

Single-Voice Separation From Monaural Recordings using robust principal component analysis

report by :

Mohsen Nabian

April 2017

- goal: separating voices from music accompaniment

- assumption of this paper:

Music accompaniment \rightarrow low rank subspace
Singing voices \rightarrow relatively sparse within songs.

- Technique :

we use RPCA (Robust principal component analysis)
which is a matrix factorization algorithm for

Solving underlying low-rank and sparse matrices
① ②

- mathematical formulation :

$$\begin{cases} M \in \mathbb{R}^{n_1 \times n_2} \rightarrow \text{original song} \\ L \in \mathbb{R}^{n_1 \times n_2} \rightarrow \text{accompaniment music} \\ S \in \mathbb{R}^{n_1 \times n_2} \rightarrow \text{singing voice} \end{cases}$$

minimize $\|L\|_* + \lambda \|S\|_1$

s.t. $L + S = M$

use $\lambda_k = \frac{k}{\sqrt{\max(n_1, n_2)}}$

with different values of k and test which k is better.

- first we compute spectrogram of music signals as matrix M , calculated from short-time-Fourier transform (STFT)
- second, use ALM, an efficient algorithm for solving SPCA problem: $L + S = |M|$
- then we have found L and S separately.
- also to obtain waveforms back, we record of L and S

the phase of original signal: $P = \text{phase}(M)$

then append the phase to matrix L and S .

$$\begin{cases} L(m, n) = L e^{jP(m, n)} \\ S(m, n) = S e^{jP(m, n)} \end{cases}$$

Time-Frequency Masking

after finding L and S , we apply binary time-frequency masking for better separation.

binary-time frequency masking M_b as follows:

$$M_b(m, n) = \begin{cases} 1 & |S(m, n)| > \text{gain} * |L(m, n)| \\ 0 & \text{otherwise} \end{cases}$$

for all $m = 1, \dots, n_1$
 $n = 1, \dots, n_2$

So \Rightarrow

with mask

$$\begin{cases} X_{\text{singing}}(m, n) = M_b(m, n) M(m, n) \\ X_{\text{music}}(m, n) = (1 - M_b(m, n)) M(m, n) \end{cases}$$

without mask

$$\begin{cases} X_{\text{singing}}(m, n) = S \\ X_{\text{music}}(m, n) = L \end{cases}$$

3/

Experiment

~~1) The effect of λ~~

- MIR-7K dataset
- 1000 song clips Karaoke + amateur singers
- sample rate 16kHz duration 4 to 13 sec.

3 sets of mixtures:

for each clip, the singing voice and the music accompanied were mixed at -5, 0, and 5 dB SNRs.

energy music is larger Same energy level energy singing is larger

evaluations: } Source to interference ratio (SIR)
 } Source to artifacts " (SAR)
 } Source to distortion " (SDR)

$$NSDR(\hat{v}, v, x) = SDR(\hat{v}, v) - SDR(x, v)$$

\downarrow \downarrow \downarrow
singing voice original clean singing voice Preprocessed mixture

$$GNSDR(\hat{v}, v, x) = \frac{\sum_{n=1}^N \omega_n NSDR(\hat{v}_n, v_n, x_n)}{\sum_{n=1}^N \omega_n}$$

length of nth song.

4/

1) effect of λ_k :

Fig 3 : $\left. \begin{array}{l} \text{with no mask} \\ \text{with ~~no~~ binary mask} \end{array} \right\}$

$$\lambda_k = \{0.1, 0.5, 1, 1.5, 2, 2.5\}$$

$$\text{SNR} = \{-5, 0, 5\}$$

in all these cases $\left\{ \begin{array}{l} \text{SDR} \\ \text{SAR} \\ \text{SIR} \end{array} \right\}$ is reported.

2) The effect of gain factor with a binary mask

Fig 4 : Cases : $\left\{ \begin{array}{l} \lambda_1 = \{0.1, 0.5, 1, 1.5, 2\} \\ \text{SNR} = \{-5, 0, 5\} \end{array} \right\}$

in all these cases $\left\{ \begin{array}{l} \text{SDR} \\ \text{SAR} \\ \text{SIR} \end{array} \right\}$ are reported



Conclusion: higher λ_1 , results in lower power sparse matrix S . \Rightarrow larger interference
fewer artifacts

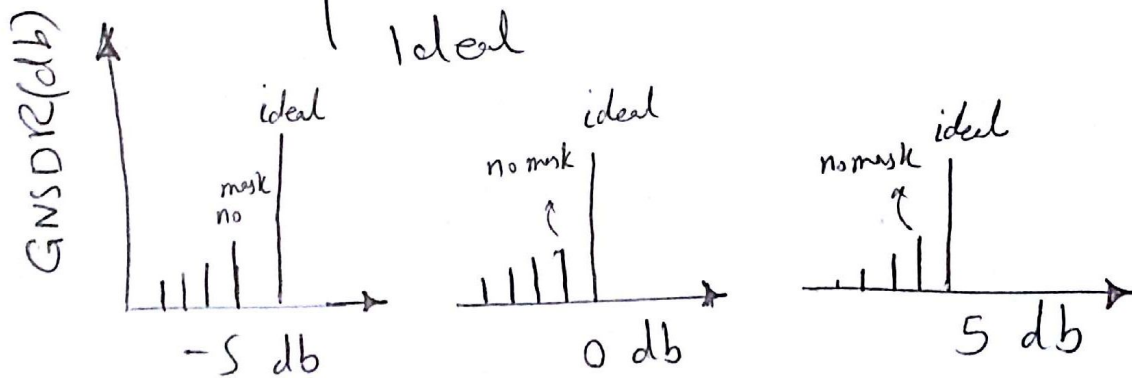
51

3) Comparison with prev. systems:

5 Cases are compared for each 3 Conditions

Conditions: $SNR = \{-5, 0, 5\}$ $\xrightarrow{\text{voice to music ratio}}$

Cases: $\left\{ \begin{array}{l} \text{Hsu algorithm} \\ \text{Rafii} \\ \text{binary mask } \lambda_1, \text{ gain}=7 \\ \text{no mask} \\ \text{Ideal} \end{array} \right.$



no mask provided the best result out of others.