Source Separation Methods for under-determined sound mixtures

Mathieu Lagrange

Analyse / Synthèse Team, IRCAM

Mathieu.lagrange@ircam.fr

ircam Centre Pompidou

ATIAM 10

Analysis of Sound Mixture

- We aim at performing
 - o Auditory Scene Analysis
 - Computationally
 - o But like human do
 - o Humans focus on one source
- Task
 - Source separation ?
 - Source classification ?
 - o Something in-between ?
 - o What then ?

Computational ASA (CASA)

- How do people analyze sound mixtures ?
 - o break mixture into small elements (in time-freq)
 - o elements are grouped in to sources using cues
 - o sources have aggregate attributes



I. Frequency Analysis (FA)

- Fourier based analysis
 - o The Short-Term Fourier Transform (STFT)
 - o By far the most widely used



(Fig. from Aphex Twin)

I. Frequency Analysis (FA)

- Perception inspired front-ends
 - Like the Correlogram
 - Designed to imitate what is known about the physiology of the inner ear
 - o Usually composed of
 - o A cascade of filterbanks
 - o Interleaved with non linear operators



(Fig. from [McDermott I I])

How to use FA for grouping ?

- Source Separation: a masking problem
- Goal: find a mask M that retrieves one source when used to filter a given time-frequency representation.



- What about the phase ?
 - Keep the one of the mixture

The Ideal Binary Mask (IBM)

- The IBM
 - o Is an "oracle" separation method, that is we know something (everything ?) we need for separating the sources.
- It provides
 - An upper bound for masking based approaches
 - o A way to understand issues with the front end
 - o Time/frequency resolution tradeoff
 - o Issues with the phase

Demonstration of the IBM



- Utterance: "That noise problem grows more annoying each day"
- Interference: Crowd noise with music (0 SNR)

2. Cues (Binaural Case)

- Have spatial location cues
 - o Termed Interchannel or Interaural
 - Phase and Intensity Differences: IPD and IID
 - Warning: profesionaly mastered audio does not preserve them.



- DUET (Degenerate Unmixing Estimation Technique) [Yilmaz&Rickard04]
 - o Histogram of IPD and IID
 - o Binary Mask created by selecting bins around histogram peaks.



(Fig. from [Yilmaz&Rickard04])

[Yilmaz&Rickard04] Ö.Yilmaz and S. Rickard. Blind Separation of Speech Mixtures via Time-Frequency Masking. IEEE Trans. on Signal Processing. Vol. 52(7), July 2004

2. Cues (Binaural Case)

- Human-assisted time-frequency masking [Vinyes06]
 - Human-assisted selection of the time-frequency bins out of the DUETlike histogram for creating the unmixing mask
 - o Implementation as a VST plugin ("Audio Scanner")



[Vinyes06] M.Vinyes, J. Bonada and A. Loscos. Demixing Commercial Music Productions via Human-Assisted Time-Frequency Masking. *120th AES convention*, Paris, France, 2006.

2. Cues (Monaural case)

- Most ASA cues can be considered
- But the most important cue is pitch



2. Cues (Monaural case)

Filterbank output

`Ideal' segmentation

Pitch candidates

Pitch tracking

Harmonic fragments



(Fig. from [Barker 11])

3. Grouping

- Bottom up approaches
 - o Statistical (Blind) approaches (NMF)
 - o Clustering approaches based on ASA cues (CASA)
- Top down approaches
 - o Model based approach
 - o Dictionary based approach
- Combination between the two
 - o Model based approach

Nonnegative Matrix Factorization (NMF)

• Given a nonnegative matrix V of dimensions FxN, NMF is the problem of finding a factorization

$V \approx WH$

- where W and H are nonnegative matrices of dimensions FxK and KxN, respectively.
- Use for transcription:
 - P. Smaragdis and J.C. Brown. Non-Negative Matrix Factorization for Polyphonic Music Transcription. Proc. IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA), New Paltz, USA, 2003.
- Use for separation:
 - o B. Wang and M. D. Plumbley. Musical Audio Stream Separation by Non-Negative Matrix Factorization. Proc. UK Digital Music Research Network (DMRN) Summer Conf., 2005.

NMF

• Along VQ, PCA or ICA, NMF provides an unsupervised linear representation of data



Mathieu Lagrange.

Statistical Tools for Audio Processing.

NMF for Vision

• By representing signals as a sum purely additive, non- negative sources, we get a parts-based representation [Lee'99]





Lee and Seung, Learning the parts of objects by nonnegative matrix factorization, Nature, 1999, 41

Update Rules for NMF

- Multiplicative (Lee & al)
 - o Minimize a cost function with positivity constraints

$$||A - B||^{2} = \sum_{ij} (A_{ij} - B_{ij})^{2}$$

o Update Rules

$$H_{a\mu} \leftarrow H_{a\mu} \frac{(W^T V)_{a\mu}}{(W^T W H)_{a\mu}} \qquad W_{ia} \leftarrow W_{ia} \frac{(V H^T)_{ia}}{(W H H^T)_{ia}}$$

- Theorem: under the update rules, the cost function is
 - o Non increasing
 - o Invariant iif @ stationary point

[Lee'01] Lee and Seung, Algorithms for Non-negativeMatrix Factorization, Nips, 2001

ICA on spectrograms



NMF for Audio



CASA

- How can we use the different cues ?
 - Earlier approach: consider the cues in sequence.
 - Sequentiality is brittle due to the propagation of errors
- All at once



Top down approaches

- Prior knowledge can be represented as an abstract model of some events of interest
 - o Recognition:
 - Example: GMM models of spoken digits like in speech recognition
 - In this case, the background can be dealt with numerous approaches
 - Noisy training
 - Multi-condition training
 - o Separation:
 - Example: separation of the singing voice in a music signal
 - Need model for
 - the singing voice
 - The music

(Fig. from [Barker 11])

GMM – Based Source Separation

• Given a mixture

$$x(n) = v(n) + m(n)$$

• Represented in the spectral domain

$$X_t(f) = V_t(f) + M_t(f)$$

• Following simple algebra



Statistical Tools for Audio Processing.

GMM – Based Source Separation



(Fig. from [Ozerov 05])

GMM – Based Source Separation



Statistical Tools for Audio Processing.

Combining Bottom-up and Top-Down

- Combining bottom up and top down approaches is
 - o the dream goal
 - o Is difficult



Combining Bottom-up and Top-Down

- One good example
 - Fragment-based spoken digit decoding
 - A simple (but terribly inefficient) implementation:



To summarize



Live coding in Matlab

- You can find the source here:
 - <u>http://recherche.ircam.fr/equipes/analyse-synthese/lagrange/teaching/atiam11/</u> <u>coursAtiam20111bm.m</u>
 - <u>http://recherche.ircam.fr/equipes/analyse-synthese/lagrange/teaching/atiam11/</u> <u>coursAtiam2011Nmf.m</u>
- You will need some external dependencies, web locations are provided in the code
- The code uses cell mode, please look at the Matlab documentation for usage

CASA for singer similarity

- Aim: discover an application of CASA for MIR
- Testbed: Music similarity by singer
 - o 2 songs are defined as similar if they have the same lead-singer
 - Evaluation metric : ranking
 - First method:
 - o Extract some features from the spectral representation of the songs
 - o Compare them
 - o Check if the closest ones are from the same singer
 - Problem: even though the lead singer is prominent, the spectral properties of the observed signal are most of the time a non linear combination of the singer and the accompaniment.
 - Question: can we use some knowledge about ASA to minimize the impact of the accompaniment ?

CASA for singer similarity

- Assumptions:
 - The accompaniment does not change throughout the song
 - The singer starts singing at about I minute
- Proposed approach
 - o Model the accompaniment as the audio signal of the beginning of the song
 - o Model the singing voice as the audio signal around 1 minute
 - o Compare songs represented as
 - o spectral features
 - o MFCC's
- Binary Masking:
 - Only consider spectral bins where amplitude of the mixture is larger than the accompaniment model.

CASA for singer similarity

- Dealing with missing data
 - Marginalization: only consider the non-zero spectral components during comparison
 - o Loose a lot of data when many zeros are present
 - o Feature representation is less powerful (can't use MFCCs)
 - o Imputation: replace zero values by default ones
 - o Can use any feature representation
 - o What are the default values to consider ?