



# A Generalized Subspace Distribution Adaptation Framework for Cross-Corpus Speech Emotion Recognition

Shaokai Li<sup>1</sup> Peng Song<sup>1\*</sup> Liang Ji<sup>1,4</sup> Yun Jin<sup>2</sup> Wenming Zheng<sup>3</sup>

<sup>1</sup>Yantai University    <sup>2</sup>Jiangsu Normal University

<sup>3</sup>The Key Laboratory of Child Development and Learning Science (Southeast University)

<sup>4</sup>The State Key Laboratory of Tibetan Intelligent Information Processing and Application

2023.04

# CONTENTS

## **01** Background

---

## **02** The challenging problem

---

## **03** The proposed method

---

## **04** Experiments

---

– 01 –

## **Background**

# 01 Background

---

**Speech Emotion Recognition (SER)** is an important research direction in affective computing, pattern recognition, signal processing and human-machine interaction. The goal of SER is to identify the emotion categories from speech signals, such as fear, anger, sadness, pleasure, and so on.



Driving assist system



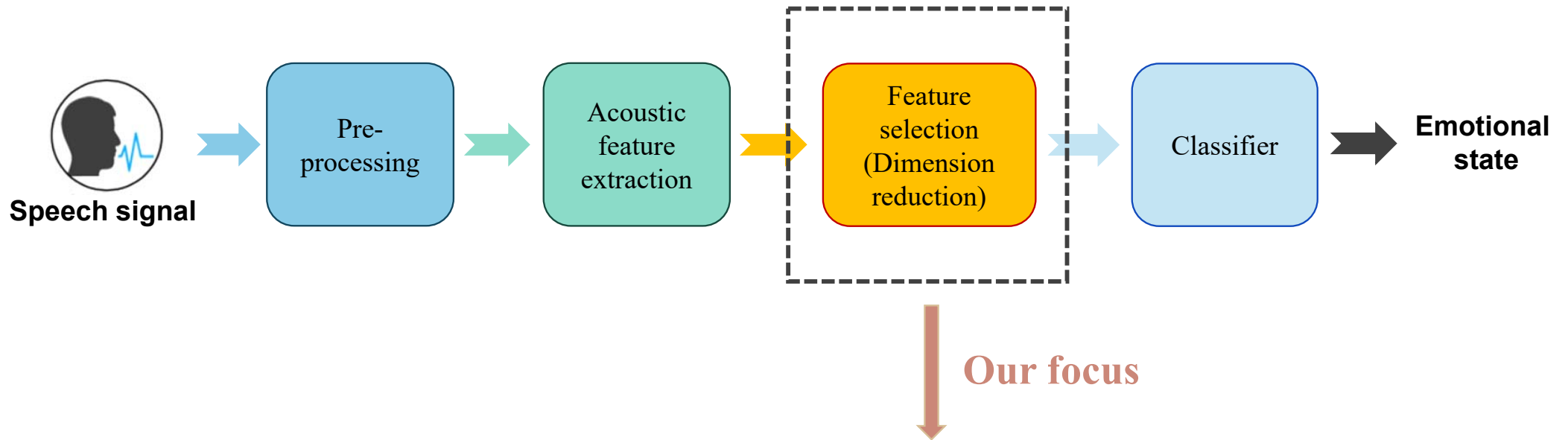
Automatic translation



Education

# 01 The process of SER

---



Learning a transfer subspace, which can obtain a common subspace to reduce the discrepancy between databases.

# 01 Traditional SER method

---

Many classification algorithms have been employed for SER, including:

- Hidden Markov model (HMM)
- Gaussian mixture model (GMM)
- Support vector machine (SVM)
- Neural network (NN)
- Deep neural network (DNN)
- Sparse representation
- Regression algorithms

– 02 –

## **The challenging problem**

## 02 The challenging problem

---

### The challenging cross-database SER problem

- **Data distribution mismatch problem:** in practical application scenarios, the speaker's gender, language, age and so on are different.
- **Insufficient labels problem:** labeling speech emotion is time-consuming, laborious, and require a large number of professionals.



## 02 Transfer learning

---

**Transfer learning:** The idea of transfer learning is to transfer the knowledge gained from one domain (source domain) to learn the knowledge of related but different domain (target domain).



We take the labeled database as the source domain and the unlabeled database as the target domain. The transfer learning can be used to solve the cross-database SER problem.

## 02 The related works

---

### Transfer learning for cross-database:

- Joint Distribution Adaptation (JDA) 2013
- Domain-adaptive least-squares regression (DaLSR) 2016
- Joint Geometrical and Statistical Alignment (JGSA) 2017
- Locality Preserving Joint Transfer (LPJT) 2019
- Transfer Sparse Discriminant Subspace Learning (TSDSL) 2019
- Deep Adaptation Networks (DAN) 2015
- Deep Subdomain Adaptation Networks (DSAN) 2021
- Deep Adaptation Regression (DAR) 2021
- Dual-level Adaptive and Discriminative (DLAD) 2022

– 03 –

**The proposed method**

## 03 The proposed method

---

### Problem Formulation:

We aim to learn a common projection subspace  $P$  by aligning the source and target distributions, where the corpus discrepancy would be well reduced. The objective function of the proposed GSDA can be formulated as follows:

$$\begin{aligned} \min_P \mathcal{F}(P, X) + \mathcal{G}(P, X) + \gamma \mathcal{S}(P) \\ \text{s.t. } P^T P = I \end{aligned}$$

The first item  $\mathcal{F}(P, X)$  is a generalized subspace learning method, in which the original feature space is projected into a low-dimensional common subspace. The second item  $\mathcal{G}(P, X)$  is the distance metric learning strategy. The third item  $\mathcal{S}(P)$  is a sparse regularization term.

## 03 The proposed method

---

**Distance Metric:**

$$\mathcal{G}(P, X) = \alpha \|V \odot S\|_1 - \beta \|V \odot D\|_1$$

**Examples of GSDA:**

1) GSDA-PCA:

$$\begin{aligned} \min_P & -\text{Tr}(P^T X X^T P) + \alpha \|V \odot S\|_1 - \beta \|V \odot D\|_1 + \gamma \|P\|_{2,1} \\ \text{s.t.} & P^T P = I \end{aligned}$$

2) GSDA-LDA:

$$\begin{aligned} \min_P & \text{Tr}(P^T (S_w - \mu S_b) P) + \alpha \|V \odot S\|_1 - \beta \|V \odot D\|_1 + \gamma \|P\|_{2,1} \\ \text{s.t.} & P^T P = I \end{aligned}$$

– 04 –

## **Experiments**

## 04 Experimental setup

---

- **Databases:** Berlin (B) , IEMOCAP (I), and CVE (C).

We select five common emotional categories, i.e., anger (AN), neutral (NE), happiness (HA), and sadness (SA), in our experiments.

- **Feature Extraction:**

**Low-level feature:** we use the openSMILE toolkit to extract the feature set of the INTERSPEECH 2010 paralinguistic challenge (1582-dimensional).

**Deep feature:** we extract the Mel spectrograms to learn 2,048-dimensional deep features by ResNet50. Specifically, given a cross-corpus task, we fine-tune a pre-trained ResNet-50 model on the source corpus, and extract 2048-dimensional deep features of the target corpus using the fine-tuned mode.

- **Classifier:** linear SVM.
- **Evaluation metric:** the weighted average recall (WAR).

# 04 Experimental results

## Results for Low-level Feature and Deep Feature

Tasks	Traditional methods		Transfer learning methods						GSDA -PCA	GSDA -LDA
	PCA	LDA	TCA	JDA	DaLSR	JGSA	LPJT	TSDSL		
B→I	44.21	40.11	46.54	46.06	49.19	45.24	45.96	49.25	<b>50.07</b>	<u>50.35</u>
B→C	45.74	39.35	49.83	48.16	48.22	49.67	46.87	<u>52.12</u>	<b>52.74</b>	51.32
I→B	30.31	40.62	55.85	53.66	49.37	59.70	<b>60.20</b>	<u>59.79</u>	57.10	59.18
I→C	35.16	30.32	43.38	47.12	51.61	46.32	<u>53.21</u>	51.19	<b>53.41</b>	48.38
C→B	56.54	56.33	59.81	62.62	49.47	58.14	59.18	63.70	<u>63.95</u>	<b>66.22</b>
C→I	44.62	32.39	47.13	48.11	49.94	47.78	46.69	46.51	<b>50.65</b>	<u>50.45</u>
Average	42.76	39.85	50.42	50.95	49.63	51.14	52.01	53.76	<b>54.65</b>	<u>54.31</u>

## Results for Deep Feature

Tasks	Traditional methods		Transfer learning methods									GSDA -PCA	GSDA -LDA
	PCA	LDA	TCA	JDA	DaLSR	JGSA	LPJT	TSDSL	DAN*	DSAN*	DAR*		
B→I	40.19	37.74	43.61	44.28	43.03	43.74	44.09	44.87	46.96	<b>47.63</b>	<u>47.23</u>	44.57	44.94
B→C	41.93	42.87	50.96	49.67	45.16	50.61	53.54	<u>55.06</u>	50.14	46.06	54.93	54.19	<b>55.80</b>
I→B	42.70	43.75	59.37	59.53	59.67	63.20	<u>66.66</u>	65.62	55.70	62.96	62.50	<b>66.75</b>	65.54
I→C	32.25	26.45	45.80	48.38	49.03	46.80	48.18	<u>49.16</u>	45.74	47.02	46.42	49.09	<b>49.23</b>
C→B	61.35	59.37	64.58	65.62	62.38	63.75	<b>69.79</b>	<u>67.71</u>	62.16	63.58	62.50	<b>69.79</b>	67.70
C→I	35.06	32.42	44.42	43.02	39.37	43.98	43.11	44.02	43.51	40.77	42.19	<b>45.39</b>	<u>44.56</u>
Average	42.24	40.43	51.45	51.75	49.77	52.01	54.22	54.40	50.70	51.33	52.62	<b>54.79</b>	<u>54.62</u>



# 04 Confusion Matrices

AN	0.75	0.01	0.22	0.02
HA	0.32	0.02	0.48	0.18
NE	0.24	0.01	0.59	0.17
SA	0.08	0.01	0.36	0.55
	AN	HA	NE	SA

(a) B→I

AN	0.83	0.15	0.02	0.00
HA	0.65	0.24	0.08	0.03
NE	0.05	0.03	0.84	0.08
SA	0.18	0.08	0.38	0.36
	AN	HA	NE	SA

(b) B→C

AN	0.95	0.00	0.05	0.00
HA	0.95	0.00	0.05	0.00
NE	0.30	0.00	0.70	0.00
SA	0.00	0.07	0.27	0.67
	AN	HA	NE	SA

(c) I→B

AN	0.85	0.05	0.10	0.00
HA	0.57	0.11	0.22	0.11
NE	0.08	0.00	0.79	0.13
SA	0.03	0.03	0.41	0.54
	AN	HA	NE	SA

(d) I→C

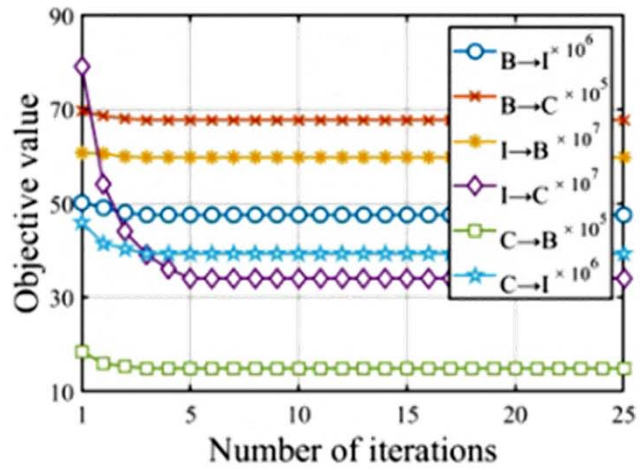
AN	0.87	0.00	0.05	0.08
HA	0.74	0.16	0.05	0.05
NE	0.13	0.00	0.87	0.00
SA	0.07	0.00	0.07	0.87
	AN	HA	NE	SA

(e) C→B

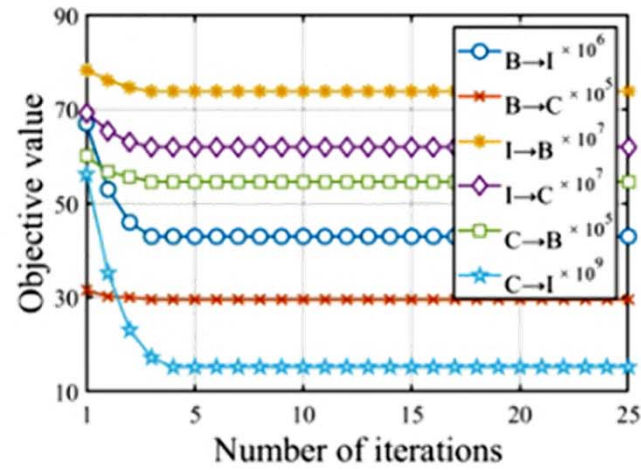
AN	0.69	0.03	0.26	0.02
HA	0.31	0.07	0.31	0.30
NE	0.17	0.01	0.57	0.25
SA	0.04	0.02	0.26	0.67
	AN	HA	NE	SA

(f) C→I

# 04 Convergence Analysis



GSDA-PCA



GSDA-LDA

# Conclusions

---

- We proposed GSDA utilizes a novel distance metric learning strategy to reduce the discrepancy between different corpora
- Extensive experimental results show that the proposed GSDA achieves superior performance than state-of-the-art compared algorithms.
- In the future, we will investigate to develop the deep transfer learning methods using the the proposed strategy to solve the cross-corpus dimensional SER problem.

A decorative graphic featuring a central white rectangle with a thin blue border. The rectangle is surrounded by four circles: a large blue circle in the top-left corner, a medium pink circle in the bottom-left corner, a medium pink circle in the top-right corner, and a large blue circle in the bottom-right corner. The circles are layered, with some overlapping the rectangle's corners.

Thank you!