

International Conference on Machine Learning and Data Engineering

Early detection of Parkinson's disease using machine learning

Aditi Govindu^a, Sushila Palwe^b^a*School of SCET, MIT WPU, Pune 411038, India*^b*School of SCET, MIT WPU, Pune 411038, India*

Abstract

Parkinson's disease (PD) is a neurodegenerative disorder affecting 60% of people over the age of 50 years. Patients with Parkinson's (PWP) face mobility challenges and speech difficulties, making physical visits for treatment and monitoring a hurdle. PD can be treated through early detection, thus enabling patients to lead a normal life. The rise of an aging population over the world emphasizes the need to detect PD early, remotely and accurately. This paper highlights the use of machine learning techniques in telemedicine to detect PD in its early stages. Research has been carried out on the MDVP audio data of 30 PWP and healthy people during training of 4 ML models. Comparison of results of classification by Support Vector Machine (SVM), Random Forest, K-Nearest Neighbors (KNN) and Logistic Regression models, yield Random Forest classifier as the ideal Machine Learning (ML) technique for detection of PD. Random Forest classifier model has a detection accuracy of 91.83% and sensitivity of 0.95. Through the findings of this paper, we aim to promote the use of ML in telemedicine, thereby providing a new lease of life to patients suffering from Parkinson's disease.

© 2023 The Authors. Published by Elsevier B.V.

This is an open access article under the CC BY-NC-ND license (<https://creativecommons.org/licenses/by-nc-nd/4.0>)

Peer-review under responsibility of the scientific committee of the International Conference on Machine Learning and Data Engineering

Keywords: Parkinson's disease (PD); MDVP dataset; telemedicine; Random Forest; SVM

1. Introduction

Parkinson's disease (PD) is a common neurological disorder impacting muscle movement in the body. It affects mobility, speech and posture leading to tremors, muscle rigidity and bradykinesia [1]. It occurs due the death of neurons, resulting in a decrease of dopamine levels in the brain. Low levels of dopamine hamper communication between synapses, causing ineffective motor functions [2]. While the progress of symptoms may vary from patient to patient, balance problems and tremors are the most prevalent side-effects of dopaminergic neuron death.

Unfortunately, there is no cure for PD, hence Patients with Parkinson's (PWP) rely on early detection and tailored treatments to slow the progress of the disease.

PD has 5 stages of progression and 90% PWP display signs of vocal cord injuries as a symptom in stage 0. Vocal impairment is not only easy to measure, but also falls under the category of telemedicine [3] or remote medicine. Patients need not travel physically to a doctor instead; they can record audio using phones and perform a simple test at home. Common voice modulation symptoms include dysphonia [4] and dysarthria [5]. Patients can be asked to hold a single vowel's pitch for as long as possible, also known as sustained phonation or running speech tests can be administered, as a realistic test of impairment. These phonation tests can be used for diagnosis of Parkinson's at stage 0.

Following early detection, doctors can cater therapeutic solutions or deep brain stimulation [6] to reactivate the dopamine producing neurons in the brain, thereby slowing the progress of PD. Owing to its complex nature, there is no cure for Parkinson's till date. However, early identification followed by right medication can reduce the tremors and imbalance symptoms in patients, enabling them to lead a normal life.

This paper focuses on early detection through audio recordings of PWP using ML techniques. This novel approach emphasizes the relevance of audio as a non-invasive biomarker [7] to detect PD. Our preliminary results show that Random Forest classifier model has an accuracy of 91.83% when trained on 22 attributes of MDVP audio data, compared to KNN, SVM and Logistic regression models. PWP suffer from mobility issues and are unable to travel for health check-ups. The proposed remote detection technique will provide a new lease of life to patients, as it classifies the severity of PD using speech data, that can be recorded on mobile phones. Our research compares and contrasts various ML models for disease classification that are not only memory efficient, but also faster compared to deep neural network learning models. We hope our promising results encourage advancements in telemedicine for PD.

1.1. Literature survey

Previous studies to predict PD have been implemented on MRI scans, gait and genetic data, but research on audio impairment for early detection is minimal. For instance, Bilal et. al. [7] studied genetic data to predict the onset of PD in senior patients with SVM model. They trained an SVM model to reach an accuracy of 0.889, while this research paper describes an improved SVM model with an accuracy of 0.9183. These results also corroborate the merits of classification of PD based on audio data, over genetic data. Raundale, Thosar and Rane [8] used keystroke data from UCI telemonitoring dataset to train a Random Forest classifier to predict the severity of PD in older patients. Cordella et. al. [9] use audio data to classify PWP, however their models are heavily reliant on MATLAB. Our research uses open-source models trained in Python, that are faster and memory efficient.

Majority of research done emphasizes the use of deep learning in PD detection, such as, Ali et. al. [10] who explain the use of ensemble deep learning models applied to phonation data, to predict the progress of Parkinson's disease. Their work lacked the use of feature selection that would improve Deep learning model (DNN) performance. Hence, this paper implements PCA on 22 attributes to select 7 major voice modalities in PD detection. Huang et. al. [11] aim to reduce PD diagnosis dependence on wearable equipment by training a traditional decision tree on 12 complex speech features of the MDVR-KCL [12] dataset. Wodzinski et. al. [13] trained a ResNet model on images of audio data, instead of training the model on the nuances of the frequency of audio. Wroge et. al [14] aimed to remove subjectivity of doctors in prediction of PD using an unbiased ML model, however their results achieved peak accuracy of 85% only.

Wang et. al. [15] implemented 12 machine learning models on 401 voice biomarkers dataset to classify patients as PD or not. They built a custom deep learning model (DEEP) with a classification accuracy of 96.45%, however the model was expensive due to large memory requirements. Alkhatib et. al. [16] implemented a linear classification model with 95% accuracy to characterize shuffling movement of PD patients. Their study focused on gait of patient and future work encouraged the use of audio and sleep data to improve the results. Ricciardi et. al [17] performed spatial-temporal analysis of brain MRI scans. They implemented decision trees, random forest and KNN to detect Mild Cognitive Impairment (MCI) in PWP. However, dataset was small and artificial data augmentation [18] was needed. A. U. Haq and colleagues [19] implemented L1-support SVM, without feature identification on vowel phonation dataset for neurological disorder patients. Their paper focused on patient age group of 46-85 years,

without considering healthy individuals in a lower age bracket. Mei et. al. [20] explain the importance of ML to detect PD, as subtle non-motor symptoms can be missed during subjective evaluation by a doctor. Their work reviews 209 studies based on dataset, ML methods and outcomes achieved.

Based on our literary review, we have implemented a PD classification model on audio data. Through our findings, we aim to contribute to the advancement of detection of PD through telemedicine. Keeping in mind, past research on biomarker data and models implemented, our research aims to explore KNN, logistic regression, random forest regression and SVM models to classify Parkinson's patient audio data. Our preliminary findings show that K nearest neighbor model is the best performing model with an accuracy of 91.83% and sensitivity of 0.95.

2. Proposed methodology

The proposed methodology collects audio data from PPMI [21] and UCI about Parkinson's patients voice modulations. Dataset contains information about jitter, shimmer and MDVP of vowel phonations. Data is preprocessed, analyzed and visualized for a thorough understanding of the attributes. Four models – Logistic regression, SVM, Random Forest Regressor and K nearest neighbors – are trained on 75% of the data. Models are trained to classify given audio data into PD or healthy, based on variations in frequency. Models are tested on 25% of the data and evaluated based on sensitivity, precision, accuracy, confusion matrix [22] and ROC-AUC score.

Figure 1 illustrates the generic process implemented. It demonstrates the stages of data ingestion from PPMI database, separation of data into testing and training sets, training of four models on data and validation of results using test data.

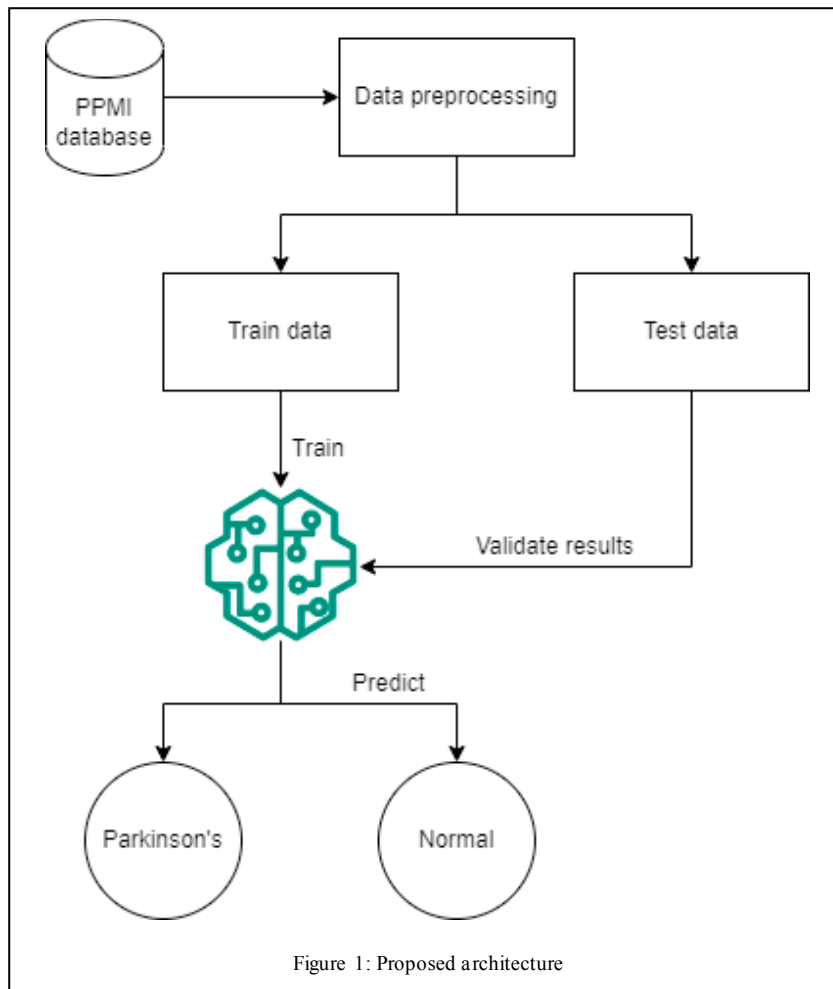


Figure 1: Proposed architecture

This research paper aims to identify the most relevant attributes in classification of PD and impact of imbalance in medical data in classification. Keeping in mind these requirements, 3 approaches have been implemented – Training on complete dataset that serves as a baseline test for PD classification, training on PCA identified attributes and training on 109 records obtained after dataset balancing. The algorithms used in each approach are described below:

Algorithm for approach 1: Models are trained on 22 attributes of data

- Collect MDVP audio data from PPPMI and UCI databases
- Perform data analysis to detect skew, imbalance and distribution of variables in data
- Scale the data to common range using Standard Scaler
- Split dataset into testing and training sets, where training data is 75% of total
- Train SVM, logistic regression, random forest and KNN models.

Algorithm for approach 2: Principal Component Analysis (PCA) is applied to identify 5 key attributes

- Collect MDVP audio data from PPPMI and UCI databases
- Perform data analysis to detect skew, imbalance and distribution of variables in data
- Scale the data to a common range using Standard Scaler
- Identify variance in every column of data and apply Principal Component Analysis (PCA) to identify 5 most relevant features to model training, out of 22 attributes.
- Split dataset into testing and training sets, where training data is 75% of total
- Retrain SVM, logistic regression, random forest and KNN models.
- Compare classification results using confusion matrix, ROC-AUC curve and accuracy

Algorithm for approach 3: Imbalance removal in dataset

- Collect MDVP audio data from PPPMI and UCI databases
- Perform data analysis to detect skew, imbalance and distribution of variables in data
- The dataset is imbalanced, with 109 records of PWP and 40 records of normal people, as illustrated in figure 2(a). The imbalance is resolved by up sampling [23] the minority class to reach 109 records each, as illustrated in figure 2(b).
- Scale the data to common range using Standard Scaler
- Split dataset into testing and training sets, where training data is 75% of total
- Retrain SVM, logistic regression, random forest and KNN models.
- Compare classification results using confusion matrix, ROC-AUC curve and accuracy

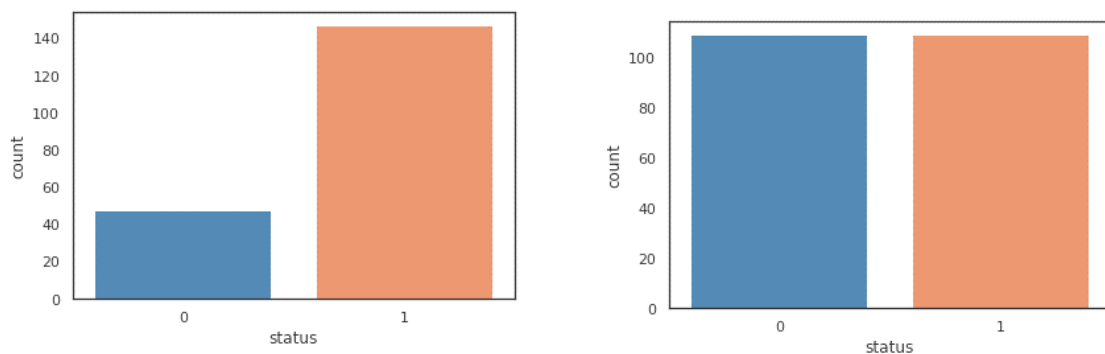


Figure 2. (a) Imbalanced data with 40 normal records; (b) Balanced data after up sampling

2.1. Dataset

Biomedical voice measurement [24] of 31 people have been gathered, where 23 patients have PD. Patients are in the age range of 46 to 85 years, while normal readings are from people of 23 years of age. An average of 6 phonation's were recorded 195 times for every person, ranging from 1 to 36 seconds in duration. The attributes of 195 records are elaborated in table 1 below:

Table 1: Dataset attributes

Attribute	Purpose
Name	Data is stored in ASCII CSV format where patient name and recording number is stored
MDVP: Fo (Hz)	Fundamental frequency of pitch period
MDVP: Fhi (Hz)	Upper limit of fundamental frequency or maximum threshold of voice modulation
MDVP: Flo (Hz)	Lower limit or minimal vocal fundamental frequency
MDVP: Jitter, Abs, RAP, PPQ, DDP	These are various Kay Pentax's multi-dimensional voice program (MDVP) measures. MDVP is a traditional measure of frequency of vibrations in vocal folds at pitch period to vibrations at start of next cycle called pitch mark [25]
Jitter and Shimmer	Measures of absolute difference between frequencies of each cycle, after normalizing the average
NHR and HNR	Signal to noise and tonal ratio measures, that indicate robustness of environment to noise
Status	0 indicates healthy person while 1 indicates PWP.
D2	Correlation dimension is used to identify dysphonia in speech using fractal objects. It is a nonlinear, dynamic attribute.
RPDE	Recurrence Period Density Entropy quantifies the extent to which signal is periodic
DFA	Detrended Fluctuation Analysis or DFA measures the extent of stochastic self-similarity of noise in speech signals.
PPE	Pitch Period entropy is used to assess abnormal variations in speech on a logarithmic scale
Spread1, spread2	Analysis of extent or range of variations in speech with respect to MDVP: Fo(Hz)

2.2. Data preprocessing

Data wrangling [26] is implemented to clean data and handle missing attributes in the dataset. Figure 3 depicts the noise to harmonic tone (NHR) ratio and harmonic tone to noise ratio (HNR) for PWP. As stages of the disease progress, noise in speech increases resulting in increased NHR, as seen in figure 3 (b). The skew in the data and low value of NHR (0.3) indicates poor voice quality.

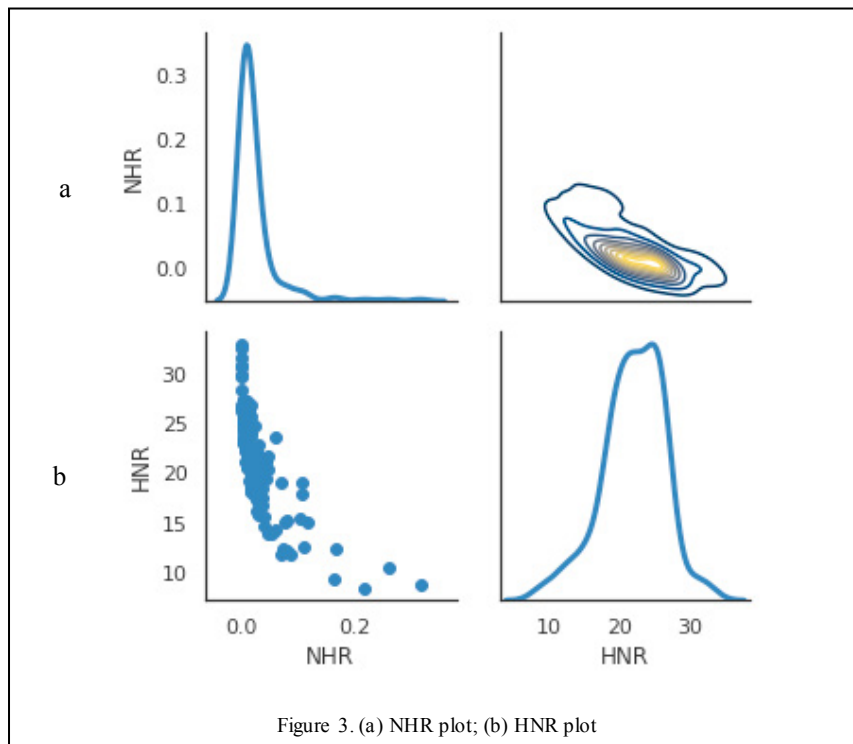


Figure 4 depicts a box plot of all 22 attributes in the dataset. It depicts the spread and skewness of data over a median quartile. Figure depicts blue as normal records and orange as PWP records. NHR data points for PWP have maximum no. of outliers, due to greater noise in speech. Similarly, HNR records have maximum data outliers below the median, for PWP records.

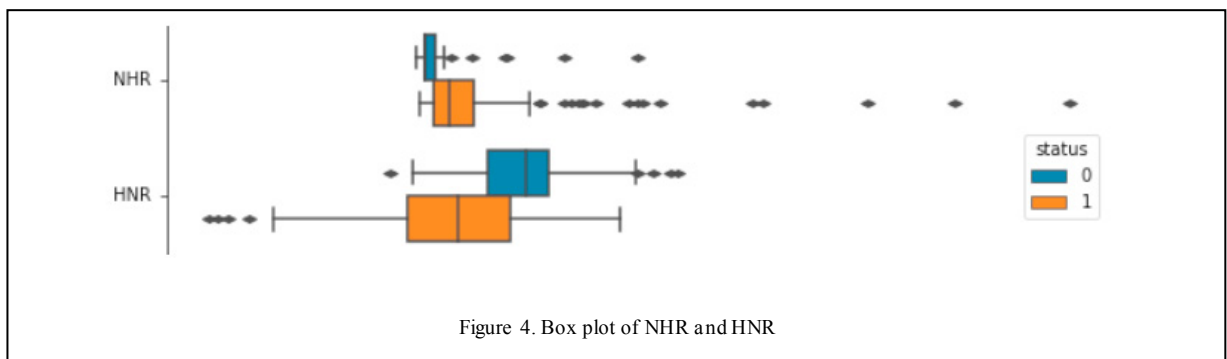
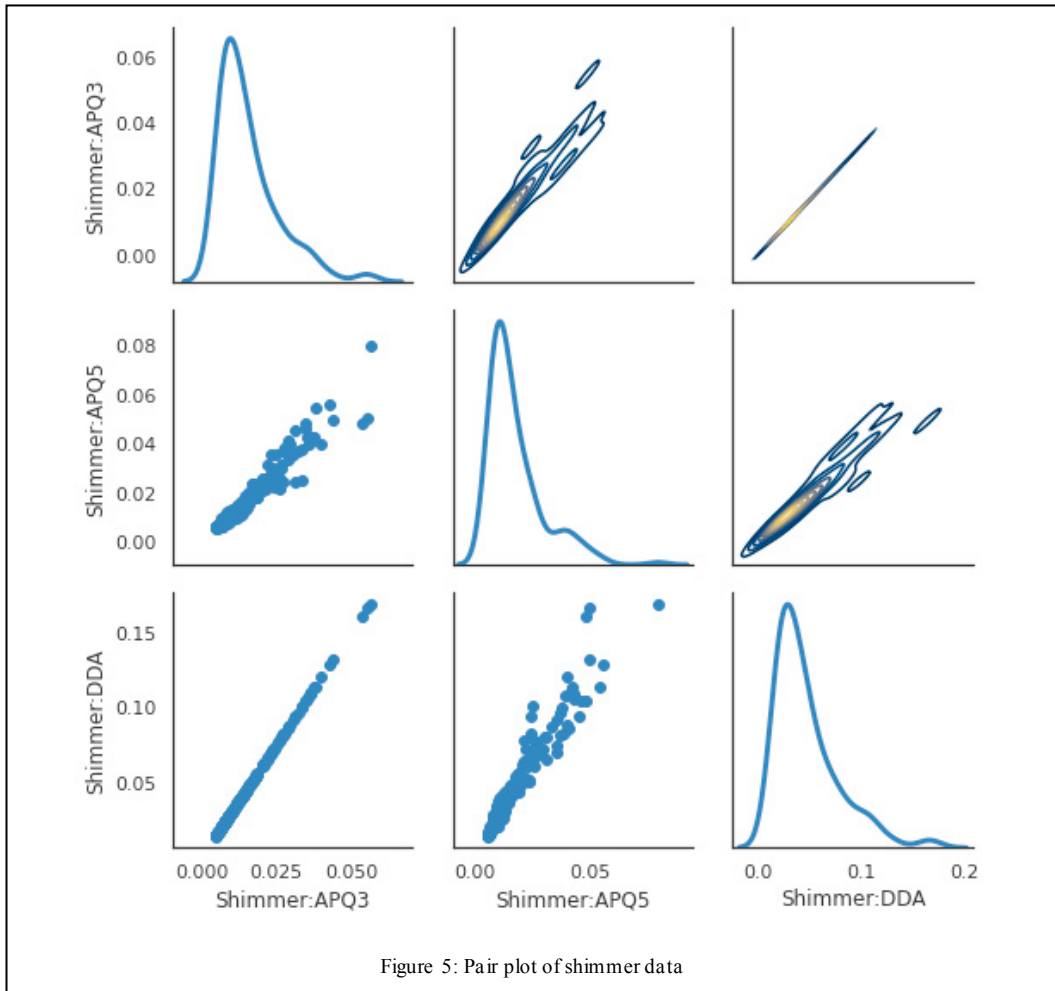


Figure 5 depicts the pair plot of shimmer data. It is used to highlight the shift in shimmer of voice for PWP, compared to healthy patients. It shows that Shimmer: APQ3 and Shimmer: DDA are linearly correlated while Shimmer: APQ5 and Shimmer: APQ3 are left skewed.



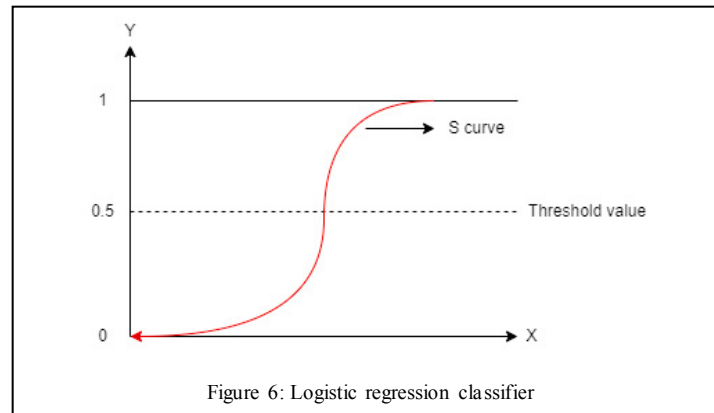
3. Model training

This research paper studies Logistic Regression, Random Forest classifier, Support Vector classifier and K nearest neighbors' models in 3 approaches:

- Complete dataset of 195 records and 22 attributes
- Dataset with 195 records and 5 attributes after Principal Component Analysis (PCA)
- Balanced dataset with 109 records and 22 attributes

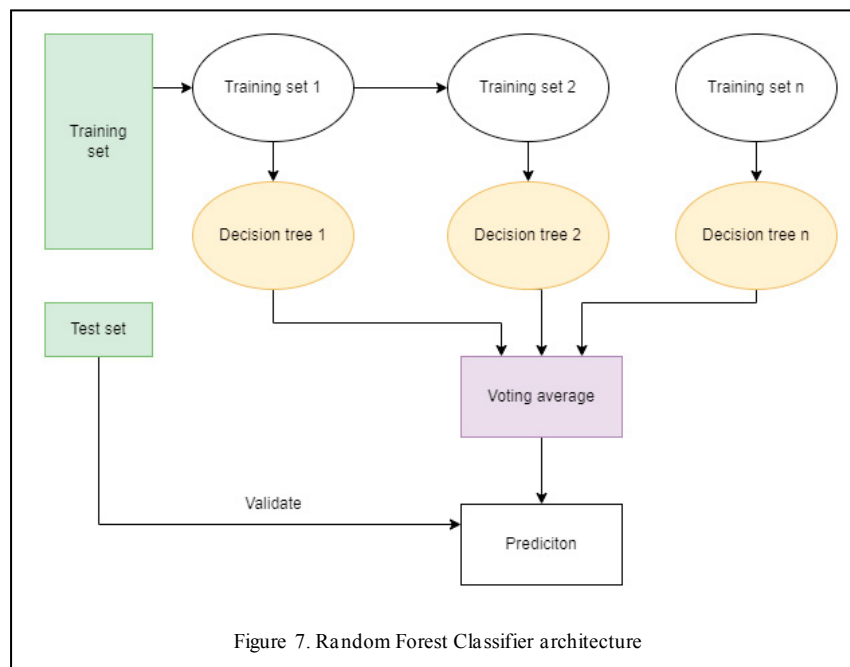
3.1. Logistic regression for classification

Logistic regression [27] is a prevalent supervised, ML algorithm that predicts categorically dependent variables using a set of independent variables. It uses 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 [28], where a line is fit to linearly predict one or more dependent variable, 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. The activation function of logistic classification has been illustrated in figure 6.



3.2. Random Forest classifier

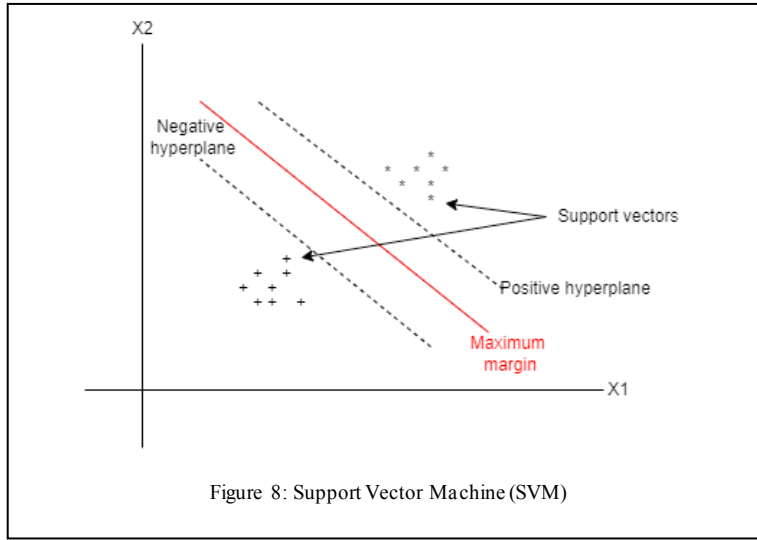
Random Forest is a supervised machine learning algorithm that is applicable to classification and regression problems. This paper implements random forest classifier [29] to train a number of decision trees on subsets of the dataset and consider the average to increase the predictive accuracy of the 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 increase, the chances of overfitting decrease. Architecture of random forest classifier used in this research paper, has been illustrated in figure 7 below.



3.3. Support Vector Machine (SVM)

Support vector machine (SVM) [30] is a supervised machine learning algorithm that creates a hyperplane to separate N features, by mapping these features to a multidimensional space. The architecture of SVM model has been illustrated in figure 8.

Since PD voice data is not linearly separable, we use an SVM kernel to transform data into higher dimensional space. SVM performs well for PD data due to memory efficiency and support vectors formed from a subset of training data points.



3.4. K nearest neighbors (KNN)

K nearest neighbours [31] (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.

4. 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.

$$Precision = \frac{TP}{TP+FP} \quad (1)$$

$$Recall = \frac{TP}{TP+FN} \quad (2)$$

$$Accuracy = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (3)$$

Receiver Operating Characteristics or ROC [32] 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.

Table 2 below depicts the results of models after approach 1 is applied, that is, models are trained on 22 attributes of MDVP dataset.

Table 2. Results of Approach 1: 22 attributes training

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 result is predicted. Every attribute is equally weighted during 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.

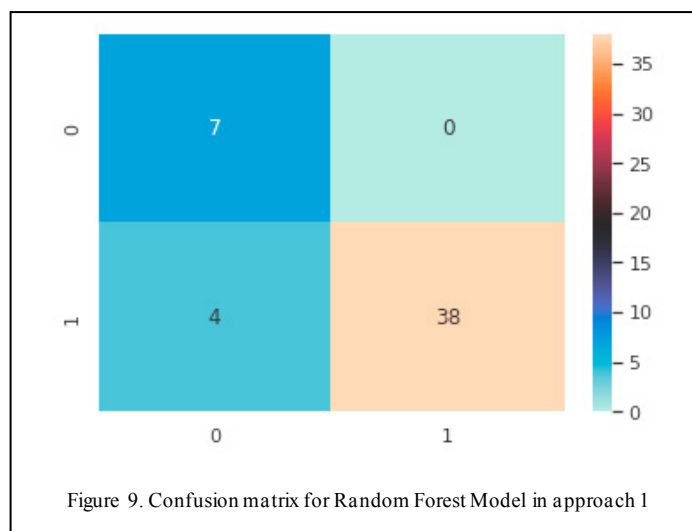


Table 3 depicts results of approach 2, after application of PCA. It yields 5 major attributes - MDVP, Shimmer, Jitter, PPE and RPDE. Upon training and evaluating models on these 5 attributes, results are:

Table 3. Results of Approach 2: 5 attributes after PCA

Metric	Logistic Regression	Random Forest	SVM	KNN
Accuracy	83.67%	83.67%	91.75 %	83.67 %
Precision	1.0	1.0	1.0	0.92
Recall	0.83	0.90	0.86	0.90
ROC AUC curve	0.636	0.818	0.727	0.779

Support Vector classifier with L1-support and linear kernel is ideal for PCA data, as it identifies ideal hyperplane 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.

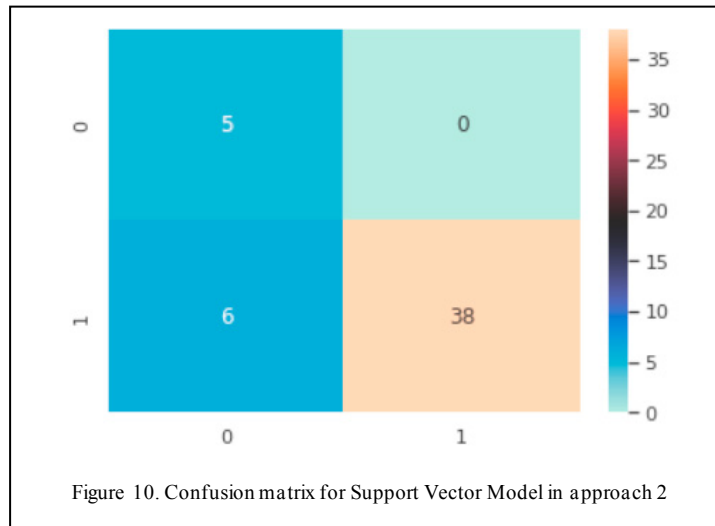
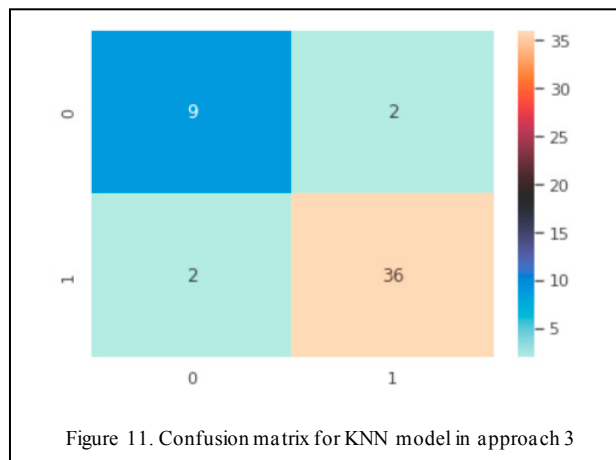


Table 4 below depicts the results for approach 3, using a balanced dataset. Models were trained on equal no. of records of normal and Parkinson's patients' data. Balancing ensures equal weightage is given to PWP and non-Parkinson's patients. The results are as follows:

Table 4. Results of Approach 3: Balanced dataset

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 neighbor'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 confusion matrix in figure 11. KNN model classifies data into 9 TN (no PD), 2 FN, 36 TP (PWP) and 2 FP.



5. 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 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 neighbor (KNN) model also performs well for balanced dataset, as classification into 2 categories without presumptions of data is favored. Thus, we recommend the use of Random Forest model to classify progress of the disease. It is a non-invasive, simple and accurate technique to provide long-term relief to PWP, globally.

In the future, we propose to use audio and REM sleep data to improve the results, as audio data alone is not a sufficient biomarker for classification of Parkinson's disease. We hope these findings encourage the use of mobile recorded audio to classify PD through telemedicine.

Acknowledgements

We express our sincere thanks and gratitude to the professors and coordinators at our esteemed university for their invaluable feedback. We also wish to acknowledge our peers, family and well-wishers who inspire and encourage us to work tirelessly.

References

- [1] Prabhavathi, K., Patil, S. (2022). "Tremors and Bradykinesia. In: Arjunan, S.P., Kumar, D.K. (eds) Techniques for Assessment of Parkinsonism for Diagnosis and Rehabilitation". *Series in BioEngineering. Springer*. 135–149 https://doi.org/10.1007/978-981-16-3056-9_9
- [2] Braak, H., Braak, E. (2000) "Pathoanatomy of Parkinson's disease" *J Neurol* **247**, 113–110. <https://doi.org/10.1007/PL00007758>
- [3] F. Amato, I. Rechichi, L. Borzi and G. Olmo, (2022), "Sleep Quality through Vocal Analysis: A Telemedicine Application," *2022 IEEE International Conference on Pervasive Computing and Communications Workshops and other Affiliated Events (PerCom Workshops)*, 706–711, doi: 10.1109/PerComWorkshops53856.2022.9767372.
- [4] Neighbors C, Song SA. "Dysphonia" (2022) *StatPearls [Internet]. Treasure Island (FL): StatPearls Publishing*.
- [5] Serge Pinto, Canan Ozsancak, Elina Tripoliti, Stéphane Thobois, Patricia Limousin-Dowsey, Pascal Auzou, "Treatments for dysarthria in Parkinson's disease", (2004) *The Lancet Neurology*, **3**(9): 547-556, ISSN 1474-4422, [https://doi.org/10.1016/S1474-4422\(04\)00854-3](https://doi.org/10.1016/S1474-4422(04)00854-3).
- [6] Nicolás G. Pozzi, Ioannis U. Isaias (2022), "Chapter 19 - Adaptive deep brain stimulation: Retuning Parkinson's disease", *Elsevier* **184**: 273–284. <https://doi.org/10.1016/B978-0-12-819410-2.00015-1>
- [7] Alatas Bilal, Moradi Shadi, Tapak Leili, Afshar Saeid (2022), "Identification of Novel Noninvasive Diagnostics Biomarkers in the Parkinson's Diseases and Improving the Disease Classification Using Support Vector Machine", *BioMed Research International, Hindawi*
- [8] P. Raundale, C. Thosar and S. Rane (2021), "Prediction of Parkinson's disease and severity of the disease using Machine Learning and Deep Learning algorithm," *2021 2nd International Conference for Emerging Technology (INCET)*, pp. 1-5, doi: 10.1109/INCET51464.2021.9456292.
- [9] F. Cordella, A. Paffi and A. Pallotti (2021) "Classification-based screening of Parkinson's disease patients through voice signal," *2021 IEEE International Symposium on Medical Measurements and Applications (MeMeA)*, pp. 1-6, doi: 10.1109/MeMeA52024.2021.9478683.
- [10] Ali, L., Chakraborty, C., He, Z. et al. (2022) "A novel sample and feature dependent ensemble approach for Parkinson's disease detection". *Neural Comput & Applic*. <https://doi.org/10.1007/s00521-022-07046-2>
- [11] F. Huang, H. Xu, T. Shen and L. Jin (2021), "Recognition of Parkinson's Disease Based on Residual Neural Network and Voice Diagnosis," *2021 IEEE 5th Information Technology, Networking, Electronic and Automation Control Conference (ITNEC)*, pp. 381-386, doi: 10.1109/ITNEC52019.2021.9586915.
- [12] D. Trivedi H. Jaeger and M. Stadtschnitzner. (2019) "Mobile Device Voice Recordings at King's College London (MDVR-KCL) from both early and advanced Parkinson's disease patients and healthy controls." <https://doi.org/10.5281/zenodo.2867216>
- [13] M. Wodzinski, A. Skalski, D. Hemmerling, J. R. Orozco-Arroyave and E. Nöth, (2019) "Deep Learning Approach to Parkinson's Disease Detection Using Voice Recordings and Convolutional Neural Network Dedicated to Image Classification," *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pp. 717-720, doi: 10.1109/EMBC.2019.8856972.
- [14] T. J. Wroge, Y. Özkanca, C. Demiroglu, D. Si, D. C. Atkins and R. H. Ghomi, (2018), "Parkinson's Disease Diagnosis Using Machine Learning and Voice", *2018 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*, pp. 1-7, doi: 10.1109/SPMB.2018.8615607.
- [15] W. Wang, J. Lee, F. Harrou and Y. Sun, "Early Detection of Parkinson's Disease Using Deep Learning and Machine Learning," in *IEEE Access*, vol. **8**, pp. 147635-147646, 2020, doi: 10.1109/ACCESS.2020.3016062.

- [16] R. Alkhatib, M. O. Diab, C. Corbier and M. E. Badaoui, "Machine Learning Algorithm for Gait Analysis and Classification on Early Detection of Parkinson," in *IEEE Sensors Letters*, **vol. 4, no. 6**, pp. 1-4, June 2020, Art no. 6000604, doi: 10.1109/LSSENS.2020.2994938.
- [17] C. Ricciardi et al., "Machine learning can detect the presence of Mild cognitive impairment in patients affected by Parkinson's Disease," *2020 IEEE International Symposium on Medical Measurements and Applications (MeMeA)*, 2020, pp. 1-6, doi: 10.1109/MeMeA49120.2020.9137301.
- [18] X. Yang, Q. Ye, G. Cai, Y. Wang and G. Cai, (2022), "PD-ResNet for Classification of Parkinson's Disease from Gait," in *IEEE Journal of Translational Engineering in Health and Medicine*, **vol. 10**, pp. 1-11, 2022, Art no. 2200111, doi: 10.1109/JTEHM.2022.3180933.
- [19] A. U. Haq et al., "Feature Selection Based on L1-Norm Support Vector Machine and Effective Recognition System for Parkinson's Disease Using Voice Recordings," in *IEEE Access*, **vol. 7**, pp. 37718-37734, 2019, doi: 10.1109/ACCESS.2019.2906350.
- [20] Mei Jie, Desrosiers Christian, Frasnelli Johannes, (2021), "Machine Learning for the Diagnosis of Parkinson's Disease: A Review of Literature", in *Frontiers in Aging Neuroscience*, **vol. 13**, doi: 10.3389/fnagi.2021.633752.
- [21] <https://www.ppmi-info.org/access-data-specimens/download-data>
- [22] Amalia Luque, Alejandro Carrasco, Alejandro Martín, Ana de las Heras, (2019), "the impact of class imbalance in classification performance metrics based on the binary confusion matrix", *Pattern Recognition*, **Volume 91**, Pages 216-231, ISSN 0031-3203, <https://doi.org/10.1016/j.patcog.2019.02.023>.
- [23] J. R. Barr, M. Sobel and T. Thatcher (2022), "Upsampling, a comparative study with new ideas," *2022 IEEE 16th International Conference on Semantic Computing (ICSC)*, pp. 318-321, doi: 10.1109/ICSC52841.2022.00059.
- [24] Little, M.A., McSharry, P.E., Roberts, S.J. et al. (2007) "Exploiting Nonlinear Recurrence and Fractal Scaling Properties for Voice Disorder Detection" *BioMed Eng OnLine* **6 (23)**. <https://doi.org/10.1186/1475-925X-6-23>
- [25] Little, Max A et al. (2009) "Suitability of dysphonia measurements for telemonitoring of Parkinson's disease." *IEEE transactions on bio-medical engineering* vol. **56 (4)**: 1015. doi:10.1109/TBME.2008.2005954
- [26] D. Barrejón, P. M. Olmos and A. Artés-Rodríguez, "Medical Data Wrangling With Sequential Variational Autoencoders," in *IEEE Journal of Biomedical and Health Informatics*, **vol. 26, no. 6**, pp. 2737-2745, June 2022, doi: 10.1109/JBHI.2021.3123839.
- [27] Y. Guan (2021), "Application of logistic regression algorithm in the diagnosis of expression disorder in Parkinson's disease," *2021 IEEE 2nd International Conference on Information Technology, Big Data and Artificial Intelligence (ICIBA)*, 2021, pp. 1117-1120, doi: 10.1109/ICIBA52610.2021.9688135.
- [28] Anbalagan, B., Karnam Anantha, S. & Kalpana, R. (2022) "Novel Approach to Prognosis Parkinson's Disease with Wireless Technology Using Resting Tremors". *Wireless Pers Commun*. <https://doi.org/10.1007/s11277-022-09694-y>
- [29] Gupta, I., Sharma, V., Kaur, S., & Singh, A. K. (2022). "PCA-RF: An Efficient Parkinson's Disease Prediction Model based on Random Forest Classification". *arXiv preprint arXiv:2203.11287*.
- [30] D. Yadav and I. Jain (2022), "Comparative Analysis of Machine Learning Algorithms for Parkinson's Disease Prediction," *2022 6th International Conference on Intelligent Computing and Control Systems (ICICCS)*, pp. 1334-1339, doi: 10.1109/ICICCS53718.2022.9788354
- [31] D. V. Rao, Y. Sucharitha, D. Venkatesh, K. Mahamthy and S. M. Yasin (2022), "Diagnosis of Parkinson's Disease using Principal Component Analysis and Machine Learning algorithms with Vocal Features," *2022 International Conference on Sustainable Computing and Data Communication Systems (ICSCDS)*, 2022, pp. 200-206, doi: 10.1109/ICSCDS53736.2022.9760962.
- [32] Jin Huang and C. X. Ling, "Using AUC and accuracy in evaluating learning algorithms," in *IEEE Transactions on Knowledge and Data Engineering*, **17 (3)**, pp. 299-310, March 2005, doi: 10.1109/TKDE.2005.50.
- [33] Aghav-Palwe S, Mishra D (2020) "Statistical tree-based feature vector for content-based image retrieval." *Int J Comput Sci Eng* 2020 <https://doi.org/10.1504/IJCSE.2020.106868>
- [34] Aghav-Palwe, S., Mishra, D. (2019). "Color Image Retrieval Using Statistically Compacted Features of DFT Transformed Color Images. In: Bhatia, S., Tiwari, S., Mishra, K., Trivedi, M. (eds) *Advances in Computer Communication and Computational Sciences. Advances in Intelligent Systems and Computing*," **vol 760**. Springer, Singapore. https://doi.org/10.1007/978-981-13-0344-9_29