

Voice Feature Extraction for Gender and Emotion Recognition

Dr. Madhu M. Nashipudimath^{1,*}, Pooja Pillai¹, Anupama Subramanian¹, Vani Nair¹, and Sarah Khalife¹

¹Department Of Computer Engineering, Pillai College of Engineering, New Panvel – 410 206

Abstract. Voice recognition plays a key function in spoken communication that facilitates identifying the emotions of a person that reflects within the voice. Gender classification through speech is a popular Human Computer Interaction (HCI) method on account that determining gender through computer is hard. This led to the development of a model for “Voice feature extraction for Emotion and Gender Recognition”. The speech signal consists of semantic information, speaker information (gender, age, emotional state), accompanied by noise. Females and males have specific vocal traits because of their acoustical and perceptual variations along with a variety of emotions which bring their own specific perceptions. In order to explore this area, feature extraction requires pre-processing of data, which is necessary for increasing the accuracy. The proposed model follows steps such as data extraction, pre-processing using Voice Activity Detector(VAD), feature extraction using Mel-Frequency Cepstral Coefficient(MFCC), feature reduction by Principal Component Analysis(PCA) and Support Vector Machine (SVM) classifier. The proposed combination of techniques produced better results which can be useful in healthcare sector, virtual assistants, security purposes and other fields related to Human Machine Interaction domain.

Keywords – Human Computer Interaction(HCI), Voice Feature Extraction, Gender Recognition, Emotion Recognition, Voice Activity Detector(VAD), Mel-Frequency Cepstrum (MFC), Mel-Frequency Cepstral Coefficient(MFCC), Principal Component Analysis(PCA), Support Vector Machine (SVM).

1 Introduction

The rapid growth of technology and increasing human demand for automation has made voice recognition systems one of the most desired software programs in various devices. Speech recognition technology converts speech into text and lets users control digital devices by speaking instead of using conventional tools such as keystrokes, buttons, keyboards etc. Some examples of such software are Google voice, digital assistants, car blue-tooth etc. Speech signals contain large amounts of information. Two such pieces of information are gender and emotion which can be distinguished relatively more easily by humans than by computers. At present, the research on voice recognition mainly focuses on the identification of single information, which is not enough to understand the true meaning of speech. Here we intend to use voice feature extraction to identify the gender and emotion of the person using SVM classifier, PCA and MFCC.

2 Literature Survey

Speech is the most fundamental way of communication. The main purpose of this project is to determine which classifiers and pre-processors will be most effective at recognising the speaker's gender and emotional state from their speech. Convolutional Neural Networks (CNN),

Support Vector Machines (SVM), Long Short Term Memory (LSTM), Deep Neural Network (DNN), and Recurrent Neural Networks (RNN) have all been employed in recent years. The review of such techniques is as follows.

2.1 Pre-processing voice signals

Speech signals from predefined datasets [1] cannot be directly passed into the feature extraction module. The input signals are first pre-processed using Voice Activity Detector (VAD) to identify active frames [12] and filter out silent frames that carry no information. The benefit of employing VAD is that, even if there is a significant pause at the start or end of an utterance, the classifier's actions will not be influenced.

2.2 Feature Extraction

For the purpose of gender and emotion recognition, it is critical to extract relevant features from speech signals. In this step the processed speech signals are transformed into a concise but logical representation which is more distinct and reliable than the actual signal. Jiang et al [5] proposed feature extraction module, a heterogeneous unification module, and a fusion network module for a hybrid neural network.

The short term power spectrum of sound is described by Mel-Frequency Cepstrum (MFC) [4], on the basis of a

*e-mail: madhu.nashipudi@yahoo.in

linear cosine transform to log power spectrum with a non-linear Mel scale of frequency. By converting the conventional frequency to Mel Scale, MFCC [11] accounts for human perception for sensitivity at acceptable frequencies. MFC is easy to implement and hence has become a widely used method for speech recognition. The accuracy of the system decreases if the sound samples used have low emotional intensity. It is observed that accuracy is increased when multiple datasets are involved.

Linear Prediction Coding (LPC) [16] approximates speech samples as a linear combination of past samples. Then, over a finite interval, a unique set of predictor coefficients can be calculated by minimising the total of the squared differences between the real speech samples and the linearly predicted samples. An automatic vowel classification system can be presented based on LPC and neural networks. Where traditional linear prediction suffers aliased auto-correlation coefficients LPC gives a very accurate estimate of speech parameters and is comparatively efficient for computation. At the same time, the performance of LPC degrades on the presence of noise in audio signals. Predictive coding [18] determines the best suitable representation of the speech using minimum mean-squared technique.

Linear Predictive Cepstral Coefficients (LPCC) [7] gives a stable representation of the input speech signal in compressed form as compared to LPC. LPCC are derived from the fourier transform of the log magnitude spectrum of LPC. The input signal is analyzed by approximating the frequency bands by removing their effects from the signal and approximates the intensity and frequency of the remaining signal. With the help of Discrete Wavelet Transform (DWT), time domain and frequency domain information of the signal can be fetched. DWT [2] decreases the quantity of signals required to recognize the emotions. Different feature extraction techniques like MFCC [10], pitch, energy, Zero Crossing Rate (ZCR) and DWT are used to extract maximum information of the speech signal and achieve better accuracy with less processing time[2].

2.3 Feature Selection

There are many features in a speech signal but not all of them are needed for implementation of the proposed system. Feature selection [1] is required to extract features from audio signals for selection of principal components, as well as to remove redundant and unused information. Xavier et al [15] used prosodic feature parameters like pitch contour, utterance timing and energy contour. Principal Component Analysis (PCA)[17] is used to find the principal components out of all available features. PCA is a statistical tool [13] which is used to convert a set of observations of correlated variables to a set of values of linearly uncorrelated variables . It also reduces the processing time since a large set of information requires more processing time.

2.4 Classification

Emotion and gender recognition is a supervised learning problem. Each pattern used for the training of the classifier

carries the correct emotion/gender class label. The most popular approaches for classifications include Bayesian learning[6], the Linear Discriminant Analysis (LDA), the Support Vector Machine (SVM) [1, 3] which is used as an extension of LDA with a high-dimensional feature space, the multi-layer Neural Network (NN) [14], and the Hidden Markov model (HMM) which captures temporal state transitions. The intensity of emotion [3] fluctuates on a voice from a low to a high level of emotion.. Hosain et al [9] used Cepstral Coefficient (CC) as voice feature and a fixed valued k-means clustered method for feature classification where value of k is determined by the number of emotional events that are evaluated in human physiology.SVM[7, 8, 11] is the most widely used classifier due to its efficiency in classifying high dimensional data where the number of features is greater than number of observations. SVMs have a major benefit over Artificial Neural Network(ANN) that, unlike ANNs, the solution to an SVM is global and exclusive. SVM has a straightforward geometric interpretation and produces a sparse solution.

3 Proposed System Model

Detecting users' emotion and gender accurately, from his/her voice input, requires complex algorithms and intricate deep learning models. To overcome this, we pre-processed the input data meticulously, followed by classification with the help of SVM, which resulted in precise results without the need of complicated Neural Networks. The proposed system shown in Fig. 1 is trained using pre-defined datasets. More emphasis is given to pre-processing of input voice signals contrary to most of the existing systems.

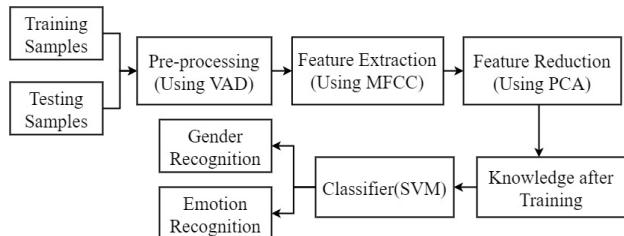


Figure 1. Proposed Model Architecture.

The signals are first pre-processed using Voice Activity Detection (VAD) which is used to determine whether the input signals contain speech or not. The next step is feature extraction. Various acoustic features are extracted using Mel- frequency Cepstral Coefficients (MFCC). The next step is to reduce the number of features. This is done using Principal Component Analysis (PCA). The set of correlated features are transformed into a new set of uncorrelated features called as principal components. Based on the data available after performing the pre-processing steps on the input signal, Support Vector Machine (SVM) classifier is trained to get accurate results.

In the following section, the methodology of Data Selection, Data Pre-processing, Feature Extraction and Data

Classification to obtain the classification results from the system are discussed in detail.

3.1 Data Selection

The training of the proposed model is employed using RAVDESS [2, 11] (Ryerson Audio-Visual Database of Emotional Speech and Song) dataset. The database of size 24.8 GB contains a total of 7356 files consisting of voice samples of 24 professional actors ,12 female and 12 male, vocalizing 2 lexically-matched statements. These are in a neutral North American accent. Speech consists of 8 emotions like happy, sad, angry, surprise, calm, disgust and fearful expressions. With additional neutral expression, every expression is produced at 2 levels of emotional intensity- normal and strong.

3.2 Data Pre-processing using Voice Activity Detection (VAD)

The input signals are pre-processed using Voice Activity Detection (VAD). It is a technique used to detect the presence or absence of speech [12] in an input signal. The unvoiced portions in a signal are removed. It is a binary decision i.e. the output can either be 0 or 1. 0 represents absence of speech whereas 1 represents presence of speech. The function $y = VAD(x)$ is used to express VAD algorithm , wherein the desired target output is:

$$y* := 0, x \text{ is not speech}; = 1, x \text{ is speech} \quad (1)$$

VAD can also be determined as a probability that an input signal contains speech or not. This is called Speech Presence Probability (SPP). The SPP is expressed as the probability which is always in the range 0 to 1. The main objective of calculating SPP is to determine whether the input signal contains speech or not. The SPP is calculated and compared with a threshold value. If the probability is less than the threshold value, then speech is absent in the input signal. A possible definition for the VAD is then

$$VAD(X) := \{0, SPP(X) < \theta; 1, SPP(X) \geq \theta\} \quad (2)$$

where θ is a scalar threshold.

3.3 Feature Extraction using Mel-Frequency Cepstral Coefficient (MFCC) and Principal Component Analysis (PCA)

3.3.1 Mel-Frequency Cepstral Coefficient (MFCC)

One of the most popular audio feature extraction methods is the Mel-Frequency Cepstral Coefficient (MFCC) [10]. It is a technique which takes voice samples as inputs and determines coefficients unique to a particular sample after processing. It provides enough frequency channels to analyze the audio. MFCC involves: framing and windowing, applying the DFT, Mel frequency warping, computing the log of the magnitude, and then applying the inverse DFT.

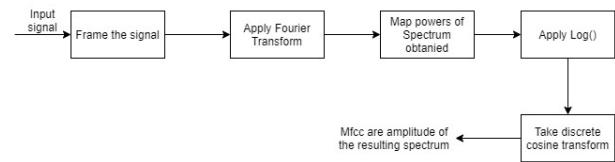


Figure 2. Flowchart of MFCC Feature Extraction process

The flow of extracting the MFCC features is shown in Fig 2.

Acoustic features are extracted using MFCC from the output of VAD i.e. the pre-processed speech signals. The Mel scale compares the perceived frequency of a sound to the measured frequency. It scales the frequency to best match what the human auditory system can hear. The following formula is used to convert a frequency measured in Hertz (f) to the Mel scale

$$Mel(f) = 2595 \log(1 + f/100) \quad (3)$$

3.3.2 Principal Component Analysis

Principal Component Analysis (PCA)[17] is applied to the output of MFCC to reduce the size of the signals. It is a statistical technique which is used to reduce the dimension of the data which is then plotted with lesser dimension compared to original data. It is used for dimensionality reduction of the extracted features without any loss of information. One set of variables is transformed into another smaller set using PCA. The newly created variable is difficult to interpret. In many implementations, PCA is used to produce information on true dimensionality of the data set. Consider there are X variables in the data set, among those all X variables will not represent the needed information. PCA transforms a set of correlated variables into a new set of uncorrelated variables called principal components. This reduces the amount of time required to train the classifier as well as the memory space required. Fig 3 illustrates the steps involved in PCA.

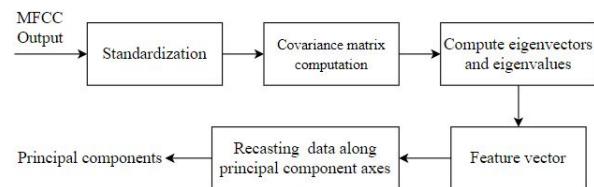


Figure 3. Steps involved in PCA

3.4 Classification using Support Vector Machine (SVM)

The classification model of gender and emotion recognition proposed here is based on the SVM classifier. SVM

[7, 8, 11] uses binary or multi-class classification. It uses hyper planes in the feature space of high dimensions. This helps to differentiate values based on a particular specification. SVM classification is identical to supervised learning which includes feature extraction to generate desired outputs. The major advantage of SVM is that it is very effortless to train. It is capable of scaling high dimensional data better than neural networks[3].

Based on the knowledge available after pre-processing the speech signals, SVM classifier is used for classification and pattern recognition. Thus, pre-defined datasets and integrated algorithms will be used to classify the gender and emotion of the speech signal. All the pre-processing steps are applied to the speech signal. This pre-processed signal along with the knowledge obtained after training is given as input to the SVM Classifier.

4 Experimental Evaluation

4.1 Data Sources

The dataset used here is Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS). It consists of 7356 files available in three modality formats Audio-only, Audio-Video and Video-only. Audio-only files are present in two forms- Speech and Song. For the implementation of this project, audio speech files are used. The audio speech database contains 1440 files recorded by 24 professional actors (12 female, 12 male), vocalizing two lexically-matched statements in a neutral North American accent. It includes various emotions such as calm, happy, sad, angry, fearful, surprise, and disgust to name a few, wherein each expression is produced at two different levels of emotional intensity (normal and strong) along with an additional neutral expression. This project focuses on more frequently observed emotions like happy, sad, angry, neutral.

4.2 Variance using Principal Components

On applying MFCC feature extraction, 180 features are extracted. In order to reduce the number of features, principal components are calculated by computing eigenvalues and eigenvectors from the covariance matrix. 55 principal components are required to achieve 95% of variance as shown in fig 4. In the graph plotted below x axis represents the 180 MFCC features extracted from the dataset and y axis represents the cumulative variance for each feature. As observed from the graph, 95% variance is achieved in first 55 features which are considered for further processing.

The number of principal components are calculated by combining computational efficiency and performance of the classifier. Two eigenvectors are chosen for illustration purpose as data is plotted using a two-dimensional scatter plot. The first two components alone contribute around 33% of information in the model as shown in Fig 5.

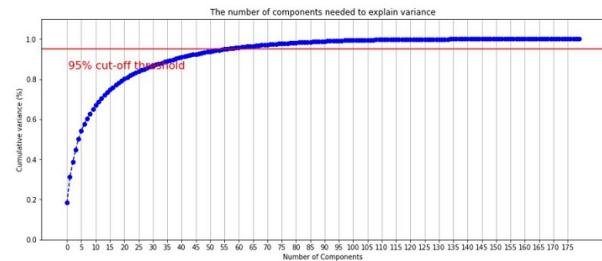


Figure 4. Plot indicating the number of components needed to illustrate variance

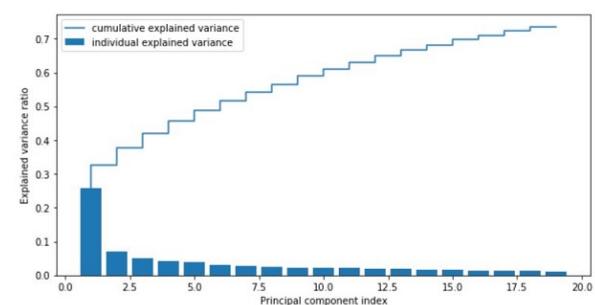


Figure 5. Plot of ‘Explained variance ratio’ vs ‘Principal component index’

4.3 Scatterplot for gender

Fig 6 shows the scatter plot for gender of the first two principal components. The x and y axes represent the first two principal components respectively. The data is divided into two clusters, each for male and female represented by distinguishing colours in the figure below depending on the first PC. Few outliers are also observed.

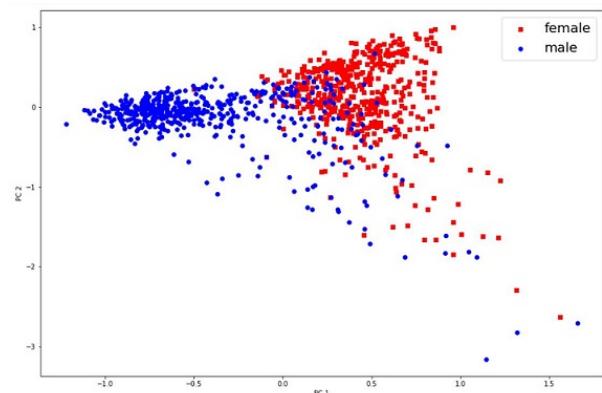


Figure 6. 2-dimensional PCA scatter plot for Gender

4.4 Scatterplot for emotion

Fig 7 shows the scatter plot for emotion of the first two principal components. The x and y axes represent the first two principal components respectively. The data is divided into four clusters, each for happy, sad, angry and neutral

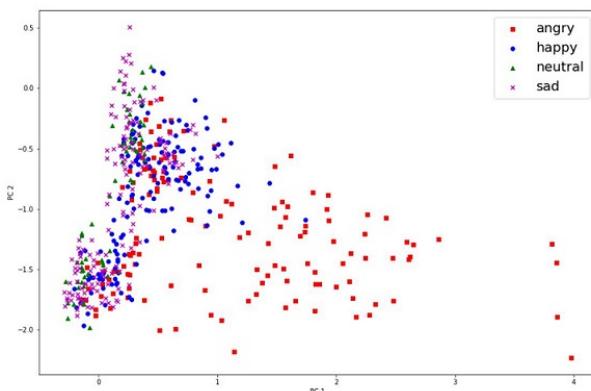


Figure 7. 2-dimensional PCA scatter plot for Emotion

Table 1. Evaluation metrics for gender recognition

Gender	Precision	Recall	F1-score	Support
Female	0.98	1.00	0.99	171
Male	1.00	0.98	0.99	183

Table 2. Evaluation metrics for emotion recognition

Emotion	Precision	Recall	F1-score	Support
Angry	0.80	0.84	0.82	51
Happy	0.66	0.66	0.66	44
Neutral	0.69	0.76	0.72	33
Sad	0.71	0.60	0.65	40

represented by distinguishing colours in the figure below. Overlapping features are observed for happy, sad and neutral emotions whereas the angry emotion is easily differentiable with few outliers.

4.5 Implementation and Evaluation Measures

The dataset is split into training and testing sets, 75% is used for training the system and 25% for testing. The system is evaluated with the help of various parameters such as Precision, Recall, F1-score and Support which is calculated for each gender and emotion individually. The results of which are displayed in Table I for gender recognition system and in Table II for emotion recognition system.

The accuracy score obtained by the gender recognition system is 98.88% where the train accuracy score is 100% and test accuracy score is 98.88%.

The accuracy score obtained by the emotion recognition system is 72.02% where the train accuracy score is 100% and test accuracy score is 72.02%.

4.6 Evaluation Results

The output of the system is analysed using a confusion matrix. In contrast to other machine learning classification metrics like “Accuracy” give less useful information, as Accuracy is simply the difference between correct predictions divided by the total number of predictions. Fig

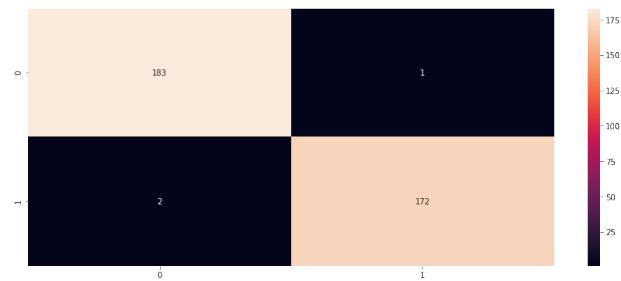


Figure 8. Confusion matrix for gender recognition system



Figure 9. Confusion matrix for emotion recognition system

8 represents the confusion matrix as obtained for the gender recognition system and Fig 9 represents the confusion matrix for emotion recognition system.

Conclusion

This proposed system uses 1440 voice samples in .wav format and produces 98.88% accuracy for gender recognition and 72.02% accuracy for emotion recognition. This improvement in accuracy is achieved due to better pre-processing of data with the help of VAD and better feature extraction by using MFCC and PCA.

Assorted evaluation parameters are considered while computing the final result. Inclusion of kernel PCA will produce more reliable emotion classification results. Additional speaker identification features along with multiple categories of emotions can be added, which may potentially improve the system's results and performance.

Acknowledgment

It is our privilege to express our sincerest regards to our supervisor Dr. Madhu M. Nashipudimath for the valuable inputs, able guidance, encouragement, whole-hearted co-operation throughout the duration of this work. We deeply express our sincere thanks to Head of the Department, Dr. Sharvari Govilkar and Principal, Dr. Sandeep M. Joshi, Pillai College of Engineering, New Panvel for support and boost.

References

- [1] Sharma, Gyanendra & Mala, Shuchi, 2020 *10th International Conference on Cloud Computing, Data Sci-*

- ence Engineering (*Confluence*), "Framework for gender recognition using voice", 32-37, (IEEE, India, 2020).
- [2] Koduru Anusha, Hima Bindu Valiveti, and Anil Kumar Budati, International Journal of Speech Technology **23**, 1 - 11 (2020).
- [3] Mr. Sundar Ka, Sadagopan E.Nb, Chandran Mc, Aswin Raja S, "Emotion Recognition Using Support Vector Machine." (2020).
- [4] Gumelar, Agustinus Bimo, et al, 2019 IEEE 7th International Conference on Serious Games and Applications for Health (SeGAH), "Human Voice Emotion Identification Using Prosodic and Spectral Feature Extraction Based on Deep Neural Networks" (IEEE, Japan, 2019).
- [5] Jiang, Wei & Wang, Zheng & Jin, Jesse & Han, Xianfeng & Li, Chunguang."Speech Emotion Recognition with Heterogeneous Feature Unification of Deep Neural Network. Sensors" (2019)
- [6] Aggarwal, Gaurav, and Rekha Vig, 2019 Amity International Conference on Artificial Intelligence (AICAI), "Acoustic Methodologies for Classifying Gender and Emotions using Machine Learning Algorithms.",(IEEE,United Arab Emirates, 2019), 672 - 677 (2019).
- [7] Jain, Manas & Narayan, Shruthi & Balaji, Pratibha & Bhowmick, Abhijit & Muthu, Rajesh. "Speech Emotion Recognition using Support Vector Machine", 45-55 (2018).
- [8] Poonam Rani, and Ms Geeta, International Journal of Electronics Engineering (ISSN: 0973-7383), **10** • Issue 2, 165-174 (2018).
- [9] Hossain, Nazia & Jahan, Rifat & Tunka, Tanjila, International Journal of Software Engineering & Applications, **9**, 37 - 44 (2018).
- [10] Kerkeni, Leila & Serrestou, Youssef & Mbarki, Mohamed & Raoof, Kosai & Mahjoub, Mohamed, 10th International Conference on Agents and Artificial Intelligence, 175 - 182 (2018).
- [11] Alshamsi, Humaid & Képuska, Veton & Alshamsi, Hazza & Meng, Hongying, 2018 9th IEEE Annual Ubiquitous Computing, Electronics Mobile Communication Conference (UEMCON), (IEEE, USA ,2018), "Automated Speech Emotion Recognition on Smart Phones", 44 - 50 .
- [12] Wang, Zhong-Qiu & Tashev, Ivan. Learning utterance-level representations for speech emotion and age/gender recognition using deep neural networks, (2017).
- [13] Sengupta, Saptarshi & Yasmin, Ghazaala & Ghosal, Dr.Arijit, International Conference on Computing, Communication and Networking Technologies (ICCCNT 2017) "Classification of Male and Female Speech Using Perceptual Features" (2017).
- [14] Pahwa, Anjali & Aggarwal, Gaurav, International Journal of Image, Graphics and Signal Processing, **8**, 17 - 25 (2016)
- [15] Xavier, Arputha rathina, International Journal of Computer Science, Engineering and Applications, **2**, 99 - 107 (2012).
- [16] Paulraj, M. P., et al. "A speech recognition system for Malaysian English pronunciation using Neural Network." (2009).
- [17] Rong, Jia & Li, Gang & Chen, Yi-Ping Phoebe, Information Processing & Management **45**, 315 - 328 (2009).
- [18] Rosenberg, Aaron & Sambur, Marvin, IEEE Transactions on Acoustics, Speech, and Signal Processing **23**, 169 - 176 (1975).