

Multimodal Parkinson's Disease Risk Assessment Application Based On Vocal and Improved Spiral Test

Chan Yi Jie Kelvin A0178430E, Gopa Sen A0178297J, Han Yuen Kwang Andy A0178227X, Lim Pier A0178254X, Teresa Cheng Siew Loon A0178510H

ABSTRACT

The paper proposes a multi-modal approach combining voice and image test for early detection of Parkinson disease (PD). Research studies done earlier have used data related to either voice or spiral drawing to detect PD. However, different people experience different symptoms and different levels of severity of PD. Hence in this paper, we propose a multi-modal approach to enhance the reliability of identifying PD patient. Additionally, we propose to implement the multi-modal approach into a touch-enabled smartphone-based application to carry out preliminary PD tests conveniently, without the need of supervision of additional medical personnel or any specialized equipment. To substantiate our idea, we have evaluated both voice and spiral test data using various machine learning models. The results based on the two types of dataset demonstrate an excellent level of accuracy for PD identification.

Pairwise correlation and k-means clustering techniques are used to extract features from the vocal dataset. In this classification problem, the highest accuracy of 95.89% is obtained using an ensemble of 3 classification models.

The Pearson's correlation is used to extract features from the image dataset. The best accuracy of 99.6% is achieved using the k-Nearest Neighbors classifier in the Dynamic Spiral Test (DST). An accuracy of 98.8% and 94.9% are achieved using the Logistic Regression classifier and the Adaptive Boosting classifier on the Static Spiral Test (SST) and Stability Test on Certain Point (STCP) respectively. A second ensemble making use of results from DST, SST, and STCP will provide the overall result of the spiral test.

The final ensemble for the application makes use of the results of the respective ensemble from the vocal and spiral test.

1. INTRODUCTION

According to statistics, there are about 60,000 Americans diagnosed with PD each year while there are more than 10 million people worldwide with PD [1]. The probability of getting Parkinson's disease increases with age with approximately four percent of the PD patients worldwide are under the age of 50 [1]. It is estimated, that the number of Parkinson's disease patients will increase from 4.1 million in 2005 to 8.7 million in 2030 [2].

There are a few measurement techniques [3] for PD, such as Unified Parkinson's Disease Rating Scale (UPDRS), the Hoehn-Yahr Scale, the Schwab and England Scale of Activities of Daily Living and the Parkinson's Disease Questionnaire 39. Among them, UPDRS is the most common technique used to follow the PD progression and evaluate the results of surgical, medical, and other interventions of the disease. Speech characteristics are key components for early-stage detection where vocal impairment was detected in approximately 90% of the patients who are in their earlier stages of the disease [4]. Hence, there is an increasing interest in building PD diagnostic and telemonitoring systems based on vocal features.

Lately, researchers found out that diagnostic tools such as biological and genetic biomarkers as well as imaging techniques have high accuracies in predicting PD too [3]. In 2017, a group of Australian scientists developed a new diagnostic method called the Spiral test [5]. The test could detect signs of the disease by analyzing the amount of time used to draw spirals, amount of pressure exerted and the characteristics of the lines.

Currently, there are more than eight apps in the market [6]. Some are speech therapy apps, while some apps track chronic conditions of Parkinson patients' movements such as tremor, balance, and gait using sensitive accelerometers (sensors) which are available in most smartphones. The patient can choose to hold the smartphone equipped with the app in their hand or apply the device to their ankle for as long as 30 seconds

to allow the accelerometer of the tablet or smartphone to record movement.

This paper, on the other hand, proposes a smartphone application to detect early onset of PD utilizing two of the early symptoms displayed by People with Parkinson's (PWP): vocal test and spiral test. The app will also enable remote collection of data from PWP to facilitate further research. Two datasets have been used in this paper - one vocal test dataset and one spiral test dataset to substantiate the suitability of these tests in the application.

The paper is organized as follows, Section 2 will present the related work on the vocal and spiral test datasets for the detection and monitoring of PD. Section 3 will provide a description of both datasets as well as the data exploration done. In addition, the proposed algorithm and its configurations are also presented in this section. Section 4 will present the classification results from the analysis done on both datasets. Lastly, the proposed smartphone application will be discussed in Section 5 and the conclusion will be presented in Section 6.

2. RELATED WORK

Vocal Data Set

Studies in the area of using vocal test to detect Parkinson disease can be categorized into two main groups: to ascertain which are the most effective vocal features [7-11] to detect Parkinson's disease, while some studies focus on enhancing classification accuracy (to distinguish between healthy and unhealthy subjects) [12 -19].

"Accurate telemonitoring of Parkinson's disease progression by noninvasive speech tests" [11] focused on features extraction using signal processing techniques applied on voice samples taken from 42 patients with early-stage PD to estimate the unified Parkinson's disease rating scale (UPDRS) using linear and nonlinear regression. Their results show an accuracy difference of about 7-point from clinical UPDRS estimations.

"Novel speech signal processing algorithms for high-accuracy classification of Parkinson's disease" [12] applied feature selection techniques using 132 dysphonia measurements on 263 samples. The team obtained 99% overall classification accuracy. In another work, "Collection and analysis of a Parkinson speech dataset with multiple types of sound recordings" [10] mutual information-based selection algorithm was

applied using a permutation test to select features and rank them based on maximum-relevance-minimum-redundancy (mRMR) into an SVM classifier. Using leave-one-subject-out (LOSO) as the cross-validation technique of their model in order to avoid bias. Their approach achieved 92.75% classification accuracy.

Spiral Test Data Set

This area of study is relatively new, it is often used to ascertain if subjects have PD [20 – 23].

In "Digitized Spiral Drawing: A Possible Biomarker for Early Parkinson's Disease" [21], the analysis was done on drawn spirals, capturing kinematic, dynamic, and spatial abnormalities. These were used to calculate indices that quantify motor performance and disability. Linear mixed effect models adjusting for age, gender and handedness were used to compare the different indices between cases and controls. Results show that spiral indices yielded good accuracy in diagnosing early PD from cross-validation studies.

"A new quantitative evaluation method of spiral drawing for patients with Parkinson's disease based on a polar coordinate system with varying origin" [20] was carried out where seven characteristic parameters were extracted from hand movement in spiral drawing experiment. The characteristic value was linearly interpolated. The team identified a few good predictors of PD in this paper.

3. MATERIALS AND METHOD

3.1 Data

Two data sets, The Parkinson Speech Dataset with Multiple Types of Sound Recordings (Sakar, Isenkul et al. 2013) and Improved Spiral Test Using Digitized Graphics Tablet for Monitoring Parkinson's Disease, were obtained from the UCI Machine Learning Repository. The findings from these data sets will substantiate the proposed application for mobile detection and monitoring of Parkinson's disease.

3.1.1 Vocal Data Set

The PD dataset [24] used in this study, comprising of 26 speech samples of 20 PD patients and 20 healthy individuals, was collected by the Department of Neurology in Cerrahpasa Faculty of Medicine of Istanbul University. The 26 speech samples consist of sustained vowels, numbers, words and short sentences which are recorded using a Trust MC-1500 microphone with a frequency range between 50Hz and 13kHz. For

the collection of the speech samples, the microphone was placed at a distance of 10cm from the subject. A description of the samples is described in the table below.

Table 1: Description of speech samples collected

Sample No.	Description
1	Sustained vowel "a"
2	Sustained vowel "o"
3	Sustained vowel "u"
4 -13	Numbers from 1 – 10
14 - 17	Rhymed Short sentences
18 - 26	Words in Turkish language

From each of the samples collected, the Praat acoustic analysis software is used to extract 26 linear and time-frequency based features with reference to previous works in this field of study [9] [25]. The features of the PD dataset used in this study are as seen in the table below.

Table 2: Time-frequency-based features extracted from speech samples

Features	Group
Jitter (local) Jitter (local, absolute) Jitter (rap) Jitter (ppq5) Jitter (ddp)	Frequency Parameters
Number of pulses Number of periods Mean period Standard dev. of period	Pulse Parameters
Shimmer (local) Shimmer (local, dB) Shimmer (apq3) Shimmer (apq5) Shimmer (apq11) Shimmer (dda)	Amplitude Parameters
Fraction of locally unvoiced frames Number of voice breaks Degree of voice breaks	Voicing Parameters
Median pitch Mean pitch Standard deviation Minimum pitch Maximum pitch	Pitch Parameters
Autocorrelation Noise-to-harmonic Harmonic-to-noise	Harmonicity Parameters

In addition to the data collected from the 40 individuals stated above, a separate dataset was collected from a test group made up of a separate 28 PD patients. For this test group, only the sustained vowels "a" and "o" samples are collected from each patient instead of the 26 samples type collected from each of the 40 individuals.

3.1.2 Image Data Set

The image data set [26] consists of 15 healthy individuals (thereafter known as the control set) and 62 individuals diagnosed with Parkinson's disease. (thereafter known as People-with-Parkinson's (PWP)). For all subjects, three types of handwriting recordings, the static spiral test (SST), dynamic spiral test (DST) and Stability Test on Certain Point (STCP) are taken.

The Static Spiral Test is frequently used for clinical research for purposes like determining motor performance (Wang et al., 2008), measuring tremor (Pullman, 1998) and diagnosing PD (Saunders et al., 2008). In this test, three wound Archimedean spirals appear on the graphics tablet and patients are asked to retrace the same spiral as much as they can with a digital pen.

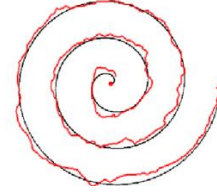


Figure 1: Static Spiral Test

The second test is the Dynamic Spiral Test which is similar to the Static Spiral Test. However, in this test, the spiral disappears and appears in certain time intervals, forcing the patient to keep the pattern in mind and continue to draw.

The third test is the Stability Test on Certain Point, in which subjects are asked to hold the digital pen on the point without the rest of the hand touching the screen in a certain time. The purpose of this test is to determine the patient's hand stability or hand tremor level.

Table 3: Attributes from Spiral Test

Attributes Collected	
X	X position (left/right) of the screen
Y	Y position (up/down) of the screen
Z	Z position of the pen (perpendicular to the screen)
Pressure	Pressure on the screen with the digital pen
Grip Angle	Individual user grip angle of the pen
Time	System time sample is recorded
Test ID	0 : Static Spiral Test (SST) 1 : Dynamic Spiral Test (DST) 2: Stability Test on Certain Point (STCP)

3.2 Data preprocessing

3.2.1 Vocal Data Set

Exploratory data analysis and specifically pairwise correlation were carried out to investigate if there were any features which are highly correlated to one another. From the investigation, it was found that there were certain features which have correlation coefficients of above 0.85 with other feature(s). The highly correlated features were thus not used in this study. In addition to the features with correlation above 0.85, Jitter (local, absolute) and Shimmer (apq11) were both removed during feature selection even though each only has a correlation of 0.74 with Jitter (ppq5) and 0.73 with Shimmer (local) respectively. It was tested and verified that the removal of Jitter (local, absolute) and Shimmer (apq11) did not have a significant impact on the results of the model.

The table below shows the features remaining after feature selection which are used in this study.

Table 4: Features remaining after feature selection

Features	Group
Jitter (ppq5)	Frequency Parameters
Number of pulses	Pulse Parameters
Standard dev. of period	
Shimmer (local)	Amplitude Parameters
Fraction of locally unvoiced frames	Voicing Parameters
Number of voice breaks	
Degree of voice breaks	
Median pitch	Pitch Parameters
Standard deviation	
Minimum pitch	
Maximum pitch	
Noise-to-harmonic	Harmonicity Parameters
Harmonic-to-noise	

Therefore, out of the initial 26 features processed, only 13 features are used for the development of the classifier using the vocal data set.

Previous studies have shown that the sustained vowels tests carry more PD discriminant information. K-means clustering, a non-supervised learning without heuristic basis, was carried out to investigate if there was a possible difference of information gain in the speech samples. K-means clustering was performed on a range of 2 to 10 clusters. A silhouette score was produced, indicating how close the features and voice samples are within each cluster. The higher the silhouette score, the more distinguished the clusters are from each other. 2 clusters produce the highest silhouette scores. Further studies using the 2 clusters showed that among the 26 speech samples, samples number 1, 2 and 3 which are the sustained vowels "a", "o" and "u", displayed a relatively high frequency in 1 cluster compared to the other, which suggest they may have a significant impact on the classification of PD.

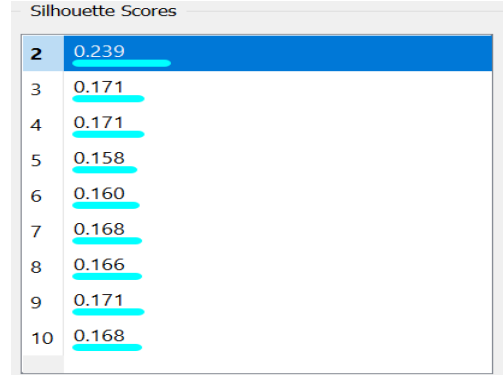


Figure 2: K-Means Clustering Silhouette Score

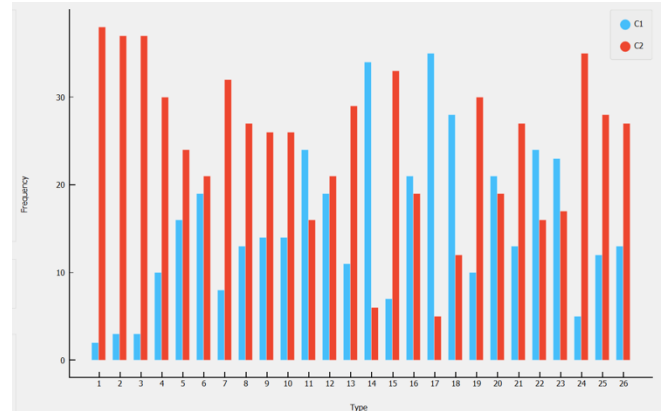


Figure 3: K-Means clustering Frequency of Speech Sample

The above exploratory study was supported by earlier work done in [9] which has shown that using only sustained vowels is sufficient to make a prediction of PD occurrence. The speech sample 1,2 and 3 was proved to show a high information. Then, we decided to only use the sustained vowels "a" and "o" samples for the study to conform to the UCI test data set. The UCI test data set only contained the "a" and "o" samples.

3.2.2 Image Dataset

The data consists of the individual samples of each test that each user of the system did. The first step was to collate all samples into one table where each sample was labeled if it belonged to a control subject or a PWP subject.

Running Pearson's Correlation on the dataset singled out the attributes "Z" and "Pressure" as very highly correlated. As more information (graduation of 1024 steps) was available in the attribute "Pressure", the attribute "Z" was removed. The dataset was then split into the 3 tests, labeled Test 0, Test 1 and Test 2. The attributes "Subject", "Timestamp" and "Test ID" were also removed for each of these 3 test datasets. Finally,

as the intended platform of the app does not have the grip angle variable in its most basic mode using a standard stylus, we left the variable “grip angle” out of the model building and analysis. Prior to the models, the attributes were scaled according to Min-Max scaling.

3.3 Methods

3.3.1 Vocal Data Set

As previously mentioned, data was collected from a total of 68 individuals out of which 28 came from a test group. For this study, data from all 68 individuals were mixed to form the overall dataset from which the train and test set were split by the 80:20 ratio. In addition, it was also found that the dataset was highly skewed towards the **PD class (“P” class)** as a result of the sampling. Thus, oversampling of the **healthy class (“H” class)** was carried out to achieve an overall “P” class to “H” class ratio of 1.3: 1 in the overall data set.

The model for the binary classifier was trained using the k-fold cross-validation method and k = 5 is used in this study. The k fold cross validation method is chosen over having a single hold-out set due to the small amount of data which may result in larger variations for performance evaluation. Another advantage of ‘k-fold cross validation’ is that it requires lower training time compared to the repeated ‘k-fold cross-validation’ and ‘leave one out cross validation’ methods while producing sufficiently accurate results in this study.

Three types of modeling techniques were built and tested for the study involving the vocal dataset as seen in the table below.

Table 5: Classifiers for vocal data set

Classifier	Description
Random Forest	The model was built using the grid search method to determine the number of variables randomly sampled as candidates at each split.
k- Nearest Neighbours	Normalization of predictors was carried out and 10 different values of k were tested and k = 9 was found to produce the best accuracy.
Support Vector Machine (SVM)	Normalization of predictors was carried out and 10 different values of regularization parameter C were tested. Radial basis kernel used.

The results from each of the classification model are further discussed in the subsequent sections.

3.3.2 Image Dataset

Previously, the data was re-grouped by the tests. Each test’s data was then split using a “Leave One Subject” out method. In our case, we left 2 subjects out per test, a PWP and a control subject to provide a balanced test analysis. Data balancing was done on the training data. For each model, 5-fold cross-validation was also performed make sure that the model did not over-fit. Finally, the model was used against the samples of the 2 test subjects to get the final classification accuracy of the model. We made use of the following models in Orange, which is based on the Python Scikit-Learn libraries for the classification of this dataset.

Table 6: Classification methods attempted for spiral image data

Classification Methods Used	
Adaboost	Estimators: 100, Algorithm: SAMME.R, Linear
k-Nearest Neighbours	4 neighbours, Mahalanobis metric, Uniform Weight
Neural Network	Layers (100,100,50), AdamOptimizer, Activation: ReLu, Alpha: 0.0001, Max iterations: 200
Simple Decision Tree	Do not split subset smaller than 5, max tree depth: 100
Naïve Bayes	-
Logistic Regression	Lasso L1 regularization
Random Forest	50 trees, limit depth of individual trees to 3
Stochastic Gradient Descent	Loss function: Hinge, L2 regularization, constant learning rate: 0.01
Support Vector Machines	Cost: 1, loss epsilon : 0.1, Cubic polynomial kernel

Interestingly, the SST, DST and STCP tests responded well to different models. For brevity, for the next section, we will only focus on the models that gave a good cross-validation and test accuracy for each spiral test.

4. MODELING

4.1 Vocal Data Set

4.1.1 Evaluation Metrics

The evaluation metrics used for the classification model are accuracy, sensitivity, and specificity. These metrics are commonly used in classification models and they are also one of the most comprehensible metrics.

The three equations below will give a simple description of these metrics - TP is the number of true positives, TN the number of true negatives, FP the number of false

positives and FN the number of false negatives. Accuracy will give the ratio of correctly classified instances to the whole instances as seen below:

$$accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (1)$$

Sensitivity will give the ratio of true positive instances to the total actual positive instances as seen below:

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

Specificity is the ratio of true negative instances to the total actual negative instances as seen below:

$$specificity = \frac{TN}{TN + FP} \quad (3)$$

The confusion matrix, which offers a simple comparison between the predicted values and their actual values, will be used to present the results of the classification models tested. The three confusion matrices seen below are the result of the test set which was not involved in the development of the training model.

Table 7: Confusion matrix for Random Forest

	Predicted "H"	Predicted "P"	Sum
Actual "H"	32	0	32
Actual "P"	3	38	41
Sum	35	38	73

Table 8: Confusion matrix for k-Nearest Neighbors

	Predicted "H"	Predicted "P"	Sum
Actual "H"	27	5	32
Actual "P"	7	34	41
Sum	34	39	73

Table 9: Confusion matrix for SVM

	Predicted "H"	Predicted "P"	Sum
Actual "H"	32	0	32
Actual "P"	4	37	41
Sum	36	37	73

Based on the results above, it is observed that the k-Nearest Neighbours and SVM models were not able to perform as well as the Random Forest model. However, it can be observed that both the Random Forest and SVM models are able to achieve zero false positive

predictions while all models exhibited false positive predictions ranging from 4 to 7 cases.

The table below gives a summary of the performance metrics of each classification model and their ensemble based on majority voting.

Table 10: Performance of all models

	Accuracy	Specificity	Sensitivity
Random Forest	95.89%	100%	91.43%
k-Nearest Neighbours	83.56%	87.18%	79.41%
SVM	93.15%	100%	86.49%
Ensemble	95.89%	100%	91.43%

In addition to the three evaluation metrics, the confidence interval for the ensemble error is calculated as follows:

$$error \pm const * \sqrt{\frac{error * (1 - error)}{n}} \quad (4)$$

where the constant is dependent on the level of confidence and n is the number of instances used for the evaluation of the model. The error translates to **0.041 +/- 0.046** with 95% confidence interval.

4.2 Image Data Set

4.2.1 Static Spiral Test

For SST, the best accuracy came from Logistic Regression classifier, giving an accuracy of **98.8%**, sensitivity of **96.1%** and specificity of **94.6%**. The error translates to **0.012 +/- 0.003** with 95% confidence interval.

Table 11: Confusion matrix for SST using Logistic Regression

	Predicted "H"	Predicted "P"	Sum
Actual "H"	2945	169	3114
Actual "P"	0	1246	1246
Sum	2945	1415	4360

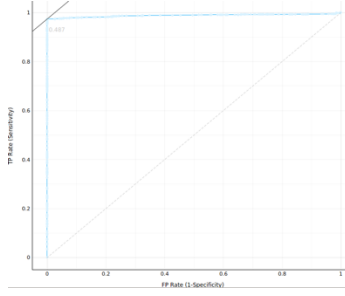


Figure 4: ROC Curve for SST with Logistic Regression

4.2.2 Dynamic Spiral Test

For DST, the k-Nearest Neighbors classifier scored at best, giving an accuracy of **99.6%**, sensitivity of **99.6%** and a specificity of **99.6%**. The error translates to **0.004 +/- 0.002** with 95% confidence interval.

Table 12: Confusion matrix for DST using k-Nearest Neighbors

	Predicted "H"	Predicted "P"	Sum
Actual "H"	2363	10	2373
Actual "P"	3	1112	1115
Sum	2366	1122	3488

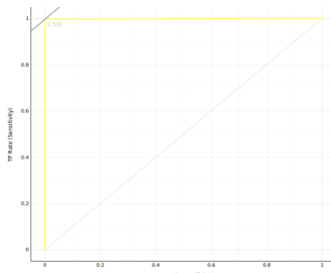


Figure 2 ROC Curve for DST with k-Nearest Neighbors

4.2.3 Stability Test on Certain Point

For STCP, the best accuracy came from the Adaptive Boosting Classifier.

$$F(x) = \text{sign}\left(\sum_{m=1}^M \theta_m f_m(x)\right),$$

where f_m stands for the m-th weak classifier and θ_m is the corresponding weight. AdaBoost is adaptive in the sense that subsequent weak learners are tweaked in favor of those instances misclassified by previous classifiers. Adaptive boosting was able to give an accuracy of **94.9%**, sensitivity of **99.9%** and specificity

of **86.9%**. The error translates to **0.051 +/- 0.007** with 95% confidence interval.

Table 13: Confusion matrix for STCP using Adaptive Boosting

	Predicted "H"	Predicted "P"	Sum
Actual "H"	1336	200	1536
Actual "P"	2	2401	2403
Sum	1338	2601	3939

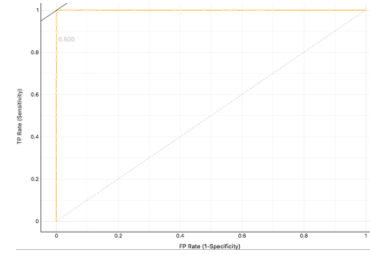


Figure 6: ROC Curve for STCP with Adaptive Boosting

4.2.4 Proposed Final Ensemble Method

For the vocal test data set, a voting ensemble which makes use of the prediction from the 3 classifier models will determine whether the patient is at risk of Parkinson's disease. As illustrated in the table below, if at least 2 of the 3 models give a positive prediction, the final ensemble will then give a positive prediction.

Table 14: Vocal tests ensemble

User	Random Forest	k-Nearest Neighbors	SVM	Ensemble (Final Prediction)
PWP	PWP	PWP	Healthy	PWP

For the spiral tests, a voting ensemble model will be used. The prediction for a single test is taken from the majority prediction of all the samples for this particular test. In this way, we would be able to use the accuracy as the percentage of confidence we have in the prediction for this particular test. The output of the ensemble is a single prediction stating whether the patient is at risk, along with the accuracy obtained for that user for test 1, 2 and 3.

Table 15: Spiral tests ensemble

User	SST (k-Nearest Neighbors)	DST (k-Nearest Neighbors)	STCP (Adaptive Boosting)	Final Prediction
PWP	PWP	Healthy	PWP	PWP

The final diagnosis will be determined by either the vocal or spiral ensemble having a positive diagnosis (ie. Having PD) This is done with the consideration that our analysis has shown that being tested positive on either

the vocal or spiral test is sufficient to diagnose the subject as suffering from PD. For clarity, the table below shows the possible combinations of results from the two tests.

Table 16: Overall Ensemble

Vocal Tests	Spiral Tests	Final Diagnosis
Negative	Negative	Negative
Positive	Negative	Positive
Negative	Positive	Positive
Positive	Positive	Positive

5. PROCESS FLOWCHART

The following flowchart summarize the steps taken from data preprocessing to final output for the vocal dataset.

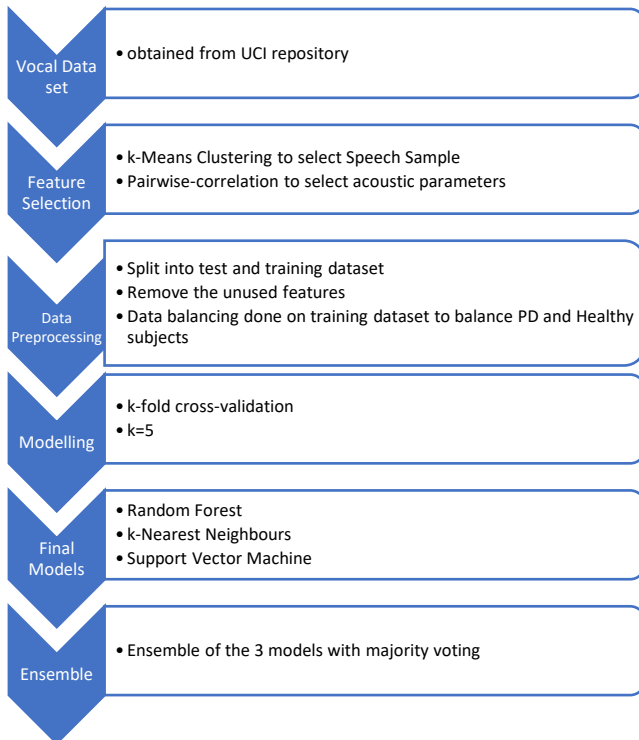


Figure 7: Flowchart of analysis done on the vocal dataset

The next diagram summarizes the steps for the analysis done on the spiral test dataset.

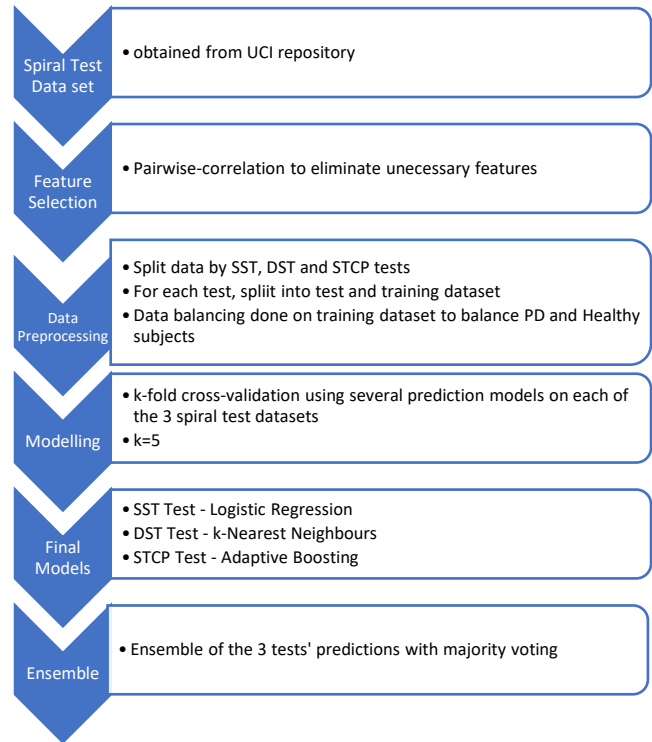


Figure 8 Flowchart of analysis done on the spiral test datasets

6. PROPOSED SMARTPHONE APPLICATION

The proposed app currently will only be on the Apple iOS platform with a 3D-Touch enabled Apple iPhone. The reason is that only the iPhone has a screen pressure sensor, which is crucial for the classification of the spiral tests. Moreover, the iOS Software Development Kit has built-in functionality to do audio recording and simple drawings, making it fast to come up with a prototype. The same app can be used on other mobile phones when the mobile phone makers build in pressure sensors into their hardware. For the vocal test, signal processing routines will have to be programmed to extract the relevant attributes from the audio clips. This can be programmed efficiently through the use of the Accelerate Framework in the iOS SDK. Apple has also recently introduced CoreML, which allows one to load Scikit-Learn based models right into the phone, this enables users to use the app offline, if required. It is not advisable to use a finger as input for

tracing the spiral. We would highly recommend the use of a standard stylus for the spiral test.

We have conducted our spiral test analysis with the assumption that the attribute “grip angle” will not be attainable by the phone hardware.

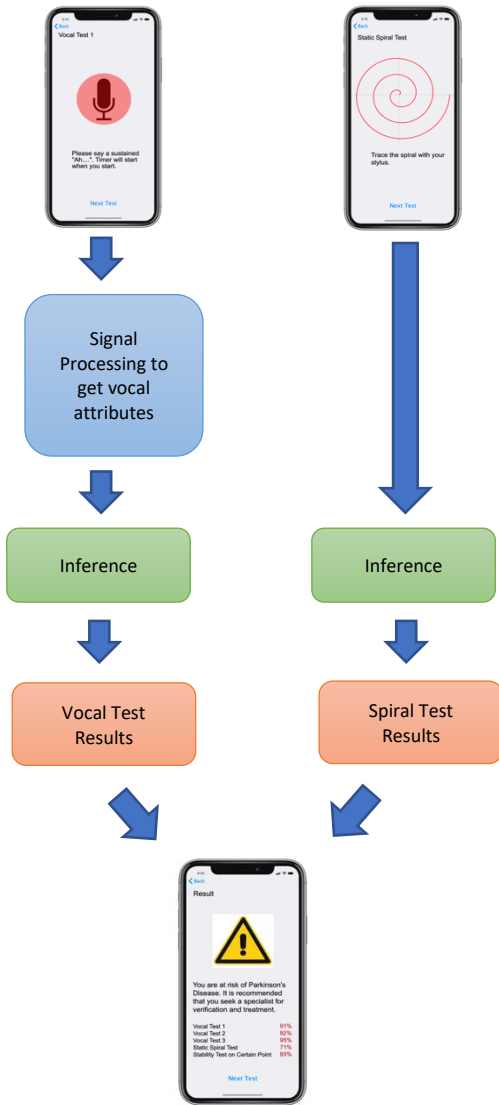


Figure 9: Proposed Operation Flow for Application

7. LIMITATIONS

The models were built using small Vocal test and Spiral test data set. Since both studies were conducted in Istanbul University, it is reasonable to deduce the

subjects were from Turkey. In addition, both data sets did not indicate the severity of the patients' Parkinson's Disease. More data is needed from different regions around the globe with patients having an early stage of Parkinson to improve the reliability and the robustness of the models.

For the proposed mobile app, the voice recording quality would likely be different compared to the Trust MC-1500 microphone employed in the test. Similarly, for the spiral test, the test done on the iPhone may vary from the graphics tablet as the iPhone screen is comparably smaller than the tablet used in the original dataset's test. Data collection has to be re-done using the app on the iPhone itself. This would improve the accuracy and sensitivity of the app.

8. CONCLUSION

The analysis of the datasets has proven that both speech and spiral test can be relied upon for Parkinson's disease detection. Early detection of PD is critical to administer early treatment and to aid them with necessary changes to adapt to the disease. This paper suggests the possibility of an early stage Parkinson's Disease mobile app using the Vocal and Spiral test. Several variables were dropped compared to existing tests and our models are have shown to have high accuracy in detecting Parkinson's disease. Furthermore, the mobile app will enable collection of data remotely from PWP, which can serve as an improved dataset for PD researchers.

The reduction from the 26 voice samplings for the Vocal test to using only the sustained vowel “a” and “o” recording will be convenient and reassuring to apprehensive users.

Due to the limitation in data, further testing is required using the iPhone on early-stage PD patients and healthy personnel for more accurate diagnosis.

Previous studies have shown early symptoms of Parkinson's disease is different for every patient – given this, using multi-modal test data involving both voice and image would provide a more effective way to detect PD for patients who may display different types of predominant symptoms. As such, we believe that having a 2-prong approach (using the two indicated test) would be more reliable in detecting early-stage Parkinson's Disease.

REFERENCES

- [1] <http://parkinson.org/Understanding-Parkinsons/Causes-and-Statistics/Statistics>
- [2] E. R. Dorsey, R. Constantinescu, etc. all, "Projected number of people with Parkinson disease in the most populous nations, 2005 through 2030", *Neurology* Jan 2007.
- [3] Michela Tinelli, "A literature review of the potential clinical and socioeconomic impact of targeting unmet needs in Parkinson's disease", *The London School of Economics and Political Science* Nov 2016
- [4] Ho AK, Iansek R, Marigliani C, Bradshaw JL, Gates S. Speech impairment in a large sample of patients with Parkinson's disease. *Behav. Neurol.* 1998;11:131–137.
- [5] Drawing a spiral test could help detect Parkinson's disease, 7 September 2017
<https://www.telegraph.co.uk/news/2017/09/07/drawing-spiral-test-could-help-detect-parkinsons-disease/>
- [6] 8 smart Parkinson's apps you need to try, 14 April 2015 <http://parkinsonslife.eu/top-apps-for-the-parkinsons-community/>
- [7] Mahnaz Behroozi and Ashkan Sami "A Multiple-Classifer Framework for Parkinson's Disease Detection Based on Various Vocal Tests" *International Journal of Telemedicine and Applications* Volume 2016, Article ID 6837498, 9 pages
- [8] Achraf Benba, Abdelilah Jilbab, Ahmed Hammouch, "Analysis of multiple types of voice recordings in cepstral domain using MFCC for discriminating between patients with Parkinson's disease and healthy people others. *PLOS One*.
- [9] 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 Transactions on Biomedical Engineering*, vol. 56, no. 4, pp. 1015–1022, 2009
- [10] B. E. Sakar, M. E. Isenkul, C. O. Sakar et al., "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, pp. 828–834, 2013
- [11] A. Tsanas, M. A. Little, P. E. McSharry, and L. O. Ramig, "Accurate telemonitoring of Parkinson's disease progression by noninvasive speech tests," *IEEE Transactions on Biomedical Engineering*, vol. 57, no. 4, pp. 884–893, 2010.
- [12] A. Tsanas, M. A. Little, P. E. McSharry, J. Spielman, and L. O. Ramig, "Novel speech signal processing algorithms for high-accuracy classification of Parkinsons disease," *IEEE Transactions on Biomedical Engineering*, vol. 59, no. 5, pp. 1264–1271, 2012.
- [13] C. O. Sakar and O. Kursun, "Telediagnosis of Parkinson's disease using measurements of dysphonia," *Journal of Medical Systems*, vol. 34, no. 4, pp. 591–599, 2010.
- [14] R. Das, "A comparison of multiple classification methods for diagnosis of Parkinson disease," *Expert Systems with Applications*, vol. 37, no. 2, pp. 1568–1572, 2010.
- [15] P.-F. Guo, P. Bhattacharya, and N. Kharm, "Advances in detecting Parkinson's disease," in *Medical Biometrics*, vol. 6165 of *Lecture Notes in Computer Science*, pp. 306–314, Springer, Berlin, Germany, 2010.
- [16] P. Luukka, "Feature selection using fuzzy entropy measures with similarity classifier," *Expert Systems with Applications*, vol. 38, no. 4, pp. 4600–4607, 2011.
- [17] D.-C. Li, C.-W. Liu, and S. C. Hu, "A fuzzy-based data transformation for feature extraction to increase classification performance with small medical data sets," *Artificial Intelligence in Medicine*, vol. 52, no. 1, pp. 45–52, 2011
- [18] A. Ozcift and A. Gulten, "Classifier ensemble construction with rotation forest to improve medical diagnosis performance of machine learning algorithms," *Computer Methods and Programs in Biomedicine*, vol. 104, no. 3, pp. 443–451, 2011.
- [19] Hu"seyin Gu"ru" ler "A novel diagnosis system for Parkinson's disease using complex-valued artificial neural network with k-means clustering feature weighting method", *Neural Comput & Applic* (2017) 28:1657–1666
- [20] Min Wang, Bei Wang, Junzhong Zoua, Masatoshi Nakamura "A new quantitative evaluation method of spiral drawing for patients with Parkinson's disease based on a polar coordinate system with varying origin" *Physica A* 391 (2012) 4377–4388
- [21] Marta San Luciano; Wang, Cuiling; Ortega, Roberto A; Yu, Qiping; Boschung, Sarah "Digitized Spiral Drawing: A Possible Biomarker for Early Parkinson's Disease" *PLoS One*; San Francisco Vol. 11, Iss. 10, (Oct 2016): e0162799.
- [22] Jonathan A., Sistia Brandon Christophea, Audrey Rakovich Seville "Computerized spiral analysis using the iPad" *Journal of Neuroscience Methods* Volume 275, 1 January 2017, Pages 50-54
- [23] André Pierre Legrand, Isabelle Rivalsb, Aliénor Richardc "New insight in spiral drawing analysis methods – Application to action tremor quantification" *Clinical Neurophysiology* Volume 128, Issue 10, October 2017, Pages 1823-1834
- [24] Parkinson Speech Dataset with Multiple Types of Sound Recordings Data Set, 12 June 2014
<https://archive.ics.uci.edu/ml/datasets/Parkinson+Speech+Dataset+with++Multiple+Types+of+Sound+Recordings#>
- [25] C. O. Sakar and O. Kursun, "Telediagnosis of Parkinson's disease using measurements of dysphonia," *J. Med. Syst.*, vol. 34, no. 4, pp. 591–599, 2010.

[26] Parkinson Disease Spiral Drawings Using
Digitized Graphics Tablet Data Set, 20 July 2017
<https://archive.ics.uci.edu/ml/datasets/Parkinson+Disease+Spiral+Drawings+Using+Digitized+Graphics+Tablet#>