Adversarial Machine Learning And Speech Emotion Recognition: Utilizing Generative Adversarial Networks For Robustness

Anonymous Author(s) Affiliation Address email

Abstract

Although deep learning has enabled unprecedented improvements in the perfor-1 2 mance of the state-of-the-art speech emotion recognition (SER) systems, recent research on adversarial examples has cast a shadow of doubt on the robustness of 3 SER systems by showing the susceptibility of deep neural networks to adversarial 4 examples that rely only on small and imperceptible perturbations. In this study, 5 we evaluate how adversarial examples can be used to attack SER systems and 6 propose the first black-box adversarial attack on SER systems. We also explore 7 potential defenses including adversarial training and generative adversarial network 8 (GAN) to enhance robustness. Experimental evaluations suggest various interesting 9 aspects of the effective utilization of adversarial examples that can be useful not 10 only for SER robustness but also other speech-based intelligent systems. 11

12 **1** Introduction

Recent progress in machine learning (ML) is reinventing the future of intelligent systems enabling plethora of speech controlled applications [1, 2]. In particular, the emotion-aware systems are on the rise. And the breakthrough in deep learning is the enabler of highly accurate and robust emotion recognition systems [3, 4].

Despite the superior performance of deep neural networks (DNNs), recent studies have shown that 17 they are highly vulnerable to the malicious attacks that use *adversarial examples*. Adversarial 18 examples are custom built by a malicious adversary through the addition of unperceived perturbation 19 with the intention of eliciting wrong responses from ML models. These adversarial examples can 20 debilitate the performance of image recognition, object detection, and speech recognition models 21 [5]. Adversarial attacks can also be used to undermine the performance of speech-based emotion 22 recognition (SER) systems [6], which is alarming due to various security-sensitive paralinguistic 23 applications of SER systems. 24

In this paper, we aim to investigate the utility of adversarial examples to achieve robustness in speech emotion classification to adversarial attacks. We consider a "black-box" attack that directly perturbs speech utterances with small and imperceptible noises. The generated adversarial examples are utilized in different schemes highlighting different trends for the robustness of SER systems. We further propose a GAN-based defense for SER systems and show that it can better withstand adversarial examples compared to the previous defense solutions such as adversarial training and random noise addition.

32 2 Background and Audio Adversarial Examples

33 Existing methods of adversarial attacks include fast gradient sign method (FGSM) [7], Jacobianbased saliency map attack (JSMA) [8], and DeepFool [9]. In another work, Carlini and Wagner 34 [10] compute the perturbation noise based on the gradient of targeted output with respect to the 35 input, which can be computed efficiently using backpropagation with the implicit assumption that 36 the attacker has complete knowledge of the network and its parameters (such methods are called 37 *white-box* attacks). While the backpropagation method, which needs to compute the derivative of each 38 layer of the network with respect to the input layers, can be efficiently applied in image recognition 39 40 due to the differentiability of all layers, the application of such methods is difficult for SER systems 41 since these systems rely on complex acoustic features of the input audio utterances—such as Mel Frequency Cepstral Coefficients (MFCCs), spectrogram, extended Geneva Minimalistic Acoustic 42 Parameter Set (eGeMAPS) [11]. The SER system's first layer therefore is the pre-processing or the 43 feature extraction layer through which there is no efficient way to compute derivative, due to which 44 gradient-based methods [8, 9, 10, 12] are not directly applicable to SER systems. 45

46 **2.1 Previous Audio Attacks**

Adversarial attacks on ML have triggered an active area of research that is focusing on understanding 47 the adversarial attack phenomenon [13] and on techniques that can make ML models robust [14]. 48 49 For speech-based systems, Carlini [5] proposed a white-box iterative optimization-based attack for DeepSpeech [15], a state-of-the-art speech-to-text model, with 100% success rate. Alzantot et al. 50 [16] proposed an adversarial attack on speech commands classification model by adding a small 51 random noise (background noise) to the audio files. They achieved 87% success without having any 52 information of the underlying model. Song et al. [17] proposed a mechanism that directly attacks the 53 microphone used for sensing voice data and showed that an adversary can exploit the microphone's 54 55 non-linearity to control the targeted device with inaudible voice commands. Gong et al. [6] presented 56 an architecture to craft adversarial examples for computational paralinguistic applications. They perturbed the raw audio file and were able to cause a significant reduction in performance. Various 57 other studies [18, 19, 20] have also presented adversarial attacks for speech recognition system. 58

Most of the previous research on targeted attacks for speech-based applications [5, 6, 16, 17, 18, 19, 59 20] has considered attacks on the model without investigating how adversarial examples may be 60 61 utilized to make the ML models more robust. Our work is different since we not only propose an 62 adversarial attack for SER system using adversarial examples but also leverage adversarial examples for making ML models more robust. We evaluate our proposed attack on two well-known emotional 63 corpora (IEMOCAP [21] and FAU-AIBO [22]) using Long Short-Term Memory (LSTM) [23], 64 a popular recurrent neural network (RNN), as the classifier. We achieved 79% success rate for 65 FAU-AIBO dataset without changing the perception of human emotion captured in an audio file. 66

3 Proposed Audio Adversarial Examples

In this work, we adopt a simple approach to craft adversarial examples by adding imperceptible 68 noise to the legitimate samples. For this, we take an audio utterance x with label y, and generate an 69 adversarial example $x' = x + \delta$ such that the SER system fails to correctly classify the given input 70 while ensuring that x and x' are very similar as perceived by humans. Previous speech-related studies 71 72 have studied different noise as adversarial noises. DolphinAttack exploits inaudible ultrasounds as 73 adversarial noise to control the victim device inconspicuously but the attack sound used was out of 74 the human perception. Similarly, Alzantot et al. [16] used random noise for creating an adversarial attack on speech recognition. It has however been empirically observed that the state-of-the-art 75 classifiers are relatively robust to random noise [13]. By considering these observations, we propose 76 a black-box attack for SER system where an adversary can add some real-world noise as adversarial 77 perturbation. We empirically show that the addition of real-world noisy speech samples can fool the 78 classier while not being perceptible to the human ear. 79

Quantification of δ : The quantification distortion caused by δ is performed through a simple rule. SER systems are designed to detect speakers' emotion independently from the background noise. In real-world scenarios, such background noise can take multiple forms such as car engine noise, passing-by-vehicle noise, or result from human discussion, music, etc. SER systems, or indeed any

Table 1: Binary class mapping of different emotions

Dataset	Positive class	Negative Class
IEMOCAP FAU-AIBO	happiness, exited, neutral neutral, motherese, and joyful	anger, sadness angry, touchy, reprimanding, and emphatic

intelligent speech-based system, have to be robust enough to tackle these sources of noise for their 84 deployment in real-world. Our aim in this study is to design adversarial example by using some of 85 these noises and their addition is done based on human perception. For this, we use three noises of 86 real-world and their addition level is based on the already existing background noises (microphone 87 noise and discussion noise) in the utterances. We estimate the existing noise in utterances using a 88 well-known technique proposed in [24] that estimate noise using spectral and log-amplitude. The 89 detected noise (N_{ex}) is used a reference for quantification of δ . For δ , we use three noise (N_{add}) (café, 90 meeting, and station) from the Demand Noise database [25] and make their the mean and variance 91 equal to the reference noise (N_{ex} as existing noise). We also use ϵ as the variation parameter to 92 further control the perturbation amplitude. In this way, the adversarial noise (N_{add}) has a very small 93 value similar to the existing noise and the adversarial example and is unrecognizable to the human 94 ear in the human perception test. This noise (N_{add}) added to utterances is multiplied with different 95 perturbation amplitude (ϵ) to generate the adversarial examples. This noise acts as the background 96 noise and does not change the emotional context of a given audio file. 97

Human Perception and Classifier Test: In order to assess the effect of added adversarial noise on the human listener, we asked five human listeners to listen to 200 adversarial examples for different perturbation amplitude (ϵ) and differentiate it from the original audio file. For the IEMOCAP and FAU-AIBO datasets, 96% and 91% of the samples were indistinguishable from the original utterances. When these examples were given to the classifier, the attack success rate was 72% and 79% for IEMOCAP and FAU-AIBO, respectively.

104 4 Experimental Setup and Results

We evaluated the generated adversarial examples using two well-known emotional corpora: IEMO-CAP and FAU-AIBO. We consider binary classification problem for both these datasets as used in [3] and [26]. Table 1 shows the considered emotions and their binary class mapping. We use the eGeMAPS features, a popular features set specifically suited for paralinguistic applications, for representing the audio samples.

Classification Model: We consider LSTM-RNN for emotion classification. LSTM is a popular RNN and widely employed in audio [27] and emotion classification [4] due to their ability to model contextual information. We find the best model structure by evaluating different number of layers. We obtained the best results with two LSTM layers, one dense layer, and softmax as the last layer. We initially used a learning rate of 0.002 to start training the model and halved this rate after every 5 epochs if performance did not improve on the test set. This process stopped when the learning rate reached below 0.00001.

Emotion Classification Results: For experimentation, we evaluated the model in a speaker inde-117 pendent scheme. IEMOCAP dataset consists of five sessions: we used four session for training and 118 one for testing, consistent with the methodology of previous works [3, 4]. In the case of FAU-AIBO, 119 we followed the speaker-independent training strategy proposed in the 2009 Interspeech Emotion 120 Challenge [26]. For emotion classification on legitimate examples, we achieved 68.35% and 56.41%121 unweighted accuracy (UA) on FAU-AIBO and IEMOCAP dataset, respectively. The results on adver-122 sarial examples are compared with these results. We generated adversarial examples with different 123 values of ϵ (0.1–2) to evaluate the performance of model with different perturbation amplitude. Figure 124 1 presents the emotion classification error on adversarial samples with different values of ϵ . 125



Figure 1: The error rate (%) with different perturbation factors for speech emotion classification for FAU-AIBO (left) and IEMOCAP (right) datasets.

126 **5** Possible Defenses

127 5.1 Training with Adversarial Examples

Adversarial training of model is considered as a possible defense to adversarial attacks when the 128 exact nature of the attack is known. Model training on the mixture of clean and adversarial examples 129 can somewhat help regularization [28]. Training on adversarial samples is different from data 130 augmentation methods that are performed based on the expected translations in test data. To the 131 best of our knowledge, adversarial training is not explored for SER systems and other speech/audio 132 classification systems. We explore this phenomenon by mixing adversarial examples with training 133 data to highlight the robustness of model against attack. We trained the model with training data 134 comprising a varying percentage of adversarial examples (10% to 100% of training data). Figure 2 135 shows the classification error rate (%) that is significantly decreased with increasing the percentage 136 of adversarial examples in training data. However, the classification error is higher (5% to 15% in 137 different scenarios) compared to the classification when the model is trained on clean utterances. 138



Figure 2: The error rate (%) with varying the percentage of adversarial samples as training data for FAU-AIBO (left) and IEMOCAP (right) datasets.

139 5.2 Training with Random Noise

Training of models by adding random noise in training data might help against adversarial attacks on the speech-based system. It is reported in [29] that the addition of a random noise layer to the neural network can prevent strong gradient-based attacks in the image domain. We evaluated this phenomenon in speech emotion classification system by adding a small random noise to overall training data and evaluated the performance against the proposed attacks. This can be noted from Table 2 that emotion classification error reduced only slightly with the addition of random noise in training data, which indicates that this strategy is not particularly effective in the SER settings.

Dataset	Adversarial Perturbations	Error (max) with adversarial attack	Error by training with random noise
FAU-AIBO	Café	56.87	54.02
	Meeting	52.58	49.24
	Station	53.57	48.51
IEMOCAP	Café	64.58	56.73
	Meeting	63.88	52.57
	Station	66.87	60.87

Table 2: Emotion classification error (%) by adding random noise in training data

147 5.3 Using Generative Adversarial Network

Generative adversarial networks (GANs) [30] are deep models that learn to generate samples, ideally indistinguishable from the real data x, that are supposed to belong to an unknown data distribution, $p_{data}(x)$. GANs consist of two networks, a generator (G) and a discriminator (D). The generator network (G) maps latent vectors from some known prior p_z to samples and discriminator tasked to differentiate between the real sample x or fake G(y). Mathematically, this is represented by the following optimization program:

$$\min_{G} \max_{D} \quad \mathbf{E}_{x}[\log(D(x))] + \mathbf{E}_{y}[\log(1 - D(G(y)))] \tag{1}$$

where G and D play this game to fool each other using this min-max optimization program. GANs 154 have already been applied for speech enhancement [31, 32] in speech recognition systems, therefore, 155 we are using GANs as a defense strategy against adversarial noises. In our case, G network is 156 tasked to remove the adversarial noise from the adversarial examples y. The G network is structured 157 like an autoencoder using LSTM layers. The LSTM based encoder-decoder architecture is well 158 suited for capturing emotions due to the demonstrated superior performance of LSTM in capturing 159 long-range context such as those present in emotions [4]. In the G network, the encoder part 160 compresses the contextual (emotional) information of the input speech features and the decoder uses 161 this representation for reconstruction. The D network follows the same encoder-decoder architecture. 162 For training G and D for different possible scenarios, we used the training data from both the datasets 163 to train the GAN. For each G step, the discriminator was updated twice. For faster convergence, we 164 pretrained the G network in each case. We train GAN using RMSProp optimizer with learning rate 165 1×10^{-4} and batch size of 32, until convergence. 166

We trained the GAN using utterances corrupted by the three adversarial noises (café, meeting, station) as noisy data and it was tasked to clean the utterances. Data cleaned by GAN (G(y)) is given to the classifier for emotion classification. Table 3 shows emotion classification results on audio utterances cleaned by GAN.

Dataset	Adversarial Perturbations	Error (max) with adversarial attack	Error by using GAN
FAU-AIBO	Café	68.82	38.31
	Meeting	62.58	36.02
	Station	66.87	35.14
IEMOCAP	Café	65.87	49.20
	Meeting	67.70	48.18
	Station	69.87	46.24

Table 3: Emotion classification error (%) by utilizing GAN for adversarial noise removal

171 6 Discussion

From the experimental evaluation for SER robustness, we have discovered that GAN-based defense against adversarial audio examples is able to better withstand adversarial examples compared to other approaches. Figure 3 shows a comparison of the different defenses using two well-known datasets: the addition of random noise in training utterances is able to slightly reduce speech emotion classification error while with adversarial training, classification error is significantly reduced; however, the best results are shown by using GAN for cleaning the utterances and then running classification on the clean utterances.

Our results highlight the power of GAN for speech emotion classification in the face of adversarial examples, which motivate further research for its utilization in other speech-based intelligent systems



Figure 3: The error rate (%) with three different approaches against adversarial examples for FAU-AIBO (left) and IEMOCAP (right) datasets.

for the minimization of adversarial perturbations. It is worth pointing out that GANs require information about the exact type and nature of adversarial examples for its training, but this is also an essential requirement for the adversarial training mechanism.

184 7 Conclusions

In this paper, we propose a black-box method to generate adversarial perturbations in audio examples 185 of speech emotion recognition system (SER) and also propose defense strategies. In particular, we 186 propose a Generative Adversarial Network (GAN)-based mechanism for enhancing the robustness of 187 SER system by first cleaning the perturbed utterances through GANs and then running a classifier on 188 it. We compared our GAN-based defense against adversarial training and the addition of random 189 noise in training examples and showed that our GAN-based defense provides consistently better 190 results in speech emotion recognition. The attack and defense that we propose can also be utilized 191 more generally for other speech-based intelligent systems. 192

193 References

- [1] Erik Cambria. Affective computing and sentiment analysis. *IEEE Intelligent Systems*, 31(2):102–107, 2016.
- [2] Soujanya Poria, Erik Cambria, Amir Hussain, and Guang-Bin Huang. Towards an intelligent
 framework for multimodal affective data analysis. *Neural Networks*, 63:104–116, 2015.
- [3] Siddique Latif, Rajib Rana, Shahzad Younis, Junaid Qadir, and Julien Epps. Transfer learning
 for improving speech emotion classification accuracy. *Proc. Interspeech 2018*, pages 257–261,
 2018.
- [4] Siddique Latif, Rajib Rana, Junaid Qadir, and Julien Epps. Variational autoencoders for learning
 latent representations of speech emotion: A preliminary study. In *Proc. Interspeech 2018*, pages
 3107–3111, 2018.
- [5] Nicholas Carlini and David Wagner. Audio adversarial examples: Targeted attacks on speechto-text. arXiv preprint arXiv:1801.01944, 2018.
- [6] Yuan Gong and Christian Poellabauer. Crafting adversarial examples for speech paralinguistics
 applications. *arXiv preprint arXiv:1711.03280*, 2017.
- [7] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples (2014). *arXiv preprint arXiv:1412.6572*.
- [8] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and
 Ananthram Swami. The limitations of deep learning in adversarial settings. In *Security and Privacy (EuroS&P), 2016 IEEE European Symposium on*, pages 372–387. IEEE, 2016.

- [9] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple
 and accurate method to fool deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2574–2582, 2016.
- [10] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In
 2017 IEEE Symposium on Security and Privacy (SP), pages 39–57. IEEE, 2017.
- [11] Florian Eyben, Klaus R Scherer, Björn W Schuller, Johan Sundberg, Elisabeth André, Carlos
 Busso, Laurence Y Devillers, Julien Epps, Petri Laukka, Shrikanth S Narayanan, et al. The
 geneva minimalistic acoustic parameter set (gemaps) for voice research and affective computing.
 IEEE Transactions on Affective Computing, 7(2):190–202, 2016.
- [12] Jiawei Su, Danilo Vasconcellos Vargas, and Sakurai Kouichi. One pixel attack for fooling deep
 neural networks. *arXiv preprint arXiv:1710.08864*, 2017.
- [13] Alhussein Fawzi, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. Robustness of
 classifiers: from adversarial to random noise. In *Advances in Neural Information Processing Systems*, pages 1632–1640, 2016.
- [14] Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier. Parse val networks: Improving robustness to adversarial examples. *arXiv preprint arXiv:1704.08847*, 2017.
- [15] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan
 Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, et al. Deep speech: Scaling up
 end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*, 2014.
- [16] Moustafa Alzantot, Bharathan Balaji, and Mani Srivastava. Did you hear that? adversarial
 examples against automatic speech recognition. *arXiv preprint arXiv:1801.00554*, 2018.
- [17] Liwei Song and Prateek Mittal. Inaudible voice commands. *arXiv preprint arXiv:1708.07238*, 2017.
- [18] Nirupam Roy, Haitham Hassanieh, and Romit Roy Choudhury. Backdoor: Making microphones
 hear inaudible sounds. In *Proceedings of the 15th Annual International Conference on Mobile Systems, Applications, and Services*, pages 2–14. ACM, 2017.
- [19] Dan Iter, Jade Huang, and Mike Jermann. Generating adversarial examples for speech recognition. http://web.stanford.edu/class/cs224s/reports/Dan_Iter.pdf, 2017.
- [20] Lea Schönherr, Katharina Kohls, Steffen Zeiler, Thorsten Holz, and Dorothea Kolossa. Adversarial attacks against automatic speech recognition systems via psychoacoustic hiding. *arXiv preprint arXiv:1808.05665*, 2018.
- [21] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim,
 Jeannette N Chang, Sungbok Lee, and Shrikanth S Narayanan. Iemocap: Interactive emotional
 dyadic motion capture database. *Language resources and evaluation*, 42(4):335, 2008.
- [22] Stefan Steidl. Automatic classification of emotion related user states in spontaneous children's
 speech. University of Erlangen-Nuremberg Erlangen, Germany, 2009.
- [23] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Yariv Ephraim and David Malah. Speech enhancement using a minimum-mean square error
 short-time spectral amplitude estimator. *IEEE Transactions on acoustics, speech, and signal processing*, 32(6):1109–1121, 1984.
- [25] Joachim Thiemann, Nobutaka Ito, and Emmanuel Vincent. The diverse environments multi channel acoustic noise database: A database of multichannel environmental noise recordings.
 The Journal of the Acoustical Society of America, 133(5):3591–3591, 2013.
- [26] Björn Schuller, Stefan Steidl, and Anton Batliner. The interspeech 2009 emotion challenge. In
 Tenth Annual Conference of the International Speech Communication Association, 2009.

- [27] Siddique Latif, Muhammad Usman, Rajib Rana, and Junaid Qadir. Phonocardiographic sensing
 using deep learning for abnormal heartbeat detection. *IEEE Sensors Journal*, 2018.
- [28] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfel low, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [29] Xuanqing Liu, Minhao Cheng, Huan Zhang, and Cho-Jui Hsieh. Towards robust neural networks
 via random self-ensemble. *arXiv preprint arXiv:1712.00673*, 2017.
- [30] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [31] Ashutosh Pandey and Deliang Wang. On adversarial training and loss functions for speech en hancement. In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing
 (ICASSP), pages 5414–5418. IEEE, 2018.
- [32] Santiago Pascual, Antonio Bonafonte, and Joan Serra. Segan: Speech enhancement generative adversarial network. *arXiv preprint arXiv:1703.09452*, 2017.