

AN EFFECTIVE WAY TO EXTRACT FEATURES FROM AUDIO

T RAMESH¹ A SRINIVASA RAO² N.TEAJA PRAKASH³ K.ANUSHA⁴

¹ASST.PROFESSOR, DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING,

²PROFESSOR, DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING,

³ASST.PROFESSOR, DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING,

⁴ASST. PROFESSOR, DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING,
^{1,2,3,4} SRI MITTAPALLI COLLEGE OF ENGINEERING

Abstract:

Expressing one's perspective to others begins with one's emotions. A person's emotional state may be conveyed via spoken or written word. When it comes to conveying feelings, audio is much more effective than text. In this study, a deep learning network is used to detect various emotions in human speech. Human emotions are often categorised as angry, silent, disgusted, afraid, happy, neutral, sad, and shocked in order to facilitate strong emotion recognition. When exposed to the disturbing effect pervasive in the environment, the acoustic characteristics determine the audio signal's reaction. Therefore, these elements play a crucial role in audio emotion recognition. We are using

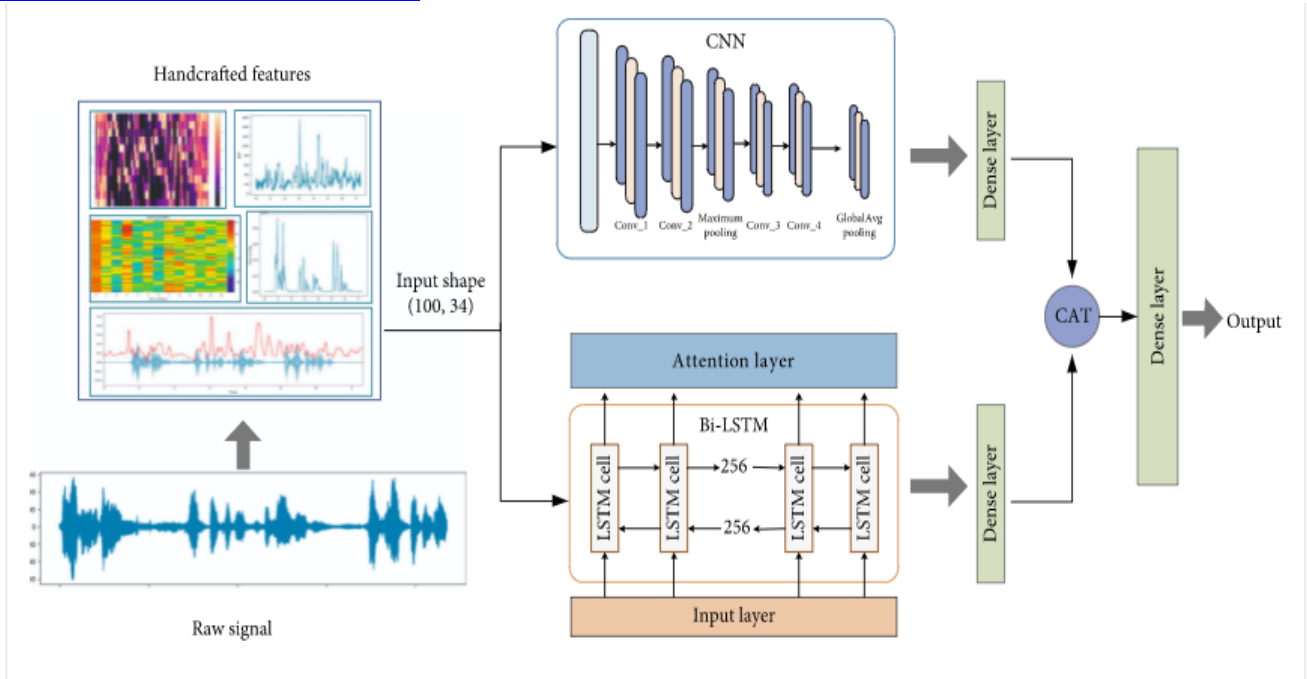
MFCC (Mel-Frequency Cepstral Coefficients) on the audio data to extract important characteristics. The most effective way to extract features from audio is MFCC. The representation of the recurrence bands in MFCC is far closer to the recurrence bands of the human auditory system. In addition, we are using increasing approaches on the audio data in order to get precise characteristics. Additionally, deep learning neural networks are entrusted with the standardised audio data in order to detect auditory emotions. The libraries Keras and TensorFlow are being used to construct deep learning neural networks.

1.Introduction

From a corporate and academic perspective, the importance of emotion recognition via audio is rapidly growing. Based on their feedback, we were able to deduce that the customer was interested in the service or product. Identifying and capturing the emotional essence of audio is crucial when dealing with consumer feedback that is presented in an auditory format. In order to enhance the administrations, such data is crucial. Emotion recognition has many purposes: it is a grammatical feature, it helps with call categorisation, and it is used in autonomous car systems. Emotion recognition relies on extracting auditory

characteristics. The acoustic characteristics of an aggravation are its sound wave progressions. In addition, the examples derived from the acoustic properties of sound will aid in perceiving the emotions. Various methods exist for extracting elements. Among them, MFCC excels in feature separation. The Mel-scale characteristics are in agreement with audible human speech. It all starts with getting the range's force value, and then using that value to generate the Cepstrum. The MFCC uses the first 12 cepstrum values taken from the Cepstrum and treats the rest as unique variations from one case to the next. The features will be separated from the audio when MFCC has standardised.

<https://ijgst.com.2023.v12.i1.pp45-53>



A portion of growth is also critical to get remarkable results. By creating synthetic data from raw data, the expansion processes are used to boost the quantity of acoustic information. Modifying the original audio data often entails injecting clamour, altering pitch, moving time, and speed. The newly synthesised data is superimposed over the existing audio data. The NumPy and librosa libraries are used to implement increase strategies in this study. With more data, the model can better forecast how people would feel. Additionally, the MFCC esteems are being scaled using Mean. Improving a model's efficiency is as simple as expanding on the audio data. The brain and other central nervous systems of living things serve as inspiration for a class of observable learning models known as artificial neural networks (ANNs). You may use them to estimate

capabilities that depend on a lot of unknowns. The "neurones" of a fake neural network are all linked and communicate with each other. Neural networks are quite adaptable when it comes to input and learning since associations have weight numerical features that may be allocated or altered to encounter. Here, we use an ANN (Artificial Neural Network) with a convolutional neural network. The convolution neural network relies on convolution layers as its fundamental building pieces. If you want to know the honest truth, convolution is just a bunch of channels applied to your input data to make starting maps. The output of the convolution layers is essentially what the actuation maps are. Using the Sequential model, we are able to classify emotions with excellent accuracy. On top of that, we are

<https://ijgst.com.2023.v12.i1.pp45-53>

improving processing.

2 Literature Survey

2.1 Perceptual Audio Features for Emotion detection:

The authors of this piece put out an alternative configuration of auditory characteristics for automated emotion recognition in audio. The characteristics are derived from the perceptual quality measures provided in the ITU BS.1387 recommendation for perceptual evaluation of audio quality. They begin with the hear-able framework's models of the outside and inside ears, and they use the masked perceptual din, which specifies broad emotion detection targets, as their feature set. Thus, the tools Gerda, openEAR toolkit, and Hidden Markov Toolkit were used to do this. As an added bonus, they made use of EMO-DB and VAM, two databases. Crema D, SAVEE, TESS, RAVDESS, and the MFCC Speech emotional database were among the five distinct databases used in the planned study, which also made use of GoogleSpec. Measures of changes of worldly envelopes, measures of sounds of the emotional contrast, the event likelihood of emotional squares, perceptual bandwidth, fractional tumult of the emotional distinction, and emotional contrast to-perceptual mask proportion are some of the features processed in basic bands based on the reference idea. To evaluate the classifier's performance, we provide a nuanced bigger part voting choice metric that fortifies the conventional larger part voting. For both "all" and "valence" tasks, EMO-DB and VAM achieve 7-11% and 7-16% improvement in emotion acknowledgement rates, respectively, when compared to state-

of-the-art frameworks like Generalised Discriminant Analysis, Hidden Markov Toolkit, and Munich Open-Source Emotion and Affect Recognition Toolkit.

2.2 The syllabus as perceptual unit in Speech Perception:

One example of an acknowledgement difficulty is speech perception. Perceivers are tasked with identifying the most accurate representation of a given voice input. In order to build up a neurological code in a pre-perceptual hear-able capability, the hear-able receptor framework transforms an audible improvement. This buffer stores the data in a pre-perceptual state for around 250 ms, when the acknowledgement cycle ought to take place, according to my regressive masking experiments and other exploratory criteria. A synthesised percept is created from the pre-perceptual image via the acknowledgement interaction. The question of what kinds of examples are helpful for speech acknowledgement arises as a result of this structure. These concrete instances of sound are referred to as perceptual units.

2.3 Emotion Detection Using MFCC and Cepstrum Features:

In the new year, Speech Emotion Recognition (SER) is the subject of a massive investigation with the overarching goal of bettering human-machine interaction. The effect of cepstral coefficients on emotion

<https://ijgst.com.2023.v12.i1.pp45-53>

recognition is carried out in this study. The same holds true for the effects of cepstrum, Mel-recurrence Cepstral Coefficients (MFCC), and synthetically widened MFCC coefficients on emotion categorisation; a comparable study is forthcoming. In comparison to Firoz Shah's earlier work, which only accurately identified four emotions in the Berlin speech corpus, our algorithm, which used a conservative element vector, shown improved acknowledgement speeds for recognising seven emotions. When compared to InmaMohino's computation, which used a component vector consisting only of synthetically extended MFCC coefficients, the suggested method achieves a much lower misclassification productivity.

3. Proposed System:

Crema D, Ravdess, Tess, SAVEE, and the MFCC Speech passionate data set are the five separate sources of audio data used to maintain the model in this article. We are using distinct growth tactics to overcome the challenge of insufficient information. Engineered data may be generated via expansion simply by changing pitch, changing speed, moving, expanding, repetition hiding, and time covering. Using the MFCC approach with 58 mel highlights, we are extracting highlights. Getting the processed data into the Convolution Neural Network (CNN) for model construction is another task. We used encoding methods to accomplish this task. In addition, after running the train test partition on the data, the Standard Scalar technique was applied to it. Following these modifications, we expand the information's components to make it feasible for demonstration. Finally, we handled all of the handled data and took care

of the model.
Data Structure: TensorFlow

In many areas of computer science, such as computer vision, imperfection location to seek after exploration, computational drug discovery, text summarisation, information retrieval, and slant examination, ML frameworks are integrated into manufacturing using TensorFlow, an interface for communicating AI calculations. Using TensorFlow as its backend, the suggested model implements a Sequential Convolutional Neural Network (CNN) with a few layers.

Keras:

Keras provides essential insights and building blocks for the rapid development and transfer of ML setups. It makes excellent use of TensorFlow's scalability and cross-stage capabilities. Layers and models are the core of Keras's information structure. All of the CNN model's layers are implemented using Keras. In this study, we used keras to run auto encoders on the data. In data processing, it helps arrange the model as a whole and converts the class vector to a binary class matrix.

Librosa:

One such Python package for analysing audio and music is Librosa. To build frameworks for music data recovery, it provides the necessary building pieces. Librosa is a fantastic Python package for working with and analysing audio. For many applications, such as differentiating between human voices and detecting unique characteristics in sounds, it is the first step towards handling large amounts of audio data. In addition, librosa is helpful for displaying waveplots and specplots, which reveal how the sound is

<https://ijgst.com.2023.v12.i1.pp45-53>

being processed. We used librosa for a variety of sound preparation tasks in this study.

Crema D:

Lots of data are lost since most sound datasets use the same amount of speakers. CREMA-D is a multi-speaker organisation. If you want to make sure your model doesn't become overfit, the CREMA-D is a great dataset to utilise for this purpose. The CREMA-D database contains details on 7,442 distinct clasps belonging to 91 different performers. The 48 male and 43 female performers who contributed clasps ranged in age from 20 to 74 years old and belonged to a wide range of racial and ethnic groups, including African Americans, Asians, Caucasians, Hispanics, and the Unspecified. A total of twelve phrases were available to the performers. One of six distinct emotions—Anger, Disgust, Fear, Happy, Neutral, or Sad—and four distinct degrees of emotion—Low, Medium, High, or Unspecified—were used to introduce the statements.

MFCCs for Speech Emotion Recognition:

An attempt to foretell the extent to which men and women would feel reliant on MFCC values informed the creation of this dataset. Using this setup (58 attributes for each emotion), we might have achieved a respectable 94% accuracy on the female emotions.

4. Experimental Results

4.1. Examination Setup

4.1.1. Data set

In the SAIL Laboratory of USC, Busso et al. collected and recorded the first IEMOCAP data set [21]. Approximately twelve hours of general media content (video, sound, text, MOCAP, etc.) make up the presentation date or the made-do. On top of that, it's one of the largest publicly available multimodal databases that welcomes enthusiastic contributions. Referring to [23], this study selects the four most prevalent emotional states—outrage, energise, indifferent, and miserable—representing more than 70% of the sample for exploratory enthusiastic classifications.

There is still a distinct gap between the presentation data and people's ordinary sentiments in everyday contact due to the exhibition data's poor believability and embellishment. Information that is ad libbed, unlike to information that is executed, is more accurate. Therefore, we use impromptu information for conversation emotion acknowledgement while considering the information's authenticity.

4.1.2. Assessment Matrices

In this article, four different types of emotions are examined using disorder grids. Each rising line represents the actual result, while each flat line represents the expected result. Any kind of emotion may be recognised and organised using a disorder grid, which is a visual tool. To evaluate the validity of the model, we combined a disarray network with weighted exactness (WA) and unweighted precision (UA) in this research.

4.2. Test Settings

4.2.1. Text Emotion Recognition Models

For the purpose of evaluating the content emotion recognition model, we constructed two models (T-LSTM and T-BL). One model that made use of the LSTM layer was the content LSTM (T-LSTM) model. The T-BL

<https://ijgst.com.2023.v12.i1.pp45-53>

model, which stands for content Bi-LSTM, made use of the bidirectional LSTM layer. For their final layout, both models used the softmax layer. A total of 256 LSTM and Bi-LSTM cells should be sufficient.

4.2.2. Discourse Emotion Recognition Models

Four models were developed for discourse emotion recognition: S-CNN, S-DNN, S-CL, and S-CBLA. Unrehearsed conversation, which includes a wealth of setting-related data, is the source of the sound information. Therefore, in our Speech-CBLA (S-CBLA) model, we used the bi-LSTM with consideration instrument and the CNN with double channel organisation. Also, we established the Speech-CNN (S-CNN), Speech-DNN (S-DNN), and Speech-CNN-LSTM (S-CL) models as a correlation. The S-CL model made use of the standard CNN and Bi-LSTM single-channel blend techniques. The S-DNN model was trained using four dense layers with 512, 256, 128 and 32 boundary weights, respectively. A max pooling layer, a global normal pooling layer, and convolutional layers were all part of Model S-CNN, which was designed to be comparable to the CNN direct in Model CBLA. The execution was configured as Leaky ReLU. The results of comparing different models' effects on accuracy are shown in Figure 5. It shows that the display becomes better for most enthusiastic classes as a result of model enhancement.

4.2.3. Multimodal Emotion Recognition Models

We used the component level methods mentioned in the previous part for the modular combination. The goal of using thick layers was to merge the highlights that were previously distinct from the top layers. Each hub in the higher layers was connected to

every hub in the thick layer. Ten24, 512, and 4 were the units that were agreed upon. For managing discourse and text-based components, the multimodal-CBLA-Bi-LSTM (M-CBLA-BL) model combined S-CBLA and T-BL. Once the indisputable level highlights from different modes were mixed, we used DNN to deduce the link between the multimodal highlights by nonlinear transformation. Finally, we grouped using the softmax layer. In addition, the M-D-BL model combined S-DNN and T-BL, whereas the M-CL-BL model combined S-CL and T-BL using multimodal-CNN-LSTM-Bi-LSTM. An approach comparable to that of model M-D-BL was used by a broad variety of other models. Regularisation was a foundational component of both the M-CL-BLR and M-CBLA-BLR models.

5 Conclusion

The proliferation of high-quality internet media has piqued the interest of many scholars in studying multimodal emotion recognition. This study proposes a multimodal emotion recognition algorithm based on conversation and text using the IEMOCAP database to address the limitations of the single mode's emotional data. In order to learn acoustic emotional highlights, the model made use of the CNN and LSTM double channel. It also used Bi-LSTM to filter out literary flourishes. The combination highlights were also familiarised with the help of a deep neural network. The model took into account both relevant and globally distributed ephemeral data in the data, used deep learning organisations, and made use of both manually constructed and significant level highlights. The model was further improved by making advantage of regularisation. Our multimodal model outperformed previous distributed multimodal models on the test datasets and

<https://ijgst.com.2023.v12.i1.pp45-53>

had a greater acknowledgement accuracy than the single modular model, according to the trial findings. Researching more effective element combination techniques to enhance the multimodal model's execution should be prioritised in the future. Additionally, we will strive to simplify the discourse modular architecture in order to get more effective sound sensation highlights.

References

1. M. Hasan, E. Rundensteiner, and E. Agu, "Automatic emotion detection in text streams by analyzing twitter data," *International Journal of Data Science and Analytics*, vol. 7, no. 1, pp. 35–51, 2019. View at: Publisher Site | Google Scholar
2. M. A. Razek and C. Frasson, "Text-based intelligent learning emotion system," *Journal of Intelligent Learning Systems and Applications*, vol. 09, no. 1, pp. 17–20, 2017. View at: Publisher Site | Google Scholar
3. C.-H. Chen, W.-P. Lee, and J.-Y. Huang, "Tracking and recognizing emotions in short text messages from online chatting services," *Information Processing & Management*, vol. 54, no. 6, pp. 1325–1344, 2018. View at: Publisher Site | Google Scholar
4. J. K. Rout, K.-K. R. Choo, A. K. Dash, S. Bakshi, S. K. Jena, and K. L. Williams, "A model for sentiment and emotion analysis of unstructured social media text," *Electronic Commerce Research*, vol. 18, no. 1, pp. 181–199, 2018. View at: Publisher Site | Google Scholar
5. K. S. Rao, S. G. Koolagudi, and R. R. Vempada, "Emotion recognition from speech using global and local prosodic features," *International Journal of Speech Technology*, vol. 16, no. 2, pp. 143–160, 2013. View at: Publisher Site | Google Scholar
6. W. Zheng, M. Xin, X. Wang, and B. Wang, "A novel speech emotion recognition method via incomplete sparse least square regression," *IEEE Signal Processing Letters*, vol. 21, no. 5, pp. 569–572, 2014. View at: Publisher Site | Google Scholar
7. S. Gupta, A. Mehra, and Vinay, "Speech emotion recognition using SVM with thresholding fusion," in *Proceedings of the 2015 2nd International Conference on Signal Processing and Integrated Networks (SPIN)*, pp. 570–574, Noida, India, February 2015. View at: Publisher Site | Google Scholar
8. J.-S. Park, J.-H. Kim, and Y.-H. Oh, "Feature vector classification based speech emotion recognition for service robots," *IEEE Transactions on Consumer Electronics*, vol. 55, no. 3, pp. 1590–1596, 2009. View at: Publisher Site | Google Scholar
9. K. Lu and Y. D. Jia, "Audio-visual emotion recognition with boosted coupled HMM," in *Proceedings of the Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, pp. 1148–1151, Tsukuba, Japan, November 2012. View at: Google Scholar

<https://ijgst.com.2023.v12.i1.pp45-53>

10. S. S. Narayanan, S. Lee, and A. Metallinou, "Audio-visual emotion recognition using Gaussian mixture models for face and voice," in *Proceedings of the 2008 Tenth IEEE International Symposium on Multimedia*, pp. 250–257, Berkeley, CA, USA, December 2008. View at: Publisher Site | Google Scholar

11. D. Li and J. Qian, "Text sentiment analysis based on long short-term memory," in *Proceedings of the 2016 First IEEE International Conference on Computer Communication and the Internet (ICCCI)*, pp. 471–475, Wuhan, China, October 2016. View at: Publisher Site | Google Scholar

12. L. L. Chao, J. H. Tao, M. H. Yang, Y. Li, and Z. Wen, "Long shot term memory recurrent neural network based on encoding method for emotion recognition in video," in *Proceedings of the 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2752–2756, Shanghai, China, March 2016. View at: Publisher Site | Google Scholar

13. J. Zhao, X. Mao, and L. Chen, "Learning deep features to recognise speech emotion using merged deep CNN," *IET Signal Processing*, vol. 12, no. 6, pp. 713–721, 2018. View at: Publisher Site | Google Scholar

14. L. Q. Cai, X. L. Liu, F. L. Chen, and M. Xiang, "Robust human action recognition based on depth motion maps and improved convolutional neural network," *Journal of Electronic Imaging*, vol. 27, no. 5, Article ID

051218, 2018. View at: Publisher Site | Google Scholar

15. T. N. Sainath, O. Vinyals, A. Senior, and H. Sak, "Convolutional, long short-term memory, fully connected deep neural networks," in *Proceedings of the 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4580–4584, Brisbane, Australia, April 2015. View at: Publisher Site | Google Scholar

16. G. Trigeorgis, F. Ringeval, R. Brückner et al., "Adieu features? End-to-end speech emotion recognition using a deep convolutional recurrent network," in *Proceedings of the 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5200–5204, Shanghai, China, March 2016. View at: Publisher Site | Google Scholar

17. S. Poria, I. Chaturvedi, E. Cambria, and A. Hussain, "Convolutional MKL based multimodal emotion recognition and sentiment analysis," in *Proceedings of the 2016 IEEE 16th International Conference on Data Mining (ICDM)*, pp. 439–448, Barcelona, Spain, December 2016. View at: Publisher Site | Google Scholar

18. P. Tzirakis, G. Trigeorgis, M. A. Nicolaou, B. W. Schuller, and S. Zafeiriou, "End-to-end multimodal emotion recognition using deep neural networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1301–1309, 2017. View at: Publisher Site | Google Scholar

<https://ijgst.com.2023.v12.i1.pp45-53>

19. L. C. Woo, K. Y. Song, J. Jeong, and W. Y. Choi, "Convolutional attention networks for multimodal emotion recognition from speech and text data," 2018, <https://arxiv.org/abs/1805.06606>. View at: Google Scholar

20. Y. Gu, S. Chen, and I. Marsic, "Deep multimodal learning for emotion recognition in spoken language," in *Proceedings of the 2018 IEEE International Conference on Acoustics, Speech and Signal Processing*, Calgary, Canada, April 2018. View at: Publisher Site | Google Scholar

21. S. Yoon, S. Byun, and K. Jung, "Multimodal speech emotion recognition

using audio and text," in *Proceedings of the 2018 IEEE Spoken Language Technology Workshop (SLT)*, pp. 112–118, Athens, Greece, December 2018. View at: Publisher Site | Google Scholar

22. Srinu, Nidamanuri, Sampathi Sivahari, and Mastan Rao Kale. "Leveraging Radial Basis Function Neural Networks for Rainfall Prediction in Andhra Pradesh." *2022 International Conference on Computer, Power and Communications (ICCPC)*. IEEE, 2022.

23. Muppavarapu, Rajasekhar, and Mastan Rao Kale. "An Effective Live Video Streaming System."