# Unsupervised Cross-Corpus Speech Emotion Recognition Using Domain-Adaptive Subspace Learning

Na Liu [1,3] , Yuan Zong [2] , Baofeng Zhang [1] , Li Liu [3] ,Jie Chen [3] , Guoying Zhao [3] , Junchao Zhu [1]

[1] *Key laboratory of computer vision and systems, Ministry of education, Tianjin University of Technology, China*

[2] *Research Center for Learning Science, Southeast University, China*

[3] *Center for Machine Vision and Signal Analysis, University of Oulu, Finland*

## Introduction

➢ *Motivation*: In this paper, we investigate unsupervised cross-corpus speech emotion recognition (SER). The training (source) and testing (target) speech signals come from two different corpora which may have different feature distributions and therefore lots of existing SER methods would not work.

➢ *Solution*: (1) Construct a label space based on the label information provided by the source speech corpora to serve as the predefined common subspace for Domain-Adaptive Subspace Learning (DoSL). (2) Learn a projection matrix which transforms the source and target speech signals from the original feature space to a common subspace.

## Method

Our DoSL aims at learning a projection matrix U to project the source speech feature matrix Xs from the original feature space to such common subspace spanned by the columns of Ls, which can be formulated as the following optimization problem:

$$\min_{\mathbf{U}} \|\mathbf{L}^s - \mathbf{U}^T \mathbf{X}^s\|_F \quad (1)$$

Enforce the projected source and target speech features share the similar distributions

$$\min_{\mathbf{U}} \left\| \frac{1}{N_s} \sum_{i=1}^{N_s} \mathbf{U}^T \mathbf{x}_i^s - \frac{1}{N_t} \sum_{i=1}^{N_t} \mathbf{U}^T \mathbf{x}_i^t \right\|^2 \quad (2)$$

Minimizing the combination of the above objective functions in Eqs. (1) and (2)

$$\min_{\mathbf{U}} \|\mathbf{L}^s - \mathbf{U}^T \mathbf{X}^s\|_F$$
$$+ \lambda_1 \left\| \frac{1}{N_s} \sum_{i=1}^{N_s} \mathbf{U}^T \mathbf{x}_i^s - \frac{1}{N_t} \sum_{i=1}^{N_t} \mathbf{U}^T \mathbf{x}_i^t \right\|^2 + \lambda_2 \|\mathbf{U}^T\|_{2,1} \quad (3)$$

### Optimization

DoSL model is solved by using inexact augmented Lagrange multiplier (IALM) method. More specifically, by introducing a auxiliary variable Q which satisfies U = Q, we convert the optimization problem of DoSL to a constrained one which can be expressed as:

$$\min_{\mathbf{U}} \|\mathbf{L}^s - \mathbf{Q}^T \mathbf{X}^s\|_F^2$$
$$+ \lambda_1 \left\| \frac{1}{N_s} \sum_{i=1}^{N_s} \mathbf{Q}^T \mathbf{x}_i^s - \frac{1}{N_t} \sum_{i=1}^{N_t} \mathbf{Q}^T \mathbf{x}_i^t \right\|^2 + \lambda_2 \|\mathbf{U}^T\|_{2,1}$$
$$\text{s.t. } \mathbf{U} = \mathbf{Q} \quad (4)$$

Subsequently, the Lagrange function of Eq. (4) can be obtained as follows:

$$L(\mathbf{U}, \mathbf{Q}, \mathbf{T}, \mu) = \|\mathbf{L}^s - \mathbf{U}^T \mathbf{X}^s\|_F + \lambda_1 \left\| \mathbf{Q}^T \overline{\mathbf{X}}_i^{st} \right\|^2$$
$$+ \lambda_2 \|\mathbf{U}^T\|_{2,1} + tr[\mathbf{T}^T(\mathbf{U} - \mathbf{Q})] + \frac{\mu}{2} \|\mathbf{U} - \mathbf{Q}\|_F^2 \quad (5)$$

where $\overline{\mathbf{X}}_i^{st} = \frac{1}{N_s} \sum_{i=1}^{N_s} \mathbf{x}_i^s - \frac{1}{N_t} \sum_{i=1}^{N_t} \mathbf{x}_i^t$ ,T is the Lagrange multiplier, and $\mu > 0$ is the regularization parameter.

Iteratively minimize the Lagrange function of Eq. (4) until convergence:

1. Update Q

$$\min_{\mathbf{Q}} \|\mathbf{L}^s - \mathbf{Q}^T \mathbf{X}^s\|_F^2 + \lambda_1 \left\| \mathbf{Q}^T \overline{\mathbf{X}}_i^{st} \right\|^2 + \lambda_2 \|\mathbf{U}^T\|_{2,1}$$
$$+ tr[\mathbf{T}^T(\mathbf{U} - \mathbf{Q})] + \frac{\mu}{2} \|\mathbf{U} - \mathbf{Q}\|_F^2$$

2. Update U:

$$\min_{\mathbf{U}} \frac{\lambda_2}{\mu} \|\mathbf{U}^T\|_{2,1} + \frac{1}{2} \left\| \mathbf{U}^T - \left(\mathbf{Q}^T - \frac{\mathbf{T}^T}{\mu}\right) \right\|_F^2$$

3. Update T and $\mu$

$$\mathbf{T} = \mathbf{T} + \mu(\mathbf{U} - \mathbf{Q}), \mu = \max(\mu_{max}, \rho\mu),$$

4. Check convergence: $\|U - Q\|_\infty < \epsilon$

$$\text{emotion\_labels} = \underset{k}{\operatorname{argmin}} \{[U_*^T X^t](k,:)\}$$

where $[U_*^T X^t](k,:j)$ means the $k^{th}$ element of the $j^{th}$ column (target speech signal) of the projected matrix $U_*^T X^t$

## Results and discussion

Results of the cross-corpus SER experiments in terms of UAR and WAR, where the common emotion states (5 classes) are Angry, Disgust, Fear, Happy and Sad.

| # | Source Corpus | Source Corpus | SVM | | KMM | | KLIEP | | uLSIF | | DALSR | | DoSL | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | UAR | WAR | UAR | WAR | UAR | WAR | UAR | WAR | UAR | WAR | UAR | WAR |
| 1 | EmoDB | eNTERFACE | 30.06₃ | 30.08₃ | 23.08₅ | 23.14₅ | 21.79₆ | 21.82₆ | 25.75₄ | 25.75₄ | 36.36₂ | 36.40₂ | **37.49₁** | **37.51₁** |
| 2 | eNTERFACE | EmoDB | 27.83₆ | 24.27₆ | 40.18₄ | 44.69₃ | 28.58₅ | 27.01₅ | 40.42₃ | 42.27₄ | **44.41₁** | **52.27₁** | 44.25₂ | 52.00₂ |
| 3 | EmoDB | AFEW4.0 | 26.07₄ | 25.99₄ | **30.39₁** | 29.78₃ | 25.47₆ | 25.57₆ | 25.75₅ | 25.93₅ | 27.51₃ | 30.19₂ | 29.10₂ | **31.00₁** |
| 4 | AFEW4.0 | EmoDB | 29.87₅ | 35.02₅ | 38.17₂ | 46.81₃ | 27.41₆ | 31.37₆ | 36.25₄ | 44.38₄ | 37.33₃ | 47.80₂ | **39.66₁** | **50.00₁** |
| 5 | eNTERFACE | AFEW4.0 | 20.80₅ | 18.39₆ | 23.79₃ | 25.72₃ | 18.66₆ | 18.60₅ | 22.61₄ | 21.21₄ | 24.67₂ | **27.70₁** | **24.83₁** | 26.20₂ |
| 6 | AFEW4.0 | eNTERFACE | 18.68₄ | 18.72₄ | 19.75₃ | 19.75₃ | 17.48₆ | 17.47₆ | 18.10₅ | 18.11₅ | **21.93₁** | **21.96₁** | 21.64₂ | 21.66₂ |

◆ It is convincing that the limited label information provided by a small number of samples in source database will lead to low recognition rate.

◆ The data imbalance between source and target databases is an important factor which will affect the cross-cropus speech emotion recognition tasks.

## Acknowledgements