

① Motivation: 应该说明 (1) 为什么要基于 spectrogram 进行 SER, (2) SER 的 discriminative features 在哪里, 从而引出文章主题 'Learning Discriminative Features'.  
 ② Challenge: 应该说明 (1) 如何设计 model architecture 以符合 spectrogram 作为输入; (2) 如何设计 model 以得到更有区分性的 features.  
 ③ Contribution: (3.1) 通过 CNN+Bi-LSTM 网络提取更丰富的 spectrogram 信息; (3.2) 通过引入 center loss 学习区分性特征.



黄健 (提取特征)

# Learning Discriminative Features from Spectrograms using Center Loss for Speech Emotion Recognition

Dongyang Dai<sup>1,2</sup>, Zhiyong Wu<sup>1,2,3</sup>, Runnan Li<sup>1,2</sup>, Xixin Wu<sup>3</sup>, Jia Jia<sup>1,2</sup>, Helen Meng<sup>1,3</sup>

<sup>1</sup>Tsinghua-CUHK Joint Research Center for Media Sciences, Technologies and Systems, Graduate School at Shenzhen, Tsinghua University, Shenzhen, China

<sup>2</sup>Tsinghua National Laboratory for Information Science and Technology (TNList), Department of Computer Science and Technology, Tsinghua University, Beijing, China

<sup>3</sup>Department of Systems Engineering and Engineering Management, The Chinese University of Hong Kong



## 1. Introduction

### Motivation

- Extract valid features from raw data for emotion recognition



(强调 spectrogram; 强调 Discriminative Features)

### Challenge

- How to design a suitable model processing directly on raw data
- Emotions are naturally ambiguous

### Contribution

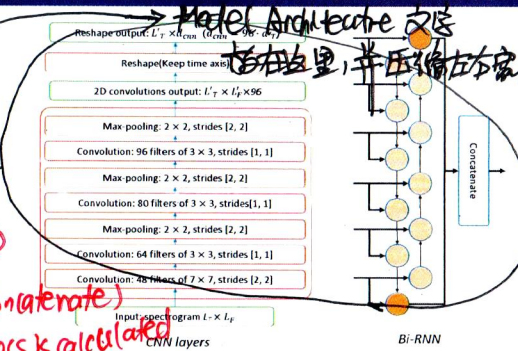
- Extract features and identify emotions directly from spectrograms
- Introduce center loss together with softmax cross-entropy loss in SER task to learn discriminative features
- Separable inter-class features
- More compact intra-class features

第三级用 'V' 号, 并与其地方年份大小一致

### Model Architecture

- Input variable length spectrograms (STFT or Mel-spectrograms)
- CNN layers extract spatial information
- Bi-RNN compresses the variable length sequence down to a fixed-length vector (强调 Bi-RNN 输出如何的 (iterate))
- FC1 output  $z \in R^d$  as the learned feature and calculate center loss according to  $z$ , from which center loss is calculated
- FC2 outputs posterior class probabilities, used to calculate softmax cross-entropy loss, from which softmax loss is computed
- Softmax Cross-entropy Loss enables the network to learn separable features
- Center Loss pulls the features belonging to the same emotion category to their center

## 2. Proposed Method



### Softmax Cross-entropy Loss

- $L_s = -\sum_{i=1}^m \omega_{y_i} \log \left( \frac{e^{W_{y_i}^T z_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T z_i + b_j}} \right)$
- $\omega_j$ : in inverse proportion to the sample number of the  $j$ -th class in training set

### Center Loss

$$L_c = \frac{1}{\sum_{i=1}^m \omega_{y_i}} \sum_{i=1}^m \omega_{y_i} \|z_i - c_{y_i}\|^2$$

$$c_j^{t+1} = \begin{cases} (1-\alpha)c_j^t + \alpha \bar{c}_j^t & \sum_{i=1}^m \delta(y_i = j) > 0 \\ c_j^t & \sum_{i=1}^m \delta(y_i = j) = 0 \end{cases}$$

$$\bar{c}_j^t = \frac{\sum_{i=1}^m \delta(y_i = j) z_i}{\sum_{i=1}^m \delta(y_i = j)}$$

① 最近介绍 center loss, 所以放在第一

② cross entropy 大家都已经清楚, 放在第二, 不用详细介绍

③ 最后说明两者结合的 joint loss

### Joint Loss

- $L = L_s + \lambda L_c$
- $\lambda$ : trades off center loss against softmax cross-entropy loss.

## 3. Experiments and Results

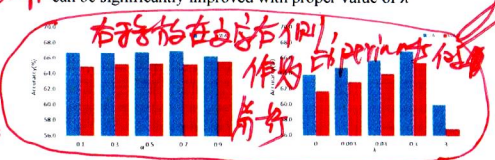
### Experimental Setup

- Data set: IEMOCAP
- Dataset: IEMOCAP
- Neutral, angry, happy, sad and excited (merged happy and excited as happy, 5531 utterances)
- 5 subsets (keep the emotion distribution); 4 subsets for training, half of the last subset as development set and half as test set
- Settings of spectrograms
- Model input: log scale STFT spectrogram or Mel-spectrogram
- Hamming window: 40ms window length and 10ms shift
- Window size: 40msec
- Window Shift: 10msec
- Sample rate: 16KHz
- DTF length: 1024
- The number of Mel bands: 128 for Mel-spectrogram
- Metrics Evaluation Metrics
- The Unweighted Accuracy (UA): the mean value of the recall for each class
- The Weighted Accuracy (WA): the number of correctly classified samples divided by the total amount of samples

5 subsets randomly divided 5531 utterances, but keeping the distribution portion of emotion categories  
 4 subsets for training, half and half ...

### Experiments

- The effect of hyperparameter  $\alpha$  and  $\lambda$  on Mel-spectrogram
- (left) fixing  $\lambda = 0.3$ , (right) fixing  $\alpha = 0.5$
- not sensitive to  $\alpha$
- can be significantly improved with proper value of  $\lambda$



### Experiments with different $\lambda$ on Mel and STFT

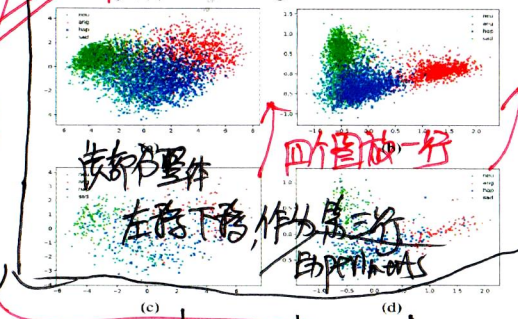
	Setting1	Setting2	Setting3	Setting4
$\lambda, \alpha$	$\lambda=0$	$\lambda=0.3, \alpha=0.5$	$\lambda=0$	$\lambda=0.3, \alpha=0.5$
Input	Mel	Mel	STFT	STFT
The UA and WA on setting 1 ~ setting 4				
	Setting1	Setting2	Setting3	Setting4
UA(%)	63.80	66.86	60.97	65.13
WA(%)	61.83	65.40	58.93	62.96

Confusion matrix on setting 1/setting 2/setting 3/setting 4(%)



### PCA embedding of feature

- (a) training set on setting 1, (b) training set on setting 2, (c) test set on setting 1, (d) test set on setting 2



## 4. Conclusion

### Conclusion

- Introducing center loss with proper  $\lambda$  could effectively improve the SER performance on both STFT spectrogram and Mel-spectrogram input
- Mel-spectrogram input, reducing the dimension based on human hearing characteristics, over-performs STFT spectrogram input
- The 2-D PCA embedding illustrates the discriminative power when using center loss, which enables the neural network to learn more effective features for SER

cf

## 5. Acknowledgment

- This work is supported by National Natural Science Foundation of China (NSFC) (61433018, 61375027), joint research fund of NSFC-RGC (Research Grant Council of Hong Kong) (61531166002, NCUHK404/15) and National Social Science Foundation of China (13&ZD189)

缩减长度  
 Experiments 部分由三行呈现结果