# VoxCeleb1: Speaker Age-Group Classification using Probabilistic Neural Network

Ameer Badr
Department of Computer Science, University of Technology, Iraq
amir.abdulbaqi@sadiq.edu.iq

Alia Abdul-Hassan
Department of Computer Science, University of Technology, Iraq
110018@uotechnology.edu.iq

**Abstract:** *The human voice speech includes essentially paralinguistic information used in many applications for voice recognition. Classifying speakers according to their age-group has been considered as a valuable tool in various applications, as issuing different levels of permission for different age-groups. In the presented research, an automatic system to classify speaker age-group without depending on the text is proposed. The Fundamental Frequency (F0), Jitter, Shimmer, and Spectral Sub-Band Centroids (SSCs) are used as a feature, while the Probabilistic Neural Network (PNN) is utilized as a classifier for the purpose of classifying the speaker utterances into eight age-groups. Experiments are carried out on VoxCeleb1 dataset to demonstrate the proposed system's performance, which is considered as the first effort of its kind. The suggested system has an overall accuracy of roughly 90.25%, and the findings reveal that it is clearly superior to a variety of base-classifiers in terms of overall accuracy.*

## 1. Introduction

Due to the several anatomical and physiological features that vary during our lives, an individual's age could be detected in their voice [11]. Speech is a typical physiological signal used in face-to-face communication and Human-Computer Interaction (HCI). Smart homes, mobile phones, and other assistant devices have pushed the development of a variety of speech applications in recent years. Speech contains paralinguistic information like emotional state, speaker identification, ethnicity, and age, along with the main linguistic information [1, 10]. This information might be beneficial for adjusting voice-operated devices to a more natural HCI for users, which is getting more popular every day. Hyper-parameters like speech synthesis speed could vary with the user's age, and multiple voice recognition models might be trained depending on age. These systems might also be utilized through call centers for speaker's classification into age groups or to profile users. Furthermore, companies engaged in targeted advertising, market research, or service customization have expressed an increased interest in these systems [23].

From various perspectives, automated speaker age-group classification is difficult. First, there is frequently a distinction between a speaker's actual age (chronological age) and their perceived age (perceptual age). Second, balancing data sets, vast age ranges and labeled age ranges are necessary when constructing a robust age classification approach. Third, a variety of factors influence voice patterns, including smoking,

mood context, and gender, to mention a few [2]. Therefore, the first systems capable of reliably recognizing the speaker's age did not emerge till early 2000s [13], and only a few attempts to design these systems have yielded positive results [23].

Although many feature groups were investigated to classify the speaker's age-group, researchers have not yet identified the best feature groups for this task. This is because most previous studies combine multiple feature groups to enhance the performance; however, classification errors from such feature groups May not be complementary. The present study aims at filling this research gap by proposing a new combination of feature groups. Furthermore, the next are some of the work's most important contributions:

1) Combining four feature groups, which are Fundamental Frequency (F0), Jitter, Shimmer, and Spectral Sub-Band Centroids (SSCs) to extract 30-dimensions from each utterance.
2) Exploring the strength of using the Probabilistic Neural Network (PNN) as a classifier in multi-class classification approach.
3) Suggesting a new age labeling to make the VoxCeleb1 dataset appropriate for the speaker's age-group classification issue.

The rest of this study is organized as follows. In section two, the suggested system-related works are provided. The suggested method is discussed in the third section. Section four displays the outcomes of experiments and simulations. Finally, part five outlines the study's future

works and conclusions.

## 2. Related Works

This work focuses mostly on speaker age-group classification; nevertheless, there have been a few previous studies that have looked into such a system.

Muller and Burkhardt [15] compared three automatic speaker age recognition systems using a combination of long-term features and short-term cepstral. A system that combines Gaussian Mixture Model (GMM) based Mel Frequency Cepstral Coefficients (MFCCs) and Support Vector Machines (SVM) based long-term voice pitch performs best. Their system classified the speaker's age into seven-class. Their results indicated a promising approach when combining the two types of features.

Yücesoy and Nabiyev [22] proposed a system that is utilized to classify speakers based on their genders and age through the use of score-level fusion of 7 sub-systems. Perceptual Linear Prediction (PLP), MFCC and prosodic features along with SVM, GMM, and GMM supervector-based SVM classifiers have been used in their system. Their system classified the speaker's age into seven-class. The best classification success achieved by their speaker's age classification system when using fusion in the classification schemes is used, in addition to the features extraction methods.

Pribil *et al*. [19] presented an experiment for automatic speaker gender and age classification by the use of GMM and five sets of features. The features they used included spectral based, formants frequency, F0, Jitter, and Shimmer. They analyzed and compared the impact of different number of mixtures and various companions of features set used in their system. Also, their system classified the speaker's age into seven-class.

Faek [7] presented an automatic gender and age recognizer from speech. As a feature, MFCC and formants were used while K-Nearest Neighbor (KNN) was used as a classifier in their speaker's age classification system. Based on the frequency range that the feature represented, a special selection of robust features is used in their work in order to improve the results. Their performance analyses conducted on database contain clean and noisy speech samples uttered in in Kurdish language. Also, their system classified the speaker's age into seven-class.

Sedaaghi [20] presented a gender and age classification in speech signals comparative study. Naïve Bayes, PNN, GMM, SVM, and KNN were used as classifiers, while spectral based, energy based, formants, and F0 were used as features extraction methods in the proposed system. The proposed system classified the speaker's age into two classes. The best classification success achieved through the use of PNN method in age classification.

## 3. The Proposed Methodology

The approach of this work, as shown in Figure 1, includes two primary stages: features extraction and age-group classification. Initially, suitable features will be recovered from each speaker's utterance. Then, the PNN is used to classify the speaker's utterance into eight age-group depending on the extracted speech features.

### 3.1. Speech Features Extraction

All classification systems' features are specified at the first stage, when the speech signal is transformed into measurable values with distinguishing properties [21]. The tonal sounds are in fact the F0 of a stationary harmonic audio signal, where tonality arranges musical scale notes. F0 is periodic waveform lowest frequency, it is also the temporal auto-correlation function first peak [21]. In general, the subjective pitch of a sound depends on its F0; each application gives for F0 a different definition. For voiced speech, the vibration rate of vocal folds is usually defined as F0 [4]. Among all methods that proposed for F0 estimation, YIN algorithm by Cheveign and Kawahara [4] presented a promising result.

Jitters and Shimmers are among the features which comprise of dynamic information about period and amplitude variations of F0. These features tend to be good for discrimination [17]. Jitter can be defined as the cycle-to-cycle variation regarding F0, put differently, it is considered as the average absolute difference between consecutive periods as indicated in Equation (1), in which Ti are the extracted F0 period lengths and N is the number of extracted F0 periods [8]. Shimmer is specified as the variability of peak-to-peak amplitude in decibels, put differently, it is the average logarithm regarding the absolute difference between the amplitudes of consecutive periods as indicated in Equation (2), in which Ai are the extracted peak-to-peak amplitude data and N is the number of extracted F0 periods [8]. The Praat voice analysis software could be used for extracting such features [9].

$$Jitter\ (absolute) = \frac{1}{N-1}\sum_{i=1}^{N-1}|T_i - T_{i+1}| \qquad (1)$$

$$Shimmer(dB) = \frac{1}{1-N}\sum_{i=1}^{N-1}\left|20log(\frac{A_{i+1}}{A_i})\right| \qquad (2)$$

The SSC feature is meant to be used in conjunction with cepstral characteristics in speech recognition. In Cepstral-based features, one of the main issues is the high sensitivity to additive noise distortion. The addition of white noise to speech signals impacts the speech power's spectrum at all frequencies, yet the impact is less recognizable in higher amplitude (formant) portions regarding the spectrum. As a result, a few formant-like features must be studied to confirm the feature's robustness. SSC features are comparable to formant frequencies and could be extracted reliably and

easily [18]. In noisy conditions, SSCs might exceed MFCCs in recognition accuracy [12, 19]. For SSCs computation, the entire frequency band (0 to Fs/2) is divided into *N* number of sub-bands, in which *Fs*

represent the speech signal sampling frequency. The first moment (centroid) regarding each one of the sub-bands is calculated after applying a filter bank to the signal power spectrum [5].
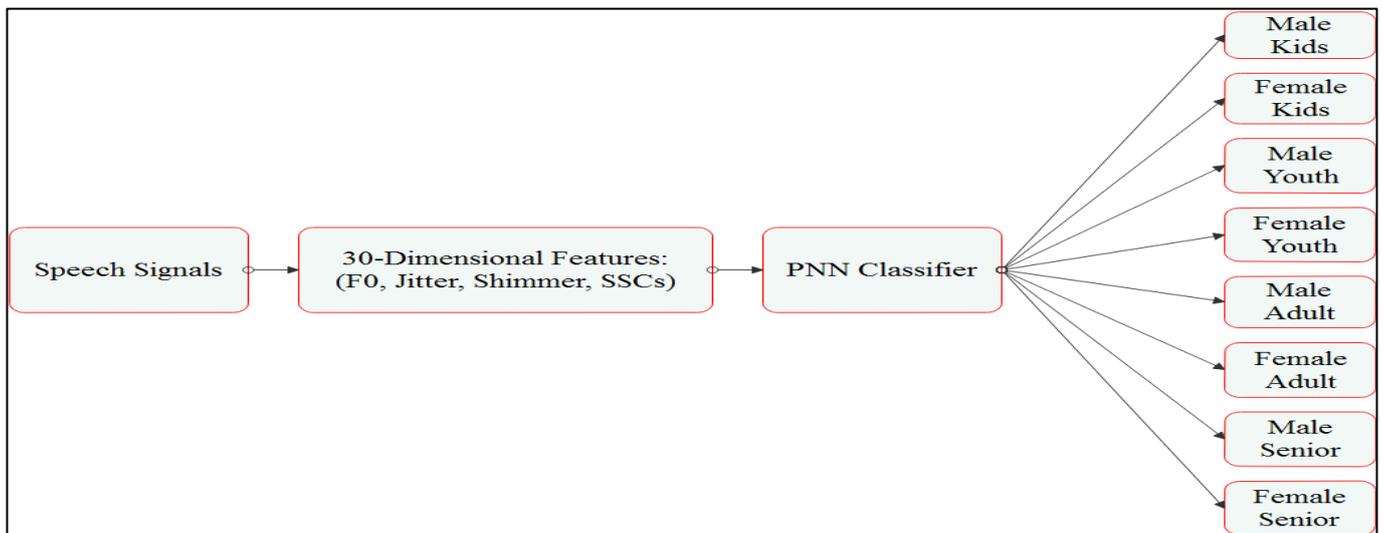


Figure 1. The proposed methodology.

SSC of the $m^{th}$ sub-band is evaluated as shown in Equation (3), in which *Fs* represents the sampling frequency, *P(f)* represents the short-time power spectrum, *ωm(f)* represents the frequency response of $m^{th}$ band pass filter and *γ* is the parameter controlling the dynamic range of the power spectrum [5].

$$C_m = \frac{\int_0^{Fs/2} f \, \omega m(f) P\gamma(f) df}{\int_0^{Fs/2} f \, \omega m(f) P\gamma(f) df} \quad (3)$$

The extracted features must be diversified in order to be informative. This research included 4 groups of features, with the recognition error from such feature groups complementing each other. To begin, each utterance of a speaker is divided into frames with a frameshift of 10 milliseconds and a window size of 25 milliseconds, ensuring that each frame contains robust data. After that, from each one of the utterance frames, a total of 4 groups of features are extracted: local Jitter (1 dimension), F0 (2 dimensions), SSC (26-dimensions), and local Shimmer (1 dimension). The extracted features at this stage have a total dimension of 30 as shown in Figure 1.

## 3.2. Age-Group Classification using The PNN

PNNs are a type of Artificial Neural Network (ANN) that was created for devising a family of Probability Density Function (PDF) estimators utilizing Parzen's technique to best approach Bayes through lowering expected risk. A PNN does not need massive back-propagation training computations. However, each one of the data patterns is represented by a unit which represents the degree of similarity between the data pattern and input [14]. In PNN, a non-parametric Parzen windows estimate approach is employed

to create the class-dependent PDF for each classification category required via Bayes theory. This allows the probability of a certain vector pattern within a specific category to be calculated. In the case when the relative frequency of each category is combined, the PNN determines the most likely category for the supplied pattern vector [3].

The PNN classifier is utilized in the proposed system for speaker age-group classification due to its strong classification ability in multi-class classification issues. The classification stage involves two steps, which are the testing step and the training step. Regarding the latter, the PNN classifier will be trained with the training part of the database to find the best model, the std parameter of the PNN is set to 25. This model is fed then to the testing step for recognizing the speaker's age into eight age-groups.

## 4. Experimental Results and Discussions

### 4.1. Datasets Description

In this study, two datasets have been used which are VoxCeleb1 [16] and CMU Kids [6] corpus in order to classify speaker utterances into eight age-group can be seen in in Table 1.

Table 1. The proposed speakers age-group.

| Class No. | Class Name | Abb. | Age Range | No. Speaker/Total Utterances |
|-----------|------------|------|-----------|------------------------------|
| 1 | Male Kids | MK | 6 – 11 | 18 / 250 |
| 2 | Female Kids | FK | 6 – 11 | 18 / 250 |
| 3 | Male Youth | MY | 15 – 29 | 10 / 250 |
| 4 | Female Youth | FY | 15 – 29 | 10 / 250 |
| 5 | Male Adult | MA | 30 – 54 | 10 / 250 |
| 6 | Female Adult | FA | 30 – 54 | 10 / 250 |
| 7 | Male Senior | MS | 55 – 90 | 10 / 250 |
| 8 | Female Senior | FS | 55 – 90 | 10 / 250 |

The VoxCeleb1 dataset contains 1,250 speakers in text-independent settings, which were acquired using automatic pipelines from open-source media. At a sample rate of 16 kHz, all utterances are encoded with a 16-bit resolution. Since it contains a variety of varied utterances from YouTube, the dataset contains a variety of background noise and utterance durations. The VocCeleb1 data set was utilized for age-group classification in this work, which is the first time it was utilized for such a system to our knowledge. The metadata was created for creating 6 classes from VoxCeleb1 dataset: female youth, male youth, female adult, male adult, female senior, male senior. To acquire enough data for training and testing the system, 10 speakers with 25 utterances were picked for each class. Table 2 shows the metadata details.

Table 2. The proposed voxceleb1 metadata.

| Age-Group | VoxCeleb1 Speakers ID |
|-----------|----------------------|
| MY | id10028, id10062, id10066, id10071, id10076, id10094, id10099, id10208, id10404, id10405 |
| FY | id10007, id10087, id10098, id10181, id10259, id10602, id10825, id11046, id11072, id11201 |
| MA | id10001, id10002, id10003, id10004, id10005, id10009, id10010, id10011, id10012, id10017 |
| FA | id10006, id10008, id10013, id10014, id10024, id10026, id10030, id10032, id10036, id10047 |
| MS | id10019, id10020, id10022, id10029, id10049, id10051, id10059, id10069, id10089, id10090 |
| FS | id10027, id10035, id10038, id10078, id10079, id10086, id10088, id11083, id10136, id10188 |

The CMU Kids Corpus consists of sentences read by children aged 6 to 11 recorded in a controlled environment. For a total of about 9 hours, the dataset consists of 52 female and 24 male speakers. In this study, 18 male speakers with total 250 utterances and 18 female speakers with total 250 utterances have been chosen.

## 4.2. Results

The optimal configuration regarding the suggested speaker age-group classification system characteristics is determined through two experiments. The first experiment involves evaluating the suggested system's performance using a variety of measures, while the second experiment includes a comparison of the proposed system with a number of base-line classification methods. All experiments are conducted on two datasets which are VoxCeleb1 and CMU-kids as seen in Table 1.

The suggested speaker age-group classification system's performance is measured in terms of precision, recall, F-score, overall accuracy, and confusion matrix in the first experiment. This experiment also shows the influence of training data size in the proposed system. Tables 3 and 4 show the results of this experiment. As shown in Table 3, the efficiency regarding the suggested system is estimated with the use of four measures: precision, recall, F-score, and

overall accuracy. All of these measures have been applied to different training data sizes; they are 80%, 66% and 50%. The overall accuracy of the proposed system was affected slightly by the size of training data, as it decreased from 90.25% when the training size was 80% to 85.20% when the training size became 50%. On the other hand, FS class achieved the highest F-score compared with the rest of the classes for 80% and 66% training sizes, while FK achieved the highest F-score for 50% training size. As shown in Table 4, the confusion matrix regarding the suggested system was estimated in terms of accuracy. The results indicate that the highest confusion ratio for MK class happened with FK class only, and same for FK. For MY and FY classes, the high confusion ratio happened with MA and FA classes respectively, and that because the neighborhood between those classes. Another remarkable case, the highest confusion ratio for MS class occurred with MA class, this may be explained as the neighborhood between these groups.

In the second experiment, a comparison will be made with respect to the overall accuracy between the proposed classifier, which is PNN, and eight base-classification methods, namely Adaboost, Logistic Regression (LR), Bagging, SVM, KNN, Gaussian Naïve Bayes (GNB), Decision Tree, and Random Forest (RF). All classification methods are conducted on the same proposed features vector. Figure 2 shows the experiment results. As seen in Figure 2, the proposed method has been compared with eight base-classifiers in terms of overall accuracy, same proposed features vector with 80% training size have been used. The proposed classification method showed superiority, and this confirms the strength of the PNN method as compared with other methods of classification.

## 5. Conclusions and Future Works

This paper proposes an automatic technique for categorizing a speaker's age into eight age groups without relying on text. To enhance system performance, four groups of features are combined, leading to 30-dimensional feature vectors for each speaker utterance. Then, the use of PNN as classification method has a vital effect on the system efficiency taking advantage of the PNN strength in multi-class classification issues. The experimental results indicate the efficiency of the suggested system with overall accuracy of 90.25% using two datasets: VoxCeleb1 and CMU Kids. For future work, a Deep Neural Networks (DNNs) may be utilized for speaker age-group classification in which the network could be fed with the same feature vectors suggested in the presented work.

Table 3. The classification report of the proposed speaker age-group classification method (%).

| Train Size | Class Name | Precision | Recall | F-Score | Overall Accuracy |
|---|---|---|---|---|---|
| **Train 80% -Test 20%** | MK | 93.02 | 83.33 | 87.91 | 90.25 |
| | FK | 83.93 | 95.92 | 89.52 | |
| | MY | 90.24 | 77.08 | 83.15 | |
| | FY | 95.56 | 91.49 | 93.48 | |
| | MA | 85.25 | 96.30 | 90.43 | |
| | FA | 89.13 | 89.13 | 89.13 | |
| | MS | 96.15 | 90.91 | 93.46 | |
| | FS | 91.07 | 96.23 | **93.58** | |
| **Train 66% -Test 33%** | MK | 89.47 | 87.18 | 88.31 | 87.27 |
| | FK | 83.91 | 92.41 | 87.95 | |
| | MY | 91.94 | 73.08 | 81.43 | |
| | FY | 93.42 | 84.52 | 88.75 | |
| | MA | 78.85 | 97.62 | 87.23 | |
| | FA | 85.71 | 83.54 | 84.62 | |
| | MS | 91.95 | 86.02 | 88.89 | |
| | FS | 86.81 | 92.94 | **89.77** | |
| **Train 50% -Test 50%** | MK | 89.19 | 85.34 | 87.22 | 85.20 |
| | FK | 85.11 | 94.49 | **89.55** | |
| | MY | 84.04 | 66.95 | 74.53 | |
| | FY | 92.92 | 83.33 | 87.87 | |
| | MA | 72.96 | 89.23 | 80.28 | |
| | FA | 86.55 | 85.12 | 85.83 | |
| | MS | 90.48 | 86.36 | 88.37 | |
| | FS | 84.67 | 89.23 | 86.89 | |

Table 4. The accuracy based confusion matrix of the proposed system (%).

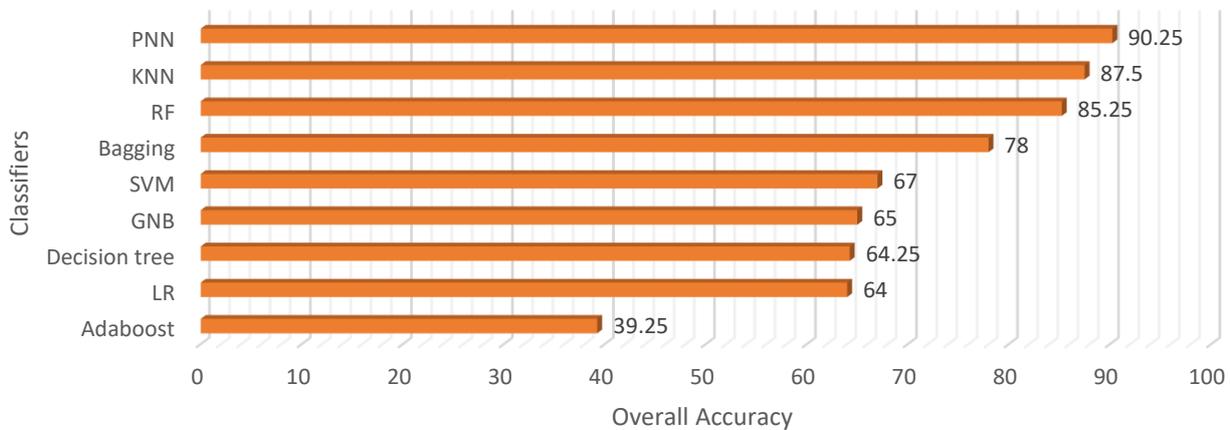| Predicted / Actual | MK | FK | MY | FY | MA | FA | MS | FS |
|---|---|---|---|---|---|---|---|---|
| **MK** | **87.91** | 15.38 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **FK** | 4.16 | **89.52** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **MY** | 0.0 | 0.0 | **83.15** | 0.0 | 9.17 | 4.16 | 4.0 | 3.84 |
| **FY** | 0.0 | 0.0 | 0.0 | **93.48** | 0.0 | 6.45 | 0.0 | 1.94 |
| **MA** | 0.0 | 0.0 | 2.10 | 0.0 | **90.43** | 0.0 | 0.0 | 1.81 |
| **FA** | 0.0 | 0.0 | 2.29 | 4.39 | 1.86 | **89.13** | 0.0 | 1.96 |
| **MS** | 2.04 | 1.80 | 4.16 | 0.0 | 7.72 | 0.0 | **93.46** | 0.0 |
| **FS** | 0.0 | 0.0 | 0.0 | **0.0** | 3.50 | 0.0 | 0.0 | **93.58** |



Figure 2. Overall accuracy based comparison between the proposed classification method and number of base-classifiers.

# References

[1]   Amich H., Mohamed M., and Zrihui M., "Multi-Level Improvement for a Transcription Generated by Automatic Speech Recognition System for Arabic," *The International Arab Journal of Information Technology*, vol. 16, no. 3, pp. 460-466, 2019.

[2] Bahari M. and Van hamme H., "Speaker Age Estimation Using Hidden Markov Model Weight Supervectors," *in Proceedings of 11th International Conference on Information Science, Signal Processing and their Applications*, Montreal, pp. 517-521, 2012.

[3] Bolat B. and Sert S., "Classification of Parkinson's Disease by Using Probabilistic Neural Networks," *in Proceedings of International Symposium on Innovations in Intelligent Systems and Applications*, Trabzon, 2009.

[4] Cheveigne A. and Kawahara H., "YIN, a Fundamental Frequency Estimator for Speech and Music," *The Journal of the Acoustical Society of America*, vol. 111, no. 4, pp. 1917-1930, 2002.

[5] Chougule S. and Chavan M., "Speaker Recognition in Mismatch Conditions: A Feature Level Approach," *International Journal of Image, Graphics and Signal Processing*, vol. 4, pp. 37-43, 2017.

[6] Eskenazi M., Mostow J., and Graff D., "The CMU Kids Corpus," in Linguistic Data Consortium, Philadelphia, USA, 1997.

[7] Faek F., "Objective Gender and Age Recognition from Speech Sentences," *The Scientific Journal of Koya University*, vol. 3, no. 2, pp. 24-29, 2015.

[8] Farrús M., Hernando J., and Ejarque P., "Jitter and Shimmer Measurements for Speaker Recognition," *in Proceedings of 8th Annual Conference of the International Speech Communication Association*, Antwerp, pp. 778-781, 2007.

[9] Feinberg R. "Parselmouth Praat Scripts in Python," OSF, 2019.

[10] Ghahremani P., Nidadavolu P., Chen N., Villalba J., Povey D., Khudanpur S., and Dehak N., "End-to-End Deep Neural Network Age Estimation," *in Proceedings of Interspeech*, Hyderabad, pp. 277-281, 2018.

[11] Grzybowska J. and Kacprzak S., "Speaker Age Classification and Regression Using I-Vectors," *in Proceedings of Interspeech*, San Francisco, 2016.

[12] Kinnunen T., Zhang B., Zhu J., and Wang Y., "Speaker Verification with Adaptive Spectral Subband Centroids," *in Proceedings of International Conference on Biometrics*, Seoul, Korea, pp. 58-66, 2007.

[13] Minematsu N., Sekiguchi M., and Hirose K., "Automatic Estimation of One's Age with His/her Speech Based Upon Acoustic Modeling Techniques of Speakers," *in Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing*, Orlando, pp. 137-140, 2002.

[14] Mohebali B., Tahmassebi A., Meyer-Baese A., and Gandomi A., "Probabilistic Neural Networks: A Brief Overview of Theory, Implementation, and Application," *Handbook of Probabilistic Models*, pp. 347-367, 2020.

[15] Muller C. and Burkhardt F., "Combining Short-term Cepstral and Long-Term Pitch Features for Automatic Recognition of Speaker Age," *in Proceedings of the 8th Annual Conference of the International Speech Communication Association (Interspeech)*, Antwerp, pp. 2277-2280, 2007.

[16] Nagraniy A., Chungy J., and Zisserman A., "VoxCeleb: a Large-Scale Speaker Identification Dataset," *in Proceedings of the Interspeech*, Stockholm, 2007.

[17] Naini A. and Homayounpour M., "Speaker Age Interval and Sex Identification Based on Jitters, Shimmers and Mean MFCC using Supervised and Unsupervised Discriminative Classification Methods," *in Proceedings of 8th International Conference on Signal Processing*, Guilin, 2006.

[18] Paliwal K., "Spectral Subband Centroid Features for Speech Recognition," *in Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP '98 (Cat. No.98CH36181)*, Seattle, pp. 617-620, 1998.

[19] Pribil J., Pribilova A., and Matousek J., "GMM-based Speaker Age and Gender Classification in Czech and Slovak," *Journal of Electrical Engineering*, vol. 68, no. 1, pp. 3-12, 2017.

[20] Sedaaghi M., "A Comparative Study of Gender and Age Classification in Speech Signals," *Iranian Journal of Electrical and Electronic Engineering*, vol. 5, no. 1, pp. 1-12, 2009.

[21] Sharma G., Umapathy K., and Krishnan S., "Trends in Audio Signal Feature Extraction Methods," *Applied Acoustics*, vol. 158, 2020. 107020,

[22] Yücesoy E. and Nabiyev V., "A New Approach with Score-Level Fusion for the Classification of A Speaker Age and Gender," *Computers and Electrical Engineering*, vol. 53, pp. 29-39, 2016.

[23] Zazo R., Nidadavolu P., Chen N., Rodriguez J., and Dehak N., "Age Estimation in Short Speech Utterances Based on LSTM Recurrent Neural Networks, *IEEE Access*, vol. 6, pp. 22524-22530, 2018.

**Ameer Badr** received the B.Sc. degree and the M.Sc. degree in System Software from Computer Science Department, University of Technology, Baghdad, Iraq, in 2014 and 2018 respectively, and the Ph.D. degree in AI from Computer Science Department, University of Technology, Baghdad, Iraq, in 2021. Currently, he is a lecturer at Imam Ja'afar Al-Sadiq University, Salahaddin, Iraq. He has authored or co-authored more than 10 refereed journal and conference papers. His research interests include AI, Machine learning, Speech processing, Speech Enhancement, Speech recognition, Speaker recognition and verification, and, Voice-based HRI.

**Alia Abdul-Hassan** received the B.Sc. degree, the M.Sc. degree and the Ph.D. degree from Computer Science Department, University of Technology, Baghdad, Iraq, in 1993, 1999 and 2004 respectively. She is working as a Dean of Computer Science Department since Feb 2019 till now .She was supervised on more than 30 M.Sc. & Ph.D. thesis in Computer Science since 2007. Her research interests include Soft computing, Green computing, AI, Data Mining, Software Engineering, Electronic Management, and Computer security.