

Joint Audio Source Separation and Diarization

J Source signals, reverberated and summed, are being recorded at I microphones.

- Separation: Recover the J original source signals!
- Diarization: Classify each source as active/inactive along time!

Standard Formulation in the STFT domain

- Separate a mixture of J sources with I microphones.
- In STFT domain the problem becomes:

$$\mathbf{x}_{f\ell} = \mathbf{A}_f \mathbf{s}_{f\ell} + \mathbf{b}_{f\ell}$$

mixture $[I \times 1]$
observed

mixing matrix $[I \times J]$
unknown!

source STFT $[J \times 1]$
unknown!

sensor noise $[I \times 1]$
unknown!

- $f=1:F$ frequency bins, $\ell=1:L$ time frames.

Modelling Diarization

- The Standard Model: $\mathbf{x}_{fl} = \mathbf{A}_f \mathbf{s}_{fl} + \mathbf{b}_{fl}$ has all sources active;
- Look between \mathbf{A}_f and \mathbf{s}_{fl} : You notice the identity matrix????

$$\mathbf{x}_{fl} = \mathbf{A}_f \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} \begin{bmatrix} s_{1,fl} \\ s_{2,fl} \\ s_{3,fl} \end{bmatrix} + \mathbf{b}_{fl}.$$

- This is a special case where the diagonal entries are fixed to 1.
- What if an entry was 0 instead, e.g.

$$\mathbf{x}_{fl} = \mathbf{A}_f \begin{bmatrix} 1 & & \\ & 1 & \\ & & 0 \end{bmatrix} \begin{bmatrix} s_{1,fl} \\ s_{2,fl} \\ s_{3,fl} \end{bmatrix} + \mathbf{b}_{fl} =$$

$$\mathbf{a}_{1,f} s_{1,fl} + \mathbf{a}_{2,f} s_{2,fl} + \cancel{\mathbf{a}_{3,f} s_{3,fl}},$$

- where $\mathbf{a}_{1,f}, \dots, \mathbf{a}_{3,f} \in \mathbb{C}^l$ are the columns of \mathbf{A}_f .

The state variable

- For $J = 3$ sources there are $N = 8$ possible matrices:

$$\mathbf{D}_1 = \begin{bmatrix} 0 & & \\ & 0 & \\ & & 0 \end{bmatrix}, \mathbf{D}_2 = \begin{bmatrix} 0 & & \\ & 0 & \\ & & 1 \end{bmatrix}, \dots, \mathbf{D}_8 = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix}.$$

- Let a categorical variable $Z_\ell = n, n \in [1, N]$ choose the \mathbf{D}_n at time frame ℓ .
- The hidden variable Z_ℓ has a temporal model on $\ell = 1, \dots, L$

$$\begin{aligned} p(Z_\ell = n | Z_{\ell-1} = r) &= T_{nr}, \\ p(Z_\ell = n) &= \lambda_n, \end{aligned}$$

with λ_n, T_{nr} prior and transition parameters.

Audio Mixture with Diarization

- We can compactify:

$$\mathbf{x}_{f\ell} = \mathbf{A}_f \mathbf{D}_{Z_\ell} \mathbf{s}_{f\ell} + \mathbf{b}_{f\ell},$$

- or probabilistically (white isotropic $\mathbf{b}_{f\ell}$):

$$p(\mathbf{x}_{f\ell} | Z_\ell = n, \mathbf{s}_{f\ell}) = \mathcal{N}_c(\mathbf{A}_f \mathbf{D}_n \mathbf{s}_{f\ell}, v_f \mathbf{I}_I),$$

where $\mathbf{A}_f \in \mathbb{C}^{I \times J}$ mixing matrix, and $v_f \in \mathbb{R}_+$ microphone noise variance, parameters to be estimated.

Associated Graphical Model

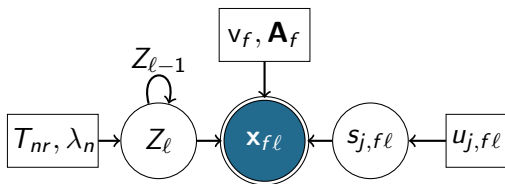


Figure: Joint Audio Source Separation and Diarization.

Associated EM algorithm

- Hidden variables:

$$\mathcal{H} = \{\mathbf{s}_{fl}, Z_l\}_{f,\ell=1}^{F,L}.$$

- Observed data:

$$\mathcal{X} = \{\mathbf{x}_{fl}\}_{f,\ell=1}^{F,L}.$$

- Joint a posteriori distribution **it factorizes!!!!!!!!!!!!!!!!!!!!!!**

$$p(\mathbf{s}_{1:F1:L}, Z_{1:L} | \mathcal{X}) \propto p(\mathbf{s}_{1:F1:L} | Z_{1:L}, \mathcal{X}) p(Z_{1:L} | \mathcal{X}).$$

- Parameters to be estimated (e.g. $J = 3, N = 8$):

$$\theta = \{\mathbf{A}_f, \mathbf{v}_f, u_{j,fl}, T_{nr}, \lambda_n\}_{f,\ell,j,n,r=1}^{F,L,J,N,N}.$$

E-Step Outline 1/2

- E-step of **Sources|Diarization**: For every value of $Z_\ell = n$:

$$p(\mathbf{s}_{f\ell} | Z_\ell = n) = \mathcal{N}_c(\boldsymbol{\mu}_{f\ell n}^s, \boldsymbol{\Sigma}_{f\ell n}^s).$$

- Closed form expressions

$$\boldsymbol{\Sigma}_{f\ell n}^s = \left[\text{diag}_J \left(\frac{1}{u_{j,f\ell}} \right) + \mathbf{D}_n \frac{\mathbf{A}_f^H \mathbf{A}_f}{v_f} \mathbf{D}_n \right]^{-1},$$

$$\boldsymbol{\mu}_{f\ell n}^s = \boldsymbol{\Sigma}_{f\ell n}^s \mathbf{D}_n \mathbf{A}_f^H \frac{\mathbf{x}_{f\ell}}{v_f}.$$

E-Step Outline 2/2

- E Step of **Diarization**: Estimate the Responsibilities of Z_ℓ :
- $\eta_{\ell,n} = p(Z_\ell = n|\mathcal{X})$ is found by solving a HMM with emission probability

$$p(\mathcal{X}|Z_\ell = n) \propto \exp \left(\sum_{f=1}^F \left(\log |\boldsymbol{\Sigma}_{f\ell n}^s| + \frac{\mathbf{x}_{f\ell}^H}{v_f} \mathbf{A}_f \mathbf{D}_n \boldsymbol{\mu}_{f\ell n}^s \right) \right),$$

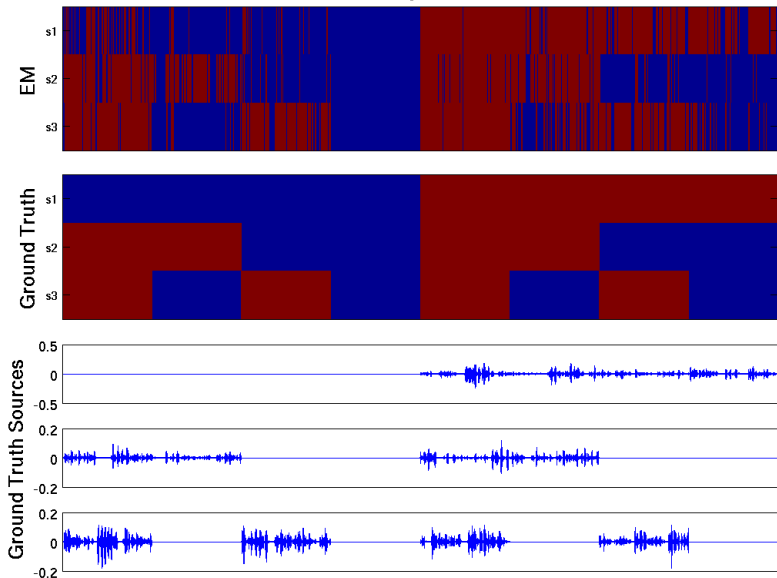
transition T_{nr} , and beginning-of-time probabilities λ_n .

- Closed form solution via the forward-backward algorithm.

M-step Outline

- The parameters $\mathbf{A}_f, \mathbf{v}_f$ updated in closed form: Typical rules for the Gaussian.
- T_{nr}, λ_n updated in closed form: provided for free by the forward-backward algorithm.
- $u_{j,f\ell}$ is composed by two sets of parameters $\{w_{fk}\}_{f,k}, \{h_{k\ell}\}_{k,\ell}$: Typical (closed form, alternating) updates for the NMF.

Estimated Activity on a Mix of $J = 3$ Sources via EM



Some Quantitative (Continuously active)

Comparisson of Separation performance (dB) on a 2×3 mix of continuously emmiting sources:

	Proposed			Ozerov & Févotte '10		
	SDR	SIR	SAR	SDR	SIR	SAR
<i>s1</i>	9.2	13.4	13.6	9.3	13.8	14.0
<i>s2</i>	7.1	15.2	13.4	7.1	14.5	14.1
<i>s3</i>	9.6	14.0	13.9	9.6	13.4	14.6

Thank you !