

DESIGN AND DEVELOPMENT OF AN ALGORITHM TO DETECT AND DIAGNOSE PARKINSON'S DISEASE

Singa Reddy Bhagya Lakshmi¹, Siri S Gowda², Trisha A³

^{1,2,3}Student, Dept of Electronics and Communication, R.V. College of Engineering, Bangalore, INDIA

¹singareddybl.ec17@rvce.edu.in, ²sirisgowda.ec17@rvce.edu.in, ³trishaa.ec17@rvce.edu.in

Abstract: Parkinson's disease is neurological disorder that affects the neurons that produce a chemical substance known as dopamine. The symptoms of this disease are tremors, stiffness all around the body, vocal impairment. If not treated early this disease might even lead to death. This disease doesn't have a treatment, but early prediction might help in the reduction of the progress of the disease. Detecting the disease is much more difficult as there is no quantitative test that can be conducted to detect the disease. Voice is one of the primary symptoms of this disease and therefore the features present in the voice can be extracted and can be used to train a model to detect whether the person is suffering from the disease or not. An ensemble learning voting classifier algorithm is designed with a performance accuracy of more than ninety percent and is used for the prediction of the disease using the vocal features extracted from the person. This algorithm is trained with the dataset which contains both normal as well as the affected person's voice features. Decision tree classifier algorithm, Logistic Regression and that of the Support Vector Machine Algorithm are used as the input for the voting classifier which is used for detecting the disease.

Keywords: Support Vector Machine, Voting Classifier, Parkinson's Disease, Ensemble learner

1. INTRODUCTION

Parkinson's disease is a neural dysfunction that influence the motion of human body. Symptoms are seen slowly, sometimes starting with a hardly noticeable shaking in one hand. Shakiness is common, but the dysfunction is also associated with a constrained movement. The face may show little or no expression in the beginning stage of the disease. The arm may not show much movement while walking. Also, the stammering may increase, the voice might become soft. The symptoms become more severe as the disease advances.

This disease causes some specific neurons in the brain to slowly die. Most of the symptoms are caused by the affinity of neurons that are responsible for producing a chemical called dopamine. An unusual brain malfunction happens which results in this symptom when this chemical's level falls. It has been estimated that there are 7 million people in India who are affected by this disease. There is currently no blood or laboratory test for this disease. The diagnosis is done by looking at the family history and reviewing the signs and symptoms shown by the patient.

This kind of diagnosis is not very efficient and sometimes this results in the ignorance of the early signs of the disease. It has been confirmed that Parkinson's Disease can affect the vocal ability of the patient. The speech of a diseased patient has change in the frequency spectrum in their voice because they lose the control of the limb, which decreases the frequency of the audio. Therefore, the main goal here is to detect the disease by considering speech as the parameter. There are many parameters of the vocal that are taken into consideration here, they are Jitter, Shimmer, Harmonic to Noise ratio and Noise to Harmonic ratio. A total of sixteen parameters are taken into consideration and the dataset that is utilized to train the model is taken from UCI Machine Learning Repository.

2. LITERATURE REVIEW

The problem of attribute election in case of Parkinson's Disease mechanized recognition using adaptive-based methods is mentioned in this paper [1], where the part to be highlighted is the accuracy over an evaluating set. All these methods increased the outcomes over the standard dataset, according to the experiments.

In this paper [2], Convolutional Neural Network were used to investigate Parkinson's Disease recognition from a vocal wave. Spectrograms and a variety of other attributes were used as input for Convolutional Neural Network. The influence of each segment on the Parkinson's Disease detection output was evaluated and compared to the decision level fusion of all segments case. One segment had a detection performance of twenty-nine percent, while the other segment had a detection performance of twenty percent. This suggests that some aspects of these recordings are more efficient than others in detecting Parkinson's Disease.

In this paper, the author gathered a wide range of vocal samples and different sounds from Person with the disease talking exercises [3].

It provides a chance to explore the validity of already present models. As a result of the analysis of the dataset, Sustained vowels were found to have more Parkinson's Disease selective data than separated words, which is inline with the findings published in the literature discussed in this paper. They found that representing a subject's samples with the mean and Standard Deviation are better than others. This method of representation seemed to be more precise.

In this paper, the author explained Parkinson's Disease using various Machine Learning and Deep Learning algorithms to differentiate people into the groups of healthy and Parkinson's Disease affected based on numerous signs in order to come up with an effective way to diagnose Parkinson's Disease [4]. The outcomes of numerous studies were contrasted using various methods, it was determined that Deep Learning is the best method for studying two main symptoms: twisted walking style and speech disability. The data was collected from the UCI Irvine Machine Learning

Repository. Two modules, Vascular Endothelial Growth Factor Spectrogram Detector using Convolutional Neural Network and Voiced disfigurement using Artificial Neural Network, have been implemented to differentiate PD patients based on the signs gait and speech disability, with an accuracy of eighty eight percent and eighty nine percent for the two modules on the testing dataset, respectively, and compared with three algorithms, Extreme Gradient Boost, Support Vector Machine and Multilayer perceptron.

Employing auto-immune voice

biomarkers as attributes, automated machine learning architectures

can diagnose and predict the disease [5]. The research presented in this paper compares the success of different Machine Learning classifiers in disease detection with noisy and multi-scaled data. Clinical level precision is possible after careful feature selection. These findings are encouraging because they could pave the way for new ways to use voice data to evaluate patient well-being and neurological diseases.

In this

paper, the authors investigated how to differentiate Person with Parkinson's disease from normal people using vowel phonations using a wide variety of old and new algorithms for testing [6]. In recent years, this binary discrimination issue has piqued interest, with the best

results showing a classification accuracy of around ninety three percent on a subset of twenty-two features. They showed that using ten defective measures, we can achieve nearly ninety nine percent accuracy. They used a speech data and added several newly suggested defects of speech steps that had never been used in this application before. They also explored Radio frequency. They see this research as a first step towards a bigger goal of developing technologies for treatment plan in Parkinson's Disease. On the basis of comprehensive Curriculum vitae tests, the algorithms used in this study tend to be very successful at distinguishing diseased from people who are not affected by it.

This research looked at a stage-variant, medication-variant category of diseased patients, suggesting that the findings apply to Parkinson's Disease patients at all stages of the disease.

[7]. Several aspects of speech deterioration have been linked to disease progression among people with Parkinson's Disease. This paper shows whether nonlinear dynamic studies reveal a connection between phonatory pathology and the seriousness of Parkinson's Disease. Overall, the findings indicate that nonlinear dynamic analysis may be a useful new approach for studying Parkinson's vocal pathology, complementing conventional voice analysis methods.

In this paper the mentioned procedure is to try and gain a comparatively better accord and consistency in the process for practical assessment of pathologic sounds [8]. Because of the wide range of methods used to determine functional effects, systematic analysis of the outcomes of voice treatments are usually constrained, if not impossible. A multidimensional set of basic computation is suggested that can be used to diagnose all common dysphonia. Perception, videostroboscopy, acoustics, aerodynamics, and subjective patient rating are the five different approaches. Instrumentation is held to a minimum, but it is considered essential for phono surgery professionals. A skilled and certified speech therapist will assist the Ear-Nose-Tongue specialists surgeon in conducting this simple series of measurements.

In this paper [9], the author obtained a number of sound tracks from various subjects uttering the vowel 'a' to improve the assessment of Parkinson's Disease. Many frames in the extracted Mel Frequency Cepstral Coefficient from different participants take up the most processing time during the classification process, preventing accurate recognition.

In this paper, numerous voice track refining algorithms were used to draw out information for Parkinson's Disease evaluation in this report, and the obtained attributes were given to the algorithms to build accurate resolution support structures. For feature extraction, Tunable Q-factor wavelet transform was applied to the voice tracks of Parkinson's Disease patients [10]. They collected voice recordings from two hundred and fifty-two people and uprooted various characteristic subcomponents from the tracks. Multiple classifiers are given to the attribute subcomponents, and the classifier prognosis are merged with ensemble learning techniques. The result of their thesis shows that Tunable Q-factor wavelet transform is the better for PD diagnosis.

3. COMPARATIVE ANALYSIS OF CLASSIFIER ALGORITHMS

The classification algorithms that are taken into consideration to study how they perform for the inputted dataset are Decision Tree Classifier, Support Vector Machine, Logistic Regression, Naïve Bayes and K Nearest Neighbor algorithms. These algorithms are designed using a software application called as Google Colaboratory. The dataset that is taken from UCI ML Repository, is first pre-processed and then is inputted into this algorithm to obtain the performance accuracy. Sixteen parameters are taken into consideration for prediction purpose. The accuracy value that was obtained is shown in the Table 1.

Table 1. Comparison of classifier algorithms with respect to accuracy

| Classifier Algorithms | Performance Accuracy (in percentage) |
|------------------------------------|--------------------------------------|
| Support Vector Machine Algorithm | 89.79 |
| Decision Tree Classifier Algorithm | 89.79 |
| Logistic Regression Algorithm | 81.86 |
| Naïve Bayes Algorithm | 63.26 |
| K Nearest Neighbor Algorithm | 80.93 |

From this it can be deduced that for this particular dataset support vector machine, decision tree classifier algorithm and logistic regression algorithm have outperformed Naïve Bayes and K Nearest Neighbor Algorithm. Therefore, these three algorithms will be used with ensemble learning voting classifier algorithm in order to create an algorithm which is much more precise in detecting the disease with the given dataset.

4. PERFORMANCE COMPARISON OF ENSEMBLE LEARNERS

Multiple algorithms are combined together and each of their decisions are taken into consideration while making the final prediction. Because a multiple combined model will definitely perform better than a single classifier algorithm. Diversification always leads

to betterment in accuracy and prediction and same goes for machine learning models. Diverse set of classifier algorithms taken into consideration to make a stronger model where the error of wrong prediction is almost zero to none. Three ensemble learners are taken into consideration for the comparison purpose.

Random forest is a bagging technique which has decision tree classifier algorithm as its base estimator. AdaBoost algorithm uses Support Vector Machine algorithm as its base learner and executes the dataset. XGBoost algorithm is also known as Extreme Gradient Boosting algorithm is one of the most efficient boosting algorithms presents at the moment. Voting classifier involves summing up of the predictions made by the weak learner and take the average of the said algorithms. For hard voting, majority of the votes is taken into consideration.

Voting, Bagging and Boosting are the three types of ensemble learners present. All the algorithm with their respective and suitable base learner is applied to the dataset and the accuracy is obtained. This accuracy value shows the performance efficiency of the said algorithm.

It is shown in the below Table 2. It can be inferred from the table that Voting classifier is the most efficient performance wise. Therefore, this algorithm along with its base learner, that is Decision Tree Classifier, Support Vector Machine and Logistic Regression is chosen for integration with the front end in order to predict whether the person is affected by the disease or not.

Table 2.
Performance efficiency of different ensemble learner algorithms

| Ensemble Learner Algorithm | Accuracy (in percentage) |
|-----------------------------|--------------------------|
| Random Forest Algorithm | 90 |
| AdaBoost Algorithm | 77.55 |
| XGBoost Algorithm | 88.59 |
| Voting Classifier Algorithm | 95 |

5. IMPLEMENTATION

The completed designing of the algorithm was done using the software application called as Google Colaborative which is provided by the technology company named Google. The algorithm was designed using three classification algorithms which were then ensemble to the voting classifier which predicts the output depending upon majority of the output given by the weak classifier algorithm. All these three classifiers are considered as weak learners. These three classifier algorithms are combined and are put with the ensemble voting classifier. The voting method that is used here is called as the hard voting.

The dataset that is obtained from UCI ML Repository has the voice measurement, a total of thirty-one people, out of which twenty-three were affected by the disease. Each column represents a feature of the voice extracted. The last column represents whether the person is affected by the disease or not. If the last column is zero, it means normal person. If the last column is one, then it means the person is affected by the disease.

Data pre-processing is done in order to scale the values of the dataset and to fill in

the missing values if any is present within the file. Importing of the required the libraries are done. Later the entire dataset is divided into testing and training set. Normally it will be

eighty percent for training set and the remaining is for the test set. The training set is inputted to the model in order to make it understand how exactly the prediction has to be done. Once the training is done, the test dataset is inputted and the output of the model, that is the predicted value is compared with the known outcome of the test dataset and the performance accuracy value of the algorithm is obtained. The confusion matrix is also obtained to understand in detail on how exactly the designed algorithm is performing with the values given. The entire methodology of algorithm designing is shown in the Figure 1.

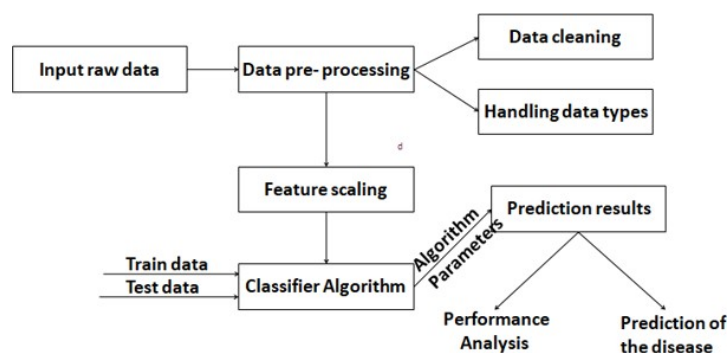


Figure 1. Methodology

This algorithm after the completed design is implemented in an application called as Spyder IDE. Later flask is used for creating a web application which can input the data. The same application is used for the integration of the front-end model and that of the web pages. The prediction done by the algorithm for the sixteen features inputted is shown in the Figure 2. Here the binary one means the person is suffering from the Parkinson's disease and binary zero means the person is normal.

```

+ Code + Text
print(Voting_Classifier.predict([[330.721,337.175,315.198,0.00341,0.000010,0.00168,0.00177,0.00503,
[0]

[ ] print(Voting_Classifier.predict([[269.335,288.624,254.342,0.01026,0.000038,0.00437,0.00509,0.01310
[0]

[ ] print(Voting_Classifier.predict([[266.719,268.513,264.069,0.00150,0.0000563,0.00078,0.00097,0.0023
[0]
  
```

Figure 2. Prediction of the disease

The values of the voice features are taken from a software called as Praat which is used for the analysis of the vocal phonetics. It is a free application, where the monotone voice is recorded multiple times and the features required for the algorithm are extracted. But before that voice feature is sampled and the external environment noises that are present are removed as much as possible. Even though one hundred of it can't be removed, it will still increase the accuracy in which the prediction is happening. Below Figure 3 shows the file which has the extracted values.

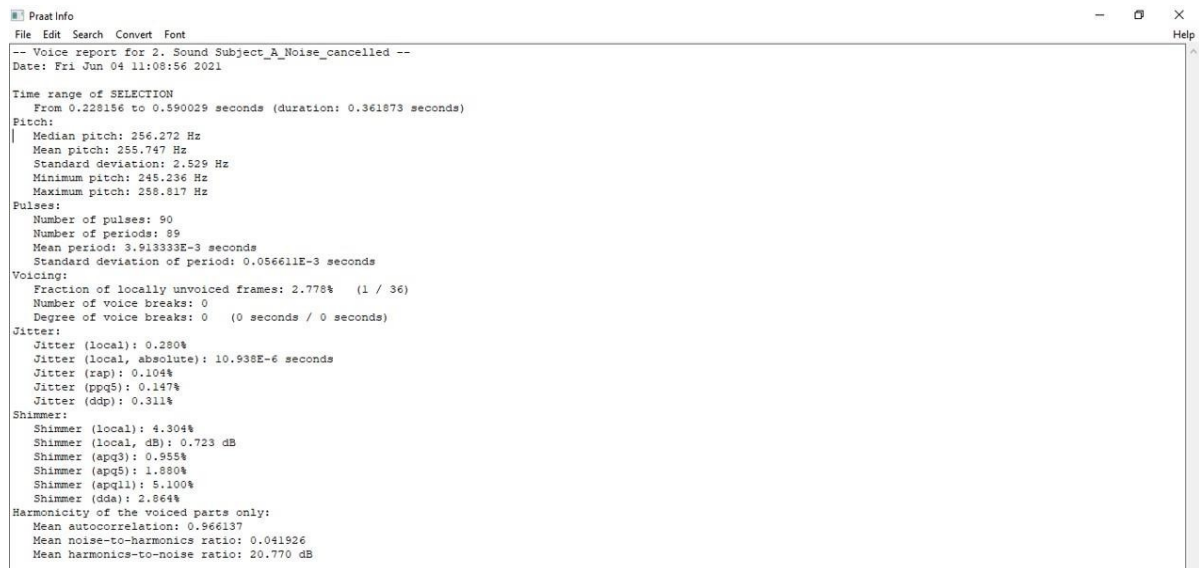


Figure3.Extractedfeaturesvalues

6. CONCLUSIONANDFUTURESCOPE

One of the earliest symptoms of Parkinson's disease is the distortion in the voice, it is much different from the normal voice. The pitch variation in the diseased person will be less when compared to the normal person. Therefore, the features present in the voice are extracted. A dataset is used to train the designed algorithm and the extracted value is used as the input for prediction. The designed voting classifier algorithm with that of the decision tree, support vector machine and logistic regression as the input has a performance accuracy of ninety-five percent. The reason to choose three as the base estimator is because they have an accuracy performance that beats the rest of the weak classifier algorithm. And therefore, the prediction will happen with little to no error.

There are multiple other features that can be taken into consideration for making the prediction. These are all called as non-linear features and their extraction requires a much more sophisticated equipment and an almost noiseless environment. A hybrid algorithm can be designed with two ensemble learners combined together. Wherein depending upon the subset of features chosen, the ensemble learner will be chosen.

REFERENCES

- [1] A. e A. Spadoto and R. C. Guido, "Improving Parkinson's Disease Identification Through Evolutionary-Based Feature Selection.", IEEE, 2011.
- [2] A. G. Evaldas Vaiciukynas, "Parkinson's Disease Detection from Speech Using Convolutional Neural Networks.", Kaunas University of Technology, Studentu 50, 51368 Kaunas, Lithuania, 2018.
- [3] M. E. I. Betul Erdogdu Sakar, "Collection and Analysis of a Parkinson Speech Dataset with Multiple Types of Sound Recordings.", IEEE Journal of Biomedical and Health Informatics, Vol. 17, No. 4, 2013.
- [4] A. J. Shivangi and A. Tripathi, "Parkinson Disease Detection Using Deep Neural Networks.", IEEE, 2019, isbn: 978-1-7281-2360-8.
- [5] C. D. Timothy J. Wroge, "Parkinson's Disease Diagnosis Using Machine Learning and Voice.", IEEE, 2018, isbn: 978-1-5386-5917-5.
- [6] A. Tsanas and P. E. McSharry, "Novel Speech Signal Processing Algorithms for High Accuracy Classification of Parkinson's Disease.", IEEE Transactions on Biomedical Engineering, vol. 59, no. 5, 2012.
- [7] M. C. Douglas A. Rahn III, "Phonatory Impairment in Parkinson's Disease: Evidence from Nonlinear Dynamic Analysis and Perturbation Analysis.", University of Wisconsin-Madison, 2007.
- [8] P. B. Philippe H. Dejonckere, "A basic protocol for functional assessment of voice pathology, especially for investigating the efficacy of (phonosurgical) treatments and evaluating new assessment techniques Guideline elaborated by the Committee on Phoniatrics of the European Laryngological Society (ELS).", Eur Arch Otorhinolaryngol, 2001.
- [9] A. J. Achraf Benba, "Voiceprint's analysis using MFCC and SVM for detecting patients with Parkinson's disease.", IEEE, 2015.
- [10] G. S. C. Okan Sakar, "A Comparative Analysis of Speech Signal Processing Algorithms for Parkinson's Disease Classification and The Use of The Tunable Q-Factor Wavelet Transform.", Applied Soft Computing (APPLSOFTCOMPUT), 2018.