



CLASSIFICATION OF CARDIOVASCULAR DISEASES USING DYSPHONIA MEASUREMENT IN SPEECH

Abdelhamid BOUROUHOU, Abdelilah JILBAB, Chafik NACIR, Ahmed HAMMOUCH

University Mohammed V, Ecole Normale Supérieure de l'Enseignement Technique, Rabat, Morocco

e-mail: abdelhamid.bourouhou@um5s.net.ma

Abstract

Cardiovascular disease is the leading cause of death worldwide. The diagnosis is made by non-invasive methods, but it is far from being comfortable, rapid, and accessible to everyone.

Speech analysis is an emerging non-invasive diagnostic tool, and a lot of researches have shown that it is efficient in speech recognition and in detecting Parkinson's disease, so can it be effective for differentiating between patients with cardiovascular disease and healthy people?

This present work answers the question posed, by collecting a database of 75 people, 35 of whom suffering from cardiovascular diseases, and 40 are healthy. We took from each one three vocal recordings of sustained vowels (aaaaa..., ooooo... and iiiiiii... ..). By measuring dysphonia in speech, we were able to extract 26 features, with which we will train three types of classifiers: the k-near-neighbor, the support vectors machine classifier, and the naive Bayes classifier.

The methods were tested for accuracy and stability, and we obtained 81% accuracy as the best result using the k-near-neighbor classifier.

Keywords: Cardiovascular disease, speech analysis, dysphonia measurement, classification methods, PCA features selection

1. INTRODUCTION

A healthy human body returns to perfect blood circulation; the engine element behind this circulation is the heart. The heart is a muscle that pumps blood throughout the body.

Cardiovascular disease is abnormalities that harm the normal functioning of the heart; they are disorders that affect the heart and blood vessels. We can cite, as examples, coronary heart disease, cerebrovascular disease, peripheral arteriopathy, rheumatic heart disease, and venous thrombosis.

The world heart federation has declared that cardiovascular disease (CVD) is responsible for 17.5 million deaths per year worldwide. This death rate has put CVD at the top of the world's death causes. The world health organization (WHO) predicts that by 2030 the number of deaths will reach 23,6 million, and 1 of 10 people aged from 30 to 70 years old will die prematurely from cardiovascular disease. On the other hand, 80% of premature deaths could be avoided or delayed [1]. The following figure [2] shows a top 10 causes of death worldwide for the year 2017, published by the Institute for Health Metrics and Evaluation (IHME).

The 80% that can be avoided, requires early detection of CVD. For this reason, we need a reliable, precise, fast, and inexpensive tool that will make a distinction between cardiovascular patients and healthy people.

In order to respond to this task, a lot of researches have been carried out. Some of them tried to develop an automatic diagnosis based on the techniques of medical imaging, echocardiography, or magnetic resonance imaging (MRI). Other researches were based on the analysis of electrocardiogram signals (ECG) [3,4], while some others chose phono-cardiogram signals (PCG) [5,6]. The starting information and methods are different, but the goal remains the same – the development of an automatic diagnostic tool. All the previously provided methods used one or more information from the clinical examination of the individual concerned.

It has been suggested that the characteristics of the voice signal are associated with a number of different disease entities, including dyslexia, attention deficit hyperactivity disorder, Parkinson's disease and other neurological disorders [7, 8, 9]. Why is not it the case with CVD?

So, we came up with the idea of achieving the same goal but using a simple information source that does not require a clinical examination. This source is the human speech. Successful research into Parkinson's detection and speech recognition has inspired us, whence the present paper.

Clinical practice uses sustained vowels to assess the quality of the voice; it is for the speaker to pronounce a vowel maintained as long as possible and at a comfortable level. [10] Healthy people can stationary produce a sustained vowel, while this is

not the case for those with vocal impairments. [11, 12, 13]

And to distinguish between healthy and sick people using speech requires finding the difference in spoken speech. Therefore, we must study the speech dysfunctions of the two teams.

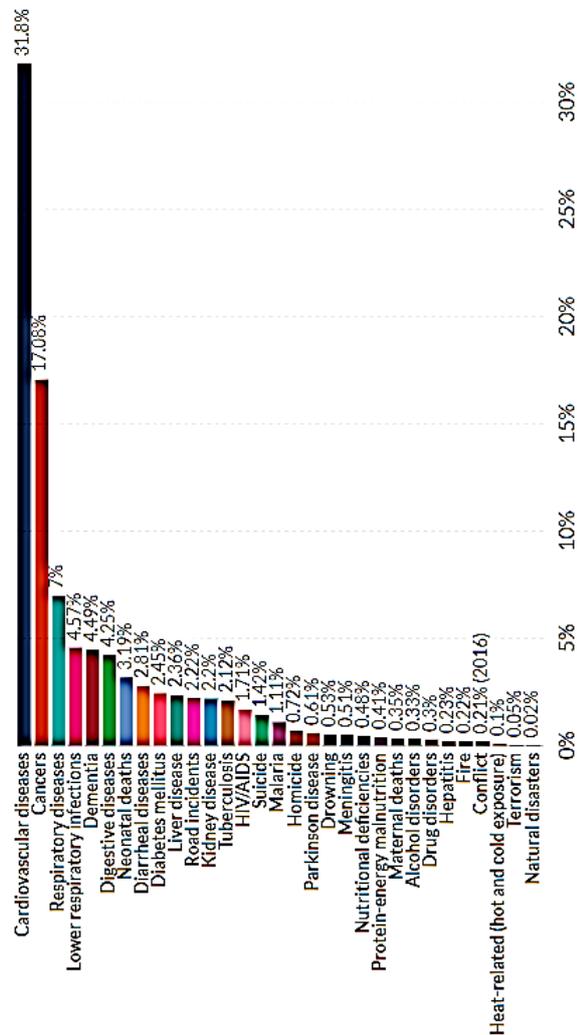


Fig. 1. Top 10 causes of death worldwide for 2017 from World Health Organization [2]

2. DATABASE DESCRIPTION

To conduct research on the differentiation between healthy people and patients with cardiovascular disease (CVD), we need a database with these two categories of people.

The first challenge encountered was collecting this database because, at first, it took a long time before convincing the director of the CHU as well as the head of the cardiology department. Afterwards, we had to wait for the hospitalization of cardiovascular patients, the end of their treatment period and the arrival of a new group of patients.

Our database consists mainly of 75 people, 40 of whom are healthy and 35 suffer from cardiovascular disease (CVD). We took 3 voice recordings for each individual pronouncing sustained vowels (aaaaa / oooooo / iiiiii

.....). The duration of each recording is 5s, and the used device is a smartphone.

The following table summarizes the entire contents of the database:

Table 1. Database description

Database			
Total person	75	Healthy	Patient
		40	35
Gender	Women	Man	
	48%	52%	
Age average	55,5 years		
	Healthy	Patient	
	55 years	56 years	
Total records	225		
Sampling frequency	44100 Hz		

3. METHODS

To distinguish between the people with (CVD) and the healthy ones using speech, we used Dysphonia measures for each of the two categories.

Before proceeding to the measures of the Dysphonia, we filtered and fragmented the recordings in order to keep only the useful part.

Then, we proceeded to the Dysphonia measurements which will be the subject of 26 characteristics extracted from speech for each recording; these built the training base for our used classifiers (k-NN, Naive Bayes, SVM). Validation was ensured by the k-folds cross-validation technique. The diagram in fig.2 presents all steps of our approach.

3.1. Signal acquisition

In this part, we used a smartphone, equipped by MEMS technology microphone, model (CMM-3729 AB-38308-TR) to obtain the records. This microphone operates in a bandwidth of [100 10,000 Hz], has 65 dBA signal to noise ratio measure at 1 kHz by 94 dBA signal, and has 0.2% total harmonic distortion measure at 1 kHz by 94 dB signal. [14]

3.2. Segmentation & filtering

After the signal acquisition, we proceed on two essential steps. The first one consists of signal segmentation because our records take between 6 seconds and 10 seconds for each one, and contains some additional useless sounds. So, we propose a segmentation for each record to keep just the 5 seconds which represents the sustained vowel.

The second step is filtering because our records are a speech (vocal production), and this vocal production is made up of sounds with very specific frequency components, often given as the fundamental frequency (F0, corresponding to the carrier signal), and the first formants (Fi, spikes in the spectral amplitude due to the resonances of the duct voice). The vowels –in our case- are very

easily characterized by their formants. The frequencies used by human speech can therefore be between 110 and 7 kHz (speech communications).

But we filter the signal to keep just the bandwidth [300-3400 Hz] because our speech records contain sustained vowels, and the chosen bandwidth includes the fundamental frequency and the first 3 formants, which is sufficient for us.

So, we applied a Butterworth high pass filter with a cutoff frequency at 300Hz and seventh order succeeded by Butterworth low pass filter with a cutoff frequency at 3400Hz, and also with seventh order. The Butterworth filter choice came from its characteristics, the cutoff frequency is the same regardless of the filter order; the response in the bandwidth is very flat, and the attenuation slope can be increased by increasing the filter order.

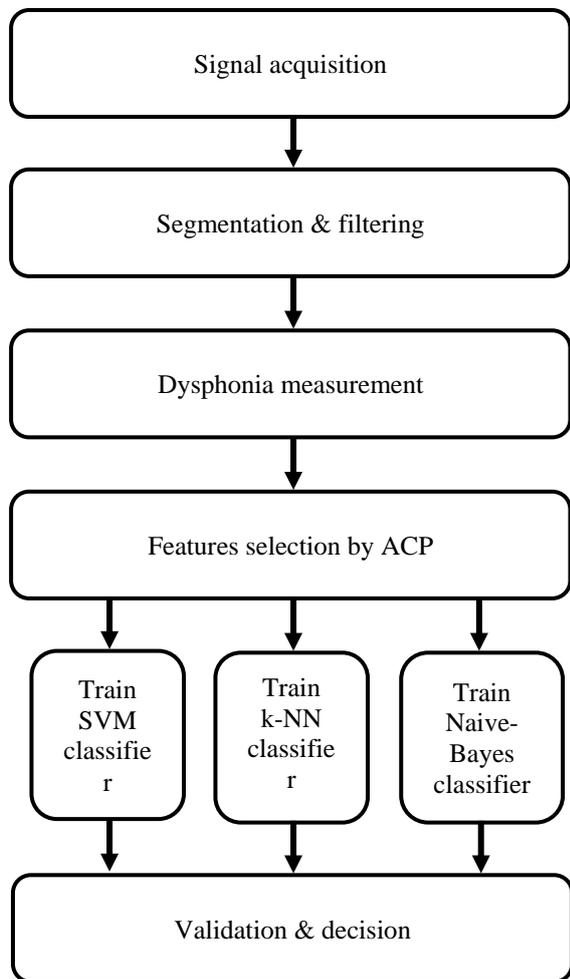


Fig. 2. Diagram of the used method

3.3. Dysphonia measurement

In this step, we try to extract 26 features from each record. These features built the training datasets and were the keys to classify entities into healthy and sick with (CVD).

The 26 features are linear and non-linear parameters; we had frequency parameters, amplitude parameters, harmonic parameters, and

some specific parameters dedicated to vocal disorder analysis.

We used the “Praat” software to extract all parameters, and the following table sums up all extracted features:

Table 2. Features extracted from voice signal

Features	Description	
Pitch	1. Median (Hz)	Depends on the number of vibrations per second produced by the vocal cords. The main acoustic correlate of tone and intonation
	2. Mean (Hz)	
	3. Standard deviation (Hz)	
	4. Maximum (Hz)	
	5. Minimum (Hz)	
Pulses	6. Number of pulses	
	7. Number of periods	
	8. Mean of periods (s: seconds)	
	9. Standard deviation of periods (s)	
Jitter	10. Absolute (s)	Fundamental frequency variation measurements
	11. Local (%)	
	12. RAP (%)	
	13. PPQ5 (%)	
Shimmer	14. DDP (%)	Amplitude variation measurements
	15. Local (%)	
	16. Local (dB)	
	17. APQ3 (%)	
	18. APQ5 (%)	
Harmonicity	19. APQ11(%)	
	20. DDA (%)	
	21. Mean autocorrelation	
Intensity	22. Mean noise to harmonics ratio	
	23. Mean harmonics to noise ratio (dB)	
	24. Mean (dB)	
	25. Maximum (dB)	
	26. Minimum (dB)	

The table below presents the used mathematical expressions to extract features:

Table 3 Mathematics formulas of the extracted features

	Mathematical expression
Jitter	$jitta = \frac{1}{N-1} \sum_{i=1}^{N-1} T_i - T_{i-1} $
	$jitt_{local} = \frac{jitta}{\frac{1}{N} \sum_{i=1}^N T_i} \times 100$

	$RAP = \frac{\frac{1}{N-1} \sum_{i=1}^{N-1} T_i - (\frac{1}{3} \sum_{n=i-1}^{i+1} T_n) }{\frac{1}{N} \sum_{i=1}^N T_i} \times 100$
	$PPQ5 = \frac{\frac{1}{N-1} \sum_{i=2}^{N-2} T_i - (\frac{1}{5} \sum_{n=i-2}^{i+2} T_n) }{\frac{1}{N} \sum_{i=1}^N T_i} \times 100$
	$DDP = \frac{\frac{1}{N-2} \sum_{i=2}^{N-1} T_i - ((T_{i+1} - T_i) - (T_i - T_{i-1})) }{\frac{1}{N} \sum_{i=1}^N T_i} \times 100$
Shimmer	$Shim = \frac{\frac{1}{N-1} \sum_{i=1}^{N-1} A_i - A_{i+1} }{\frac{1}{N} \sum_{i=1}^N A_i} \times 100$
	$ShdB = \frac{1}{N-1} \sum_{i=1}^{N-1} \left 20 * \log \left(\frac{A_{i+1}}{A_i} \right) \right $
	$APQ3 = \frac{\frac{1}{N-1} \sum_{i=1}^{N-1} A_i - (\frac{1}{3} \sum_{n=i-1}^{i+1} A_n) }{\frac{1}{N} \sum_{i=1}^N A_i} \times 100$
	$APQ5 = \frac{\frac{1}{N-1} \sum_{i=2}^{N-2} A_i - (\frac{1}{5} \sum_{n=i-2}^{i+2} A_n) }{\frac{1}{N} \sum_{i=1}^N A_i} \times 100$
	$APQ11 = \frac{\frac{1}{N-1} \sum_{i=5}^{N-5} A_i - (\frac{1}{11} \sum_{n=i-5}^{i+5} A_n) }{\frac{1}{N} \sum_{i=1}^N A_i} \times 100$
	$DDA = \left[\frac{1}{N-2} \sum_{i=2}^{N-1} A_i - ((A_{i+1} - A_i) - (A_i - A_{i-1})) \right] \times 100$
Harmonicity	$NHR = 10 \times \log_{10} \frac{AC_V(0) - AC_V(T)}{AC_V(T)}$
	$HNR = 10 \times \log_{10} \frac{AC_V(T)}{AC_V(0) - AC_V(T)}$

Where, T_i represents the length of period, N represents the number of periods and A_i represents the Peak-to-peak amplitude.

3.4. Features selection by PCA

Now, we have a training matrix where columns represent 26 extracted features and each row is an observation (a record). To get reliable results in classification, we proceeded into a features selection technique, firstly, to reduce the matrix and, secondly, to increase the accuracy of discrimination.

The used technique is the Principal Component Analysis (PCA); it's the mostly used statistical technique in data analysis and data compression. The main idea consists of a data projection from the original space of D variables to a subset

characterized by d variables uncorrelated; this subset contains the principal components and conserves the information contained in the original space.

In our case, we used an algorithm of PCA selection to reduce features number from 26 to 23. The choice of 23 is not random, but it is determined by running the PCA algorithm many times changing the number of principal components to keep in each time. Then, we compared the classification result for all the tries. We noted that the PCA features selection reduces the number of features by computing automatically new “ n ” features (n can be a number from 1 to 26) from the 26 original features. In our case the best projection was for 23 new features computed automatically from the original 26 features because it provides the best result.

3.5. Training & Validation of Classifiers

Finally, we had to train and validate our classifiers. We used 3 types of classifiers that are the k near neighbor, the support vector machine, and the Naive Bayes. All classifiers are subjected to the same cross-validation technique.

3.5.1. Training phase

We have 75 people from whom we took 3 records each, so we collected 225 records. Then we extracted from each record 26 features which was reduced to 23 by the PCA features selection technique. Thus, we built a training dataset (225 rows = 225 observations; 23 columns = 23 features) for our used classifiers.

3.5.2. The k -NN classifier

It is one of the most used classifiers in machine learning; it is meant to determine the k which represents the number of neighbors to take into consideration in the classification. A new entity will be classified as the same class as the majority of neighbors.

The nearest neighbor from the entity to classify depends on the distance which can be Euclidean, Cosine, Minkowski, or other types of distance. So, the number of neighbor k and the type of distance must be chosen carefully because it influences the accuracy of classification.

We applied an optimization algorithm with the aim to select the optimum parameters to train our k -NN classifier. And the chosen parameters were Cosine distance and $k=5$ as a number of the neighbor.

3.5.3. The SVM classifier

It is a technique intended for discrimination and regression problems. Generally, this type of classifier is used for two-class discrimination problems but can be extended for multiclass problems.

The concept consists in building a decision limit named hyperplane separator in the features space,

maximizing the margin between samples from two different classes. These hyperplane separators are the support vectors, which determine the class of the new entity.

To calculate the support vectors, the SVM classifier algorithm uses different kernels, and each kernel has specific proprieties. Examples of these kernels are the Gaussian kernel, the RBF kernel, and the Linear kernel. The choice of the kernel is crucial for prediction accuracy, so we run an optimization algorithm to get the most appropriate SVM classifier proprieties, and the result was the Gaussian kernel.

3.5.4. The Naïve Bayes classifier

The Naïve Bayes classifier is a type of classifier probabilistic based on the theorem of Bayes; it is simple and belongs to the linear classifier family.

The prediction accuracy result for this classifier depends on the kernel choice, Gaussian kernel, Triangular kernel, Box kernel, or Epanechnikov kernel. The difference between all these kernels is determined by the formulas used in the algorithm. The best result using the Naïve Bayes classifier was for the Triangular kernel.

3.5.5. Validation phase

To validate our classifiers, we proceeded into k-folds cross-validation. This type of cross-validation consists of a random subdivision of the database by “k”, keep one of “k” subsets for validation, and training the classifier by all the rest “k-1” subsets. We repeated the operation “k” times until all the subsets may be used one time as a validation set; we calculated the performance score each time. The mean of the “k” squared errors average is finally calculated to estimate the prediction error. In our case the k=5.

After building the cross-validation model, the judgment consists of calculating three essential parameters – the accuracy, the sensitivity, and the specificity; their formulas are given below:

$$Accuracy(\%) = \frac{TP + TN}{TP + FP + TN + FN} \times 100\#(1)$$

$$Sensitivity(\%) = \frac{TP}{TP + FN} \times 100\#(2)$$

$$Specificity(\%) = \frac{TN}{TN + FP} \times 100\#(3)$$

4. RESULTS & DISCUSSION

In this section, we will present the findings. Fig.3 below presents two extracted features, from people with CVD and from healthy ones plotted together to note the difference between our two classes of people.

In the first part of the figure, we plotted the standard deviation of the pitch for records from the people with CVD and for the healthy one records.

In the second part of the figure, we have the standard deviation of the period for the two classes.

We can see that the area occupation for the people with CVD is more important and differs significantly from that of the healthy people. So, we come to conclude that the distinction can be done using speech. Table 4 describes the reached results for a first classification try, using all 26 features extracted from records to train classifiers.

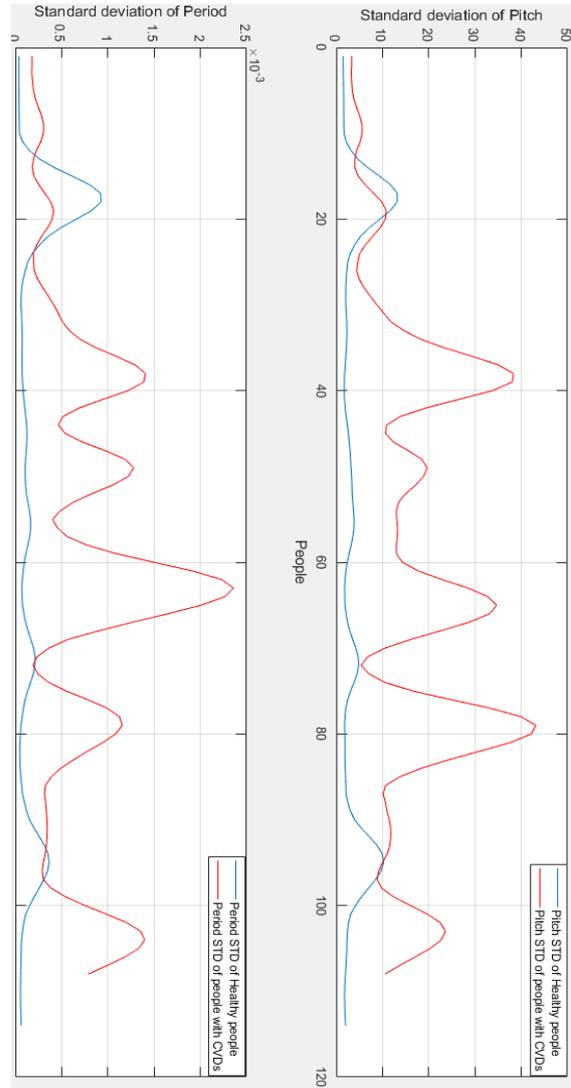


Fig. 3. Comparison between healthy people & people with CVD

Table 4 Classifiers result without application of features selection method

Classifier	Accuracy (%)	Sensitivity (%)	Specificity (%)
K- Near Neighbor	78.46	80.07	76.85
Support Vectors Machine	70.54	69.08	72.00
Naïve Bayes	63.88	65.45	62.31

We have proceeded to a first classification try, by using all the extracted features, which provides us a training matrix 26×225 , where 225 is the number of observations. We have subjected this matrix to different classifiers algorithms. The best provided result was reached by the KNN classifier that means it is the most adapted classifier for our database.

Fig.4 below shows the confusion matrix of the k-near-neighbor classifier; this confusion matrix specifies that our KNN classifier can detect the true positive feature by 80.7%, and the true negative one by only 76.85%. These results reveal the relation supposed to be holding between speech and CVD, and also that we have to improve the classification accuracy.

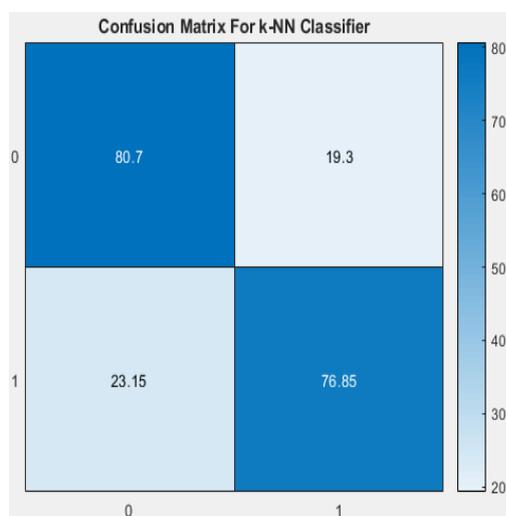


Fig. 4. Confusion matrix for the k-NN Classifier before features selection method

To increase the classification accuracy, we have proceeded to a PCA features selection technique. Table 5 describes the results for our second try of classification; this time we trained our classifiers by the matrix generated by the PCA algorithm, which contains 23 principal components instead of our 26 extracted features.

Table 5 Classifiers result after application of features selection method "PCA"

Classifier	Accuracy (%)	Sensitivity (%)	Specificity (%)
K- Near Neighbor	81.51	82.46	80.56
Support Vectors Machine	75.28	78.50	72.06
Naïve Bayes	70.45	70.00	70.09

The used features selection technique has improved all the performances of our used classifiers, and the KNN classifier is still the best. Fig.5 shows the confusion matrix of the k-near-

neighbor classifier after PCA technique of features selection.

We can also plot the Receiver Operating Characteristic of our three classifiers to compare the results achieved. Fig. 6 below describes that comparison.

As shown in the figure, the area under the plot of the KNN classifier is larger than those of the SVM and Naïve Bayes classifiers; therefore, we conclude that the KNN classifier is better in this case.

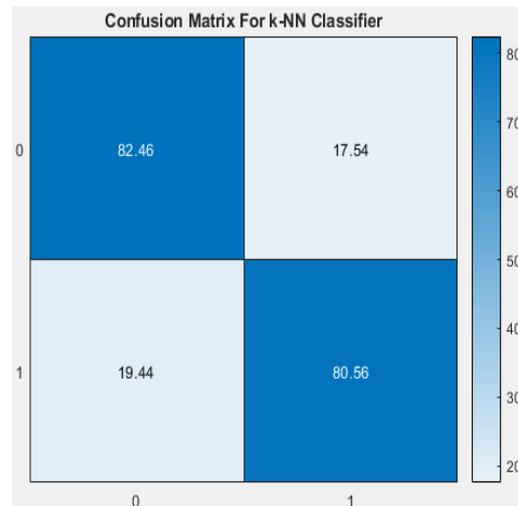


Fig. 5 Confusion matrix for the k-NN Classifier after using "PCA" algorithm

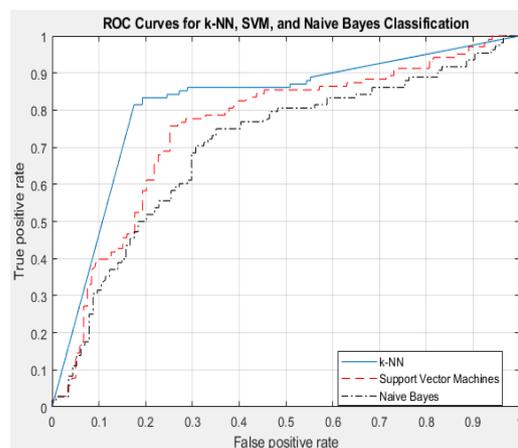


Fig. 6. ROC Curve for all classifiers result after PCA method

5. CONCLUSION

The CVD is still the cause number one of death worldwide, but we can avoid 80% of death if we detect earlier people with CVD, and to ensure detection we have to make the assessment and diagnostic precise, fast, and inexpensive.

To achieve this, we used the speech as a tool of diagnostic, to distinguish between people with CVD and healthy ones. So, we have collected a database that contains multiple records from different people who pronounce sustained vowels

/a/, /o/ and /i/. Then, we have extracted 26 voice features from each record.

To improve the assessment of CVD, we have reduced the training matrix which contains 26 features and 225 observations, to another matrix 23*222 by the PCA feature selection technique.

We have used 3 classifiers, the K-near-neighbor, the support vectors machine, and the Naive Bayes. Choosing the K- near neighbor and cosine distance made the KNN classifier the most successful classifier. The best classification accuracy result was 81.50%.

REFERENCES

- <https://www.euro.who.int/en/health-topics/noncommunicable-diseases/cardiovascular-diseases/data-and-statistics>
- <https://ourworldindata.org/what-does-the-world-die-from>
- Rawther NN, Cheriyan J. Detection and classification of cardiac arrhythmias based on ECG and PCG using temporal and wavelet features. IJARCCCE. 2015; 4(4).
- Bouguila Z, Moukadem A, Dieterlen A, Ahmed Benyahia A, Hajjam A, Talha S, Andres E. Autonomous cardiac diagnostic based on synchronized ECG and PCG signal. In: 7th International Joint Conference on Biomedical Engineering Systems and Technologies—ESEO, Angers. 2014
- Ghassemian H, Kenari AR. Early detection of pediatric heart disease by automated spectral analysis of phonocardiogram in children. J. Inf. Syst. Telecommun. 2015; 3(2): 66–75.
- Nabih-Ali M, El-Dahshan E-SA, Yahia AS. Heart diseases diagnosis using intelligent algorithm based on PCGsignal analysis. Circuits Syst. 2017; 8(7):184–190.
- Levanon Y, Lossos-Shifrin L. Inventors; Google, assignee. Method and system for diagnosing pathological phenomenon using a voice signal 2008. US patent 7,398,213 B1.
- Bonneh YS, Levanon Y, Dean-Pardo O, Lossos L, Adini Y. Abnormal speech spectrum and increased pitch variability in young autistic children. Front Hum Neurosci. 2011;4:237.
- Uma Rani K, Holi MS. Automatic detection of neurological disordered voices using mel cepstral coefficients and neural networks. In: 2013 IEEE Point-of-Care Healthcare Technologies (PHT) 2013:76-79. Bangalore, India, 2013.
- Titze Ingo. Phonation into a straw as a voice building exercise. Journal of Singing. 2000; 57: 27-28.
- Cnockaert L, Schoentgen J, Auzou P, Ozsancak C, Defebvre L, Grenez F. Low-frequency vocal modulations in vowels produced by Parkinsonian subjects, Speech Communication. 2008;50(4):288-300. <https://doi.org/10.1016/j.specom.2007.10.003>
- Little MA, McSharry PE, Hunter EJ, Spielman J, Ramig, LO. Suitability of dysphonia measurements for telemonitoring of Parkinson's Disease. IEEE Transactions on Biomedical Engineering. 2009;56(4):1015-1022. <https://doi.org/10.1109/TBME.2008.2005954>
- Tsanas A, Little MA, McSharry PE, Spielman J, Ramig LO. Novel speech signal processing algorithms for high-accuracy classification of Parkinson's disease. IEEE Trans Biomed Eng. 2012;59(5):1264-1271. <https://doi.org/10.1109/TBME.2012.2183367>
- <https://www.cuidevices.com/product/resource/cmm-3729ab-38308-tr.pdf>
- Bourouhou A, Jilbab A, Nacir C, Hammouch A. Detection and localization algorithm of the S1 and S2 heart sounds. 2017 International Conference on Electrical and Information Technologies (ICEIT), Rabat. 2017:1-4 <https://doi.org/10.1109/EITech.2017.8255217>
- Bourouhou A, Jilbab A, Nacir C, Hammouch A. Comparison of classification methods to detect the Parkinson disease. 2016 International Conference on Electrical and Information Technologies (ICEIT), Tangiers, 2016:421-424. <https://doi.org/10.1109/EITech.2016.7519634>
- Bourouhou A, Jilbab A, Nacir C, Hammouch A. Heart Sounds classification for a medical diagnostic assistance. International Journal of Online and Biomedical Engineering (iJOE) 2019; 15(11): 88–103.

Received 2020-09-27

Accepted 2021-01-19

Available online 2021-01-25



Abdelhamid BOUROUHOU

was born in Rabat, Morocco on December 26th, 1989. Received the Master degree in Electrical Engineering from ENSET, Rabat Mohammed V University, Morocco, in 2014 he is a research student of Sciences and Technologies of the Engineer in ENSIAS, Research Laboratory in Electrical Engineering LRGE, Research Team in Computer and Telecommunication ERIT at ENSET, Mohammed V University, Rabat, Morocco. His interests are in sounds classification for medical diagnostic assistance.



Abdelilah JILBAB

Professor at ENSET Rabat, Morocco; he graduated in electronic and industrial computer aggregation in 1995. Since 2003, he is a member of the laboratory LRIT (Unit associated with the CNRST, FSR, Mohammed V University, Rabat, Morocco). He acquired his PhD in Computer and Telecommunication from Mohammed V-Agdal University, Rabat, Morocco in 2009. His domains of interest include signal processing and embedded systems.



NACIR Chafik

Teacher Researcher in Mathematics. Former Head of the Department of Mathematics and Computer Science. Former member of the Scientific Commission ENSET of Rabat

Morocco.



Ahmed HAMMOUCH received the master degree and the PhD in Automatic, Electrical, Electronic by the Haute Alsace University of Mulhouse (France) in 1993 and the PhD in Signal and Image Processing by the Mohammed V University of Rabat in 2004. From 1993 to 2013 he was a professor in the Mohammed V

University in Morocco. Since 2009 he manages the Research Laboratory in Electronic Engineering. He is an author of several papers in international journals and conferences. His domains of interest include multimedia data processing and telecommunications. He is with National Center for Scientific and Technical Research in Rabat.