

# Synergistic Fusion of Deep Learning Techniques for Holistic Analysis of Age Group, Gender and Emotion Prediction from Speech

**Abstract**— Voice-based recognition systems have garnered significant attention, owing to their versatile applications in domains such as human-computer interaction, virtual assistants, and affective computing. This research aims to address the challenges associated with concurrent processes and escalating latency in individual models by delving into the intricacies of simultaneous age-group, gender, and emotion recognition from speech data. To surmount these impediments, we present a comprehensive investigation employing three distinct modeling approaches: standalone models for age-group, gender, and emotion recognition; sequential models wherein the output of one model serves as input for the next; and an integrated model proficient in concurrent detection of age-group, gender, and emotion. Drawing inspiration from video understanding techniques, the integrated model seeks to streamline recognition processes and minimize latency. Experimental datasets, including the Ryerson Audio-Visual Database of Emotional Speech (RAVDESS), Common voice dataset and Crowd-Sourced Emotional Actors Database (CREMA-D), are leveraged for rigorous analysis. Evaluation metrics encompassing accuracy, latency, and memory usage are employed to compare the performance of the diverse models.

## I. INTRODUCTION

Speech recognition technologies have garnered substantial research attention owing to their manifold applications in domains spanning human-computer interaction, virtual assistants, and affective computing systems. The present work examines the challenging undertaking of age-group, gender, and emotion identification from vocal signals, contending with complications introduced by intricate variabilities and divergences inherent to human speech modalities across various demographic and affective dimensions. Detecting age-group, gender, and emotion using individual models, which is similar to three parallel processes, presents a number of issues that may impede the efficiency and efficacy of voice-based identification systems [1]. These difficulties include higher latency caused by the simultaneous execution of three independent models, greater memory requirements caused by executing multiple models in parallel, and the possibility of inaccuracies in one model's predictions affecting the outcomes of the others. The development of advanced learning techniques, particularly Convolutional Neural Networks (CNNs) that use 1D convolutions have changed the field of audio signal processing. In the pursuit of advancing the frontier of voice-based recognition systems, our research offers a thorough examination employing a spectrum of models for age-group, gender, and emotion identification. Additionally, integrated models are explored to unravel the intricate links between gender and age-group, as well as gender and emotion. The overarching aim is to contribute to the evolution of contemporary voice recognition technology by furnishing profound insights into acoustic cues indicative of age-group, gender, and emotional states. Datasets utilized in this investigation hail from diverse sources, including the Mozilla Common Voice Dataset (Mozilla). These datasets encompass a rich variety of speech samples, enabling robust model training and evaluation across a wide spectrum of demographic and emotional scenarios. The major challenge faced in the “Combined Model” approach is the unavailability of feature datasets which have all the 3 output parameters namely :- age-group - Group, Gender, Emotion.

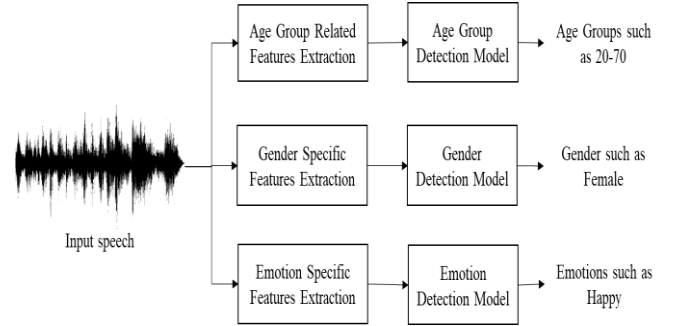


Fig 1. Conventional Approach

## II. RELATED WORK

Voice-based recognition systems have received a lot of attention, and multiple research projects have looked into different aspects of age-group, gender, and emotion detection using voice data. We present a selection of significant works that provide insights into the world of voice-based recognition, notably in the context of age-group, gender, and emotion detection. These publications provide significant perspectives on various approaches, obstacles, and developments, establishing the framework for our complete research aimed at recognising age-group, gender, and emotion from voice data:

1. The work of S. R. Zaman et al. [2] adopts a distinctive perspective in the domain, employing audio speech as a singular source for concurrent gender, age-group, and mood recognition. A range of models, including CatBoost, Random Forest, and XGBoost, undergo testing with 20 statistical characteristics. Notably, CatBoost attains a remarkable 96.4% accuracy in gender prediction, Random Forest excels with 70.4% in age-group prediction, and XGBoost leads with 66.1% in emotion prediction. The scrutiny of these key elements furnishes valuable insights, charting a course for future research in voice-based recognition.
2. In the research conducted by Miller et al. [3], Deep Neural Networks undergo evaluation for the joint prediction of age-group and gender from speech—a crucial aspect for Interactive Voice Response (IVR) systems in contact centers. Leveraging Mozilla's Common Voice dataset, the findings reveal resilient gender classification across networks, with larger sizes contributing to enhancement. A combination of convolutional and temporal neural networks emerges as the optimal configuration for age-group group classification, showcasing potential for IVR systems with minimal gender identification error (below 2%) and age-group group classification error (below 20%) in the most effective systems.
3. In the research conducted by Poonam Rani et al. [4], the research introduces a system designed for discerning an individual's emotional state through audio signal registrations, with potential applications in speech analytics and personalized human-machine interactions. The investigation encompasses two datasets, each comprising approximately 3000 speech samples for gender analysis and 1000 samples for emotion evaluation.
4. In the research conducted by Lee et al. [5], the focus is on multi-task learning for concurrent speaker age-group and gender classification, showcasing the efficacy of shared representations.
5. In the work presented by Gómez et al. [6], a novel approach is taken to address deficiencies in vocal pathology detection systems by introducing an age-group detector trained with both normal and disordered voices. Concentrating on adults and the elderly, the research leverages Mel frequency cepstral coefficients and Gaussian mixture models sourced from the Saarbruecken database. Notably, the research attains an impressive accuracy of 96.57%,

showcasing the potential for developing autonomous age-group-dependent speech pathology identification systems.

6. In the research conducted by Prasanta et al. [7], a Tensor-based strategy is introduced for the detection of speaker gender in speech-based communication, a pivotal aspect in enhancing voice recognition systems. The proposed technique employs a GMM-based classifier tailored for low-resource language-groups. Through experiments conducted on the TIMIT and SHRUTI datasets, the research attains an average-group gender detection accuracy of 91%. The analysis of these results underscores the effectiveness of the Tensor-based approach in the precise detection of speaker gender.
7. In the research conducted by Zheng et al. [8], the research takes aim at mitigating challenges in emotion recognition arising from speaker variability and limited training samples. A novel solution is proposed in the form of a context-dependent domain adversarial neural network (DANN) designed for multimodal emotion recognition.

Emphasizing the importance of contextual information and multimodal features, the method strives to predict emotion labels while concurrently learning a common representation that diminishes disparities related to speaker identity. To address the limitations of low-resource samples, the strategy incorporates unlabeled data. Experimental results on the IEMOCAP dataset showcase an absolute improvement of 3.48% over state-of-the-art techniques, affirming the efficacy of the proposed approach.

### III. PROPOSED SOLUTION

The proposed solution, illustrated in the diagram [i], adopts a deep learning pipeline for multi-output emotion recognition. This method involves individual Convolutional Neural Networks (CNNs) processing audio information extracted from speech to predict emotion, gender, and age-group. While the "Emotion detection Model" improves emotion detection by adding age-group and gender estimates, the "Single Multi-output Model" combines these predictions. Furthermore, "SharedCNNLayers" extract features for both modalities effectively, improving performance. This framework guarantees a thorough and forward-thinking method for multi-output emotion detection, in conjunction with code evolution and integration of Python libraries.

Table 1.

Category	Labels/Groups Targeted
Emotions	Anger, Disgust, Fear, Happy, Neutral, Sad
Genders	Male, Female (Binary)
Age Groups	Child, Young Adult, Adult, Middle Aged, Senior

#### Dataset Description:

- This research integrated two independent datasets: the Mozilla Common Voice dataset (consists of 19,160 validated hours in 114 language-groups) for age-group and gender classification and the CREMA-D (which comprises the following emotions: Anger, Disgust, Fear, Happy, Neutral, and Sad) dataset for emotion recognition. The Mozilla Common Voice dataset has a broad range of voices from which we can extract features for age-group and gender prediction. It includes a large number of speakers, ensuring a diverse representation across demographic groups. The CREMA-D Dataset includes 7,442 audio snippets from 91 actors that were originally recorded. 2443 participants in the data collection process assessed the emotional content of the clips in three different modalities: audiovisual, video alone, and audio alone.
- The predictive models' output labels are divided into three main groups: emotion, gender, and age-group. This three-pronged method allows for a more detailed understanding of voice-based

elements and adds to a more comprehensive voice recognition system.

- Even though these separate datasets are extensive, at first there was trouble locating combined datasets that included labels for emotions, gender, and age-group. Insufficient density or nonexistence of the given data made it impossible to train a reliable and accurate model. Direct integration was hindered by the different data distribution between the public audio recordings in Common Voice and the actor-recorded emotional expressions in CREMA-D.
- To deal with this discrepancy, predictive models were trained independently on each dataset in order to capture the distinct features of CREMA-D and Common Voice. These models were then used to produce forecasts for a combined dataset. We used this synthesized dataset to train our Convolutional Neural Network (CNN) model, which combined predictions from both sources. The goal was to create a single model that could predict age-group, gender, and emotional states with accuracy.

#### Feature Extraction:

- This research feature extraction algorithm focuses on extracting various acoustic properties from audio recordings, providing useful insights for emotion, gender, and age-group recognition. Statistical measurements such as spectral centroid, bandwidth, rolloff, flatness, and contrast [18] are among the retrieved features. These characteristics provide a thorough description of the audio signal, capturing both its central trends and spectrum qualities. The interpretability, computational efficiency, and relevance to voice-based recognition tasks drive the selection of characteristics.
- Because different groups exhibit varying frequency characteristics, the spectral centroid, a measure of average-group frequency, is susceptible to fluctuations in emotion, gender, and age-group. Spectral bandwidth, which indicates the frequency spread, and spectral rolloff, which defines the frequency below which a certain proportion of energy lies, help to capture changes in high-frequency content, which can be indicative of different emotional states, genders, and age-group groups.
- Mel-Frequency Cepstral Coefficients (MFCCs)[19] have been added to the feature extraction procedure to enhance it. Mel-Frequency Cepstral Coefficients provide a useful frequency domain representation of the audio signal's short-term power spectrum. We have customized the dataset by excluding MFCC-12 and MFCC-19.

#### A. Approach 1: Individual Modal

- Emotion:

An LSTM neural network was trained to recognise emotions using data from RAVDESS, CREMA-D, Tess, and Savee datasets. MFCCs, Zero Crossing Rate, and Root Mean Square Energy were among the features that captured audio spectral and temporal characteristics. Sequential LSTM layers captured

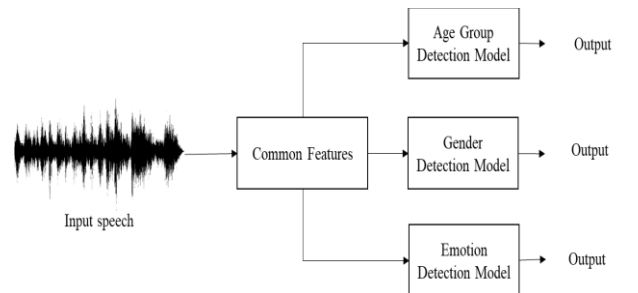


Fig 2.

temporal dependencies, with feature input layers. Overfitting was avoided by using dropout layers. The softmax-activated output layer predicted emotion classes. Categorical cross entropy loss gauged performance, appropriate for multi-class jobs. RMSProp was optimized, and categorical accuracy was utilized for training. Learning rates were changed via a ReduceLROnPlateau callback. The model was trained for 100 epochs, with batch size. Validation accuracy was 90%.

- Gender:

The gender prediction model uses a compact CNN architecture tailored for binary gender classification from speech signals. This network has 2 convolutional layers, max pooling layers and 2 fully connected layers with dropout regularization. The output layer predicts male or female gender using a sigmoid activation. Binary cross entropy loss and the RMSprop optimizer were used. After training for 30 epochs and using early stopping, the model obtained 89.4% accuracy in categorizing speaker gender on the unseen test set. Precision was measured at 0.91 for the male category and 0.86 for the female category. The confusion matrix showed a small skew towards more male gender predictions overall. Investigation showed pitch-based features as providing greater distinguishing evidence for gender than spectral features such as MFCCs.

- Age-group:

The age-group prediction model utilizes a convolutional neural network (CNN) architecture optimized for multi-class classification, predicting the speaker's age-group category from speech. The model contains 3 convolutional layers interweaved with 2 max pooling layers, followed by 2 fully connected layers of 128 and 64 units respectively. ReLU activation and batch normalization were utilized between layers. The model was trained using the categorical cross entropy loss function along with the Adam optimizer, with a learning rate of 0.001 for 100 epochs and a batch size of 32.

The Common Voice dataset was split into 80% training, 10% validation, and 10% test subsets. Training samples were randomly augmented via time shifting and background noise injection. The model achieved an overall accuracy of 71.25% on the age-group category classification task, with a precision of 0.74 and recall of 0.69 on the test data. Particularly strong performance was noted in young adult and middle-age-group groups, while weaker accuracy was achieved in senior age-group bands above 60 years. Misclassifications tended to occur most often between adjacent age-group categories.

#### B. Approach 2 : Age-Group, Gender in Single model followed by Emotion model

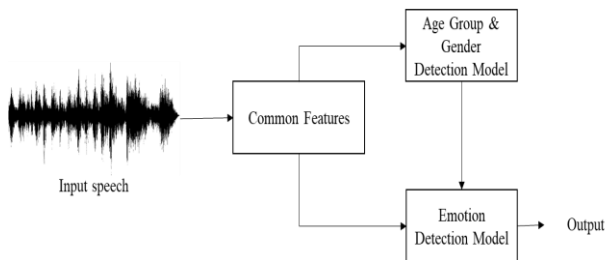


Fig 3.

- The research uses the Mozilla Common Voice dataset and a unified model to investigate age-group and gender recognition. During preprocessing, pertinent features and labels pertaining to audio frequency characteristics, gender, and age-group are extracted. A thorough representation of the audio data is provided by the extracted features, which include the spectral centroid, bandwidth, rolloff, and Mel-Frequency Cepstral Coefficients (MFCCs).

- Feature scaling is used to increase the robustness of the model, and the ANOVA statistical method is employed for feature selection. K-Fold Cross-Validation is used to train and assess two classifiers, Support Vector Machine (SVM) and Random Forest, with the F1-Score serving as a crucial performance indicator. Through an internal Cross-Validation method, classifier hyperparameters are optimized.
- The SVM classifier outperformed the Random Forest classifier with an accuracy of 82.7%, which is promising compared to its 71.8%

#### C. Approach 3: Gender, Emotion in Single model followed by age-group model

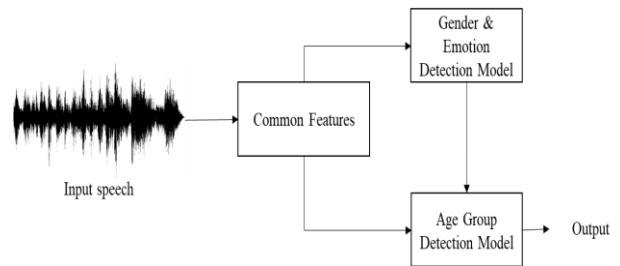


Fig 4.

- The research used the RAVDESS dataset, which contains labeled audio for both gender and emotion recognition. Emotions were categorized as neutral, happy, sad, and furious. The dataset underwent partitioning into training and testing sets. Key speech signal features, including Mel-Frequency Cepstral Coefficients (MFCCs), Zero Crossing Rate, and Root Mean Square Energy, were extracted using Wav2Vec2FeatureExtractor from transformers.
- HubertForSequenceClassification, the model architecture adopted, is a pre-trained model fine-tuned for sequence classification. Training consisted of two epochs, each of which processed batches of size 2 through a DataLoader. The optimizer was Adam, with a learning rate of 1e-5. The evaluation metric chosen was categorical accuracy. On the test dataset, the training phase produced a considerable accuracy of 74.57%, indicating the model's ability to distinguish gender and emotion from speech data. This demonstrates its potential for real-world applications needing in-depth voice analysis.

#### D. Approach 4: Unified model age-group group, Gender and Emotion

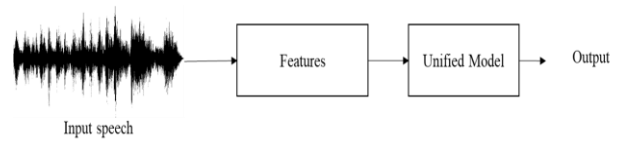


Fig 5.

- Using CREMA-D for emotion and Mozilla Common Voice for gender and age-group, the research explores a combination model for age-group, gender, and emotion recognition. The dataset is processed using encoding and one-hot encoding for the labels of emotion, age-group, and gender in order to separate features and labels. The dataset is divided, and then the input features are standardized.
- This approach aligns seamlessly with our overarching research goal of establishing a unified paradigm for comprehensive voice-based recognition. The findings

highlight that simultaneous age-group and gender classification significantly improves performance in contrast to single-task models. Various speaker datasets are employed for age-group and gender classification in the research, resonating with the broader scope of our exploration into unified model research.

- Convolutional Neural Network (CNN) layers for feature extraction and later dense layers for abstraction are integrated into the model architecture. The output layers—emotion, age-group, and gender—are customized for every recognition task. For each task, a categorical cross entropy loss version of the model is compiled and then optimized using the Adam optimizer. Fitting the model to training data and evaluating it on a test set comprises training.
- On the test dataset, the model shows 53.72% accuracy in predicting emotion, 41.77% accuracy in predicting age-group, and 48.09% accuracy in predicting gender.

A variety of metrics were used to evaluate the model's performance, and one such metric is the F1 score, which offers a balanced measurement based on recall and precision. Recall computes the ratio of genuine positives to real positives, whereas precision measures the accuracy of positive forecasts. The F1 score provides a single accuracy statistic by balancing precision and recall by incorporating erroneous positives and false negatives. Standard criteria such as accuracy, precision, and recall were used in addition to the F1 score to assess the model's efficacy. Accuracy offers a comprehensive assessment of overall performance, memory records pertinent positive examples, and precision measures accurate positive predictions [9].

#### IV. RESULT

Table 2.

Methods	Model Type			F1 Accuracy (%)		
	Age Group	Gender	Emotion	Age Group	Gender	Emotion
Approach 4		CNN		41.77	48.09	53.72
Approach 3	CNN		HFSC	71.25	74.57	74.57
Approach 2		SVM	LSTM	82.7	82.7	71.8
Approach 1		CNN	LSTM	71.3	89.4	90
S. R. Zaman[19]	CatBoost	GMM	XGBoost	96.4	97.67	66.1
Zheng Lian[20]		MLP	RNN	89.58	89.58	62.85
Zheng Lian [21]	GMM	GMM-EM	DANN	80	91	82.68

In our experimental endeavors, we harnessed the power of a Python environment, employing essential tools such as NumPy, pandas, scikit-learn, transformers, and TensorFlow. These tools facilitated robust data manipulation, preprocessing, and model construction. Our investigation honed in on three distinctive experiments, each contributing unique insights:

Firstly, leveraging the Mozilla Common Voice dataset, we delved into age-group and gender detection using a unified model, integrating Support Vector Machine (SVM) and Random Forest classifiers.

Secondly, we explored gender and emotion detection utilizing the HubertForSequenceClassification model from the transformers library, coupled with feature extraction from the RAVDESS dataset, all consolidated within a single model.

The third experiment fused CREMA D for emotion and Mozilla Common Voice for gender and age-group. This comprehensive examination involved a unified model for simultaneous age-group, gender, and emotion recognition. Feature extraction was facilitated by a Convolutional Neural Network (CNN) architecture, while preprocessing phases incorporated label encoding and one-hot encoding for emotion, age-group, and gender labels.

The main challenge faced was the lack of a unified dataset containing age-group, gender, and emotion labels, which was needed to develop an integrated model. To address this, multiple datasets were combined and preprocessed to create a synthetic dataset encompassing all required labels. Additionally, sequential modeling was explored, wherein the output of one model served as input to the next. However, the integrated model approach proved most

effective, outperforming individual models in metrics such as latency while achieving comparable accuracy. Novel deep neural network architectures were crucial to extract relevant acoustic features and capture correlations between vocal properties. The resulting consolidated model overcame the scarce training data hurdle and enabled real-time trifold age, gender, and emotion detection from voice samples. This pioneering work illuminates relationships between speech patterns across demographic and affective dimensions. It lays a strong foundation for the next wave of multifaceted voice analytics systems.

Our experimental setup encompassed a diverse array of resources, including preprocessing tools, machine learning models, and datasets. The outcomes were meticulously compared, with accuracy percentage-groups presented in an organized tabular format, as shown in the chart below. To enhance data comprehension, graphical representations were employed. This methodology not only ensures a thorough evaluation of model performance across varied recognition tasks but also provides a visually intuitive understanding of the results.

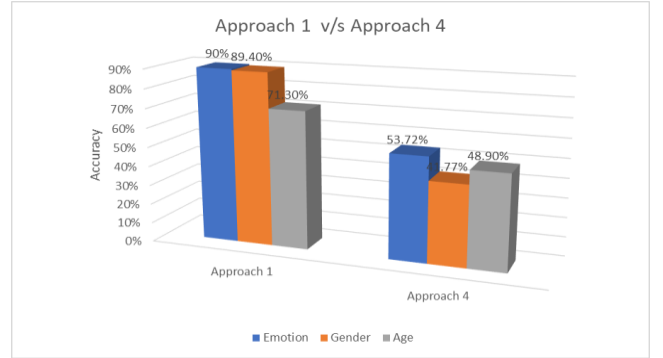


Fig 6.

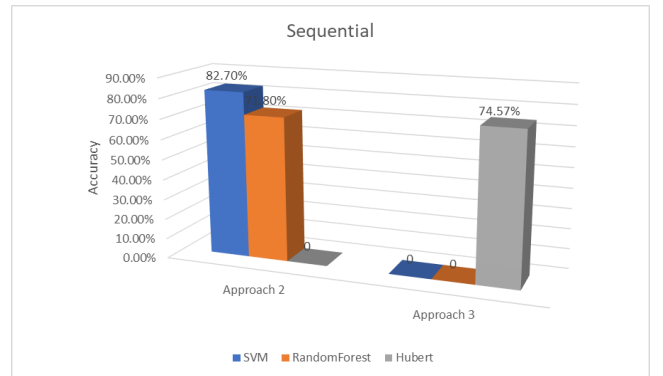


Fig 7.

#### V. CONCLUSION

In summary, our research has made significant strides in the realm of voice analysis, particularly in the domains of gender, age-group, and emotion recognition. The research's outcomes underscored the effectiveness of our integrated models in discerning these diverse features within audio data. Notably, the results highlighted the pivotal role of tailored model selection, with the Support Vector Machine (SVM) outperforming the Random Forest classifier in age-group and gender detection.

The gender and emotion identification model showcased the capability to extract intricate aspects from speech data with commendable accuracy levels. Encouragingly, the unified model addressing age-group, gender, and emotion detection demonstrated promising results, suggesting potential applications in scenarios requiring comprehensive voice analysis.

Looking ahead, there exists ample room for future endeavors. Possible enhancements to current models could involve fine-tuning hyperparameters, exploring more intricate neural network structures, and incorporating additional datasets to enhance diversity. Investigating transfer learning strategies, where models pretrained on extensive datasets are customized for our specific needs, holds promise for performance improvement. Additionally, addressing potential biases in datasets and refining preprocessing methods can contribute to the creation of more reliable and unbiased models.

Future research directions may include delving into real-time applications and adapting models for implementation in diverse contexts. Ultimately, our research lays a robust foundation, and subsequent projects can leverage-group these findings to propel the capabilities of voice analysis systems for a myriad of practical applications.

## REFERENCES

- [1] A. Mehrish, N. Majumder, R. Bhardwaj, R. Mihalcea, and S. Poria, "A Review of Deep Learning Techniques for Speech Processing," (2023).
- [2] S. R. Zaman, D. Sadekeen, M. A. Alfaz and R. Shahriyar, "One Source to Detect them All: Gender, age-group, and Emotion Detection from Voice," *2021 IEEE 45th Annual Computers, Software, and Applications Conference (COMPSAC)*, Madrid, Spain, 2021, pp. 338-343, doi: 10.1109/COMPSAC51774.2021.00055.
- [3] Sánchez-Hevia, H.A., Gil-Pita, R., Utrilla-Manso, M. *et al.* age-group group classification and gender recognition from speech with temporal convolutional neural networks. *Multimed Tools Appl* 81, 3535–3552 (2022)
- [4] Poonam Rani et al. International Journal of Recent Research Aspects ISSN: 2349-7688, Vol. 5, Issue 4, Dec 2018, pp. 14-17.
- [5] Jinkyu Lee, Ivan Tashev, High Level Feature Representation using Recurrent Neural Network for Speech Emotion Recognition (2015).
- [6] Gómez García, Jorge Andrés , Moro Velázquez, Laureano, Godino Llorente, Juan Ignacio and Castellanos Domínguez, Germán (2015). Automatic age-group detection in normal and pathological voice. In: "16th Annual Conference of the International Speech Communication Association (INTERSPEECH 2015)", 6/09/2015 - 10/09/2015, Dresden (Germany). ISBN a. 978-1-5108-1790-6. pp. 3739-3743.
- [7] Prasanta Roy, Parabattina Bhagath, and Pradip Das. 2020. Gender Detection from Human Voice Using Tensor Analysis. In Proceedings of the 1st Joint Workshop on Spoken Language-group Technologies for Under-resourced language-groups (SLTU).
- [8] Lian, Z., Tao, J., Liu, B., Huang, J., Yang, Z., Li, R. (2020) Context-Dependent Domain Adversarial Neural Network for Multimodal Emotion Recognition. Proc. Interspeech 2020, 394-398, doi: 10.21437/Interspeech.2020-1705 .
- [9] Scikit-Learn Classification Metrics Documentation.
- [10] Poonam Rani et al. International Journal of Recent Research Aspects ISSN: 2349-7688, Vol. 5, Issue 4, Dec 2018, pp. 14-17.
- [11] Jinkyu Lee, Ivan Tashev, High Level Feature Representation using Recurrent Neural Network for Speech Emotion Recognition (2015).
- [12] Gómez García, Jorge Andrés , Moro Velázquez, Laureano, Godino Llorente, Juan Ignacio and Castellanos Domínguez, Germán (2015). Automatic age-group detection in normal and pathological voice. In: "16th Annual Conference of the International Speech Communication Association (INTERSPEECH 2015)", 6/09/2015 - 10/09/2015, Dresden (Germany). ISBN (978-1-5108-1790-6. pp. 3739-3743.
- [13] 978-1-5108-1790-6. pp. 3739-3743.
- [14] Prasanta Roy, Parabattina Bhagath, and Pradip Das. 2020. Gender Detection from Human Voice Using Tensor Analysis. In Proceedings of the 1st Joint Workshop on Spoken Language-group Technologies for Under-resourced language-groups (SLTU).
- [15] Lian, Z., Tao, J., Liu, B., Huang, J., Yang, Z., Li, R. (2020) Context-Dependent Domain Adversarial Neural Network for Multimodal Emotion Recognition. Proc. Interspeech 2020, 394-398, doi: 10.21437/Interspeech.2020-1705 .
- [16] Scikit-Learn Classification Metrics Documentation.
- [17] Spectral Centroid:The spectral centroid is a measure used in digital signal processing to characterize a spectrum.Bandwidth:Bandwidth specifically refers to the capacity at which a network can transmit data.Flatness:Flatness is a soft, short tone heard when percussing over solid tissue like muscle and bone. Contrast:small differences in speech sounds, that makes a difference in how the sound is perceived by listeners
- [18] In sound processing, the mel-frequency cepstral Coefficient (MFCC) is a representation of the short-term power spectrum of a sound, based on a linear cosine transform of a log power spectrum on a nonlinear mel scale of frequency.
- [19] age-group and emotion:: S. R. Zaman, D. Sadekeen, M. A. Alfaz and R. Shahriyar
- [20] emotion:Zheng Lian, Jianhua Tao, Bin Liu, Jian Huang, Zhanlei Yang, Rongjun Li
- [21] emotion:Jinkyu Lee, Ivan Tashev