# PARKINSON'S DISEASE PREDICTION USING XGBOOST, RANDOM FOREST AND CATBOOST

Umang Tiwari

CSE department,

MAIT,

Rohini, Delhi

umangtiwari2604@gmail.com

Himani Sheoran

CSE department,

MAIT,

Rohini, Delhi

sheoran.himani@gmail.com

Udit Jain

CSE department,

MAIT,

Rohini, Delhi

jain30udit@gmail.com

Zameer Fatima,

Assistant Professor,

MAIT, Rohini, Delhi

zameerfatima@mait.ac.in

**ABSTRACT:**

We are predicting Parkinson's disease with the assistance of voice Dataset efee which allows to treat the human beings in early stages. Parkinson's disease is a neurological sickness that ends in shaking and issue in strolling, balance, and coordination. In worst instances, sufferers have super problem on foot or status even they're not capable to live by themselves and require a wheelchair to move around an assistance is wished in all each day activities.

Besides motor signs and symptoms, the character may additionally see, pay attention, or revel in things that are not real (hallucinations), or agree with matters

that aren't genuine (delusions). Parkinson's disease patients normally have a low-quantity voice with a monotone high-quality.The speech sample of Parkinson's affected person is often produced in quick bursts with beside the silences between words and long pauses earlier than initiating speech. The voice dataset have the functions like MDVP:Fo(Hz)- average vocal fundamental frequency, MDVP:Fhi(Hz)- most vocal essential frequency, jitter, Simmerand so forth.

First, we balanced the data using SMOTE (synthetic Minority Oversampling approach) and then train and test one-of-a-kind model like Random Forest, Cat boost, XGboost and tuned the hyperparameter with the assist of GridSearchCV.

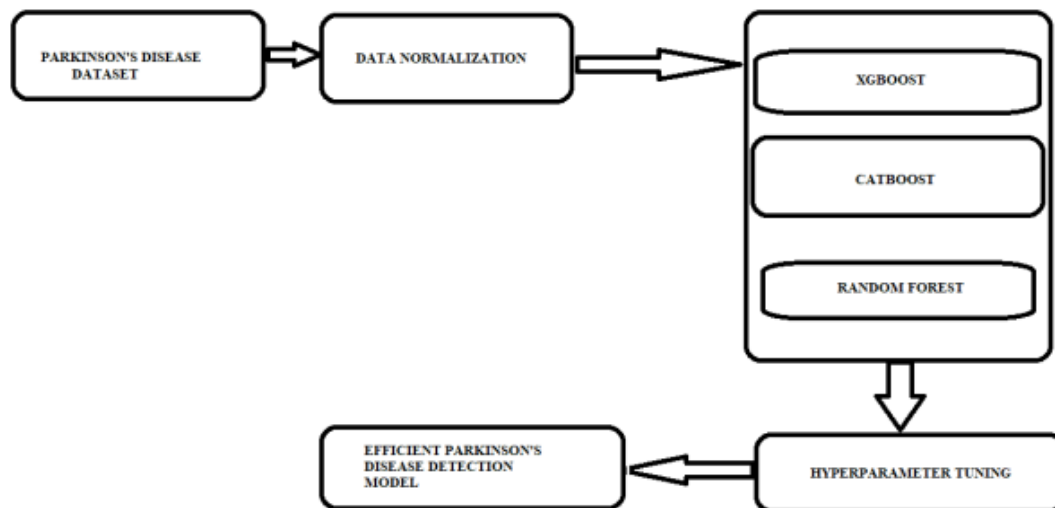On assessment we located that Cat boost showing the higher accuracy – 96.6%

and high Matthews Correlation Coefficient (MCC) – 91.43% among these kinds of fashions.

## INTRODUCTION:

Parkinson's ailment is a frightened sickness that influences motion. every so often symptoms starting with a barely substantial tremor in only one hand.

Tremors on this sickness are not unusual, however additionally usually reasons stiffness or slowing of movement. In the early stages of Parkinson's sickness, someone may additionally display very little expression or fingers won't swing even as walking and speech may come to be gentle or slurred. In Parkinson's disorder, sure nerve cells (called as neurons) inside the mind gradually spoil down or die. Most of the symptoms are because of a lack of neurons that produce a chemical messenger to your mind called dopamine. When dopamine tiers starts decrease, it reasons strange brain interest, leading to impaired movement and other signs of Parkinson's ailment. Parkinson's sickness causes vocal impairment that consequences speech, motor abilities, and different functions. Subsequently, on this paper, there is a try and explore a higher machine gaining knowledge of primarily based version for an early detection of Parkinson's ailment from the voice samples. From the above, it is able to be discovered that various device mastering strategies had been carried out in latest research works over voice primarily based Parkinson's disorder detection but it can be discovered that in none of these works the SMOTE (Artificial Minority Oversampling method) is used motive it is imbalanced dataset first we need to balance it after which train our version so the model can teach at the minority elegance and are expecting it more appropriately.

## METHODOLOGY:



Dataset element:

Dataset is collected from UCI internet site. This dataset has 195 precise values and 24 columns.

Matrix column entries (attributes):

call - ASCII problem name and recording quantity

MDVP:Fo(Hz) - average vocal essential frequency

MDVP:Fhi(Hz) - most vocal essential frequency

MDVP:Flo(Hz) - minimum vocal essential frequency

MDVP:Jitter(%),MDVP:Jitter(Abs),MDVP:RAP,MDVP:PPQ,Jitter:DDP– Severalmeasures of variation in fundamental frequency

MDVP:Shimmer,MDVP:Shimmer(dB),Shimmer:APQ3,Shimmer:APQ5,MDVP:APQ,Shimmer:DDA – numerous measures of variation in amplitude

NHR,HNR - measures of ratio of noise to tonal components within the voice

repute - fitness status of the situation (one) - Parkinson's, (zero) - healthy

RPDE,D2 - two nonlinear dynamical complexity measures

DFA - sign fractal scaling exponent

spread1,spread2,PPE - 3 nonlinear measures of essential frequency variation

records Preprocessing:

This step includes two process which is Normalization and balancing the dataset that is give an explanation for in element

below:

**Normalization:**

Normalization is a method that's applied as a part of information preparation . The want of normalization is to trade the values of numeric columns within the dataset to a common scale, without changing variations inside the levels of values.

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

where xnew is a selected function represented with the aid of a column in the dataset, x is a value of this column.

The minimum cost of the column is represented as xmin and the most value of the column is xmax.

**Stability Dataset:**
Balancing of the dataset is required while there is a class in minority which made the dataset made the dataset biased towards the alternative elegance.

We balanced our facts the use of SMOTE (synthetic Minority Oversampling method) with the assist of imbalance library.

## MODELLING AND ANALYSIS:

**Models:**
Following are the fashions that we used:

### XGBOOST:

XGBoost is a gradient boosting library.
It facilitates to implements gadget studying algorithms under the Gradient Boosting framework.
XGBoost is a parallel tree boosting which solves many machine studying problems in a fast and easy way.
The same code runs on distributed environment and solves many systems getting to know problems.

### CATBOOST:

CatBoost is currently an open-sourced system studying algorithms evolved by way of Yandex.

It reduces time spent on Hyperparameter tuning, due to this fact CatBoost presents awesome results with default parameters.

It enables to improve your schooling results that allows you to use non-numeric factors, instead of having to pre-system your records or spend effort and time turning it to numbers.

### RANDOM FOREST:

Random Forest is the ensemble approach that work on the big numbers of selection tree. every character tree inside the random offers a category prediction and the class with the maximum votes will become our version's prediction.

**Hyperparameter Tuning:**

The purpose of hyperparameter tuning is to get the satisfactory possible parameter for our model.

We did Hyperparameter tuning with the assist of GridSearchCV cause it searches for best set of hyperparameter from a grid of hyperparameter values.

**MATHEW CORRELATION COFFECIENT (MCC):**

The Matthews correlation coefficient (MCC) or phi coefficient is used to measure of the fine of binary classifications, brought by biochemist Brian W. Matthews in 1975.

The variety of values of MCC lie between -1 to +1. MCC takes all of the four fee of values of confusion matrix into consideration.

If the MCC value is close to 1 way that each lessons are predicted nicely.

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$$

in which TP is the actual effective, TN is the real negative, FP is the false nice and FN is the false poor

**RESULTS AND DISCUSSIONS:**

On assessment we found that CatBoost algorithm model has the higher accuracy-96.6% and MCC-91.4% as compared to the models like XGBoost and Random Forest.

Following are the accuracy and MCC value of different algorithm.

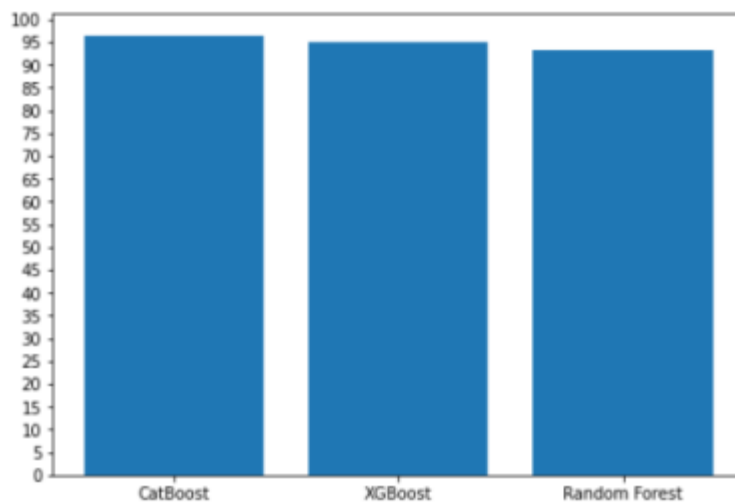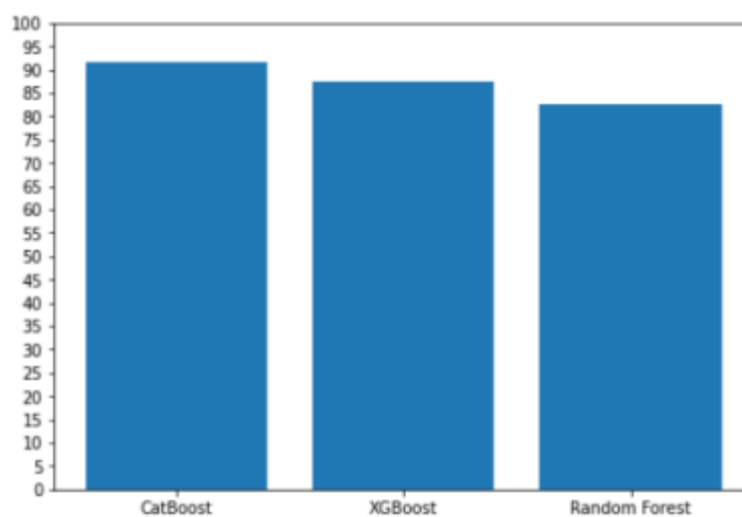| Algorithm | MCC Value | Accuracy |
|-----------|-----------|----------|
| Catboost | 91.4 | 96.6 |
| Xgboost | 87.4 | 94.9 |
| Random Forest | 82.4 | 93.2 |

**Fig.-1:** Accuracy comparison



**Fig.-2:** MCC comparison

## CONCLUSION:

Early detection of Parkinson's disease is very helpful to ensure timely treatment of the patients before it reaches to severe stage.

From this study we analyzed the different machine learning algorithms like CatBoost, XGBoost and Random Forest and got an efficient Parkinson's Disease prediction model with highaccuracy-96.6% and high MCC-91.4% which will help to predict Parkinson before getting it to worst stage.

**REFERENCES:**

[1] M. A. Little, P. E. McSharry, E. J. Hunter, J. Spielman, and L. O. Ramig, "Suitability of dysphonia measurements for telemonitoring of Parkinson's disease," IEEE Trans. Biomed. Eng., vol. 56, no. 4, pp. 1010–1022, 2009.

[2] Bhattacharya, I., & Bhatia, M. P. S. (2010, September).SVM classification to distinguish Parkinson disease patients. In Proceedings of the 1st Amrita ACM-W Celebration on Women in Computing in India (p. 14).ACM.

[3] Leandro A. Passos, Parkinson Disease Identification using Residual Networks and Optimum-Path Forest, SACI 2018, IEEE 12th International Symposium on Applied Computational Intelligence and Informatics, May 17-19, TimiÃ°oara, Romania.

[4] Deepak Gupta, Optimized cuttlefish algorithm for diagnosis of Parkinsons disease, Cognitive Systems Research, Volume 52, December 2018, Pages 36-48https://doi.org/10.1016/j.cogsys.2018.06.006.

[5] Mathur, R., Pathak, V., & Bandil, D. (2019). Parkinson Disease Prediction Using Machine Learning Algorithm. InEmerging Trends in Expert Applications and Security (pp.357-363). Springer, Singapore.

[6] Benba, A., Jilbab, A., Hammouch, A., & Sandabad, S.(2015, March). Voiceprints analysis using MFCC and SVM for detecting patients with Parkinson's disease. In2015 International conference on electrical and information technologies (ICEIT) (pp. 300-304). IEEE