

## Objective:

Develop a classification model to predict Parkinson's disease using the given patient's biomedical voice measurements.

## Context:

Parkinson's Disease (PD) is a degenerative neurological disorder marked by decreased dopamine levels in the brain. It manifests itself through a deterioration of movement, including the presence of tremors and stiffness. There is commonly a marked effect on speech, including dysarthria (difficulty articulating sounds), hypophonia (lowered volume), and monotone (reduced pitch range). Additionally, cognitive impairments and changes in mood can occur, and risk of dementia is increased. Traditional diagnosis of Parkinson's Disease involves a clinician taking a neurological history of the patient and observing motor skills in various situations. Since there is no definitive laboratory test to diagnose PD, diagnosis is often difficult, particularly in the early stages when motor effects are not yet severe. Monitoring progression of the disease over time requires repeated clinic visits by the patient. An effective screening process, particularly one that doesn't require a clinic visit, would be beneficial. Since PD patients exhibit characteristic vocal features, voice recordings are a useful and non-invasive tool for diagnosis. If machine learning algorithms could be applied to a voice recording dataset to accurately diagnosis PD, this would be an effective screening step prior to an appointment with a clinician.

## Data Description:

This dataset is composed of a range of biomedical voice measurements from 31 people, 23 with Parkinson's disease (PD). Each column in the table is a particular voice measure, and each row corresponds one of 195 voice recording from these individuals ("name" column). The main aim of the data is to discriminate healthy people from those with PD, according to "status" column which is set to 0 for healthy and 1 for PD.

The data is in ASCII CSV format. The rows of the CSV file contain an instance corresponding to one voice recording. There are around six recordings per patient, the name of the patient is identified in the first column.

name - ASCII subject name and recording number  
MDVP:Fo(Hz) - Average vocal fundamental frequency

MDVP:Fhi(Hz) - Maximum vocal fundamental frequency  
MDVP:Flo(Hz) - Minimum vocal fundamental frequency  
MDVP:Jitter(%),MDVP:Jitter(Abs),MDVP:RAP,MDVP:PPQ,Jitter:DDP - Several measures of variation in fundamental frequency  
MDVP:Shimmer,MDVP:Shimmer(dB),Shimmer:APQ3,Shimmer:APQ5,MDVP:APQ,Shimmer:DDA - Several measures of variation in amplitude  
NHR,HNR - Two measures of ratio of noise to tonal components in the voice  
status - Health status of the subject (one) - Parkinson's, (zero) - healthy  
RPDE,D2 - Two nonlinear dynamical complexity measures  
DFA - Signal fractal scaling exponent  
spread1,spread2,PPE - Three nonlinear measures of fundamental frequency variation

### Steps and Milestones (100%):

- Setup Environment and Load Necessary Packages (5%)
- Data Preparation (40%)
  - Loading Data <sup>(5%)</sup>
  - Cleaning Data <sup>(10%)</sup>
  - Data Representation & Feature Engineering (If Any) <sup>(15%)</sup>
  - Creating Train and Validation Set <sup>(10%)</sup>
- Model Creation (30%)
  - Write & Configure Model <sup>(10%)</sup>
  - Compile Model <sup>(10%)</sup>
  - Build Model & Checking Summary <sup>(10%)</sup>
- Training and Evaluation (25%)
  - Run Multiple Experiments <sup>(10%)</sup>
  - Reason & Visualize Model Performance <sup>(5%)</sup>
  - Evaluate Model on Test Set <sup>(10%)</sup>

### Learning Outcomes:

- Predictive Analytics
- Ensemble Classifiers – Random Forests
- Decision Tree Classifier
- Fine-tuning Model with Grid Search
- Data Preparation
- Feature Engineering
- Visualization