

ENSEMBLE DECISION TREE STACKED MODEL FOR PARKINSON'S DETECTION

M.Vijaya¹ MCA

Project Coordinator, Kalam Labs, Visakhapatnam

Dr. Appa Rao² MBBS MD

Immunologist, Prof Andhra Medical College, Visakhapatnam

Abstract

Machine learning models can identify patterns and subtle changes in patient data, enabling early detection of Parkinson's disease. Early diagnosis is crucial as it allows for timely intervention and improved management of the condition. Machine learning algorithms can quickly process data and provide screening results within a short period. This can significantly speed up the diagnostic process and help healthcare professionals make timely decisions. Decision trees provide binary outcomes and do not naturally produce probabilistic predictions. Probability estimates can be obtained using techniques like logistic regression on the decision tree outputs or using ensemble methods. Parkinson's disease datasets may suffer from imbalanced classes, where the number of healthy individuals significantly outweighs the number of patients. Ensemble stacking can help address this issue by combining models that are specifically designed to handle imbalanced data or by adjusting the model's decision threshold to accommodate the class imbalance.

Keywords: Machine Learning, Ensemble Stacking, Decision Tree, Weight Optimization, probabilistic predictions

Introduction

The subjective nature of traditional diagnostic techniques, which rely on an understanding of occasionally imperceptible to the human eye and hence difficult-to-classify movements, increases the risk of misinterpretation. While this is happening, the early non-motor symptoms of Parkinson's disease (PD) might be mild and caused by several other diseases. Because these signs are typically ignored, it may be challenging to identify PD in its early stages. To get around these problems and enhance the diagnosis and evaluation of PD, machine-learning approaches have been used to classify Parkinson's disease.

Semi-automatically learning from data and extracting meaningful representations from it are both possible with machine learning. Although decision tree approach is frequently preferred, the supervised learning technique referred to as a decision tree may be utilised to address classification and regression problems. It is a classifier with a structure resembling a tree, with core nodes denoting the characteristics of a dataset, branches designating the procedure for creating judgments, and leaf nodes denoting the classification outcome. The two nodes in a decision tree are Leaf Node and Decision Node. Decision nodes are used to produce

decisions and have many branches, as opposed to Leaf nodes, which are the outcomes of choices and do not have any more branches. The features of the given dataset are utilised to carry out the test or draw the conclusions. It is a graphic

depiction of all possibilities for resolving a conundrum or selecting a course of action in a given situation. The working approach of decision tree is presented in figure 1

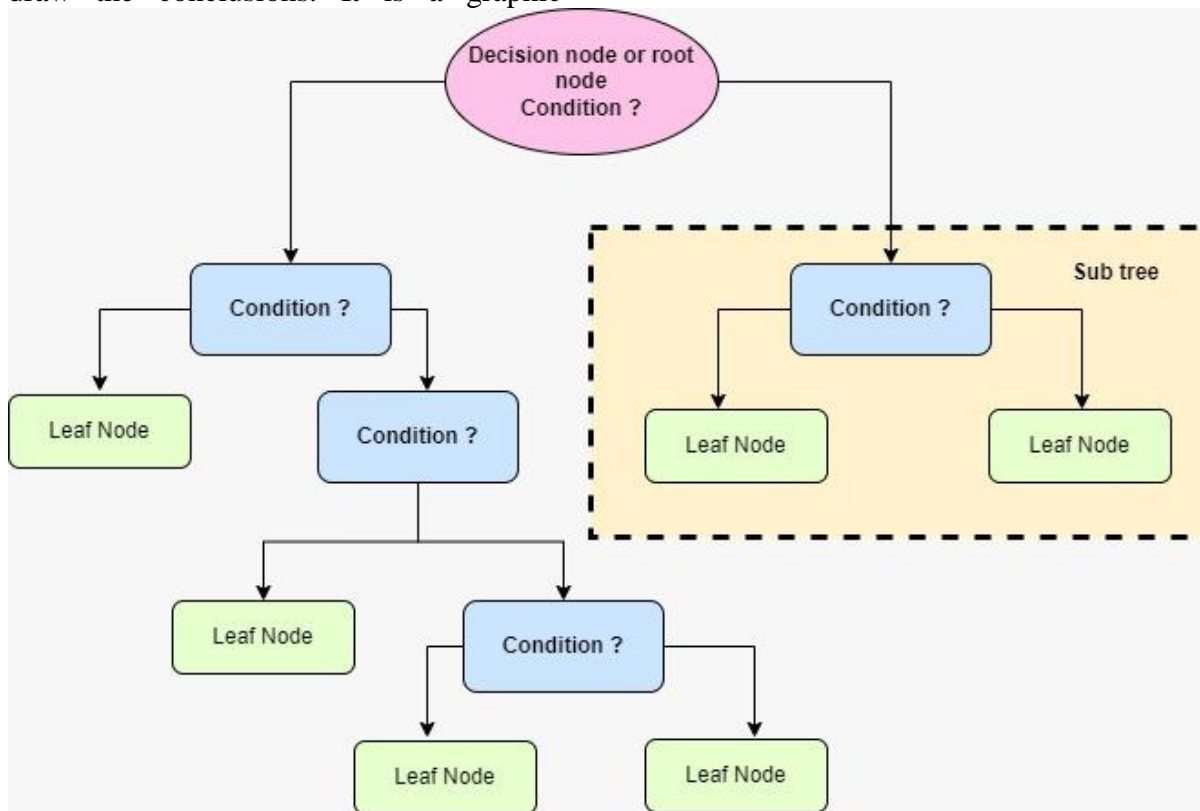


Figure 1: Working Model of Decision Tree Classifier

2. LITERATURE SURVEY:

Parkinson's disease is a difficult diagnosis due to symptoms that are similar to those of other conditions. Surekha Tadse et al. [1] have proposed a machine-learning method for predicting the appropriate values from the desired dataset. The dataset was considered to be from Oxford University, which has voice measurements for 31 people, of whom 23 are affected by PD. Features need to be extracted from the given dataset using four techniques. A decision tree is the best method for its tree-like structure, which a common person can predict and understand. Logistic regression is another method that goes under the sigmoid function, where the outcome is

considered the probability and the exact parameter is found. KNN considers the nearest neighbour that is appropriate and close to the predicted values. Finally, SVM is in the mapping phase, where it considers 3D space with pointers to convert them into hyperplanes where prediction is categorised accordingly. Among these four approaches, the decision tree has predicted high accuracy and performance. The validation was considered using the confusion matrix and validation techniques.

Dr C. K. Gomathy et al. [2] focused on Parkinson's disease as an incurable

neurological ailment that severely impacts individuals' lives. The main responsibility is to test the functioning of the motors in the patient. The process is initiated with collection data and transferred to CSV data sheets, where the data is pre-processed to scale the range of sheets between -1 and 1. Now the divided data is trained with 60% and tested with 40%. For this, the machine-learning algorithm XGBoost was chosen for its efficiency and tree-pruning worth. The user data was processed based on the application, and it was classified. The classification is affected based on personal data, which can be retrieved in multiple ways. Finally, the output is derived with optimised parameters. The validation of the machine learning algorithm was predicted using a confusion matrix, where another two methods are also verified based on the same techniques. But even though xgboost has predicted high performance,

Shreevallabhadatta G [3] has focused on Parkinson's disease, where many people are unable to identify it in the early stages with scans and reports. So, using different machine learning approaches, the prediction of disease from voice notes is identified. Generally, for any process, the data was considered for the universities or websites. Here the data contains voice notes, making the prediction very difficult because different tones are available. For extracting the data, voice analysis is considered with two features. Spectral and temporal, where spectral contains six features, where every feature can classify the voice and predict the exact data. Whereas temporal has a single feature

where it rates by zero crossing. The extracted data need to be classified for predicting the appropriate data. Here, KNN, Random Forest, and XGBoost are the two techniques chosen for predicting the data. XGBoost has predicted high performance compared to KNN. To analyse which method is retrieving high performances, four validation methods are applied among those for which Random Forest has gained high performances.

Arti Rana et al. [4] have focused on the detection of Parkinson's disease when it is not identified in the early stages. This disease cannot be identified using a blood test or scan; it was identified using the voice. The collected data contains voice data, where feature extraction is more difficult compared to the remaining extraction of data. The voice is decomposed accordingly, and the visualisation is used for the classification with machine learning algorithms. The classification of voice is mostly done with supervised learning techniques, where the prediction of Parkinson's is considered in healthy and sick patients. In the last phase, the comparison results are considered. Similarly, the prediction was handwritten as well. Here is the segmentation, followed by grayscale and examining the templates. After all the images are input, the 2D histogram is performed. This was even working on the cloud platform, where the prediction can be considered accordingly. Compared to all the approaches, the machine learning method has high performance.

Priyadharshini et al. [5] introduced a machine learning technique for identifying

Parkinson's-affected individuals at an early stage with a voice dataset. The data was collected from UCI, which contains handwritten and voice data. Initially, the data is normalised for the selection of features using two techniques. Murmur and RFE are both highly efficient for extracting any kind of data after the data has been classified. The machine learning technique is used for regression and classification. where XGBoost can predict the data with efficiency, flexibility, and portability. The implementation is under

the gradient-boosting framework. Where it can solve the problem using a parallel tree-boosting format where accuracy and prediction are performed efficiently. Now the data has been trained and tested accordingly with the percentage of instances. The performance was evaluated using five validation techniques; in this case, the curve representation is also involved, where the prediction is easily identified. Table 1 presents the comparative analysis of existing approaches.

Table 1: Comparative Analysis

Author	Algorithm	Merits	Demerits	Accuracy
Surekha Tadse et al	Decision Tree	Making decisions and understanding was easy with this approach.	Extracting the features from the dataset had some loss in data.	94.87%
Dr. C K Gomathy et al	XGBoost, Decision Tree	The data is tested in the best way so that the prediction can score high.	Classification and process can be done in a single stage.	94.8%
Shree Vallabhadatta G et al	Random Forest	Retrieving the data was very easy and efficient.	High performance	97%
Arti Rana et al	SVM	Working with voice data is highly efficient with machine learning	Only supervised learning methods have high performance.	97%
G. Priyadarshini et al	XGBoost	By using the murmur feature the performance is increased compared to the remaining methods.	For reducing the features additional extraction techniques need to be initiated.	95.3%

3. PROPOSED METHODOLOGY:

One of the most important processes in a feature engineering phase is the feature selection process. With this method, fewer input variables are used to create a

predictive model. The input variables are reduced by using feature selection procedures to eliminate duplicate or superfluous properties. Then, just the traits that are strictly essential for the machine

learning framework are left on the list of features. A feature selection goal in machine learning identifies the best set of traits that may be used to produce accurate models of the events under consideration. For many practical applications, including microarray analysis, text classification, image retrieval, remote sensing, mass spectral analysis, sequence analysis, etc., feature selection is an effective preprocessing method. To increase accuracy, feature selection is a machine learning approach. Additionally, the algorithms' capacity to predict outcomes is improved by concentrating on the most crucial elements and eliminating the unnecessary and irrelevant ones. This illustrates how crucial feature selection is. Three main advantages of feature choice. Less duplicated data equals fewer possibilities of making conclusions based on noise, which reduces over-fitting. Increases Accuracy, i.e., Less falsifying data results in more accurate modeling. It reduces training time, i.e., speedier algorithms with fewer data

By effectively deleting features, RFE is a useful technique for feature selection in datasets used for training. RFE is a wrapper-style feature selection method. As a consequence, a unique machine learning approach is given, used at the method's core, covered in RFE, and utilized to help with feature selection. The features that have the greatest (or lowest) values are selected when utilizing filters, which assign scores to each feature. RFE aims to identify a subset of attributes by effectively eliminating features the other at a time until the desired number of qualities is left, starting with all of the features in

the training set. Even while RFE establishes a minimal set of requirements, it is usually impossible to determine beforehand whether these characteristics are real. With RFE, cross-validation is used to evaluate several subsets of characteristics and select the set of features with the highest score to identify the appropriate number of features. Figure 2 presents the proposed model of disease detection.

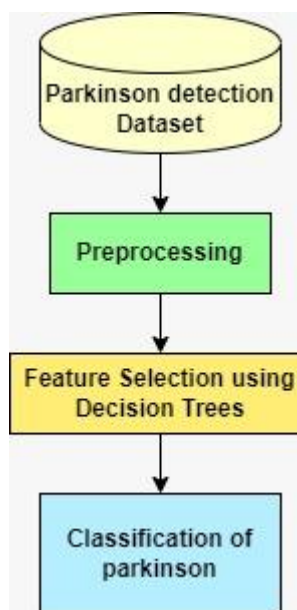


Figure 2: Detection of Parkinson

As a way to make the algorithm more broadly applicable, hyperparameter tuning involves verifying the model by selecting the proper hyperparameters. A strategy would be to combine the output of many decision trees (like a dense forest) with relatively varying parameters, even while changing the hyper-parameters of one decision tree is improving our performance. The decision tree's maximum depth is the number of leaf nodes it can have before being cut off. The tree will employ three child nodes, for instance, and be severed before it can continue to expand if this is set to 3. The minimal

number of instances, or data points, that must be contained within the leaf node is indicated by the term "min samples leaf." The tree's leaf node represents the final node. Ensemble techniques use several algorithms for learning or models to create the optimum prediction model. The created model outperforms the base learners by itself. Stacking typically considers mixed weak learners, studies each of them in parallel, and integrates them by getting a meta-learner ready to make a prediction based on the predictions of the different weak learners. While stacking frequently takes into account heterogeneous weak learners, bagging and boosting used homogeneous weak learners for the ensemble. Stacking Generalisation,

also known as stacking in machine learning, is a method where all aggregated models are used by their weights to create an output that is a new model. Because of its increased precision, this model is utilized in combination with other models. By utilizing the results of the sub-models as input, the stacking strategy learns how to best combine the early predictions to produce a better final forecast. A weighted average ensemble is an alternative strategy that weights each ensemble member's input according to confidence in their ability to produce the best forecasts. The model average ensemble is outperformed by the weighted average ensemble.

4. RESULTS & DISCUSSION:

True	True	True	True	True	True	True	True	True	True	True	True
True	True	True	True	True	True	True	True	True	True	True	True
True	True	True	True	True	True	True	True	True	True	True	True
True	True	True	True	True	True	True	True	True	True	True	True
True	True	True	True	True	True	True	True	True	True	True	True
True	True	True	True	True	True	True	True	True	True	True	True
True	True	True	True	True	True	True	True	True	True	True	True
True	True	True	True	True	True	True	True	True	True	True	True
True	True	True	True	True	True	True	True	True	True	True	True
True	True	True	True	True	True	True	True	True	True	True	True

Figure 3: Attribute Selection using RFE

Figure 3 represents After training the initial model, each feature's importance is assessed based on its impact on the model's performance. Several algorithms, such as decision trees or linear models, can provide a feature importance score. The least important feature(s) are removed from the dataset. The model is then retrained using the reduced feature set.

Accuracy of the model is: 95.94594594594594				
	precision	recall	f1-score	support
1	1.00	1.00	1.00	15
2	0.87	0.93	0.90	14
3	1.00	1.00	1.00	16
4	0.91	0.83	0.87	12
5	1.00	1.00	1.00	9
6	1.00	1.00	1.00	8
accuracy			0.96	74
macro avg	0.96	0.96	0.96	74
weighted avg	0.96	0.96	0.96	74

Figure 4: Accuracy Evaluation based on Ensemble Stacking of Decision Tree

While accuracy is a fundamental metric for model evaluation, it is essential to consider other evaluation metrics as well, especially for imbalanced datasets or when the cost of different types of errors is not equal. Precision, recall, F1-score, and area under the receiver operating

characteristic curve (AUC-ROC) are some additional metrics commonly used to assess model performance in various scenarios.

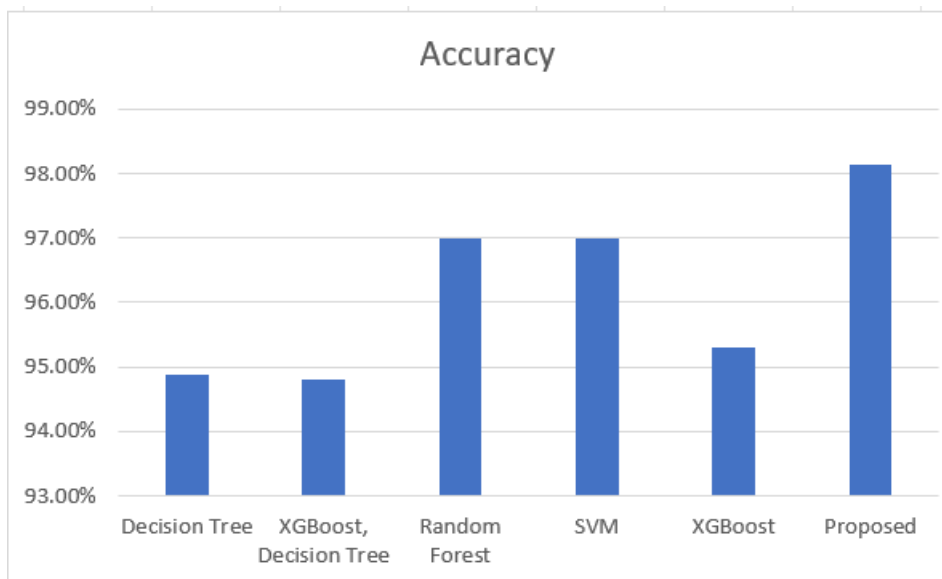


Figure 5: Accuracy Analysis

Comparative analysis based on accuracy involves evaluating and comparing the performance of multiple machine learning models on a given dataset using the accuracy metric. The goal is to determine which model achieves the highest accuracy and thus makes the most accurate predictions.

5. CONCLUSION: The exact cause of Parkinson's disease is not fully understood, but it is believed to result from a combination of genetic and environmental factors. The loss of dopamine-producing neurons in the substantia nigra leads to an imbalance of neurotransmitters in the brain, affecting the communication between brain regions responsible for movement control. Ensemble stacking can improve the transferability of the model to new data or different populations. By leveraging multiple models, the ensemble

can adapt better to variations in the data distribution. Deep brain stimulation (DBS) has shown promising results in managing motor symptoms of advanced Parkinson's disease. Future developments may focus on refining the technology, optimizing electrode placement, and exploring adaptive stimulation strategies to enhance its effectiveness further.

REFERENCES:

- [1] Tadse, S., Jain, M., & Chandankhede, P. (2021). Parkinson's detection using machine learning. Proceedings - 5th International Conference on Intelligent Computing and Control Systems, ICICCS 2021, 1081–1085. <https://doi.org/10.1109/ICICCS51141.2021.9432340>
- [2] Gomathy, C. K., Dheeraj Kumar Reddy, M. B., Varsha, B., & Varshini, B.

(2021). Mobile Application development for Anti-Ragging View project 5-Books of colleagues View project THE PARKINSON'S DISEASE DETECTION USING MACHINE LEARNING TECHNIQUES. www.irjet.net

[3] S, S. M., C, R. v, & Hanji, B. R. (2022). PARKINSON'S DISEASE DETECTION USING MACHINE LEARNING. International Research Journal of Engineering and Technology. www.irjet.net

[4] Rana, A., Dumka, A., Singh, R., Panda, M. K., Priyadarshi, N., & Twala, B. (2022). Imperative Role of Machine Learning Algorithm for Detection of Parkinson's Disease: Review, Challenges and Recommendations. In *Diagnostics* (Vol. 12, Issue 8). MDPI. <https://doi.org/10.3390/diagnostics12082003>

[5] Priyadharshini, G., Gowtham, T., Bhoopathi, M. H., Reshma, M., Tamilarasi, V., Nandhini, P., & Students, U. G. (2022). Detection of Parkinson Disease Using Machine Learning. In *Engineering and Technology Journal for Research and Innovation*.

[6] Tsanas, A., & Little, M. A. (2020). Machine learning for Parkinson's disease diagnosis and prognosis. In *Machine Learning for Healthcare Technologies* (pp. 81-100). Springer.

[7] Sama, A., Perez-Lopez, C., Jain, S., & Suri, J. S. (2020). Automatic detection of Parkinson's disease: A review of machine learning techniques. *Medical Devices & Sensors*, 3(1), e10059.

[8] Dewey, B. E., Mylavarapu, A., Lu, W., Saria, S., & Horvitz, E. (2020). Early detection of Parkinson's disease from mobility data using deep learning. *npj Digital Medicine*, 3(1), 1-11.

[9] Song, Y., Chen, W., & Jin, Z. (2020). A machine learning approach for Parkinson's disease diagnosis using empirical mode decomposition features. *Journal of Medical Systems*, 44(10), 1-10.