



**GHENT  
UNIVERSITY**

# Sound separation

using ad-hoc distributed microphone arrays

Martijn Meeldijk — [martijn.meeldijk@ugent.be](mailto:martijn.meeldijk@ugent.be)

# OVERVIEW

- 1 INTRODUCTION: AD-HOC MICROPHONE ARRAYS
- 2 CLUSTERING SOURCE DOMINATED MICROPHONES
- 3 FEDERATED LEARNING
- 4 COHERENCE-BASED CLUSTERING
- 5 QUESTIONS?

# Introduction: Ad-hoc microphone arrays

# AD-HOC MICROPHONE ARRAYS

Also known as acoustic sensor networks (ASN):

- Several microphones
- Unknown locations in a room

Increasing interest due to:

- Declining cost of receivers
- Edge devices such as laptops, phones, tablets

# CHALLENGES

Often different nodes are only weakly synchronized

- Algorithm should be able to deal with this

Limited bandwidth

- Try to exchange little information
- Limit amount of computations

# PROBLEMS

## Clustering

- Divide receivers into clusters dominated by a source
- Optional extra cluster for background noise or interference

## Classification

- E.g. speech, music, noise, ...

## Separation

- Enhancing the signal for a given source

# Clustering source dominated microphones



# FUZZY CLUSTERING

Instead of hard clustering

- Assign fuzzy membership values (FMV)
- Indicates how much a microphone belongs to a cluster

# FUZZY CLUSTERING

## STEPS

Extract Mod-MFCC features

- Feature vector based on modulations on Mel-Frequency cepstral coefficients
  - Processed in blocks of 4s
- Averaged over time, so behaves more stationary compared to short-time audio features.

# FUZZY CLUSTERING

## STEPS

Assign fuzzy-membership values using the Fuzzy C-means Algorithm by evaluating least-squared error functional

$$J_m = \sum_{d=1}^D \sum_{n=1}^N (\mu_{nd})^\alpha ||\mathbf{v}_d - \mathbf{u}_n||_\beta^2$$

- $N$  clusters
  - $N$  cluster centers  $\mathbf{u}_n = (\mathbf{u}_{n,1}, \mathbf{u}_{n,2}, \dots, \mathbf{u}_{n,A})^T$
- $D$  observations  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_D\}$
- Weighting exponent  $\alpha$  and distance norm  $|||_\beta$  with weighting matrix  $\beta$

# FUZZY CLUSTERING

## STEPS

Detecting the background cluster

- An extra cluster is added for microphones where background noise is dominant
- Microphones assigned to a source cluster should have the background cluster as second highest FMV

# FUZZY CLUSTERING

## BACKGROUND CLUSTER FMV

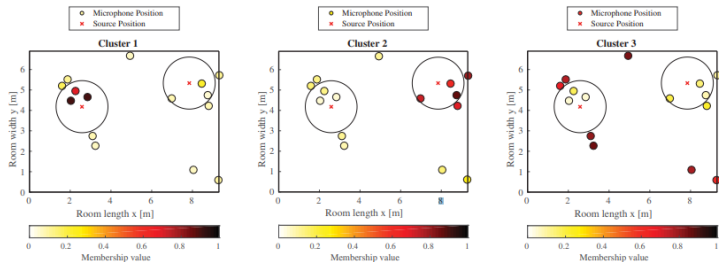
Mic. d	Cluster n		
	1	2	3
1	0.1	0.3	<b>0.6</b>
2	<b>0.5</b>	0.3	0.2
3	<b>0.7</b>	0.2	0.1
4	0.25	0.2	<b>0.55</b>
5	0.25	<b>0.6</b>	0.15
6	0.15	<b>0.65</b>	0.2
7	0.05	0.15	<b>0.8</b>

Figure: FMV's for each microphone

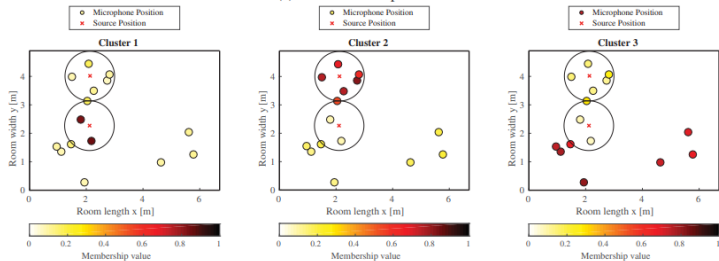
	Cluster n		
	1	2	3
Mics. in Cluster 1	<b>0.6</b>	0.25	0.15
Mics. in Cluster 2	0.2	<b>0.63</b>	0.17
Mics. in Cluster 3	0.13	0.22	<b>0.65</b>

Figure: Average FMV's per cluster

# PERFORMANCE



(a) Distant source positions.



(b) Close source positions.

# FUZZY CLUSTERING

## SOURCE SEPARATION

- 1 Estimate relative time-differences-of-arrival (TDOA)
- 2 Combine with FMV in beamforming stage
- 3 Apply cluster-related spectral masks to the output of the beamformers

# ESTIMATE RELATIVE TIME-DIFFERENCES-OF-ARRIVAL (TDOA)

- Select reference microphone (highest FMV)
- Initial source signal estimation
  - Calculate spectral mask for reference microphone

$$\mathcal{M}_n(k, b) = \begin{cases} 1 & |X_{R_n}(k - b)| > \frac{1}{B} \sum |X_{R_j}(k, b)|, \\ & j = 1, \dots, N \text{ and } j \neq n \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

- Apply to all microphones in cluster
- Correlation analysis with respect to the reference microphone to compute TDOAs



# CLUSTERING-STEERED BEAMFORMING

$$\hat{s}_{n, \text{w-DSB}}(l) = \sum_{i_n} w_{n,i_n} x_{i_n}(l + D_{i_n})$$

- $l$ : discrete time index
- $D_{i_n}$ : relative TDOA's
- $w_{n,i_n}$ : the weights allocated to each microphone  $i_n$  of cluster  $n$
- These are proportional to the FMV

By using a weighted combination of the microphone signals, we can make a better DSB with a better output SNR than uniform weighting.

# APPLY CLUSTER-RELATED SPECTRAL MASKS

## TO THE OUTPUT OF THE BEAMFORMERS

- Compute a post-filtering mask, similar to what we did before

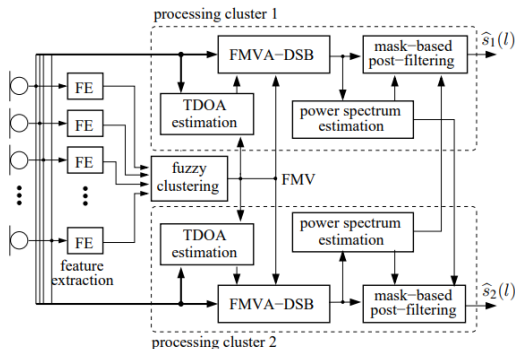


$$\mathcal{M}_n(k, b) = \begin{cases} 1 & |\hat{S}_{n,\text{FMVA-DSB}}(k, b)| > \frac{1}{B} \sum_{j=1, \dots, N, j \neq n} |\hat{S}_{j,\text{FMVA-DSB}}(k, b)|, \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

- Yields final enhanced estimate of the source signal

# OVERVIEW

(FOR TWO SOURCES)



This approach yields a slightly better performance than simple DSB.

# Federated Learning

# FEDERATED LEARNING

Is a machine learning technique that that trains a model distributed over multiple decentralized devices. Data is not shared between nodes. Instead, a three step process is followed:

- 1 The clients synchronize with the server by downloading the latest model parameters
  - This is a column vector  $\theta^\tau$
- 2 Each client independently improves its own parameters  $\theta_i^\tau$  with stochastic gradient descent on their own data
- 3 Each client uploads its model parameter updates  $\Delta\theta_i^\tau$  to the server, where they are aggregated.

$$\theta^{\tau+1} = \sum_{i=1}^M \frac{|D_i|}{|D|} \Delta\theta_i^\tau \quad (3)$$

- $M$ : number of clients
- $|D_i|$ : the cardinality of the dataset of the  $i$ -th client
- $|D|$ : the cardinality of the total dataset

# CLUSTERED FEDERATED LEARNING

## Problem

- Doesn't work if the different clients' data comes from different (incongruent) distributions.
- There's no single set of parameter updates that can optimally minimize the loss of all clients at the same time.

# CLUSTERED FEDERATED LEARNING

## Problem

- Doesn't work if the different clients' data comes from different (incongruent) distributions.
- There's no single set of parameter updates that can optimally minimize the loss of all clients at the same time.

## Solution

Cluster clients with a similar distribution and train a separate server model for each cluster.

# CLUSTERED FEDERATED LEARNING

- Calculate cosine similarity between update vectors of clients  $i$  and  $j$ , represented in matrix  $A$ , where:

$$a_{i,j} = \frac{\langle \Delta\theta_i, \Delta\theta_j \rangle}{||\Delta\theta_i, \Delta\theta_j||} \quad (4)$$

- $\langle \cdot \rangle$  denotes the inner product
- $||\cdot||$  denotes the  $L_2$  norm
- Cluster hierachically using bi-partitioning
  - Resulting two clusters  $c_1$  and  $c_2$  derived such that

$$\max_{\forall i \in c_1, k \in c_2} (a_{i,k}) < \min(\min_{\forall i,j \in c_1} (a_{i,j}), \min_{\forall k,l \in c_2} (a_{k,l})) \quad (5)$$

- maximum cross-cluster cosine similarity is always smaller than the minimum of either intra-cluster cosine similarities



# CLUSTERED FEDERATED LEARNING

- Verification using the mean and maximum Euclidian norms of weight update vectors

$$\Delta\bar{\theta}_c = \left\| \frac{1}{|c|} \sum_{i \in c} \Delta\theta_i \right\| \text{ and } \Delta\hat{\theta}_c = \max_{i \in c} (\|\Delta\theta_i\|) \quad (6)$$

- If server reaches stationary solution but clients are still converging
- $\Delta\bar{\theta}_c$  will be low and  $\Delta\hat{\theta}_c$  high, this prompts bi-partitioning

# MEMBERSHIP VALUES

To assess contribution of each node to its respective cluster. For two clusters:

- Calculate average intra-cluster similarities for each client  $i$ , arranged in vector  $\mathbf{q}$

$$q_i = \frac{1}{|c_x| - 1} \sum_{j \in c_x/i} a_{i,j} \quad (7)$$

- Calculate average cross-cluster similarities for each client  $i$ , arranged in vector  $\mathbf{r}$

$$r_i = \frac{1}{|c_y|} \sum_{k \in c_y} a_{i,k} \quad (8)$$

- $\forall i \in c_x$  and  $(c_x, c_y) \in \{(c_1, c_2), (c_1, c_1)\}$ , where  $|\cdot|$  denotes the cardinality

## MEMBERSHIP VALUES

- Now apply min-max normalization to obtain vector  $\mathbf{p}$  with:

$$p_i = \lambda \frac{q_i - \min(\mathbf{q})}{\max(\mathbf{q}) - \min(\mathbf{q})} + (1 - \lambda) \frac{r_i - \min(\mathbf{r})}{\max(\mathbf{r}) - \min(\mathbf{r})} \quad (9)$$

- Select the node with the smallest  $p_i$  as reference node
- The membership value of a node is the cosine similarity between this node and the reference node

$$\mu_i = a_{i, \arg \min(p_j)}, \forall i, j \text{ and } c_x \in \{c_1, c_2\} \quad (10)$$

- Now, apply min-max normalization again to the vector  $\mu$

# CLUSTER MEMBERSHIP VALUES

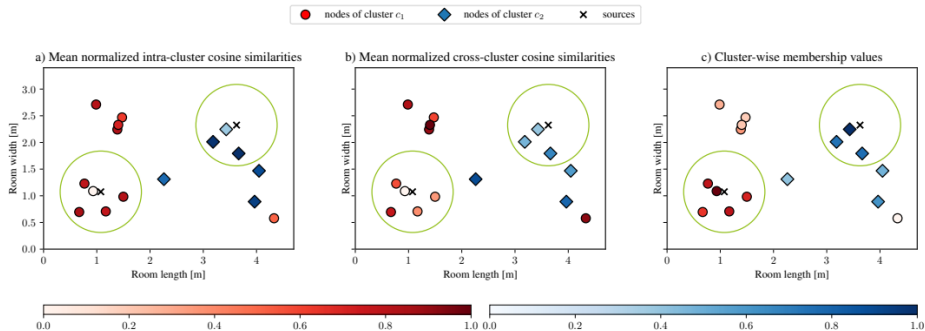


Figure: FMV's for unsupervised CFL

# Coherence-based clustering

# COHERENCE-BASED CLUSTERING

## WITH NON-NEGATIVE MATRIX FACTORIZATION

- A signal consisting of clean speech and interference:

$$x_m(t) = s_m(t) + v_m(t) \quad (11)$$

- We'll consider the signal as a frame of observation samples in vector form

$$\begin{aligned} \mathbf{x}_m(t) &= [x_m(t)x_m(t-1)\cdots x_m(t-T+1)]^T \\ &= \mathbf{s}_m(t) + \mathbf{v}_m(t) \end{aligned}$$

- With  $T$  denoting frame size

# COHERENCE-BASED CLUSTERING

## WITH NON-NEGATIVE MATRIX FACTORIZATION

The proposed clustering algorithm works in two steps:

- 1 Compute the magnitude squared coherence between the microphone observations to measure their degree of linear dependency
- 2 Apply NMF to the output of the previous step to find the optimal clustering

## MAGNITUDE-SQUARED COHERENCE

- We can measure the coherence between two signals  $x(t)$  and  $y(t)$  by computing the FFT of the signals.
- Then, we measure the coherence as a function of the center frequency of the filter.

$$\Gamma_{xy}(f) = \frac{|S_{xy}(f)|^2}{S_{xx}(f)S_{yy}(f)} \quad (12)$$

- $S_{xy}(f)$  is the cross-spectral density, which is really a spectral correlation density. This tells how the two signals are correlated in the time domain at a certain frequency and can be computed as:

$$S_{xy}(f) = \sum_{k=1-T}^{T-1} R_{xy}(k) e^{-i2\pi fk} \quad (13)$$



## MAGNITUDE-SQUARED COHERENCE

$R_{xy}(k)$  from the previous slide is the cross-correlation function between  $x(t)$  and  $y(t)$  and can be estimated by:

$$R_{xy}(k) = \begin{cases} \frac{1}{T} \sum_0^{T-1-k} x(t)y(t+k) & k = 0, \dots, T-1 \\ R_{xy}(-k) & k = -(T-1), \dots, -1 \end{cases} \quad (14)$$

## MAGNITUDE-SQUARED COHERENCE

Now we'll make a matrix consisting of all the correlations between the different microphone signals. We'd like to give the same weight to all frequency bins, regardless of their power. So we'll use the following coherence metric:

$$C_{xy} = \frac{\sum_{f=0}^F \Gamma_{xy}(f)}{F} \in [0, 1] \quad (15)$$

- With  $F$  denoting the number of frequency bins

Arranging all of the coherence measures between observations into a matrix, we obtain  $\mathbf{C}$

$$\mathbf{C} = \begin{bmatrix} 1 & \cdots & \cdots & C_{1M} \\ C_{12} & 1 & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ C_{1M} & C_{2M} & \cdots & 1 \end{bmatrix} \quad (16)$$

# NON-NEGATIVE MATRIX FACTORIZATION

- Every value  $C_{ij}$  contains the degree of correlation between the  $i$ -th and  $j$ -th observation. This means that microphones close to a certain source will be correlated.
- The matrix  $\mathbf{C}$  can be modelled as:

$$\mathbf{C} = \mathbf{B}\mathbf{B}^T \odot (\mathbf{1} - \mathbf{I}) + \mathbf{I} \quad (17)$$

- $\mathbf{B} \in \mathbb{R}^{M \times K}$  is the cluster matrix, where  $K$  is the amount of speakers (the amount of clusters)
- Because  $\mathbf{C}$  is symmetric, we model it as  $\mathbf{B}\mathbf{B}^T$
- $\odot$  is the element-wise product
- $\mathbf{I}$  is the identity matrix
- $\mathbf{1}$  is the all-ones matrix
- The latter two are introduced because the main diagonal of  $\mathbf{C}$  does not provide any relevant information in the learning process of  $\mathbf{B}$

## NON-NEGATIVE MATRIX FACTORIZATION

- Based on Euclidian divergence,  $\mathbf{B}$  is estimated with multiplicative update rules. It is first initialized with random values in the algorithm:

$$\mathbf{B} \leftarrow \mathbf{B} \odot \frac{(\mathbf{C} \odot (\mathbf{1} - \mathbf{I}))\mathbf{B}}{(\mathbf{B}\mathbf{B}^T \odot (\mathbf{1} - \mathbf{I}))\mathbf{B}} \quad (18)$$

- Now each column of  $\mathbf{B}$  contains the contribution of a microphone to each cluster. We can obtain the clustering result with:

$$\gamma_m = \{j \in [1, K] : B_{mj} \geq B_{mk}, \forall k \in [1, K]\} \quad (19)$$

- $\gamma_m$  denotes the cluster assigned to the  $m$ -th microphone. This is simply the largest value of column  $m$ .

For those wondering, the variance ratio criterion (VRC) is used to obtain the optimal number of clusters

# NON-NEGATIVE MATRIX FACTORIZATION

## RESULTS

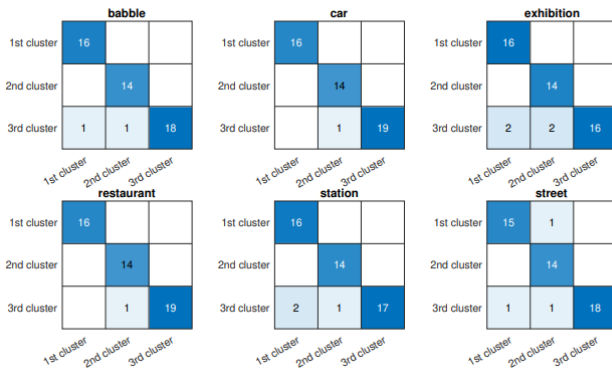
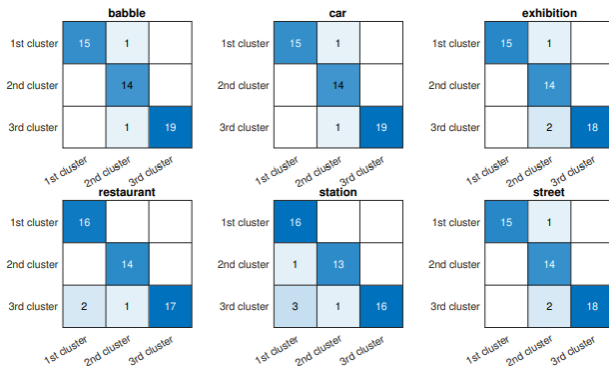


Figure: Confusion matrix for a simulation with 3 clusters and 50 microphones. For reverberation time  $T_{60} = 200ms$

# NON-NEGATIVE MATRIX FACTORIZATION

## RESULTS



**Figure:** Confusion matrix for a simulation with 3 clusters and 50 microphones. For reverberation time  $T_{60} = 400ms$

# NON-NEGATIVE MATRIX FACTORIZATION

## CONCLUSION

- Allows clustering without any prior knowledge of the acoustic scene
- Improved clustering performance in comparison to other methods

# OVERVIEW

- 1 INTRODUCTION: AD-HOC MICROPHONE ARRAYS
- 2 CLUSTERING SOURCE DOMINATED MICROPHONES
- 3 FEDERATED LEARNING
- 4 COHERENCE-BASED CLUSTERING
- 5 QUESTIONS?



# Questions?

# Sound separation

using ad-hoc distributed microphone arrays

Martijn Meeldijk — [martijn.meeldijk@ugent.be](mailto:martijn.meeldijk@ugent.be)