

Speech Emotion Recognition

Subjects: **Computer Science, Artificial Intelligence**

Contributor: Konstantinos Mountzouris , Isidoros Perikos , Ioannis Hatzilygeroudis

Speech is the most natural way of human communication. Affective computing systems based on speech play an important role in promoting human–computer interaction, and emotion recognition is the first step. Due to the lack of a precise definition of emotion and the inclusive and complex influence of emotion generation and expression, accurately recognizing speech emotions is still difficult. Speech emotion recognition (SER) is an important problem that is receiving increasing interest from researchers due to its numerous applications, such as e-learning, clinical trials, audio monitoring/surveillance, lie detection, entertainment, video games, and call centers.

speech emotion recognition

deep learning

deep belief network

1. Introduction

Machine learning (ML) is a revolutionary method in which we feed a machine an adequate amount of data, and the machine will use the experience gained from the data to improve its own algorithm and process data better in the future [1]. One of the most significant approaches in machine learning is the use of neural networks (NNs). Neural networks are networks of interconnected nodes and are loosely modeled towards the way the human brain processes information. Neural networks store data, learn from it, and improve their abilities to sort new data. For example, a neural network with the task of identifying dogs can be fed a set of characteristic values extracted from various images of dogs tagged with the type of dog. Over time, it will learn what kind of image corresponds to what kind of dog. The machine therefore learns from experience and improves itself. Deep learning (DL) is a predominant machine learning approach, where neural networks are enabled at nodes to remember past values with many layers that are trained using massive amounts of data. In deep learning, the sprawling artificial neural network is fed representations of raw data (e.g., raw image representations) and not given any other instructions. This means that, in contrast to other machine learning approaches, it determines the important characteristics and purpose of the data itself while storing it as experience. In other words, according to studies, deep neural networks (DNNs) can solve the data representation problem through learning a series of task-specific transformations [2]. The network layers extract abstract representations and filter out the irrelevant information, which leads to a more accurate classification and better generalization. Temporal models were also proposed for modelling sequential data with mid- to long-term dependencies. Deep learning models are currently used to solve problems such as face recognition, voice recognition, image recognition, computational vision, and speech emotion recognition. One of the main advantages of deep learning techniques over other machine learning techniques is the automatic selection of features, which could, for example, be applied to important features inherent in audio files that have a special emotion in the task of recognizing speech emotions.

When it comes to recognizing emotion through speech, deep learning models, such as convolutional neural networks (CNN), deep neural networks (DNNs), deep belief networks (DBNs), etc., approach the detection of high-level features for better accuracy compared to hand-made low-level features. Furthermore, the use of deep neural networks enhances the computational complexity of the entire model. However, according to Mustaqeem and Kwon [3] there are still many challenges in recognizing emotion from speech, such as the fact that the current CNN architectures have not shown significant improvement in speech accuracy and complexity in speech signal processing, or the fact that the use of recurrent neural networks (RNNs) and long short-term memory neurons (LSTMs) is useful for training sequential data, but they are difficult to train effectively and are computationally more complex.

Due to the above issues and challenges, researchers implement and test, apart from single instances, CNN- and LSTM-based architectures integrated with an attention mechanism, aiming to explore the impact of the attention mechanism on their performance. Overall, the CNN architecture, with the addition of an attention mechanism, achieved a better performance. The voice characteristics of the speakers are extracted in the form of Mel Frequency Cepstral Coefficients (MFCCs), with the help of the Librosa library.

2. Speech Emotion Recognition

There is a huge research interest, and several works attempt to perform emotion detection from speech [4][5]. Various works study the way that emotions can be automatically identified and accurately recognized in speech data [6][7]. In this regard, deep learning techniques have achieved breakthrough performance in recent years, and as a result, have been thoroughly examined by the research community [8][9]. Many existing studies in the literature have focused on improving and extending deep learning techniques [10].

In the work presented in [11], the authors present a new random deep belief network (RDBN) method for speech emotion recognition, which consists of a random subspace, DBN and SVM in the context of ensemble learning. It first extracts the low-level characteristics of the input speech signal and then applies them to the construction of many random sub-intervals. Second, it creates many different sub-intervals. In addition, the DBN continues to use the stochastic gradient descent method to optimize the parameters. To solve the problem, a random space is applied for the training of the basic classifiers as a whole, where the same classification method is used. The best accuracy achieved is 82.32% on the Emo-DB database, 48.5% on the CASIA database, 48.5% on the FAU database, and 53.60% on the SAVIEE database.

In the work presented in [12], the authors introduce a method for identifying speech emotions using a spectrogram and convolutional neural network (CNN). The proposed model consists of three convolution layers and three fully connected layers, which extract distinctive features from spectrograph images and predictions for the seven emotions of the Emo-DB Database. Layer C1 has 120 cores (11×11) applied at a rate of four pixels. The ReLU acts as an activation function instead of the standard sigmoid functions that improve the efficiency of the educational process. Layer C2 has 256 cores of size 5×5 and is applied to the input with one step. Similarly, C3 has 384 cores of size 3×3 . Each of these convolution layers are followed by ReLUs. Layer C3 is followed by three

FC layers that have 2048, 2048, and 7 nodes, respectively. More than 3000 spectrograms were generated from all the audio files in the dataset. Overall, the proposed method achieved 84.3% accuracy.

In [13], the authors present two convolutional neural networks with a long-short memory network (CNN-LSTM), one one-dimensional (1D) and one two-dimensional (2D), stacking four designed local features learning blocks (LFBL). The 1D CNN-LSTM network is intended to recognize the feeling of speaking from raw audio clips, while the 2D CNN-LSTM network focuses on learning high-level capabilities from log-Mel spectrograms. The experimental study was conducted on the Berlin Emo-DB and IEMOCAP databases. The 1D CNN LSTM network achieved 92.34% and 86.73% recognition accuracy on the speaker-dependent and speaker-independent EmoDB databases, respectively, and also delivered 67.92% and 79.72% recognition accuracy on the IEMOCAP speaker-dependent and speaker-independent databases, respectively. The 2D CNN LSTM network achieved 95.33% and 95.89% recognition accuracy on the speaker-dependent and speaker-independent Emo-DB databases, respectively, and delivered 89.16% and 85.58% recognition accuracy on the IEMOCAP speaker-dependent and speaker-dependent experiment databases, respectively.

In the work presented in [14], the authors proposed a new approach to the multimodal recognition of emotions from simple speech and text data. The attention network implemented consists of three separate convolutional neural networks (CNNs), two for extracting features from speech spectrograms and word integration sequences and one for the emotion classifier. The CNN outputs from word integration and spectrograms are used to calculate an attention matrix to represent the correlation between word integration and spectrograms in relation to emotion signaling. To evaluate the model, they used audio and text data from the CMU-Multimodal Opinion Sentiment and Emotion Intensity (CMU-MOSEI) dataset. The dataset is organized by video IDs and corresponding segments with six emotion and sentiment labels. The video IDs are then further split into segments. The training set consisted of 3303 video IDs and 23,453 segments, while the validation set consisted of 300 non-overlapping video IDs and 1834 segments. The total accuracy of the proposed method was 83.11%.

In [15], the authors present three methods based on CNNs in combination with extensive features, a CNN + RNN and ResNet, respectively. The authors investigate different types of features as the end-to-end frame input, including primary wave data, the Q-transform constant spectrogram (CQT), and the Fourier transform short-term spectrogram (STFT). In this way, the authors create multiple data samples with slightly modified speed ratios, which helps them achieve significant improvements and handle the overfitting issue in the framework from end to end. For their experiments, they used the EmotAsS dataset. The CNN + RNN model achieved the best performance (45.12%) with data balancing; the CNN model in combination with features showed a performance of 34.33% with data balancing, while the ResNet model achieved a performance of 37.78%.

In the work presented in [16], the authors propose a new architecture called attention-based 3-dimensional convolutional recurrent neural networks (3-D ACRNN) for recognizing emotion from speech, combining CRNN with an attention mechanism, because they hypothesized that calculating delta and delta-deltas for individual functions not only retains effective emotional information but also reduces the effect of emotionally unrelated factors, leading to a reduction in misclassification. First, the CNN 3-D is applied to the entire logarithmic-Mel spectrogram, which

has been compiled into a patch that contains only multiple frames. The attention layer then takes a sequence of high-level attributes as the input to generate expression-level attributes. The authors evaluated the model using the Berlin Emotional Speech Database (Emo-DB) and IEMOCAP database. From the ten speakers, for each evaluation, they selected eight as the training data, one as the validation data, and the rest as the test data. The method achieved an accuracy of 64.74% on IEMOCAP database and 82.82% on Emo-DB.

In the work presented in [17], the authors propose an attention-pooling representation learning method for recognizing emotions from speech (SER). Emotional representation is learned from end to end by applying a deep convolutional neural network (CNN) directly to speech spectrograms extracted from speech. Compared to existing aggregation methods, such as max pooling and average pooling, the proposed attention pooling can effectively integrate bottom-up class-agnostic attention maps and top-down class-specific attention maps. Given an expression, they segment it into 2 s sections for training and use an overlay of 1 s to allow them to receive more training data. Each section corresponds to the same tag with the corresponding expression. They used a 1×1 convolutional layer after Conv5 to create a top-down attention map and used another 1×1 convolutional layer to create bottom-up attention maps. The IEMOCAP improvised dataset was used, and the accuracy achieved by the proposed method was 71.75% for WA and 68.06% for UA.

In [18], the authors explore how to take full advantage of low-level and high-level audio features taken from different aspects and how to take full advantage of DNN's ability to merge multiple information to achieve better classification performance. For this reason, they proposed a hybrid platform consisting of three units, namely, a features extraction unit, a heterogeneous unification unit, and a fusion network unit. Besides low-level acoustic features, such as IS10, MFCCs, and eGemaps, that are extracted, high-level acoustic feature presentations named SoundNet bottleneck feature and VGGish bottleneck feature are considered for speech emotion recognition tasks. The heterogeneous integration unit is a Denoising AutoEncoder (DAE), which is a multilayer feed-forward neural network and is introduced in order to convert the heterogeneous space of various features into a unified representation space by deploying this unsupervised feature learning technique. The fusion network module is utilized to capture the associations between those unified joint features for emotion recognition tasks and is constructed as a four-layer neural network, containing one input layer and three hidden layers. They evaluated the model using the IEMOCAP database, and the proposed method improved the recognition performance, reaching an accuracy of 64%.

In the work presented in [19], the authors propose a platform that, at the training layer, has three main stages, such as verbal/non-verbal audio segmentation, the integration of feature extraction, and the construction of an emotion model. The verbal sections were used to train the CNN-based emotion model to derive emotion features, while the non-verbal sections were used to train the CNN audio model to extract audio features. The CNN's combined features are used as the input to the LSTM-based sequence-to-sequence emotion recognition model. Here, the sequence-to-sequence model based on the LSTM with an attention mechanism was selected for emotion recognition. The LSTM and the attention mechanism for developing a sequence emotion recognition model contained a bidirectional LSTM (Bi-LSTM) as the coder for the attention mechanism and an unidirectional LSTM as

the decoder for the emotional sequence output. They evaluated the model using the NTHU-NTUA Chinese interactive multimodal emotion corpus (NNIME); the proposed method achieved a 52.0% accuracy.

The work presented in [20] introduces a model that includes one-dimensional convolutional layers combined with dropout, batch-normalization, and activation layers. The first layer of their CNN receives 193×1 number arrays as the input data. The initial layer is composed of 256 filters with a kernel size of 5×5 and a stride of 1. After that, batch normalization is applied, and its output is activated by a rectifier linear units layer (ReLU). The next convolutional layer, consisting of 128 filters with the same kernel size and stride, receives the output of a previous input layer. The final convolutional layer, with the same parameters, is followed by the flattening layer and dropout layer, with a rate of 0.2. Their model was tested in the Berlin (EMO-DB), IEMOCAP, and RAVDESS databases and obtained 71.61% for the RAVDESS with eight classes, 86.1% for EMO-DB with 535 samples in seven classes, 95.71% for EMO-DB with 520 samples in seven classes, and 64.3% for IEMOCAP with four classes on speaker-independent audio classification tasks.

Attention-oriented parallel convolutional neural network encoders that capture the essential features required for emotion classification are introduced in [21]. The authors extracted and encoded features such as paralinguistic information and speech spectrogram data, and distinct CNN architectures were designed for each type of feature, and those encoded features were subsequently passed through attention mechanisms to enhance their representations before undergoing classification. Empirical evaluations were carried out on the EMO-DB and IEMOCAP open datasets, and the proposed model achieved a weighted accuracy (WA) of 71.8% and an unweighted accuracy (UA) of 70.9%. Furthermore, with the IEMOCAP dataset, the model yielded WA and UA recognition rates of 72.4% and 71.1%, respectively.

The authors in [22] present a work on enhancing the overall generalization performance and accuracy of SER with a balanced augmented sampling technique on spectrograms that aims to address the imbalance in sample distribution among emotional categories. A deep neural network is utilized, comprising the combination of a convolutional neural network (CNN) and an attention-based bidirectional long short-term memory network (ABLSTM) for feature extraction. Multitask learning is incorporated to enhance the deep neural network's performance. The methodology is assessed on the IEMOCAP and MSP-IMPROV databases, yielding a weighted average recall and unweighted average recall of 70.27% and 66.27% on the IEMOCAP database, respectively, while on the MSP-IMPROV database, the approach achieves 60.90% and 61.83%, respectively.

In the work presented in [23], the authors introduce a method to enhance SER performance by viewing Mel Frequency Cepstral Coefficients (MFCC), which accelerates the learning process while maintaining a high level of accuracy. The authors employ a supervised learning model, specifically a functional support vector machine (SVM), directly on the MFCCs represented as functional data. This enables the utilization of complete functional information, resulting in more precise emotion recognition. The authors' method demonstrates competitive results in terms of accuracy, underscoring its effectiveness in emotion recognition as well as reducing learning time, making it computationally efficient and practical for real-world applications.

A framework called Convolutional Auto-Encoder and Adversarial Domain Adaptation (CAEADA) for cross-corpus SER is introduced in [24]. The CAEADA framework starts by creating a one-dimensional convolutional auto-encoder (1D-CAE) for feature processing. This 1D-CAE is designed to capture correlations among adjacent one-dimensional statistical features, and the feature representation is enhanced through an encoder–decoder-style architecture. Following this, the adversarial domain adaptation (ADA) module works to reduce the differences in feature distributions between the source and target domains by confusing a domain discriminator. Specifically, it employs the maximum mean discrepancy (MMD) method to achieve effective feature transformation. The evaluation results demonstrate that the method performs quite satisfactorily on SER tasks.

In the work presented in [25], the authors present an attention-based dense long short-term memory (LSTM) approach for speech emotion recognition. The authors integrate LSTM networks, which are well-suited for handling time series data such as speech, with attention-based dense connections. This entails the incorporation of weight coefficients into skip connections for each layer, enabling the differentiation of emotional information across layers and preventing interference from redundant information in the lower layers with valuable information from upper layers. The experiments showcase an improvement in recognition performance by 12% and 7% on the eINTERFACE and IEMOCAP datasets, respectively.

References

1. Ajuzieogu, U. The Role of AI in Modern Computing and Education; Lulu Publisher: Morrisville, NC, USA, 2019; ISBN 978-0-359-72121-4.
2. Jalal, M.A.; Loweimi, E.; Moore, R.K.; Hain, T. Learning Temporal Clusters Using Capsule Routing for Speech Emo-tion Recognition. In Proceedings of the Interspeech 2019, Graz, Austria, 15–19 September 2019; pp. 1701–1705.
3. Mustaqeem; Kwon, S. A CNN-Assisted Enhanced Audio Signal Processing for Speech Emotion Recognition. *Sensors* 2020, 20, 183.
4. Singh, Y.B.; Goel, S. A systematic literature review of speech emotion recognition approaches. *Neurocomputing* 2022, 492, 245–263.
5. Wani, T.M.; Gunawan, T.S.; Qadri, S.A.A.; Kartiwi, M.; Ambikairajah, E. A Comprehensive Review of Speech Emotion Recognition Systems. *IEEE Access* 2021, 9, 47795–47814.
6. Yadav, S.P.; Zaidi, S.; Mishra, A.; Yadav, V. Survey on Machine Learning in Speech Emotion Recognition and Vision Systems Using a Recurrent Neural Network (RNN). *Arch. Comput. Methods Eng.* 2022, 29, 1753–1770.
7. Lieskovská, E.; Jakubec, M.; Jarina, R.; Chmulík, M. A review on speech emotion recognition using deep learning and attention mechanism. *Electronics* 2021, 10, 1163.

8. Khalil, R.A.; Jones, E.; Babar, M.I.; Jan, T.; Zafar, M.H.; Alhussain, T. Speech Emotion Recognition Using Deep Learning Techniques: A Review. *IEEE Access* 2019, **7**, 117327–117345.
9. Abbaschian, B.J.; Sierra-Sosa, D.; Elmaghraby, A. Deep Learning Techniques for Speech Emotion Recognition, from Databases to Models. *Sensors* 2021, **21**, 1249.
10. de Lope, J.; and Graña, M. An ongoing review of speech emotion recognition. *Neurocomputing* 2023, **528**, 1–11.
11. Wen, G.; Li, H.; Huang, J.; Li, D.; Xun, E. Random Deep Belief Networks for Recognizing Emotions from Speech Signals. *Comput. Intell. Neurosci.* 2017, **2017**, 1945630.
12. Badshah, A.M.; Ahmad, J.; Rahim, N.; Baik, S.W. Speech Emotion Recognition from Spectrograms with Deep Convolutional Neural Network. In Proceedings of the 2017 International Conference on Platform Technology and Service (PlatCon), Busan, Republic of Korea, 13–15 February 2017; pp. 1–5.
13. Zhao, J.; Mao, X.; Chen, L. Speech emotion recognition using deep 1D & 2D CNN LSTM networks. *Biomed. Signal Process. Control* 2019, **47**, 312–323.
14. Lee, C.; Song, K.Y.; Jeong, J.; Choi, W.Y. Convolutional Attention Networks for Multimodal Emotion Recognition from Speech and Text Data. *arXiv* 2019, arXiv:1805.06606.
15. Tang, D.; Zeng, J.; Li, M. An End-to-End Deep Learning Framework for Speech Emotion Recognition of Atypical Individuals. In Proceedings of the Interspeech 2018, Hyderabad, India, 2–6 September 2018; pp. 162–166.
16. Chen, M.; He, X.; Yang, J.; Zhang, H. 3-D Convolutional Recurrent Neural Networks With Attention Model for Speech Emotion Recognition. *IEEE Signal Process. Lett.* 2018, **25**, 1440–1444.
17. Li, P.; Song, Y.; McLoughlin, I.; Guo, W.; Dai, L. An Attention Pooling Based Representation Learning Method for Speech Emotion Recognition. In Proceedings of the INTERSPEECH 2018, Hyderabad, India, 2–6 September 2018.
18. Jiang, W.; Wang, Z.; Jin, J.S.; Han, X.; Li, C. Speech Emotion Recognition with Heterogeneous Feature Unification of Deep Neural Network. *Sensors* 2019, **19**, 2730.
19. Huang, K.; Wu, C.; Hong, Q.; Su, M.; Chen, Y. Speech Emotion Recognition Using Deep Neural Network Considering Verbal and Nonverbal Speech Sounds. In Proceedings of the ICASSP 2019—2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Brighton, UK, 12–17 May 2019; pp. 5866–5870.
20. Issa, D.; Demirci, M.; Yazici, A. Speech emotion recognition with deep convolutional neural networks. *Biomed. Signal Process. Control* 2020, **59**, 101894.

21. Makhmudov, F.; Kutlimuratov, A.; Akhmedov, F.; Abdallah, M.S.; Cho, Y.-I. Modeling Speech Emotion Recognition via Attention-Oriented Parallel CNN Encoders. *Electronics* **2022**, *11*, 4047.
22. Liu, Z.-T.; Han, M.-T.; Wu, B.-H.; Rehman, A. Speech emotion recognition based on convolutional neural network with attention-based bidirectional long short-term memory network and multi-task learning. *Appl. Acoust.* **2023**, *202*, 109178.
23. Saumard, M. Enhancing Speech Emotions Recognition Using Multivariate Functional Data Analysis. *Big Data Cogn. Comput.* **2023**, *7*, 146.
24. Wang, Y.; Fu, H.; Tao, H.; Yang, J.; Ge, H.; Xie, Y. Convolutional Auto-Encoder and Adversarial Domain Adaptation for Cross-Corpus Speech Emotion Recognition. *IEICE Trans. Inf. Syst.* **2022**, *105*, 1803–1806.
25. Xie, Y.; Liang, R.; Liang, Z.; Zhao, L. Attention-Based Dense LSTM for Speech Emotion Recognition. *IEICE Trans. Inf. Syst.* **2019**, *102*, 1426–1429.

Retrieved from <https://encyclopedia.pub/entry/history/show/122576>