-level programming language for general-purpose programming. Created by Guido van Rossum and first released in 1991, Python has a design philosophy that emphasises code readability, notably using significant whitespace.

Python features a dynamic type system and automatic memory management. It supports multiple programming paradigms, including object-oriented, imperative, functional, and procedural, and has a large and comprehensive standard library.

- Python is Interpreted – Python is processed at runtime by the interpreter. You do not need to compile your program before executing it. This is similar to PERL and PHP.
- Python is Interactive – You can sit at a Python prompt and interact with the interpreter directly to write your programs.

Python also acknowledges that speed of development is important. Readable and terse code is part of this, and so is access to powerful constructs that avoid tedious repetition of code. Maintainability also ties into this may be an all but useless metric, but it does say something about how much code you have to scan, read, and/or understand to troubleshoot problems or tweak behaviours. This speed of development, the ease with which a programmer of other languages can pick up basic Python skills, and the huge standard library are key to another area where Python excels. All its tools have been quick to implement and saved a lot of time, and several of them have later been patched and updated by people with no Python background - without breaking.

## 6.2.2 Modules Used in Project

**Streamlit:**

Purpose: used for creating the web interface of the application

Key Functions: st.title(), st.text_input(), st.button(), st.write(), st.cache(), st.error()

**Requests (requests):**

Purpose: used for making HTTP requests to fetch content from URL.

Key Functions: requests.get(), response.raise_for _status()

**Testing data samples:**

Purpose: Used to test the machine model by using some sample data with real time readings.

# 7. SYSTEM TESTING

The purpose of testing is to discover errors. Testing is the process of trying to discover every conceivable fault or weakness in a work product. It provides a way to check the functionality of components, sub-assemblies, assemblies and/or a finished product. It is the process of exercising software with the intent of ensuring that the Software system meets its requirements and user expectations and does not fail in an unacceptable manner. There are various types of tests. Each test type addresses a specific testing requirement. System testing is a type of software testing that evaluates the overall functionality and performance of a complete and fully integrated software solution. It tests if the system meets the specified requirements and if it is suitable for delivery to the end-users. This type of testing is performed after the integration testing and before the acceptance testing.

## 7.1. Testing Knowledge Required:

- ### Machine Learning Testing Knowledge

    Understanding of overfitting, underfitting, and the need for cross-validation to ensure that the model generalises well to new, unseen data. It is the process of evaluating and validating the performance of machine learning models to ensure their correctness, accuracy and robustness. Unlike traditional software testing, which mainly focuses on code functionality, machine learning testing includes additional layers due to the inherent complexity of machine learning models. It ensures that machine learning models perform as intended, providing reliable results and adhering to industry standards.

- ### Domain-specific Knowledge

    Awareness of Parkinson's disease symptoms and medical data interpretation to validate the accuracy of predictions and feature extraction. Domain-specific knowledge is specialised knowledge or expertise in a specific field or domain. Domain-specific knowledge is also known as content knowledge. It's considered a key component of competence, expertise, academic achievement, and other cognitive learning outcomes. Domain-specific knowledge is different from general knowledge, which is also known as domain-independent knowledge. For example, a software engineer might have general knowledge of computer programming, but also have domain knowledge about developing programs for a specific industry.

- ### Data Integrity and Preprocessing Testing

    Skills in testing data pipelines for completeness, consistency, and ensuring no data leaks between training and testing phases. The integrity of data analysis is

highly dependent on the quality of data preprocessing. Data preprocessing requires the usability and interpretability of data, laying the groundwork for accurate machine learning and AI models. Data integrity refers to the accuracy, consistency, and completeness of data throughout its lifecycle. It's a critically important aspect of systems which process or store data because it protects against data loss and data leaks. Maintaining the integrity of the data over time and across formats is a continuous process involving various processes, rules, and standards.

## ● Model Evaluation Metrics Testing

Familiarity with performance metrics like accuracy, precision, recall, F1-score, and ROC-AUC to assess the model's efficacy. Evaluation matrices are important as they help: To assess the performance of a model: Evaluation metrics provide a quantitative measure of how well a model performs on a given task. Model evaluation metrics are used to assess the performance of a machine learning model, and are important for several reasons. Some metrics used in model evaluation include accuracy, precision, sensitivity, specificity, mean average precision(mAP).

## 7.2. Types of Testing used for the system:

## ● User Acceptance Testing (UAT):

Testing the system with domain experts or clinicians to ensure the model's predictions are clinically relevant and interpretable.

UAT will involve medical professionals testing the system to ensure that the model's predictions align with their medical expertise. The goal is to verify that the system's output is not only accurate but also interpretable and actionable in a clinical setting.

UAT will focus on evaluating the system's user interface (UI) and overall usability. Clinicians will test whether the system is intuitive and easy to navigate, ensuring that they can input patient data and receive predictions without technical difficulty.

UAT will simulate real-world clinical scenarios where users input patient data and receive diagnostic predictions. This will help ensure that the model functions effectively in practical settings, under varying conditions such as incomplete data or noisy input.

- 
  **Validation Testing** validation testing can be implemented using **k-fold cross-validation**. This method splits the dataset into **k subsets (folds)**. The model is trained on **k-1 folds** and tested on the **remaining fold**. This process is repeated **k times**, with each fold being used as the test set once. The results from each iteration are averaged to get a more reliable performance metric.

  Cross-validation helps detect **overfitting**, where the model performs well on the training data but fails to predict accurately on new, unseen data. By validating the model on different folds, you can gauge how well it will generalise.

  Validation testing allows you to tune your model's **hyperparameters** (like regularisation strength, learning rate, etc.). By using validation data, you can experiment with different settings and pick the one that maximises performance without compromising on generalisation.

  By using validation testing, you can compare different machine learning models (e.g., SVM, Random Forest, Neural Networks) and choose the one that performs the best in terms of validation metrics like accuracy, precision, recall, and F1-score.

- **Performance Testing**

  Stress and load testing to assess the system's ability to handle large datasets in real-time scenarios. The model focuses on evaluating how well the model performs under various conditions and its efficiency in handling data.

  Speed and Latency is important when the system needs to provide real-time or near-real-time predictions, especially if it is integrated into a clinical setting or connected to wearable devices for monitoring symptoms.

  If the system needs to process multiple inputs simultaneously (e.g., in a hospital environment with many patients), throughput testing ensures the model can handle such load efficiently.

  Resource utilisation is particularly useful for optimising the model for deployment on different platforms (e.g., cloud vs. on-premise). Lower resource consumption may also mean that your system can run on smaller devices or at a lower cost.

- 
  ## Security Testing

  Since the system handles sensitive patient data, security testing ensures that data privacy is maintained. This includes testing encryption methods used for storing and transmitting patient information to prevent unauthorised access.

  Security testing verifies that only authorised users (e.g., healthcare professionals) have access to the system. This involves testing authentication mechanisms, such as password protection, multi-factor authentication, and role-based access control (RBAC), ensuring that users are restricted to their appropriate roles and permissions.

  The system needs to comply with healthcare regulations such as HIPAA (Health Insurance Portability and Accountability Act) or GDPR (General Data Protection Regulation). Security testing verifies that the system meets these legal requirements for handling medical data.

- **End-to-End Testing**

  The test begins by simulating real-world data inputs such as patient symptoms, sensor data, or clinical records. This step ensures the system can properly handle and process various types of input, whether from files, IoT devices, or databases.

  The system's preprocessing pipeline is tested to confirm that data cleaning, normalisation, and feature extraction are functioning as expected. This ensures the data is properly prepared before being fed into the machine learning model.

  After predictions are made, the system's post-processing steps, such as formatting results for display or further analysis, are validated. This includes checking how results are presented to clinicians, ensuring they are understandable and clinically meaningful.

  If the system has an interface for users (e.g., a web or mobile app), the test ensures that users can input data, run the model, and view results seamlessly. This step ensures the system works as expected for the end-user without glitches or delays.

  If the system incorporates feedback for continuous learning or updates (e.g., from clinicians confirming a diagnosis), End-to-End testing ensures that this

- 

feedback loop is functioning correctly, allowing the system to improve over time.

## Unit Testing

Unit testing ensures that each part of the system functions as intended in isolation, helping to catch errors early in development. By testing individual components, you can guarantee that bugs in one part of the system (e.g., preprocessing) don't affect the rest, ensuring a reliable and accurate Parkinson's detection model.

Unit tests validate that these metrics are correctly calculated based on predicted and actual values. For instance, feeding in sample predictions should produce known precision or recall values, and the test verifies the accuracy of these calculations.

Model training unit is used to ensure the model is being trained properly with the correct data and that it returns consistent results under similar conditions. This can include checking if the loss decreases during training and if the model's accuracy improves.

## ● Integration testing

We need to ensure that raw data from patients (e.g., sensor readings, medical records) is correctly processed (cleaning, normalisation) and that the extracted features (e.g., motor symptoms, tremor analysis) are accurately passed to the next stage. Integration testing checks that the preprocessed data flows smoothly into the feature extraction algorithms without data loss or transformation errors.

Once features are extracted, they need to be fed into the machine learning model for training or prediction. Integration testing verifies that the extracted features are correctly formatted, labelled, and compatible with the model's input requirements. It ensures that no errors occur when moving from feature extraction to model training, such as mismatched dimensions or missing values.

- 
  After the model generates predictions (e.g., Parkinson's or non-Parkinson's), integration testing ensures the results are correctly passed to the user interface or reporting system. This step ensures that the model's output is correctly interpreted and displayed, making sure predictions are accurate and understandable for the end-user, such as clinicians.

## 7.3. Test Cases

**Test case 1: Display Login**



**Test case 2: Invalid login credentials**



43

**Test case**

### 3: Invalid login credentials



**Test case 4: Login Successful**

**Test case**

### 5: Home Page



**Test case 6: Home Page**

**Test case 7:**

# Predict Parkinson's Disease

Enter MDVP:Fo(Hz)

| 0.00 | − + |

Enter MDVP:Fhi(Hz)

| 0.00 | − + |

Enter MDVP:Flo(Hz)

| 0.00 | − + |

Enter MDVP:Jitter(%)

| 0.00 | − + |

Enter MDVP:Jitter(Abs)

| 0.00 | − + |

Enter MDVP:RAP

| 0.00 | − + |

Enter MDVP:PPQ

| 0.00 | − + |

Enter Jitter:DDP

| 0.00 | − + |

Enter MDVP:Shimmer

| 0.00 | − + |

Enter MDVP:Shimmer(dB)

46

**Test case 8:**

**Prediction model**

Enter DFA

0.78 − +

Enter spread1

-6.00 − +

Enter spread2

0.19 − +

Enter D2

0.19 − +

Enter PPE

2.53 − +

Predict

Positive, Parkinson's Found - Patient is UNHEALTHY

**Test case 9:**

**Prediction model**

Enter DFA

| 0.81 | − + |
| --- | --- |

Enter spread1

| -4.81 | − + |
| --- | --- |

Enter spread2

| 0.26 | − + |
| --- | --- |

Enter D2

| 2.30 | − + |
| --- | --- |

Enter PPE

| 0.28 | − + |
| --- | --- |

Predict

Negative, No Parkinson's Found - Patient is HEALTHY

# 8. CONCLUSION

## 1. Project Summary:

The objective of this project, "**Next-Generation Parkinson's Detection: Leveraging Machine Learning Innovations**", was to develop a system that can assist in the early and accurate detection of Parkinson's disease. Leveraging machine learning models trained on various features, such as motor symptoms and medical records, the system was designed to provide an automated, scalable solution for detecting Parkinson's. Throughout this project, we focused on improving detection accuracy and ensuring that the system is reliable and adaptable to real-world clinical settings.

The project involved several stages, starting from data collection, preprocessing, and feature extraction, to the training and evaluation of machine learning models. Multiple machine learning algorithms were explored, and various testing methodologies were applied to ensure that the system meets its objectives. In this process, the integration of multiple components, such as data input, feature extraction, and model output, was successfully achieved and validated through rigorous testing protocols, including **integration testing**.

## 2. Accomplishments:

Several key accomplishments were made during the development of this project:

**Model Accuracy and Precision:** The system achieved notable accuracy in predicting Parkinson's disease. Through thorough model evaluation and the use of metrics like precision, recall, and F1-score, we ensured that the system balances between sensitivity (detecting true positives) and specificity (minimising false positives).

**Data Preprocessing Efficiency:** One of the critical factors for achieving good model performance was efficient data preprocessing. We ensured that data was cleaned, normalised, and structured properly, addressing common issues such as missing values and data imbalance.

**User Interface Integration:** A user-friendly interface was designed, allowing clinicians or healthcare providers to easily input data and receive accurate predictions from the model. The simplicity of the UI makes the system practical for non-technical users in medical environments.

**Robust Testing Procedures:** The system underwent extensive testing, including positive and negative test cases, ensuring that both valid and invalid inputs were properly handled. In particular, **integration testing** ensured smooth communication between different system components, such as feature extraction and prediction generation, making the system robust and reliable.

## 3. Challenges Faced:

- **Data Imbalance:** One of the major issues faced was the imbalance in the dataset, where there were more instances of healthy subjects compared to those with Parkinson's disease. This required the use of techniques like oversampling, undersampling, and synthetic data generation to ensure the model could perform well on both classes.

- **Feature Overlap and Variability:** Some features used for Parkinson's detection (e.g., tremors, gait) are not exclusive to the disease and can overlap with other conditions. This made it challenging to isolate truly predictive features. We mitigated this by using feature selection techniques and domain-specific knowledge to refine the feature set.

- **Generalisation of the Model:** Ensuring that the model generalises well to unseen data was another challenge. Techniques like cross-validation and regularisation were employed to avoid overfitting and improve the model's ability to predict accurately on new data.

## 4. Future Aspects:

- **Real-Time Monitoring:** Future versions of the system could incorporate real-time monitoring of patients using wearable devices, allowing for continuous tracking of symptoms. This would improve early detection and offer personalised healthcare recommendations.

- **Enhanced Data Collection:** Incorporating more comprehensive datasets, such as those involving genetic data or advanced imaging techniques (e.g., MRI, PET scans), could significantly improve the accuracy of predictions. This would also allow the system to detect more subtle signs of Parkinson's disease.

- **Deep Learning Models:** Currently, machine learning models such as SVMs, Random Forest, and Logistic Regression have been explored. Future work could investigate the application of deep learning techniques like Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) for more complex feature extraction and classification tasks.

- **Explainable AI:** To improve adoption in clinical environments, implementing explainable AI (XAI) techniques would provide transparency into how predictions are made, helping clinicians understand and trust the model's decisions.
- **Expanded User Base:** The system could be expanded to include additional stakeholders, such as caregivers and patients themselves, providing a platform for monitoring symptom progression and adjusting treatment plans accordingly.

## 5. Conclusion:

In conclusion, **Next-Generation Parkinson's Detection: Leveraging Machine Learning Innovations** the project represents a significant step towards automating and improving the early detection of Parkinson's disease. By combining state-of-the-art machine learning techniques with domain-specific insights, the system provides a valuable tool for healthcare professionals. While challenges such as data imbalance and feature overlap were encountered, the system demonstrates strong potential for practical deployment in medical settings.

Moreover, the rigorous testing protocols, particularly **integration testing**, ensured that the various components of the system worked harmoniously, allowing for a reliable and scalable solution. Moving forward, there is ample opportunity to enhance the system with real-time data integration, deeper models, and explainable AI, which could revolutionise the way Parkinson's disease is detected and managed.

Ultimately, this project underscores the transformative power of machine learning in Health care, and with continued development, it can play a pivotal role in improving patient outcomes through earlier detection and more accurate diagnoses.

# 9. FUTURE ENHANCEMENTS

Future enhancements for the "Parkinson's Detection using Machine Learning" system could significantly improve its functionality, accuracy, and utility. One of the primary future directions would be integrating advanced data sources such as genomic data, medical imaging (MRI, CT scans), and voice analysis, alongside real-time data from wearable devices.

This would allow for more comprehensive patient profiles and improve the system's diagnostic capabilities. Another key enhancement would be shifting towards real-time monitoring and predictive analytics, which would enable the system to track a patient's condition continuously and potentially detect early onset or progression of Parkinson's at earlier stages. Through time-series analysis, the system could issue predictive alerts, allowing healthcare providers to intervene early. Additionally, the system could be expanded to incorporate personalised treatment recommendations. By analysing a patient's history, response to medication, and other clinical data, the model could suggest optimised treatment plans tailored to each patient, helping clinicians make informed decisions.
It could also monitor treatment efficacy over time, making adjustments based on real-time data from patients' wearable devices. Improving the system's user interface and accessibility is another important enhancement. The inclusion of voice command support and a mobile app would make it more user-friendly for both clinicians and patients, offering a seamless experience for interacting with the model. Multilingual support would also expand the system's usability across different regions.

Moreover, the system could incorporate Explainable AI (XAI) methods to improve transparency and trust in the model's predictions, which is critical in medical settings. Visualisation tools could help clinicians understand how the system arrived at a particular diagnosis or prediction by highlighting key features such as movement patterns or voice changes. Another exciting enhancement would be adapting the system to diagnose and manage other neurological disorders, such as Alzheimer's, ALS, and MS, leveraging the overlap in symptoms and diagnostic patterns.

This multimodal disease prediction capability could significantly broaden the system's application in healthcare. Additionally, improving security and privacy measures is critical, particularly in the handling of sensitive patient data. Future versions could implement blockchain for secure and auditable data storage, ensuring compliance with healthcare regulations like HIPAA and GDPR. Advanced encryption techniques would further safeguard patient data during transmission and storage. Federated learning could also be employed to allow the model to learn from decentralised data without compromising patient privacy.

# 10. REFERENCES

**Analysis of Tremors in Parkinson's Disease Using Accelerometer:** Niya Romy Markose, Priscilla Dinkar Moyya, and Mythili Asaithambi, Analysis of tremors in Parkinson's Disease using accelerometer, IEEE, July 2021**.**

**Building a Machine Learning Framework to Remotely Assess Parkinson's Disease Using Smartphones:** Oliver Y. Chen, Florian Lipsmeier, Huy Phan, and John Prince, IEEE, May 2020.
**Deep Convolutional Neural Network for Parkinson's Disease Based on Handwriting Screening:** Mohamed Shaban, IEEE, Sep 2020.

**Detection of Parkinson's Disease through Smell Signatures:** Shrinidhi Kulkarni, Neenu George Kalayil, Jinu James, Sneha Parsewar and Revati Shriram**,** IEEE, Sep 2020.

**Effectiveness Analysis of Bio-electric Stimulation Therapy to Parkinson's Disease via Discrete Fourier Transform Approach:** Yuxin Lin, Bingo Wing-Kuen Ling, Nuo Xu, Ringo Lam, and Charlotte Ho, IEEE, May 2021:

**Efficient Pre-diagnosis Approach for Parkinson's Disease with Machine Learning:** Yuqi Qiu, IEEE, Nov 2020.

**Feature Selection-Based Twin-Support Vector Machine for the Diagnosis of Parkinson's Disease:** Surendra Bikram, Thapa, Surabhi Adhikari, Awishkar Ghimire and Anshuman Aditya, IEEE, June 2021.

**Prediction of Parkinson's disease and severity of the disease using Machine Learning and Deep Learning algorithm:** Pooja Raundale, Chetan Thosar and Shardul Rane, IEEE, Aug 2021.

**Parkinson's Disease Detection from Spiral and Wave Drawings Using Convolutional Neural Networks A Multistage Classifier Approach:** Sabyasachi Chakraborty, Satyabrata Aich, Jong-Seong-Sim, Eunyoung Han, Jinse Park and Hee-Cheol Kim, IEEE, June 2020.

**A Speech-Based System for Parkinson's Disease Analysis and Monitoring:**
Daniel Palacios-Alonso, Guillermo Melendez-Morales, Agustin Lopez-Arribas, Carlos Lazaro-Carrascosa, Andres Gomez-Rodellar and Pedro Gomez-Vilda, MonParLoc, IEEE, Oct 2020