

VOICE ANALYSIS FOR DETECTING PERSONS WITH PARKINSON'S DISEASE USING PLP AND VQ

¹ACHRAF BENBA, ²ABDELILAH JILBAB, ²AHMED HAMMOUCH

^{1,2}Laboratoire de Recherche en Génie Electrique, Ecole Normale Supérieure de l'Enseignement Technique,

^{1,2}Mohammed V University, Rabat, Morocco

E-mail: ¹achraf.benba@um5s.net.ma,

ABSTRACT

In order to improve the assessment of speech disorders in the context of Parkinson's disease, we have used 34 voice recordings of sustained vowel / a /, from 34 subjects including 17 patients with Parkinson's disease. We subsequently extracted from 1 to 20 coefficients of the Perceptual linear prediction (PLP) from each subject. The frames of the PLP were compressed using vector quantization, with six Codebook sizes. We used Leave One Subject Out validation scheme known as (LOSO) and the Support Vector Machines (SVM) classifier with its different types of kernels, (i.e.; RBF, Linear). After viewing the variety of obtained results, we proceeded to a bench of 100 trials. The best average result obtained was 75.79%, and the maximum result obtained was 91.17% using the codebook size of 1.

Keywords: *Voice analysis, Parkinson's disease, Perceptual Linear Prediction, Vector quantization. Leave One Subject Out, Support Vector Machines.*

1. INTRODUCTION

The assessment of the quality of speech, and the identification of the causes of its degradation based on phonological and acoustic features have become major concerns of clinicians and speech pathologists. They have become more attentive to any external techniques or methods to their domain, which might provide them additional information for the diagnosis and the assessment of neurological diseases including Parkinson's disease (PD). During its course, Parkinson's disease causes different symptoms and influences the system which controls the execution of learned motor plans such as walking, talking or completing other simple tasks [1] [2] [3]. Parkinson's disease generally affects people whose age is over 50 years and causes voice deterioration in around 90% of patients [4]. For these segments of patients, physical visits for diagnosis, monitoring and treatment are not practical [5] [6].

In the case of the assessment of speech disorders in Parkinson's patients, clinicians and the speech pathologists have adopted subjective methods based on acoustic features to distinguish different disease states. In order to develop more objective assessments, recent studies use measurements of speech quality in time, spectral

and cepstral domains [7] to detect voice disorders in the context of Parkinson's disease. These measurements includes fundamental frequency of the oscillation of vocal folds (F0), absolute sound pressure level, jitter which represents pitch perturbations, shimmer which represents amplitude perturbations, and harmonicity which represents the degree of acoustic periodicity [1] [8] [9].

In this study we used the dataset which was send by Mr. M. Erdem Isenkul [5] from the Department of Computer Engineering at Istanbul University, Istanbul, Turkey. In their work [5], they analyzed multiple types of sound recordings collected from people with Parkinson's disease. The extracted time-frequency acoustic features were fed into SVM and k-NN classifiers for PD diagnosis by using a leave-one-subject-out (LOSO) cross-validation scheme and summarized Leave-One-Out (s-LOO).

However, using these methods with the extracted time-frequency acoustic features does not appear to be an efficient method to distinguish people with Parkinson's disease from healthy subjects. When they used LOSO validation scheme, they got 55.50% as classification accuracy. And in their study Betul Erdogdu Sakar et al, argued that the reason for using s-LOO, is to reduce the effect of variations between different voice

samples of a subject [5]. They got 77.70% as a classification accuracy using the s-LOO method and 85% as maximum classification accuracy after 1000 runs.

In this study we focused on the measurements in cepstral domain by applying Perceptual Linear Predictive cepstral coefficients which have been traditionally used in speaker identification applications and was first proposed by Hynek Hermansky [10]. We have extracted PLP coefficients from the speech signals provided in a database and used vector quantization for feature compression. We then used the LOSO validation scheme with SVM for feature classification in order to discriminate Parkinsonian patients from healthy individuals.

This paper is organized as follows: the voice recordings database is described in section II. The PLP processes and vector quantization are presented successively in section III and IV. The methodology of this research is presented in section V. The results and discussion are presented in Section VI and conclusion in Section VII.

2. DATA ACQUISITION

Dysarthria is the set of speech disorders related with disturbances of muscular control of the speech organs. Dysarthria includes all malfunctions associated to breathing, phonation, articulation, nasalization and prosody. These indications can be measured and detected by analyzing various cues of voice. The data collected in the context of this study (figure 1) belongs to 17 patients with Parkinson's disease (6 female, 11 male) and 17 healthy subjects (8 female, 9 male). Voice signals were recorded via a standard microphone at a

sampling frequency of 44,100 Hz using a 16-bit sound card in a desktop computer. The microphone

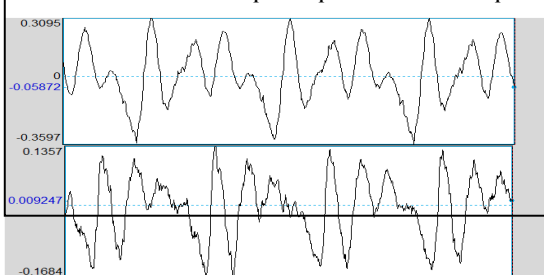


Figure 1: Waveform of a voice sample belonging to healthy individual (top) and Parkinsonian patient (bottom). The horizontal axis represents time and the vertical axis represents the amplitude. This figure was captured using Praat Software

was placed at a 15 cm distant from individuals and they were asked to say sustained vowel /a/ at a comfortable level. All voice recordings were made in mono-channel mode and saved in WAVE format; acoustic analyses were done on these voice recordings. All the voice samples were collected by Mr. M. Erdem Isenkul from the Department of Computer Engineering at Istanbul University, Istanbul, Turkey.

3. THE PLP PRECESSES

Our first aim was to transform the speech waveform to some type of parametric representation for further analysis and processing [13]. The speech signal is a slow time varying signal which is called quasi-stationary [13]. When it is observed over a short period of time, it appears fairly stable [13]. However, over a long period of time, the speech signal changes its waveform. Therefore, it should be characterized by doing short-time spectral analysis [13]. The process of computing the PLP is shown in Figure 1 and described in the next paragraphs.

3.1 Spectral Analysis

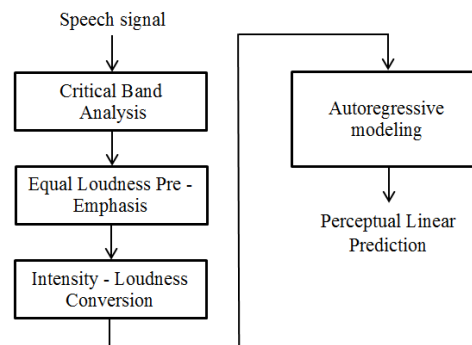


Figure 2: Block diagram of Perceptual Linear Prediction coefficients (PLP)

Since the speech signal is a real signal, it is finite in time; thus, a processing is only possible on finite number of samples [14]. To this end, the first step of PLP process is to weight the speech segment by Hamming window [10]. The aim is to reduce signal discontinuities, and make the ends smooth enough to connect with the beginnings [14]. This was done by using Hamming window to taper the signal to zero in the beginning and in the end of each frame, by applying the following formula to the samples [10]:

$$W(n) = \left\{ 0,54 - 0,46 \cdot \cos\left(\frac{2\pi n}{N-1}\right) \right\} \quad (1)$$

where N is the length of the Hamming window, with a length about 20 ms.

The next processing step consists on converting each frame of N samples from time domain into frequency domain by applying the Fast Fourier Transform (FFT) [13]. We used the FFT for the reason that it is a fast algorithm to implement the Discrete Fourier Transform (DFT) [13]. As known, the DFT is defined on the set of N samples (S_n) as follow [13]:

$$S_n = \sum_{k=0}^{N-1} s_k e^{-2\pi jkn/N}, n = 0,1,2,\dots, N-1 \quad (2)$$

The short-term power spectrum is computed by adding the square of the real and imaginary components of short-term speech spectrum, as follow [10]:

$$P(\omega) = \text{Re}[S(\omega)]^2 + \text{Im}[S(\omega)]^2 \quad (3)$$

3.2 Critical Band Analysis

The short-term power spectrum $P(\omega)$ is warped along its frequency axis ω where ($\omega=2\pi f$), into Bark frequency Ω by using the following equation [10]:

$$\Omega(\omega) = 6 \ln \left\{ \frac{\omega}{1200 \pi} + \sqrt{\left[\left(\frac{\omega}{1200 \pi} \right)^2 + 1 \right]} \right\} \quad (4)$$

$$\Omega(f) = 6 \ln \left\{ \frac{f}{600} + \sqrt{\left[\left(\frac{f}{600} \right)^2 + 1 \right]} \right\} \quad (5)$$

$$\Omega(f) = 6 \sinh^{-1} \left(\frac{f}{600} \right) \quad (6)$$

where ω is the angular frequency in [rad/s], and f is the frequency in [Hz]. The aim of the next step, is to convolve the resulting warped power with the power spectrum of the simulated critical-band masking curve $\Psi(\Omega)$ approximated by Hynek Hermansky [10] as follow:

$$\Psi(\Omega) = \begin{cases} 0 & \text{for } \Omega < -1,3 \\ 10^{2,5(\Omega+0,5)} & \text{for } -1,3 \leq \Omega \leq -0,5 \\ 1 & \text{for } -0,5 < \Omega < 0,5 \\ 10^{-1,0(\Omega-0,5)} & \text{for } 0,5 \leq \Omega \leq 2,5 \\ 0 & \text{for } \Omega > 2,5 \end{cases} \quad (7)$$

It is a rather curd approximation of the shape of auditory filters.

The samples of the critical-band power spectrum are produced by doing the discrete convolution of $\Psi(\Omega)$ with $P(\omega)$ by applying the following equation [10]:

$$\theta(\Omega_i) = \sum_{\Omega=-1,3}^{2,5} P(\Omega - \Omega_i) \Psi(\Omega) \quad (8)$$

The convolution between the relatively broad critical-band masking curve $\Psi(\Omega)$ and the short-term power spectrum $P(\omega)$, reduces the spectral resolution of $\theta(\Omega)$ in comparison with the original $P(\omega)$ [10].

3.3 Equal-loudness Preemphasis

The next step in this process is to preemphasis the samples $\Theta[\Omega(\omega)]$ using the simulated equal-loudness curve, by applying the following equation [10]:

$$\Xi[\Omega(\omega)] = E(\omega) \times \Theta[\Omega(\omega)] \quad (9)$$

here, $E(\omega)$ is an approximation to the non-equal sensitivity of human ear perception at different frequencies. The practical approximation used in this research was adopted by Hynek Hermansky [10] and was first proposed by Makhol and Cosell [15] and given by the following equation:

$$E(\omega) = \frac{(\omega^2 + 56,8 \times 10^6) \omega^4}{(\omega^2 + 6,3 \times 10^6)^2 \times (\omega^2 + 0,38 \times 10^9)} \quad (10)$$

$$E(f) = \left[\frac{f^2}{f^2 + 1,6 \times 10^5} \right]^2 \times \left[\frac{f^2 + 1,44 \times 10^6}{f^2 + 9,6 \times 10^6} \right] \quad (11)$$

3.4 Intensity-loudness Power Law

The last operation before the all-pole modeling is the cubic-root amplitude compression. The following equation approximates the power law of human hearing and simulates the non-linear relation between the intensity of sound and its perceived loudness [10]:

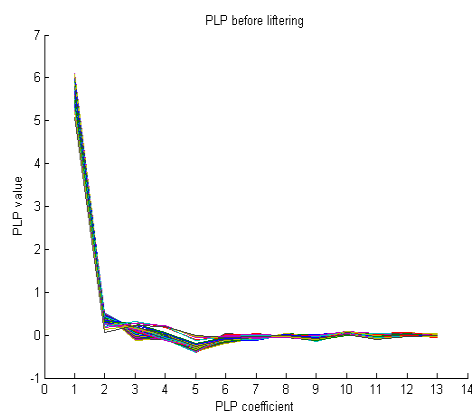


Figure 3: PLP coefficients of PD subject before liftering

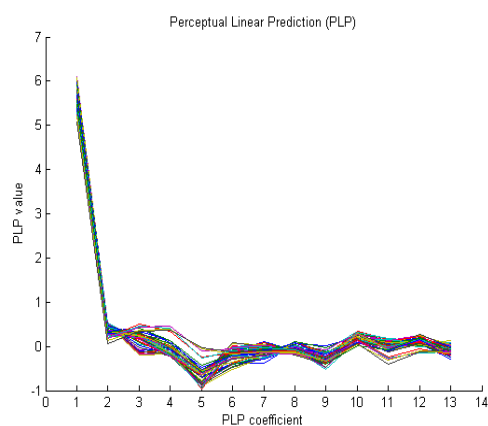


Figure 4: PLP coefficients of PD subject after liftering

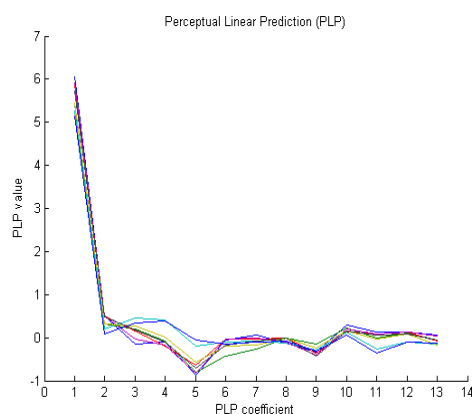


Figure 5: PLP coefficients of PD subject using a codebook size of 8

$$\Phi(\Omega) = \Xi(\Omega)^{0.33} \quad (12)$$

3.5 Autoregressive Modeling

In the final step of the Perceptual Linear Prediction process, $\Phi(\Omega)$ is approximated by the spectrum of an all-pole model using the autocorrelation method of all-pole spectral modeling, this technique is called Linear Prediction (LP) [10] [16], in which the signal spectrum is modeled by an all-pole spectrum. In this study we used the Linear Predictive Coefficient (LPC) analysis to compute the autoregressive model from spectral magnitude samples. The autoregressive coefficients are transformed to cepstral coefficients of the all-pole model; this was done by converting the LPC of 'n' coefficients into frames of cepstra [10].

3.6 Liftering

The principal advantage of cepstral coefficients is that they are uncorrelated [14]. However, the problem with them is that the higher order cepstra are quite small [14], as shown in Figure 3. For this purpose, it is essential to re-scale the cepstral coefficients to have quite similar magnitudes (Figure 4) [14]. This is done by liftering the cepstra according to the following formula [14]:

$$c'_n = \left(1 + \frac{L}{2} \cdot \sin\left(\frac{\pi \cdot n}{L}\right) \right) \cdot c_n \quad (13)$$

Where L is the Cepstral sine lifter parameter. In this work, we used ($L=0.6$).

4. VECTOR QUANTIZATION

Vector quantization (VQ) is a compression method with data losses [17]. The basic idea of this method is to take a large number of feature vectors and reduce it to a smaller group of feature vectors, which represent the centers of gravity of the distribution.

The VQ technique consists of extracting a small number of the most representative features to characterize different individuals.

Here VQ is used to reduce the number of frames of the coefficients of the PLP in order to have only the most significant vectors which represent the center of gravity of the distribution of other frames of the PLP coefficients. In this study, we have made tests using codebook sizes of 1, 2, 4, 8, 16 and 32 [18]. Figure 5 represents the first 13 PLP coefficients extracted from one patient with Parkinson's disease, with data compression using a codebook size of 8, which gives us the 8 most significant

frames representing the center of gravity of the distribution of other frames of the PLP coefficients.

5. METHODOLOGY

The first step in this study was to build a database containing voice recordings of patients with Parkinson's disease and normal individuals. Ultimately, we were able to collect 17 voice recordings from both groups. This gave us 34 voice samples [19]. These recordings were made through a standard microphone at a sampling rate of 44100 Hz. All participants were asked to pronounce the vowel / a / at a comfortable level.

All the algorithms were executed on a desktop computer with a Core (TM) i3-2120 CPU and a processing speed of 3.30 GHz.

We then extracted from each voice sample, cepstral coefficients of the Perceptual Linear Prediction. The number of PLP coefficients extracted ranged from 1 to 20. We proceeded in this way to get the optimal coefficient number needed for the best diagnostic accuracy.

The PLP coefficients extracted from each sample contains a large number of frames which require extensive processing time for classification and prevents making the correct diagnostic decision. To overcome this problem, and reduce the processing time, we used a method of lossy data compression known as vector quantization (VQ). The detailed description of this method has been made in section IV. As we know, VQ compresses the frames according to the number of Codebooks. In this paper we have used six codebooks of size 1, 2, 4, 8, 16 and 32. We applied this method over 20 PLP coefficients that have already been extracted from each voice sample, and which contains from 1 to 20 coefficients per subject. This makes a total of 120 extraction operations per subject ($6 * 20$).

To train and validate our classifier, we used a method of classification called Leave One Subject Out, that is, we left out all the compressed frames of the PLP coefficients of one individual to be used for validation as if it were an unseen individual, and trained a classifier on the rest of the compressed frames of other subjects [6]. We used the Leave One Subject Out method of classification iteratively for each coefficient per subject until all 20 coefficients per subject for the six different codebooks size. In this paper, we used the SVM classifier with its different types of kernels, i.e.; RBF, and Linear.

During the test section, we noticed that the obtained results when using a Codebook size of 1 are not stable. Unlike the other Codebook sizes,

namely 2, 4, 8, 16 and 32, the compression of the PLP frames using a codebook size of 1, did not always give the same location of the centroids of the clusters forming the compressed PLP coefficients. Therefore, every time we redid the same test, we will not get the same classification results. To assess how the results change, we used a test-bed of 100 times. This test-bed allows us to obtain the minimum, maximum and mean value of the diagnostic results for Parkinson's disease [19].

We made a test-bed of 100 times, for the codebook size of 1, and 5 times for the codebook size of 2, 4, 8 and 16, and only one time for the codebook size of 32. As already mentioned, the obtained results using a codebook size higher than 1 are stable, nonetheless we did the test-bed 5 times on the others to get an idea of the variation in execution time and to be sure that the results remained the same.

6. RESULTS AND DISCUSSION

In their study, Betul Erdogdu Sakar et al [5] used a classification with Leave-One-Subject-Out validation scheme, in which all the 26 voice samples of one individual were left out to be used for validation as if it was an unseen individual, and the rest of the samples are used for training [5].

According to their method, if the majority of the voice samples of a test individual are classified as unhealthy, then the individual is classified as positive [5]. In their study, they presented another classification with Summarized Leave-One-Out (s-LOO). The aim of using this method was to compare the success of conventional Leave-One-Subject-Out validation with an unbiased Leave-One-Out [5]. In this method, the feature values of the 26 voice samples of each individual are summarized using central tendency and dispersion metrics such as mean, median, trimmed mean (10% and 25% removed), standard deviation, interquartile range, mean absolute deviation, and a novel form of dataset consisting of N samples is formed where N is the number of individuals [5]. The purpose of summarizing the voice samples of individuals is to minimize the effect of variations between different voice samples of a subject [5]. The best classification accuracy achieved in their research was 77.50% using s-LOO with linear kernel of SVM and the best results of 1000 runs of selecting a random voice samples from each individual was 85% with the same SVM kernel [5].

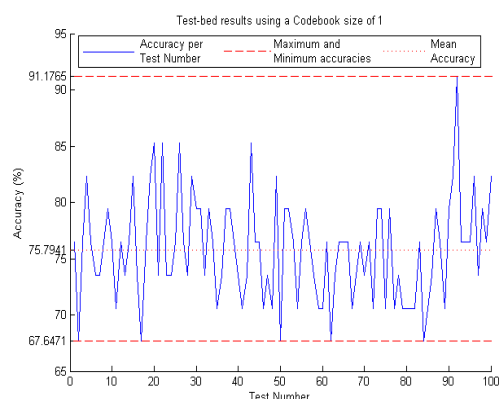


Figure 6: The test-bed results using the codebook size of 1 with minimum, maximum and mean classification accuracies

Table 1: classification results for different codebook sizes.

Codebook sizes	Max accuracy (%)	Min accuracy (%)	Mean accuracy (%)
1	91.1765	67.6471	75.7941
2	75	75	75
4	68.3824	68.3824	68.3824
8	63.6029	63.6029	63.6029
16	62.3162	62.3162	62.3162
32	61.8566	61.8566	61.8566

Table 2: Execution time of the classification program for different sizes of the codebook

Codebook sizes	Max Time (second)	Min Time (second)	Mean Time (second)
1	94.1673	75.5406	78.7189
2	82.0656	81.5608	81.7872
4	130.6943	130.2041	130.4668
8	315.4144	312.1552	314.1749
16	3.1511e+03	3.1383e+03	3.1445e+03
32	4.9475e+04	4.9475e+04	4.9475e+04

Based on our results, it is clear that using higher codebook size, decreases the accuracy of diagnosis (Table I), and the time required for processing becomes longer (Table II).

The extracted PLP coefficients from each subject contains in addition to the number of coefficients used, many frames with different values. The use of a large number of frames leads

us to a diversity of results, often very close to the extracted values from other subjects (PD and Normal) [19]. This similarity of results between different individuals prevents making the correct diagnostic decisions [19]. By way to explanation; assuming that the frames are in the form of points distributed in space, increasing the number of frames leads to interference between these points [19]. Therefore, the task of the classifier, to find a hyper plane able to separate perfectly the two groups of subjects (namely patients with Parkinson's disease and healthy individuals), will be very difficult, if not impossible [19].

As can be seen in Table II for a single test using a codebook size of 1, we need 78.71 seconds. Each time we increase the size of the codebook, the processing time for classification becomes longer. For a single test using a codebook size of 16 we need about 53 minutes and with a codebook size of 32 we need about 14 hours. For this size, it is not practical to apply a test-bed if we already know that the results will remain the same even after 100 trials.

The test-bed accuracy results using the codebook size of 1 are represented in Figure 6. A maximum classification accuracy of 91.17% was achieved using a codebook size of 1 as shown in Table I, by linear kernel SVMs. As seen from Table I, the best mean classification accuracy of 75.79% was achieved using a codebook size of 1.

7. CONCLUSION

Dysarthria symptoms accompanying Parkinson's disease do not appear abruptly. It is a slow process whose early stages may go unobserved. To improve the assessment of Parkinson's disease we collected a variety of voice samples from different subjects during the pronunciation of sustained vowel /a/. The extracted PLP coefficients from different participants contain many frames which take maximum processing time in the classification section, and prevent making accurate diagnosis. For this reason we have compressed the extracted PLP coefficients using vector quantization with different codebook sizes.

After doing the tests we noticed that the obtained results using a codebook size of 1 were not stable. To assess on how the results change, we proceeded to a bench of 100 trials. The compression of the frames of the PLP coefficients using Vector Quantization with the codebook size of 1 has shown to be a good parameter for the detection of voice disorder in Parkinson's disease, showing a mean

classification accuracy of 75.79% and a maximum classification accuracy of 91.17%.

8. ACKNOWLEDGEMENT

The authors would like to thank Mr. Erdem Isenkul from Department of Computer Engineering at Istanbul University, Mr Thomas R. Przybeck and Mr. Daniel Wood from United States Peace Corps Volunteers (Morocco 2013-2015), and all of the participants involved in this study.

REFERENCES:

- [1] Little, M. A., McSharry, P. E., Hunter, E. J., Spielman, J., & Ramig, L. O. (2009). Suitability of dysphonia measurements for telemonitoring of Parkinson's disease. *Biomedical Engineering, IEEE Transactions on*, 56(4), 1015-1022.
- [2] Ishihara, L., and C. Brayne. "A systematic review of depression and mental illness preceding Parkinson's disease." *Acta Neurologica Scandinavica* 113.4 (2006): 211-220.
- [3] Jankovic, Joseph. "Parkinson's disease: clinical features and diagnosis." *Journal of Neurology, Neurosurgery & Psychiatry* 79.4 (2008): 368-376.
- [4] S. B. O'Sullivan, T. J. Schmitz, "Parkinson disease," *Physical Rehabilitation, 5th ed.* Philadelphia, PA, USA: F. A. Davis Company, 2007, pp. 856–894.2007, pp. 856–894.
- [5] Huse, Daniel M., et al. "Burden of illness in Parkinson's disease." *Movement disorders* 20.11 (2005): 1449-1454.
- [6] Sakar, Betul Erdogdu, et al. "Collection and Analysis of a Parkinson Speech Dataset With Multiple Types of Sound Recordings." *Biomedical and Health Informatics, IEEE Journal of* 17.4 (2013): 828-834.
- [7] Uma Rani, K., and Mallikarjun S. Holi. "Automatic detection of neurological disordered voices using mel cepstral coefficients and neural networks." *Point-of-Care Healthcare Technologies (PHT), 2013 IEEE. IEEE*, 2013.
- [8] Little, Max A., et al. "Exploiting nonlinear recurrence and fractal scaling properties for voice disorder detection." *BioMedical Engineering OnLine* 6.1 (2007): 23.
- [9] Rahn III, Douglas A., et al. "Phonatory impairment in Parkinson's disease: evidence from nonlinear dynamic analysis and perturbation analysis." *Journal of Voice* 21.1 (2007): 64-71.
- [10] Hermansky, Hynek. "Perceptual linear predictive (PLP) analysis of speech." *the Journal of the Acoustical Society of America* 87.4 (1990): 1738-1752.
- [11] R. Frail, JI. Godino-Llorente, N. Saenz-Lechon, V. Osma-Ruiz, C. Fredouille, "MFCC-based remote pathology detection on speech transmitted through the telephone channel," *Proc Biosignals*, Porto, 2009.
- [12] Jafari, A, "Classification of Parkinson's disease patients using nonlinear phonetic features and Mel-frequency cepstral analysis," *Biomedical Engineering: Applications, Basis and Communications* 25.04 (2013).
- [13] Ch. S. Kumar, P. R. Mallikarjuna, "Design of an automatic speaker recognition system using MFCC, Vector Quantization and LBG algorithm," *International Journal on Computer Science and Engineering*, Vol. 3, no. 8, 2011.
- [14] S. Young, G. Evermann, T. Hain, D. Kershaw, X. Liu, G. Moore, J. Odell, D. Ollason, D. Povey, V. Valtchev, P. Woodland, "The HTK Book (for HTK Version 3.4)," Copyright. 2001-2006, Cambridge University Engineering Department.
- [15] Makhoul, John, and Lynn Cosell. "LPCW: An LPC vocoder with linear predictive spectral warping." *Acoustics, Speech, and Signal Processing, IEEE International Conference on ICASSP'76.* Vol. 1. IEEE, 1976.
- [16] Makhoul, John. "Spectral linear prediction: properties and applications." *Acoustics, Speech and Signal Processing, IEEE Transactions on* 23.3 (1975): 283-296.
- [17] Kapoor, Tripti, and R. K. Sharma. "Parkinson's disease Diagnosis using Mel-frequency Cepstral Coefficients and Vector Quantization." *International Journal of Computer Applications* 14.3 (2011): 43-46.
- [18] J. Martinez, H. Perez, E. Escamilla, M. M. Suzuki, "Speaker recognition using mel frequency cepstral coefficients (MFCC) and Vector Quantization (VQ) techniques," *IEEE Electrical Communications and Computers*, Cholula, Puebla, Feb 2012, pp. 248-251.
- [19] Achraf BENBA, Abdelilah JILBAB and Ahmed HAMMOUCH. "Voice analysis for detecting persons with Parkinson's disease using MFCC and VQ." *The 2014 International Conference on Circuits, Systems and Signal Processing, 2014.*
- [20] Achraf BENBA, Abdelilah JILBAB and Ahmed HAMMOUCH. "Hybridization of best



acoustic cues for detecting persons with Parkinson's disease," *2nd World conference on complex system (WCCS'14), IEEE, 2014.*