(O)

2019　2021

Construction of a computational model to deal with the cocktail-party problem for intelligent speech interface

Lu, Xugang

3,300,000

1.

1

2.

In cocktail party scenarios, many information need to be explored in order to identify different speech (or sound) sources. Under this project, we have the following contributions: 1. For identifying speech source, who is speaking (speaker information) is one of the most important information. Besides developing speaker embedding system, we proposed a coupling of generative and discriminative learning for speaker recognition. Our framework showed a large improvement compared with state of the art models. 2. Concerning speech source recording environments may change in different domains, we proposed a new distance metric for unsupervised domain adaptation technique. We have tested the proposed adaptation algorithm on both speaker and language recognition tasks, and obtained promising improvement when speech recording environments are changed.

様　式　Ｃ－１９、Ｆ－１９－１、Ｚ－１９（共通）

## １．研究開始当初の背景

For most speech technology application systems, when they are applied in cocktail-party scenarios, i.e., real acoustic environments with mixed sound sources, the performance is drastically degraded. The reason is that the conventional computational models (no matter with or without deep learning algorithms) used in those applications take all entangled sound sources without actively selecting one of (or some of) them for processing and recognition. Parsing mixed sound sources in cocktail-party scenarios is a very important "intelligent" function of human beings during speech communication. In this study, our ambition is to construct a new computational model to realize this "intelligent" function with selective attention for speech interface. The computational model should integrate both bottom-up sound saliency detection and top-down selective attention processing to dynamically parse the incoming mixed sound sources. Among this computational framework, integrating prior knowledge is one of the most effective ways for parsing the mixed speech sources. Among the available prior information, speaker information and language information are important factors that could be accurately identified before speech source separation.

## ２．研究の目的

The speaker and language identification algorithms could be directly applied for our purpose. However, there are several problems remained in real applications: 1. For short utterance, it is difficult to obtain a satisfied performance due to limited information (speaker information for speaker recognition task, and language information for language recognition task). 2. For a model trained with a certain domain data, the performance will degrade drastically when the domain is different in testing (e.g., real recording environments are changed frequently). The purpose of this study is to deal with the two problems for real application purpose.

## ３．研究の方法

For the first problem we mentioned above, besides adopting an embedding method (speaker or language embedding), we further propose to integrate a generative model and a discriminative model for improving the performance. As the generative model focuses on class-conditional feature distributions while the discriminative model focuses on classification boundaries, the generative model could have a good generalization for short utterances (but less discriminative power), the discriminative model has high discriminative capacity (but less generalization ability to short utterances). Fig. 1 shows the two different focuses of the two types of models. In this figure, only two classes are showed. By coupling generative model in a discriminative neural network learning framework, we could combine both the advantages of generative and discriminative models to constrain large model variation (due to large feature variation of short
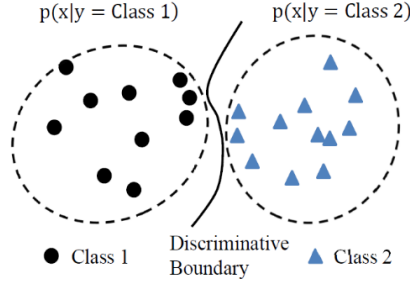
utterances).



Fig. 1. Generative model learning focuses on class conditional feature distributions (dashed-circles of feature distribution shapes), and discriminative model learning emphasizes the class discriminative boundary (solid curve).

For the second problem, we adopt an adaptation method to deal with the cross-domain problem. The problem can be well described as: source domain data set $D_s = \{(x_i^s, y_i^s)\}_{i=1,2,\ldots,N}$, and target domain data set $D_t = \{(x_i^t, y_i^t)\}_{i=1,2,\ldots,M}$. Due to domain changes (e.g. recording channels), there exists domain discrepancy, i.e., $p_s(x,y) \neq p_s(x,y)$ . The purpose for domain adaptation is to reduce this discrepancy to make $p_s(x,y) \approx p_s(x,y)$ . Based on Bayesian theory, $p_s(x,y) = p_s(y|x)p_s(x)$ and $p_t(x,y) = p_t(y|x)p_t(x)$, we need to approximate the two terms as: $p_s(y|x) \approx p_t(y|x)$ and $p_s(x) \approx p_t(x)$. For finding a latent transformed space $z = \varphi(x)$, the approximation will be $p_s(y|z) \approx p_t(y|z)$ and $p_s(z) \approx p_t(z)$.

４．研究成果

For dealing with the first problem, we proposed coupling a generative model with a discriminative learning framework, and applied to improve speaker recognition performance. The proposed model framework is showed in Fig. 2.
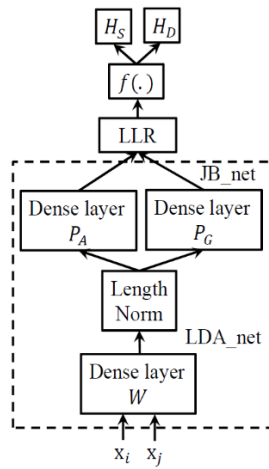


Fig. 2 The proposed two-branch Siamese neural network with coupling of the generative joint Bayesian model structure.

In this Fig. 2, the model framework was adopted for speaker verification task with two hypothesis labels $H_S$ and $H_D$ as the two compared utterances are from the same speaker and different speakers, respectively. LLR means log-likelihood ratio score, and JB_net as joint Bayesian model (as generative model), LDA_net as linear discriminative net for dimensional reduction. Dense layers were used to fit the functions of model parameters (coupling to the generative model). And the input feature vectors $x_i$ and $x_j$ are two compared vectors representing two utterances. For discriminative training, we further proposed an objective function based on false alarm and miss measure metrics which are used in detection tasks. The idea was illustrated in Fig. 3.
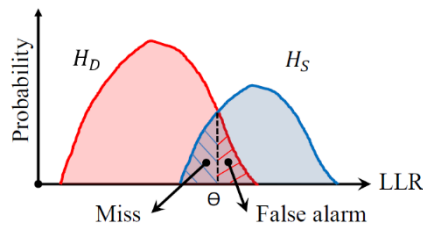


Fig. 3 The LLR distribution for $H_S$ and $H_D$ conditions.

Based on the proposed framework, the two hypothesis distributions ($H_S$ and $H_D$) were further separated as showed in Fig. 4. And the speaker verification experiments confirmed the improved performance.
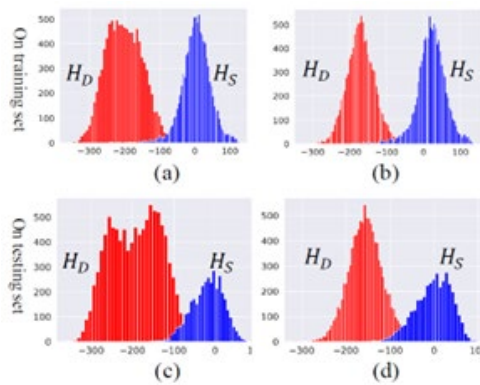


Fig 4. LLR distributions in $H_S$ and $H_D$ spaces: the first row (a, b) and second row (c, d) are for the training and test sets, respectively; the left column (a, c) for setting with generative model parameters learned based on the EM algorithm, and the right column (b, d) for setting with discriminatively trained parameters after initializing with generative model parameters.

For dealing with the cross-domain problem, our proposed model framework is showed in Fig. 5. In this framework, $x^t$ and $x^s$ represent feature vectors from target and source domains, with a transform (X-vector extraction and a neural dense layer transform, and length normalization), we obtained latent feature representations as $z^t$ and $z^s$, and finally obtained their classification

labels $y^t$ and $y^s$. Based on our theoretical analysis, we optimized the model parameters for approximation of joint distributions for cross-domain adaptation.
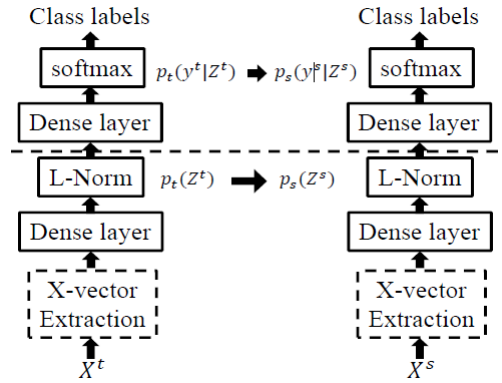


Fig. 5 Joint distribution adaptation for language recognition.

Based on the proposed adaptation algorithm, the domain discrepancy was reduced. An example of the feature distributions of two languages before and after adaptation was showed in Fig. 6. From this figure, we can see that after adaptation, the distributions of training and test sets have large overlap.
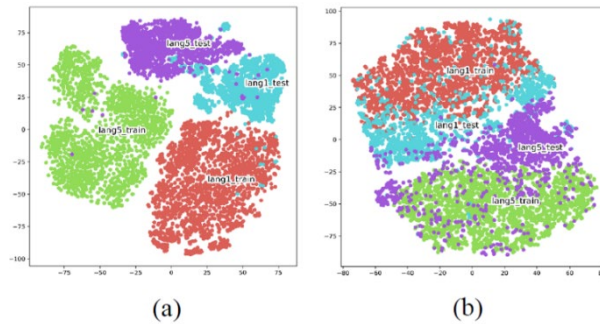


(a)                    (b)

Fig. 6 Language cluster distributions based on the TSNE for cross-domain test: before adaptation (a), and after adaptation (b).

| | |
|---|---|
| Xugang Lu, Peng Shen, Sheng Li, Yu Tsao, Hisashi Kawai | 29 |
| Coupling a Generative Model With a Discriminative Learning Framework for Speaker Verification | 2021 |
| IEEE/ACM Transactions on Audio, Speech, and Language Processing | 3631-3641 |
| DOI<br>10.1109/TASLP.2021.3129360 | |
| | |

| |
|---|
| Xugang Lu, Peng Shen, Yu Tsao, Hisashi Kawai |
| Class-Wise Centroid Distance Metric Learning for Acoustic Event Detection |
| Interspeech 2019 |
| 2019 |

| |
|---|
| Peng Shen, Xugang Lu, Komei Sugiura, Sheng Li, Hisashi Kawai |
| Compensation on x-vector for short utterance spoken language identification |
| Odyssey 2020 The Speaker and Language Recognition Workshop |
| 2020 |

| |
|---|
| Xugang Lu, Peng Shen, Yu Tsao, Hisashi Kawai |
| UNSUPERVISED NEURAL ADAPTATION MODEL BASED ON OPTIMAL TRANSPORT FOR SPOKEN LANGUAGE IDENTIFICATION |
| ICASSP2021 |
| 2021 |

Xugang Lu, Peng Shen, Sheng Li, Yu Tsao, Hisashi Kawai

Siamese Neural Network with Joint Bayesian Model Structure for Speaker Verification

O

| | | |
|---|---|---|
| | | |

O

| | |
|---|---|
| | |

Xugang Lu, Peng Shen, Sheng Li, Yu Tsao, Hisashi Kawai