Model selection for deep audio source separation via clustering analysis

Alisa Liu, Prem Seetharaman, Bryan Pardo DCASE 2020



Northwestern University

Audio source separation



We can automatically distinguish different sound sources in an auditory scene...

Audio source separation



We can automatically distinguish different sound sources in an auditory scene...

Audio source separation



... Moreover, we can easily handle different types of auditory environments!

The problem

Source separation models are trained to separate on domain-specific data, and do not generalize across domains

Motivating question

Given an **audio mixture** whose source domain is unknown, can we **automatically select** the best model for the mixture?



We develop a confidence measure for systems that perform clustering-based separation...



We develop a confidence measure for systems that perform clustering-based separation...



Confidence = .2

We develop a confidence measure for systems that perform clustering-based separation...



Confidence = .2

We develop a confidence measure for systems that perform clustering-based separation... to automatically select the model output with the best predicted separation quality



Confidence = .2



Map each time-frequency point...



... to points in an embedding space...



... so we can cluster the embeddings...



Intuition for confidence measure



Key insight: The distribution of embedded TF points is predictive of the performance of the algorithm





For a point x_i in cluster C_k Intracluster distance



For a point x_i in cluster C_k Intracluster distance



For a point x_i in cluster C_k Intracluster distance



For a point x_i in cluster C_k Intracluster distance

 $a(x_i) =$















For a point x_i in cluster C_k Combine and scale:

$$s(x_i) = \frac{b(x_i) - a(x_i)}{\max\left\{a(x_i), b(x_i)\right\}}$$



For a point x_i in cluster C_k Combine and scale:

$$s(x_i) = \frac{b(x_i) - a(x_i)}{\max\left\{a(x_i), b(x_i)\right\}}$$

Average across a sample of 1000 TF bins in the top 1% loudest bins

$$S(X) = \sum_{x_i \in X} s(x_i)$$











For a point x_i in cluster C_k

 $p(x_i) =$





For a point x_i in cluster C_k

$$p(x_i) = \frac{K\left(\max_{k \in [0,...,K]} \gamma_{ik}\right) - 1}{K - 1}$$

Average across the top 1% loudest bins

$$P(X) = \sum_{x_i \in X} p(x_i)$$

Confidence measure

Combine silhouette score and posterior strength in a **product** to obtain the overall confidence measure for the mixture:

confidence measure

$$C(X) = S(X)P(X)$$

silhouette posterior score strength

Experimental Design



For each domain, train a model with the deep clustering objective

- 2 BLSTM layers
- 300 hidden units in both directions

2D embedding visualization of a music mixture



Evaluation of confidence measure

Strong correlations between the confidence measure and ground-truth separation quality



Approach	Speech	Music	Environ.
Oracle ensemble			
Confidence ensemble	1		
Random ensemble]		
Speech model]		
Music model			
Environ. model]		

Approach	Speech	Music	Environ.
Oracle ensemble			
Confidence ensemble	Î		
Random ensemble]		
Speech model]		
Music model			
Environ. model]		

Approach	Speech	Music	Environ.
Oracle ensemble			
Confidence ensemble	Ī		
Random ensemble	Ī		
Speech model	I		
Music model			
Environ. model	Ī		

Compare different methods of choosing the appropriate model given an audio mixture

Approach	Speech	Music	Environ.
Oracle ensemble	8.3	6.5	12.2
Confidence ensemble			
Random ensemble	1		
Speech model]		
Music model			
Environ. model]		

Oracle: select model with best ground-truth performance

Compare different methods of choosing the appropriate model given an audio mixture

Approach	Speech	Music	Environ.
Oracle ensemble	8.3	6.5	12.2
Confidence ensemble	7.6	6.4	10.5
Random ensemble		-	
Speech model]		
Music model			
Environ. model]		

Confidence: select model with highest confidence

Compare different methods of choosing the appropriate model given an audio mixture

Approach	Speech	Music	Environ.
Oracle ensemble	8.3	6.5	12.2
Confidence ensemble	7.6	6.4	10.5
Random ensemble	4.8	4.2	2.8
Speech model		-	
Music model			
Environ. model]		

Random: select model randomly with equal probability

Compare different methods of choosing the appropriate model given an audio mixture

Approach	Speech	Music	Environ.
Oracle ensemble	8.3	6.5	12.2
Confidence ensemble	7.6	6.4	10.5
Random ensemble	4.8	4.2	2.8
Speech model	8.2	2.0	3.0
Music model	1.4	6.5	2.5
Environ. model	2.1	1.7	11.9

Domain-specific models: applied without switching

Compare different methods of choosing the appropriate model given an audio mixture

Approach	Speech	Music	Environ.
Oracle ensemble	8.3	6.5	12.2
Confidence ensemble	7.6	6.4	10.5
Random ensemble	4.8	4.2	2.8
Speech model	8.2	2.0	3.0
Music model	1.4	6.5	2.5
Environ. model	2.1	1.7	11.9

Domain-specific models: applied without switching

Compare different methods of choosing the appropriate model given an audio mixture

Approach	Speech	Music	Environ.
Oracle ensemble	8.3	6.5	12.2
Confidence ensemble	7.6	6.4	10.5
Random ensemble	4.8	4.2	2.8
Speech model	8.2	2.0	3.0
Music model	1.4	6.5	2.5
Environ. model	2.1	1.7	11.9

Domain-specific models: applied without switching

Approach	Speech	Music	Environ.
Oracle ensemble	8.3	6.5	12.2
Confidence ensemble	7.6	6.4	10.5
Random ensemble	4.8	4.2	2.8
Speech model	8.2	2.0	3.0
Music model	1.4	6.5	2.5
Environ. model	2.1	1.7	11.9

Approach	Speech	Music	Environ.
Oracle ensemble	8.3	6.5	12.2
Confidence ensemble	7.6	6.4	10.5
Random ensemble	4.8	4.2	2.8
Speech model	8.2	2.0	3.0
Music model	1.4	6.5	2.5
Environ. model	2.1	1.7	11.9

Approach	Speech	Music	Environ.
Oracle ensemble	8.3	6.5	12.2
Confidence ensemble	7.6	6.4	10.5
Random ensemble	4.8	4.2	2.8
Speech model	8.2	2.0	3.0
Music model	1.4	6.5	2.5
Environ. model	2.1	1.7	11.9

Conclusion

- Confidence measure effectively estimates the performance of clustering-based source separation algorithms
- Apply the confidence measure to effectively select the appropriate model for a given mixture