

SINGLE CHANNEL SPEECH-MUSIC SEPARATION USING MATCHING PURSUIT AND SPECTRAL MASKS

Emad M. Grais and Hakan Erdogan

Faculty of Engineering and Natural Sciences,
Sabanci University, Orhanli Tuzla, 34956, Istanbul.
Email: grais@su.sabanciuniv.edu, haerdogan@sabanciuniv.edu

ABSTRACT

A single-channel speech music separation algorithm based on matching pursuit (MP) with multiple dictionaries and spectral masks is proposed in this work. A training data for speech and music signals is used to build two sets of magnitude spectral vectors of each source signal. These vectors' sets are called dictionaries, and the vectors are called atoms. Matching pursuit is used to sparsely decompose the magnitude spectrum of the observed mixed signal as a nonnegative weighted linear combination of the best atoms in the two dictionaries that match the mixed signal structure. The weighted sum of the resulting decomposition terms that include atoms from the speech dictionary is used as an initial estimate of the speech signal contribution in the mixed signal, and the weighted sum of the remaining terms for the music signal contribution. The initial estimate of each source is used to build a spectral mask that is used to reconstruct the source signals. Experimental results show that integrating MP with spectral mask gives good separation results.

Index Terms— Source separation, single channel source separation, speech music separation, speech processing, and Matching pursuit.

1. INTRODUCTION

Single channel audio source separation is an interesting and hard research problem. There are various studies focusing on separating source signals from a mixture of multiple speakers [1], multiple musical instruments [2], speech and music [3], or speech and noise [4]. In this paper, we focus on separating speech signals from background music signals.

Single channel source separation aims to separate the original source signals from only one observed mixture of these source signals. Since only single mixed signal of these source signals is available, the separation techniques usually rely on prior knowledge which is training data of each original source signal that are in the mixture. These training data are used to build a representative model for each source that can capture the characteristics of these source signals. These models can be probabilistic distributions like Gaussian mixture models (GMM) [5], or HMMs [1]. It can be also a model of a set of trained basis vectors [4]. The idea of decomposing the mixed signal with a set of trained bases for each source signal was used in previous works, and the estimate of each source signal is found by grouping the decomposition results that are related to each source bases. In [4], Non-negative matrix factorization (NMF) was used to train a set of basis vectors from a training data of each source signal, then NMF is used with these trained bases to decompose the mixed signals. In [6, 3], learned codebooks or dictionaries

that consist of a number of code vectors for the source signals were used. Different optimization algorithms were used to find the best combination of codebooks' vectors that can explain the observed mixed signal. Usually the codebooks are learned using any clustering technique like k-means or using dictionary learning algorithms to optimally represent each source.

In this paper, matching pursuit and spectral masks are used to separate a mixture of speech and music signals. Given sufficient training data from each source signal we propose that, the best model of these data is the data themselves as claimed in [7]. These training data for each source are used to build two dictionaries of speech and music magnitude spectral (MS) vectors or "atoms". After observing the mixed signal, matching pursuit is used to sparsely decompose the MS of the mixed signal as a weighted linear combination of the dictionaries' atoms for each source. The weighted sum of the terms that include atoms from speech dictionary in the decomposition result is used as an initial estimate of the MS of the speech signal. The weighted sum of the remaining decomposition terms that include atoms from music dictionary is used as an initial estimate of the MS of the music signal. These initial estimates of the sources are used to build different spectral masks. These masks are used to find a new estimate of each source signal from the mixed signal. The novelty in this work is in many aspects. First, we build source specific dictionaries instead of using dictionaries of Gabor atoms that represent every source as in [8]. Second, our separation algorithm works in the frequency domain rather than in the time domain as in [8], because spectral power or magnitude is considered to be a more powerful representation of audio signals than the time domain representation. Third, we consider the linkage and smoothness between the consequent frames in our decomposition algorithm, that is instead of decomposing one MS frame every time we stack a number of consequent frames in one super vector, then every super vector is shifted by one frame at a time. This gives us a chance to decompose every frame multiple times with different neighbor frames, and then average the results. Fourth, combining the spectral mask with MP to get a better separation of the mixed signal.

The rest of this paper is organized as follows: In section 2, a mathematical description of the problem is given. In section 3, a brief explanation about matching pursuit and how we built our dictionaries is described. In section 4, an explanation of applying matching pursuit with different spectral masks in source separation is given. The remaining sections are for the results of our experiments and our observations.

2. PROBLEM FORMULATION

Single channel speech-music separation problem can be formulated as follows: Assume we have a single observed signal $y(t)$, which is the mixture of two sources speech $x(t)$ and music $m(t)$. The source separation problem aims to find estimates for $x(t)$ and $m(t)$ from $y(t)$. The framework in this paper is in the short time Fourier transform (STFT) domain. Let $Y(t, f)$ be the STFT of $y(t)$, where t represents the frame index and f is the frequency-index. Due to linearity of the STFT, we have:

$$Y(t, f) = X(t, f) + M(t, f), \quad (1)$$

$$|Y(t, f)| e^{j\phi_Y(t, f)} = |X(t, f)| e^{j\phi_X(t, f)} + |M(t, f)| e^{j\phi_M(t, f)}. \quad (2)$$

In this work, we assume the sources have the same phase angle as the mixed signal for every frame, that is $\phi_Y(t, f) = \phi_M(t, f) = \phi_X(t, f)$. This assumption was shown to yield good results in earlier work. Hence, we can write the magnitude spectrum of the measured signal as the sum of source signals' magnitude spectra.

$$M_y(t, f) = M_x(t, f) + M_m(t, f). \quad (3)$$

Here $M_x(t, f)$ and $M_m(t, f)$ are unknown magnitude spectra, and need to be estimated using observed mixed signal and training speech and music signals. The magnitude spectrum for the observed signal $M_y(t, f)$ is obtained by taking the magnitude of the DFT of the windowed signal.

3. SIGNAL DECOMPOSITION USING MATCHING PURSUIT

Assume a set of training data for speech and music signal is available. A speech dictionary D_s ¹ which contains normalized MS vectors or atoms g_s of the training speech signal is built, and a music dictionary D_m with MS atoms g_m is also built using the music training data. The main idea in this work is to sparsely decompose the MS of the mixed signal y as a weighted linear combination of the speech atoms g_s and the music atoms g_m according to the contribution of every source in the mixed signal.

$$y \approx \underbrace{\sum_{D_s} c_s g_s}_{\text{speech part}} + \underbrace{\sum_{D_m} c_m g_m}_{\text{music part}}, \quad (4)$$

in matrix form

$$y \approx D_s c_s + D_m c_m. \quad (5)$$

The coefficients' vectors c_s and c_m of the speech and music parts respectively are unknowns and need to be calculated and enforced to be sparse. Therefore, we need to solve the following L_0 norm problem:

$$c_s, c_m = \arg \min_{\hat{c}_s, \hat{c}_m} (\|\hat{c}_s\|_0 + \|\hat{c}_m\|_0). \quad (6)$$

Subject to

$$\begin{aligned} y &= D_s c_s + D_m c_m, \\ D_s c_s &\geq 0, \\ D_m c_m &\geq 0. \end{aligned}$$

Where $\|c\|_0$ counts the number of non-zero entries in vector c . The exact solution to (6) is not easy to find. However, a simple iterative and greedy algorithm like matching pursuit can find good approximations.

¹The notations here are as follows: bold capital letters are for matrices, bold small letters are for vectors others are for scalars.

3.1. Matching pursuit

Matching pursuit [9] is an algorithm that approximates any signal or vector by decomposing it into a weighted linear combination of a set of basis elements called atoms. MP iteratively picks out the best atoms that can match the structure of the signal from a dictionary which contains a huge number of atoms in a greedy fashion. Let the dictionary $D = [D_s, D_m]$, an atom $g_n \in D$, and given a signal y in \mathbb{R}^d space, MP tries to find a good approximation \hat{y} as a linear combination of N atoms g_n selected from dictionary D .

$$\hat{y} = \sum_{n=0}^{N-1} c_n g_n \quad g_n \in D, \quad (7)$$

such that

$$\left\| y - \sum_{n=0}^{N-1} c_n g_n \right\|_2 < \epsilon.$$

For sparse solution $N \ll d$. Let $R^0 y = y$, the vector $R^0 y$ can be decomposed in the first iteration into

$$R^0 y = \langle g_0, R^0 y \rangle g_0 + R^1 y. \quad (8)$$

Where $R^1 y$ is the residual vector after approximating y in g_0 's direction, $\langle g_0, R^0 y \rangle$ is the dot product, and g_0 is atom chosen to maximize the correlation with $R^0 y$ as

$$g_0 = \arg \max_{g \in D} |\langle g, R^0 y \rangle|.$$

The residue $R^1 y$ in the second iteration is decomposed by projecting it on the atom that best matches it from D , and a new residue $R^2 y$ is obtained. This procedure is repeated until the norm of the residual is less than a predefined threshold or when the maximum number of iterations is achieved. After N iterations, the following approximation is constructed:

$$y = \sum_{n=0}^{N-1} \langle g_n, R^n y \rangle g_n + R^N y, \quad (9)$$

and the norm of the final residual is

$$\|R^N y\|^2 = \|y\|^2 - \sum_{n=0}^{N-1} \langle g_n, R^n y \rangle^2.$$

Since we are working with atoms and mixed signals vectors in the MS domain, we should make sure that the residual vector's components after every iteration $n \in [0, N-1]$ are nonnegative.

$$R^n y = [R^n y]_+.$$

Where $[R^n y]_+ = \max[R^n y, 0]$. One of the advantages of matching pursuit is the complete freedom in designing the dictionary, but the performance highly depends on the structure of the dictionary.

3.2. Building the dictionaries

The magnitude spectrograms of the whole available training data for speech and music signals are found. We build two different dictionary matrices one for each source. Every column or atom g_s in the speech dictionary matrix D_s is formed by stacking its corresponding frame with its "2L" surrounding frames from the speech spectrogram in one super vector. Therefore, every column vector in D_s is

built by stacking “ $2L + 1$ ” neighbor frames from the speech spectrogram in one vector. For example, the column number l in the speech dictionary matrix D_s is

$$\mathbf{g}_s(l) = \left[\mathbf{m}_s^T(l-L), \dots, \mathbf{m}_s^T(l), \dots, \mathbf{m}_s^T(l+L) \right]^T.$$

Where $\mathbf{m}_s(l)$ is the magnitude spectrum column vector corresponding to frame number l of the training speech signal spectrogram. A mirror imaging at the edges of the spectrogram is performed. Any training frame which has energy less than a threshold is removed. The music dictionary matrix D_m is built from the music spectrogram in the same way. Then all atoms in both dictionaries are normalized. The two dictionaries matrices D_s and D_m are used to build one big dictionary matrix $D = [D_s, D_m]$.

4. SIGNAL SEPARATION

In this section, we explain the proposed method for separating the mixed signal. We use matching pursuit and multi-dictionaries to decompose the mixed signal. Then the initial estimates of the underlying source signals in the mixed signal are found by grouping every source’s atoms from the decomposition result, and the summation for every group is taken separately. These initial estimates are used to build different masks that define at each time-frequency component the ratio of every source signal in the mixed signal. The mask is used to scale the mixed signal in the STFT domain to find a smooth estimate of every source signal in the mixture.

4.1. Signal decomposition using matching pursuit

After observing the mixed signal, the spectrogram of the mixed signal is computed. A matrix Y of the observed mixed MS’s frames is built. Every column of this matrix is formed by stacking its corresponding frame with its surrounding “ $2L$ ” frames from the mixed signal spectrogram in one super vector. This is analogous to the way we constructed the dictionaries’ matrices in section 3.2. For example, the column number l of the matrix Y is

$$\mathbf{y}(l) = \left[\mathbf{m}_y^T(l-L), \dots, \mathbf{m}_y^T(l), \dots, \mathbf{m}_y^T(l+L) \right]^T,$$

where $\mathbf{m}_y(l)$ is the frame number l of the observed mixed signal spectrogram. We need to decompose each column vector \mathbf{y} in the matrix Y with the best atoms that match its structure from the dictionary matrix D . We build a matrix X for the initial estimated separated speech MSs and a matrix M for the initial estimated separated music MSs.

$$Y \approx X + M.$$

For every column \mathbf{y} , \mathbf{x} , and \mathbf{m} in the matrices Y , X , and M the separation algorithm works as follows:

1. $n = 0$, $\mathbf{x} = \mathbf{m} = \mathbf{0}$, and let $\mathbf{R}^0 \mathbf{y} = \mathbf{y}$.
2. Project $\mathbf{R}^n \mathbf{y}$ at iteration n into the dictionary $D = [D_s, D_m]$.
3. Find the atom $\mathbf{g}_n \in D$ that gives maximum dot product

$$\mathbf{g}_n^* = \arg \max_{\mathbf{g} \in D} (\mathbf{R}^n \mathbf{y})^T \mathbf{g},$$

$$c_n^* = (\mathbf{R}^n \mathbf{y})^T \mathbf{g}_n^*.$$

4. If $\mathbf{g}_n^* \in D_s$

$$\mathbf{x} = \mathbf{x} + c_n^* \mathbf{g}_n^*.$$

Elseif $\mathbf{g}_n^* \in D_m$

$$\mathbf{m} = \mathbf{m} + c_n^* \mathbf{g}_n^*.$$

5. Remove \mathbf{g}_n^* from the dictionary D .

6. $\mathbf{R}^{n+1} \mathbf{y} = [\mathbf{R}^n \mathbf{y} - c_n^* \mathbf{g}_n^*]_+$.

7. If $\frac{\|\mathbf{R}^{n+1} \mathbf{y}\|_2^2}{\|\mathbf{R}^0 \mathbf{y}\|_2^2} > \epsilon$ and $n < N - 1$

$n = n + 1$; go to step 2.

else stop.

Where N is the maximum allowed number of iterations.

As shown in step 4 in the previous algorithm, the initial estimates for every speech and music MS vectors are found by finding the weighted sum of the decomposition terms that include atoms from speech dictionary and music dictionary respectively as follows:

For speech part, we got

$$\mathbf{x} \approx \sum_k c_s(k) \mathbf{g}_s(k) \quad \mathbf{g}_s \in D_s, \quad (10)$$

and for music part

$$\mathbf{m} \approx \sum_j c_m(j) \mathbf{g}_m(j) \quad \mathbf{g}_m \in D_m, \quad (11)$$

where $\mathbf{g}_s(k)$ is the best atom from the speech dictionary D_s that match the structure of the mixed signal vector \mathbf{y} at iteration k and $c_s(k) = \langle \mathbf{g}_s(k), \mathbf{R}^k \mathbf{y} \rangle$ is its weight, and $\mathbf{g}_m(j)$ is the best atom from the music dictionary D_m that match \mathbf{y} at iteration $j \neq k$ and $c_m(j)$ is its weight. Applying the previous procedures on all vectors in matrix Y we get the matrix X for the initial estimated separated speech MSs with \mathbf{x} in its columns, and the matrix M for initial estimated separated music MSs with \mathbf{m} in its columns. For example, the column number l in the matrix X is

$$\mathbf{x}(l) = \left[\widetilde{\mathbf{m}}_x^T(l-L), \dots, \widetilde{\mathbf{m}}_x^T(l), \dots, \widetilde{\mathbf{m}}_x^T(l+L) \right]^T,$$

where $\widetilde{\mathbf{m}}_x(l)$ is the frame number l of the initial estimated spectrogram of the separated speech signal from the frame $\mathbf{m}_y(l)$ of the mixed signal spectrogram. Notice that every frame $\widetilde{\mathbf{m}}_x(l)$ is differently estimated $2L + 1$ times in $2L + 1$ different columns \mathbf{x} in matrix X . We find the overall initial estimated spectrograms’ frames $\widetilde{\mathbf{m}}_x$ and $\widetilde{\mathbf{m}}_m$ for each speech and music frame by averaging their corresponding frames in the $2L + 1$ neighbor columns in the matrices X and M respectively.

4.2. Source signals reconstruction and masks

We can directly use the initial estimate spectrograms of the speech and music signals that are found in section 4.1 as the final estimate of every source, but the two estimated spectra $\widehat{M}_x(t, f)$ and $\widehat{M}_m(t, f)$ may not sum up to the mixture $M_y(t, f)$. We usually get nonzero decomposition error. Thus, MP gives us an approximation:

$$M_y(t, f) \approx \widehat{M}_x(t, f) + \widehat{M}_m(t, f).$$

Assuming noise is negligible in our mixed signal, the spectrogram of the source signals’ sum should be directly equal to the spectrogram of the mixed signal. To make the error zero, we use the initial estimate $\widehat{M}_x(t, f)$ and $\widehat{M}_m(t, f)$ to build a mask as follows:

$$H(t, f) = \frac{\widehat{M}_x^p(t, f)}{\widehat{M}_x^p(t, f) + \widehat{M}_m^p(t, f)}, \quad (12)$$

where $p > 0$ is a parameter. Notice that elements of $\mathbf{H} \in (0, 1)$. Using different p values leads to different kinds of masks. When $p = 2$ the mask $H(t, f)$ is a Wiener filter. The value of p controls the saturation level of the ratio in (12). When $p > 1$, the larger component will dominate more in the mixture. At $p = \infty$, we achieve a binary mask (hard mask) which will choose the larger source component as the only component. These masks will scale every frequency component in the observed mixed signal with a ratio that explains how much each source contributes in the mixed signal such that

$$\hat{X}(t, f) = H(t, f)Y(t, f), \quad (13)$$

where $\hat{X}(t, f)$ is the final STFT estimate of the speech signal. After finding the contribution of the speech signal in the mixed signal, the estimated speech signal $\hat{x}(t)$ can be found by using inverse STFT of $\hat{X}(t, f)$.

5. EXPERIMENTS AND DISCUSSION

The proposed algorithm is applied on simulated mixtures of speech and music data at 16kHz sampling rate. For training speech data, 540 short utterances from a single speaker was used. We left out 20 utterances for testing. For piano music data, piano music from piano society web site [10] was downloaded. We used 38 pieces from different composers but from a single artist for training and left out one piece for the testing stage. The MS dictionaries for speech and music data were trained using the STFT, a Hamming window was used, and the FFT was taken at 512 points. We form the dictionaries' atoms as we mentioned before, we find MS from every FFT frame by using the first 257 points only since the remaining points are the conjugate of the first 257 points. Then we concatenated every five ($L = 2$) MS frames in one column vector with size (5×257) as we have mentioned in section 3.2. Each vector in the speech and music dictionaries is in 1285 dimensions (5×257). The test data was formed by adding random portions of the test music file to the 20 speech utterance files at different speech to music ratio (SMR) values in dB. The audio power levels of each file were found using the "audio voltmeter" program from the G.191 ITU-T STL software suite [11]. For each SMR value, we obtained 20 test utterances this way.

The reason for working with training and testing vectors of five MS frames at a time is that we got remarkable improvement by working that way rather than working with a single frame. It is obvious that this will slow down the separation algorithm. We tried to work with concatenating ten frames but the improvement was not noticeable compared to concatenating five frames, so we worked with five frames at a time for memory capacity and speed.

Performance measurement of the separation algorithms was done using metrics introduced in [12]. Projection of the predicted signal onto the original speech signal is termed as the target signal. Source distortion ratio (SDR) is defined as the ratio of the target energy to all errors in the reconstruction.

Table 1 shows the performance results for the estimated speech signal by using MP without masks and the performance of using MP with different kinds of masks. In the case of using MP without masks, we directly use the initial spectrogram estimate of speech signal that was found in section 4.1 as the final estimate of the speech source. We get better results in the case of using MP with spectral masks. However, we should not use the hard mask in the case of low speech to music ratio. The table also shows the fact that, when the speech signal dominates more in the mixed signal, then it is better to use the mask with larger p .

6. CONCLUSION

In this work, we studied single channel speech-music separation using matching pursuit with multiple dictionaries and spectral masks. We built dictionaries from training speech and music "piano" signals. The dictionaries, matching pursuit, and spectral masks were used to separate the mixed signal. To speed up the proposed algorithm and get the ability to use bigger dictionaries, our next work will focus on combining matching pursuit with kd-tree or working with tree based matching pursuit [13].

Table 1. Source/Distortion Ratio SDR in dB for speech signal of using matching pursuit(MP) with different spectral masks.

| SMR dB | No mask | p=1 | p=2 | p=3 | Hard mask |
|--------|---------|--------------|--------------|--------------|-----------|
| -5 | 2.86 | 3.38 | 3.23 | 2.90 | 2.39 |
| 0 | 7.13 | 7.92 | 7.80 | 7.56 | 6.97 |
| 5 | 9.83 | 10.99 | 10.90 | 10.68 | 10.17 |
| 10 | 13.59 | 15.74 | 16.00 | 15.87 | 15.45 |
| 15 | 14.72 | 17.53 | 17.89 | 17.75 | 17.30 |
| 20 | 16.32 | 20.77 | 21.99 | 22.05 | 21.79 |

7. REFERENCES

- [1] J. R. Hershey, S. J. Rennie, P. A. Olsen, and T. T. Kristjansson, "Super-human multi-talker speech recognition: A graphical modeling approach," *Computer Speech and Language*, vol. 24, no. 1, pp. 45–66, Jan. 2010.
- [2] T. Virtanen, "Monaural sound source separation by non-negative matrix factorization with temporal continuity and sparseness criteria," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 15, pp. 1066–1074, Mar. 2007.
- [3] Hakan Erdogan and Emad M. Grais, "Semi-blind speech-music separation using sparsity and continuity priors," in *International Conference on pattern recognition (ICPR)*, 2010.
- [4] K. W. Wilson, B. Raj, P. Smaragdis, and A. Divakaran, "Speech denoising using nonnegative matrix factorization with priors," in *Proc. of ICASSP*, 2008.
- [5] Aarthi M. Reddy and Bhiksha. Raj, "Soft Mask Methods for single-channel speaker separation," *IEEE Trans. of Audio, Speech, and Language Processing*, vol. 15, Aug. 2007.
- [6] L. Benaroya, L. M. Donagh, F. Bimbot, and R. Gribonval, "Non negative sparse representation for Wiener based source separation with a single sensor," in *Proc. IEEE Conf. Acoust. Speech Sig. Proc.*, 2003.
- [7] Smaragdis P., M. Shashanka, and B. Raj, "A sparse non-parametric approach for single channel separation of known sounds," in *Neural Information Processing Systems, Vancouver, BC, Canada*, Dec. 2009.
- [8] N. Cho, Yu. Shiu, and Jay. Kuo, "Audio source separation with matching pursuit and content-adaptive dictionaries (mp-cad)," *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, pp. 287–290, Oct. 2007.
- [9] S.G. Mallat and Z. Zhang, "Matching pursuit with time frequency dictionaries," *IEEE Trans. Signal Process.*, vol. 41, pp. 3397–3415, Dec. 1993.
- [10] URL, "http://pianosociety.com," 2009.
- [11] URL, "http://www.itu.int/rec/T-REC-G.191/en," 2009.
- [12] E. Vincent, R. Gribonval, and C. Fevotte, "Performance measurement in blind audio source separation," *IEEE Tr. Acoust. Sp. Sig. Proc.*, vol. 14, no. 4, pp. 1462–69, July 2006.
- [13] P. Jost, P. Vanderghyest, and P. Frossard, "Tree-Based pursuit: Algorithm and properties," *IEEE Trans. Signal Process.*, vol. 54, pp. 4685–4697, Dec. 2006.