

PAPER • OPEN ACCESS

Application of blind source separation in sound source separation

To cite this article: Jiarui Xu 2019 *J. Phys.: Conf. Ser.* **1345** 032006

View the [article online](#) for updates and enhancements.



IOP | ebooks™

Bringing you innovative digital publishing with leading voices
to create your essential collection of books in STEM research.

Start exploring the **collection** - download the first chapter of
every title for free.

Application of blind source separation in sound source separation

Jiarui Xu

School of Electrical Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, 710049, P.R. China

* Author's e-mail: xvjiarui@126.com

Abstract. The classic method for solving the cocktail party problem is by utilizing Independent Component Analysis (ICA) method to separate different sounds. Since ICA method deals with each frequency point of the audio signal individually, there is a classic sorting problem. However, for Independent Vector Analysis (IVA), which makes use of the correlations between different frequency points of the same sound source, the audio of the same sound source can be isolated at one time without sorting problems. For music signal, due to the rhythm of music, there may be a strong correlation between different frequency points of the same sound source, so IVA method could be used for signal separation. This paper will discuss the application of IVA method in music sound separation and its application in music related speech recognition system from three aspects: basic principle, code implementation and performance analysis.

1. Introduction

This paper briefly describes the problems existing in speech cocktail party, introduces ICA method of blind source signal separation [1][2], briefly describes the disadvantages of ICA, and further introduces IVA method [3][4], then briefly describes the application of blind source separation in sound signal, wireless signal, brain electrical signal, and emphatically describes its application requirements and application ideas in music separation.

At a noisy cocktail party, there are many different sources of sound at the same time: the sound of multiple people talking at the same time, the clatter of cutlery, the sound of music, and so on. It's pretty easy for us as humans to figure out how to pick out a particular person's voice at a cocktail party. But for computers, there are still a lot of tricky problems to solve when it comes to breaking up an audio signal into different speech sources. This is the famous "cocktail party" problem.

If multiple microphones are used to collect these sound signals, the machine only knows the received mixed signal and has no other information about the signal source. We classify the problem of only separating different sources from the received mixed signal without knowing the signal source information as 'blind source separation' [5]. The typical method in blind source separation is Independent Component Analysis (ICA). However, ICA only processes each frequency point of the received signal separately, so there is a classical sorting problem between the separated signals. The Independent Vector Analysis (IVA) method can make use of the correlation between different frequency points of each signal source, so the sorting problem in ICA can be avoided [6]. In particular, for sound signals with strong rhythm such as music signal, there is a strong correlation between different frequency points of the same signal, so it is very suitable for IVA method for separation.



Firstly, this paper introduces the principle and mathematical derivation of two blind source separation methods: ICA and IVA. On this basis, a simulation scenario is built to simulate the separation of two music signals by IVA.

2. Blind source separation

Cocktail party problem is a classic cocktail party problem. Assuming that there are n individuals in a party, they can speak at the same time. We also have n acoustic receivers placed in some corners of the room to record sound. After the banquet, we got a set of data from n microphones:

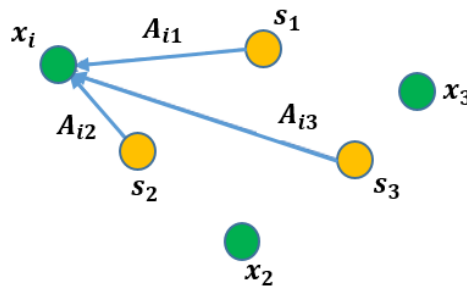
$$\left\{ \mathbf{x}^{(i)} \left(x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)} \right), i = 1, \dots, m \right\}; i \text{ represents the sampling order, which means that there are } m$$

Set of samples, and each set of samples is n -dimensional. Our goal is to identify the signals that each person speaks from the m samples alone.

So let's refine the second problem. There are n sources $\mathbf{s} = (s_1, s_2, \dots, s_n)^T, s \in \mathbb{R}^n$. Each dimension is a person's voice signal. Each person's voice signal is independent. A is an unknown mixture matrix (mixing matrix), used for composite superimposed signal \mathbf{s} , then

$$\mathbf{X} = A\mathbf{s} \quad (1)$$

\mathbf{x} here is not a vector, it's a matrix. Each of these column vectors is $\mathbf{x}^{(i)}, \mathbf{x}^{(i)} = A\mathbf{s}^{(i)}$ so let's graph that



Each component of $\mathbf{x}^{(i)}$ is represented linearly by the component of $\mathbf{s}^{(i)}$, A and \mathbf{s} are unknown, \mathbf{x} is known, and we have to figure out how to derive \mathbf{s} from \mathbf{x} , A process known as blind signal separation.

Let $\mathbf{W} = A^{-1}$, then $\mathbf{s}^{(i)} = A^{-1}\mathbf{x}^{(i)} = \mathbf{W}\mathbf{x}^{(i)}$, change \mathbf{W} as:

$$\mathbf{W} = \begin{bmatrix} w_1^T \\ \vdots \\ w_n^T \end{bmatrix} \quad (2)$$

Among them, $w_i \in \mathbb{R}^n$, just going to write w_i as a row vector. then:

$$s_j^{(i)} = w_j^T \mathbf{x}^{(i)} \quad (3)$$

2.1. Independent component analysis (ICA)

The principle and method of independent component analysis are introduced.

The independent component analysis method assumes that the multidimensional observation signal is composed of several statistically independent components, which are usually called independent sources. From the perspective of statistical analysis and information processing, if there is a physical separation or no information interaction between the two sources, then the two sources are considered to be independent of each other. In the real environment, the imagination of independent sources is common, such as the speech signals emitted by different people in the cocktail party problem, the brain electrical signals generated by different activation areas in the cerebral cortex, and the electromagnetic signals emitted by different antennas. Therefore, the assumption of "independence" has a good universality. In addition, the research object of independent component analysis is random

variable and random signal with non-gauss distribution, which makes the theory of independent component analysis more universal because naturally generated signals are mostly non-gauss. For example, the speech signal generally satisfies the super-gauss distribution, and the image signal is sub-gauss distribution. It is because of the above basic characteristics of independent components that it is regarded as a significant development in the fields of classical statistical analysis and multidimensional signal processing.

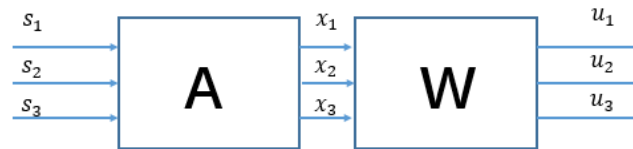


Figure 1 linear ICA model

The basic ICA model is a linear time-invariant instantaneous hybrid model shown in figure 1. Obviously, the linear time-invariant hybrid model is the nonlinear time-variant hybrid model. Instantaneous mixing refers to the delay and convolution effect of the mixed system on the source signal. In practical applications, nonlinear problems of mixed models, time-varying and convolutional effects in signal transmission, and time-varying problems of model parameters all exist to varying degrees. Therefore, in the field of ICA research, nonlinear ICA, convolutional mixed ICA and time-varying mixed ICA also appear. These problems have aroused people's attention since the early days of ICA and have been the hot issues in ICA research in recent years.

ICA algorithm is attributed to Bell and Sejnowski, where maximum likelihood estimation is used to explain the algorithm, and Infomax principal, a complex information maximization method, is used in the original paper.

We assume that each s_i has a probability density of p_s , and then the joint distribution of the original signal at a given moment is

$$p_s(s) = \prod_{i=1}^n p_s(s_i) \quad (4)$$

This formula represents the assumption that each person emits an individual sound signal. With $p(s)$, we can find $p(x)$:

$$p_x(x) = p_s(Wx) |W| = |W| \prod_{i=1}^n p_s(w_i^T x) \quad (5)$$

In the above equation, the left-hand side is the probability of x (n-dimensional vector) of each sampled signal, and the right-hand side is $|W|$ times the product of the probabilities of each original signal.

As mentioned earlier, we cannot find W and s without prior knowledge. So, we need to know $p_s(s_i)$. We're going to pick a probability density function and assign it to s , but we can't pick a gauss density function. In probability theory we know that the density function $p(x)$ is derived by the derivative of the cumulative distribution function (CDF) $F(x)$. $F(x)$ must satisfy two properties: monotonically increasing, and the range is in $[0, 1]$. We found that the sigmoid function works well, for the domain from minus infinity to infinity, for the range from zero to one, slowly increasing. We assume that the cumulative distribution function of s conforms to the sigmoid function:

$$g(s) = \frac{1}{1 + e^{-s}} \quad (6)$$

After the derivation:

$$\begin{aligned}
 p_s(s) = g'(s) &= \frac{e^{-s}}{(1+e^{-s})^2} = \frac{e^{-s}+1-1}{(1+e^{-s})^2} = \frac{e^{-s}+1}{(1+e^{-s})^2} - \frac{1}{(1+e^{-s})^2} \\
 &= \frac{1}{1+e^{-s}} \left(1 - \frac{1}{1+e^{-s}} \right) = g(s)[1-g(s)]
 \end{aligned} \tag{7}$$

This is the density function of s , where s is a real number.

If we knew the distribution function of s in advance, we would not assume it, but in the absence of such a function, sigmoid would work well for most problems.

Because in the above equation is a symmetric function, so $E[s]=0$ (the mean of s is 0), so $E[x]=E[As]=0$, and the mean of x is also 0.

We have known $p_s(s)$, and We need to know W . Given the training sample after sampling $\{x^{(i)}(x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}), i=1, \dots, m\}$ Using the probability density function of x obtained previously; the logarithmic likelihood of samples is estimated as follows:

$$L(W) = \sum_{i=1}^m \left(\sum_{j=1}^n \log p(x) + \log |W| \right) = \sum_{i=1}^m \left(\sum_{j=1}^n \log g'(w_j^T x^{(i)}) + \log |W| \right) \tag{8}$$

And then we're going to take the derivative of W , and one of the things that is involved here is taking the derivative of the determinant $|W|$, which is matrix calculus.

Take the derivative of A determinant, let's say that A is an n by n matrix, and we know that the determinant depends on the cofactors,

$$|A| = \sum_{i=1}^n (-1)^{i+j} A_{ij} |A_{i,\setminus j}|, \quad (\text{for any } j \in 1, \dots, n) \tag{9}$$

In it, $A_{i,\setminus j}$ means the sub-matrix after removing the i row and the j column, so the derivative of $A_{k,l}$ is:

$$\frac{\partial}{\partial A_{k,l}} |A| = \frac{\partial}{\partial A_{k,l}} \sum_{i=1}^n (-1)^{i+j} A_{ij} |A_{i,\setminus j}| = (-1)^{i+j} |A_{i,\setminus j}| = (\text{adj}(A))_{l,k} \tag{10}$$

$\text{adj}(A)$ means the same with A^* as what we learned in linear algebra, so:

$$\nabla_A |A| = (\text{adj}(A))^T = |A| A^{-T} \tag{11}$$

Then:

$$\nabla_w |W| = (\text{adj}(W))^T = |W| W^{-T} \tag{12}$$

In addition, the calculation method of the derivative of $\log g'(s)$ is:

$$\begin{aligned}
 \frac{\partial}{\partial s} \log g'(s) &= \frac{\partial}{\partial s} \log (g(s)[1-g(s)]) = \frac{\partial}{\partial s} \log g(s) + \frac{\partial}{\partial s} \log (1-g(s)) \\
 &= \frac{1}{g(s)} g'(s) - \frac{1}{1-g(s)} g'(s) = 1 - 2g(s)
 \end{aligned} \tag{13}$$

$$\begin{aligned}
 \frac{\partial}{\partial w} \log g'(s) &= \frac{\partial}{\partial s} \log g'(s) \cdot \frac{\partial s}{\partial w} = [1 - 2g(s)] \cdot [x^{(i)}]^T \\
 \frac{\partial s}{\partial w} &= \frac{\partial (w_j^T x^{(i)})}{\partial w} = \frac{\partial ([x^{(i)}]^T w_j)}{\partial w} = [x^{(i)}]^T
 \end{aligned} \tag{14}$$

$$\text{In it, } \frac{\partial s}{\partial w} = \frac{\partial (w_j^T x^{(i)})}{\partial w} = \frac{\partial \left(\left[x^{(i)} \right]^T w_j \right)}{\partial w} = \left[x^{(i)} \right]^T \quad (15)$$

The resulting derivative is the following,

$$W = W + \alpha \begin{pmatrix} 1 - 2g(w_1^T x^{(i)}) \\ 1 - 2g(w_2^T x^{(i)}) \\ \vdots \\ 1 - 2g(w_n^T x^{(i)}) \end{pmatrix} \left[x^{(i)} \right]^T + (W^T)^{-1} \quad (16)$$

Where, α is the gradient rise rate, which is artificially specified.

After iteration and get W , $s^{(i)} = Wx^{(i)}$ can be obtained to restore the original signal.

Note: when we calculate the maximum likelihood estimate, we assume that $x^{(i)}$ and $x^{(j)}$ are independent, but this assumption is not true for speech signals or other time-dependent properties (such as temperature). However, when there is enough data, the assumption of independence has little effect on the effect. Meanwhile, if the sample is disrupted in advance and the random gradient ascending algorithm is run, the convergence speed can be accelerated.

So just to review the cocktail party problem, s is the signal that people send, it's a continuous value, s is different at different points in time, and each person sends signals that are independent of each other (s_i and s_j independent of each other). The cumulative probability distribution function of s is the sigmoid function, but everyone who emits the sound signal conforms to this distribution. A (inverse of W) represents the change in s relative to x , which is the result of the change in s and A .

2.1.1. Uncertainty of ICA

Since w and s are uncertain, it is impossible to determine these two related parameters without prior knowledge. For example, the formula s is equal to $w x$. When w doubles, s just doubles, and the equation still satisfies, so there's no unique s . At the same time, if the numbers of people are shuffled into another order, such as the blue nodes in the figure above are changed into 3, 2, 1. Then only the column vector order of A needs to be changed, so s cannot be determined separately. These two cases are called original signal uncertainty. For an array matrix P , the receiving signal x can be expressed as $x = AP^{-1}Ps$.

The elements of Ps are the same as the elements of s but in A different order, and the mixed matrix AP^{-1} is the same as the elements of A but in A different order.

There's also a kind of ICA that's not going to be gauss. Let's say that only the sound signal of the two people is in line with the multi-value normal distribution, $s \sim N(0, I)$, and the probability density function of the unit matrix of 2×2 is not the same as the mean of the mean 0, which is the peak of the elliptic (see the multi-value gauss distribution). Because x is equal to As , so x is also gauss, the mean is 0, the covariance is $E[xx^T] = E[Ass^T A^T] = AA^T$. Let R be orthogonal matrices ($RR^T = R^T R = I$); $A' = AR$, $A' = AR$. If I replace A as A' ; then $X' = A's$. The s distribution doesn't change, so X' is still the mean of 0, the covariance is

$E[x'(x')^T] = E[A'ss^T (A')^T] = E[ARss^T (AR)^T] = ARR^T A^T = AA^T$. So whether the mixing matrix is A or A' , the distribution of X is the same, so you can't determine the mixing matrix, and you can't determine the original signal.

2.2. Independent vector analysis (IVA)

IVA method can make use of and guarantee the correlation between different frequency points in the frequency domain of each signal source, because it can avoid the output signal sequencing problem caused by ICA method. In order to achieve this correlation, the prior probability distribution of the source cannot be just a single variable generalized gauss distribution as in ICA, and the multivariate generalized gauss distribution is needed here.

It is assumed that the observation signal and the estimated signal can be expressed as:

$$\mathbf{x}(k) = \mathbf{H}(k)\mathbf{s}(k) \quad (17)$$

$$\hat{\mathbf{s}}(k) = \mathbf{W}(k)\mathbf{x}(k) \quad (18)$$

Where, $\mathbf{x}(k) = [x_1(k), x_2(k), \dots, x_m(k)]^T$ is the observation signal, $\mathbf{s}(k) = [s_1(k), s_2(k), \dots, s_n(k)]^T$ is the sound source signal. $\hat{\mathbf{s}}(k) = [\hat{s}_1(k), \hat{s}_2(k), \dots, \hat{s}_n(k)]^T$ is an estimated signal.

This is all in the frequency domain, where k is the frequency point $k = 1, 2, \dots, K$ is the KTH frequency point, K is the number of frequency points; M is the number of microphones and n is the number of sound sources. $\mathbf{H}(k)$ is the confusion matrix, dimension is $m \times n$, $\mathbf{W}(k)$ is the solution matrix, and dimension is $n \times m$. Here, we discuss the case where the number of microphones equals the number of sound sources, namely: $m = n$.

By KullbackLeibler Divergence criterion, the calculation of vector $\hat{\mathbf{s}}$ probability density $p(\hat{\mathbf{s}})$ with marginal distribution of each component of the $p(\hat{s}_n)$ the distance between as follows:

$$\begin{aligned} J &= KL(p(\hat{s}_1 \dots \hat{s}_n) \parallel \prod q(\hat{s}_i)) \\ &= \int p(\hat{s}_1 \dots \hat{s}_n) \log \frac{p(\hat{s}_1 \dots \hat{s}_n)}{\prod q(\hat{s}_i)} d\hat{s}_1 \dots d\hat{s}_n \\ &= \text{const} - \sum_{k=1}^K \log |\det(\mathbf{W}(k))| - \sum_{i=1}^n E[\log q(\hat{s}_i)] \end{aligned} \quad (19)$$

To make use of correlation between different frequency points of the same source, we assume that the distribution of the sound source satisfies the multivariable generalized Gaussian distribution, namely:

$$q(\hat{s}_i) \propto \exp\left(-\left(\frac{(s_i - \mu_i)^\dagger \sum_{i=1}^{-1} (s_i - \mu_i)}{\alpha}\right)^\beta\right) \quad (20)$$

For speech, it belongs to super gauss distribution, namely: $0 < \beta < 1$, where: μ_i and Σ_i are respectively the mean vector and variance vector of the i th sound source.

Because of the orthogonality of Fourier transform, that is, multiple sound sources are independent of each other. Therefore, in order to facilitate analysis and derivation without affecting the results, we assume that the mean value is zero and the variance matrix is a diagonal matrix.

Therefore, the estimation performance of the i sound source can be expressed as the following nonlinear function:

$$\begin{aligned}\varphi^{(k)}(\hat{s}_i(1), \dots, \hat{s}_i(k)) &= -\frac{\partial \log(q(\hat{s}_i(1), \dots, \hat{s}_i(k)))}{\partial \hat{s}_i(k)} \\ &= \frac{2\beta \hat{s}_i(k)}{\left(\sum_{k=1}^K |\hat{s}_i(k)|^2\right)^{1-\beta}} = \frac{2\beta \hat{s}_i(k)}{\left(\left(\sum_{k=1}^K |\hat{s}_i(k)|^2\right)^2\right)^{\frac{1-\beta}{2}}}\end{aligned}\quad (21)$$

Expand the denominator in the above formula to get the following:

$$\left(\sum_{k=1}^K |\hat{s}_i(k)|^2\right)^{\frac{1-\beta}{2}} = \left(\sum_{k=1}^K |\hat{s}_i(k)|^4 + \sum_{a \neq b} c_{ab} |\hat{s}_i(a)|^2 |\hat{s}_i(b)|^2\right)^{\frac{1-\beta}{2}} \quad (22)$$

The denominator in the above formula can be expanded as follows: in order to retain the correlation between different frequency points, that is, to retain the cross terms of $a \neq b$

Then $\frac{1-\beta}{2}$ must be odd

$$\frac{1-\beta}{2} = \frac{1}{2I+1} \Rightarrow \beta = \frac{1}{3} \quad (23)$$

Where, I is a positive integer. Thus, the nonlinear kernel function is:

$$\varphi^{(k)}(\hat{s}_i(1), \dots, \hat{s}_i(k)) = \frac{2\hat{s}_i(k)}{3\sqrt[3]{\left(\sum_{k=1}^K |\hat{s}_i(k)|^2\right)^2}} \quad (24)$$

Similar to the performance function of ICA, the performance function of IVA can be obtained as follows:

$$J = \sum_{i=1}^n (E[F(\sum_{k=1}^K |\hat{s}_i(k)|^2)] - \sum_{k=1}^K \lambda_i^{(k)} (\mathbf{w}_i(k)^\dagger \mathbf{w}_i(k) - 1)) \quad (25)$$

Where, \mathbf{w}_i^\dagger is the i row of solving mixed matrix \mathbf{W} , λ_i is the i Lagrange multiplier, k is the frequency point, and $E[\cdot]$ represents the expectation. $F(\cdot)$ is a nonlinear function. By taking the derivative of the above equation to the solution matrix \mathbf{W} , we can get:

$$\begin{aligned}\mathbf{w}_i(k) &\leftarrow E[F'(\sum_{k=1}^K |\hat{s}_i(k)|^2) + |\hat{s}_i(k)|^2 F''(\sum_{k=1}^K |\hat{s}_i(k)|^2)] \mathbf{w}_i(k) \\ &\quad - E[(\hat{s}_i(k)) * F'(\sum_{k=1}^K |\hat{s}_i(k)|^2) \mathbf{x}^k]\end{aligned}\quad (26)$$

Where, $F'(\cdot)$ and $F''(\cdot)$ are the first and second derivatives of $F(\cdot)$.

Assume the nonlinear function is:

$$F(\sum_{k=1}^K |\hat{s}_i(k)|^2) = \left(\sum_{k=1}^K |\hat{s}_i(k)|^2\right)^{\frac{1}{3}} \quad (27)$$

Then the first derivative is:

$$F'(\sum_{k=1}^K |\hat{s}_i(k)|^2) = \frac{2}{3\sqrt[3]{\left(\sum_{k=1}^K |\hat{s}_i(k)|^2\right)^2}} \quad (28)$$

Then the first derivative is: this is the nonlinear kernel function shown in equation 1-24.

3. Simulated analysis

Due to the music law of music, there is a strong correlation between different frequency points of the same sound source, so the application of IVA method should be able to effectively separate the sound source. Suppose there are two people singing different songs now, the girl sings Chen's "small half" and the boy sings Eason's "four seasons". This problem is also considered as microphone array beamforming [7]. In order to separate two kinds of music, at least two microphones are needed to collect sound signals. The specific scenario is shown in the figure below.

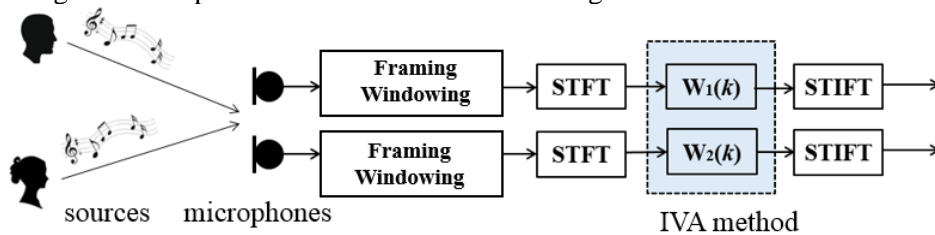
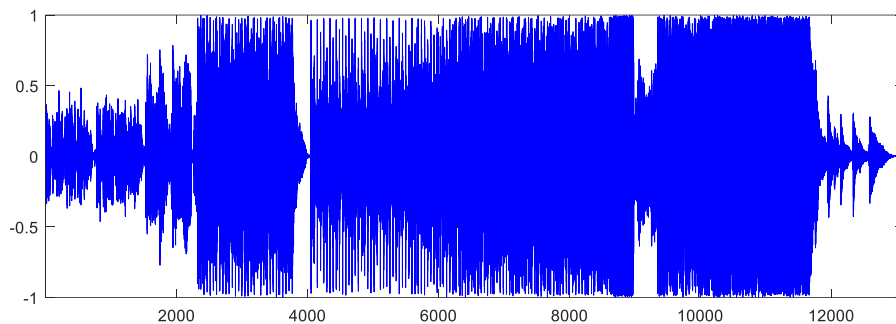


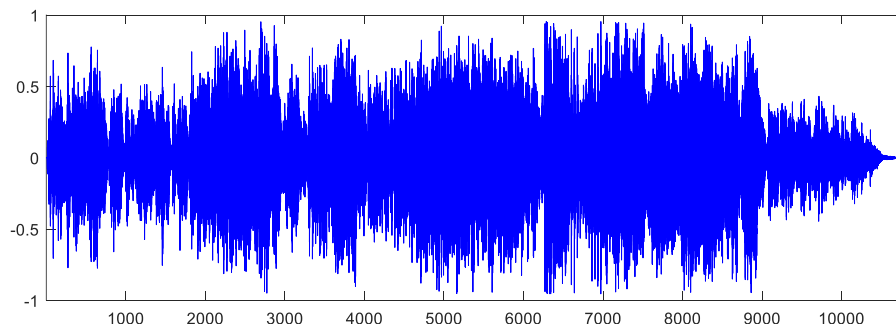
Figure 2 simulation scenario

The specific signal processing process is that two microphones collect sound signals respectively, and the sound signals collected by each microphone are the superposition of two kinds of music. Signal for each receiving all the way, in accordance with the first 10 milliseconds framing a frame, and 50% overlap between frames, after each frame for smooth add a window, USES a hanning window, to add a window by short-time Fourier transform of signal time domain signal transformation to the frequency domain, then the two way independent vector analysis method is used for source separation, after the separation by short-time Fourier inverse transformation to the frequency domain signal transformation to the time domain signal of time domain can be separated.

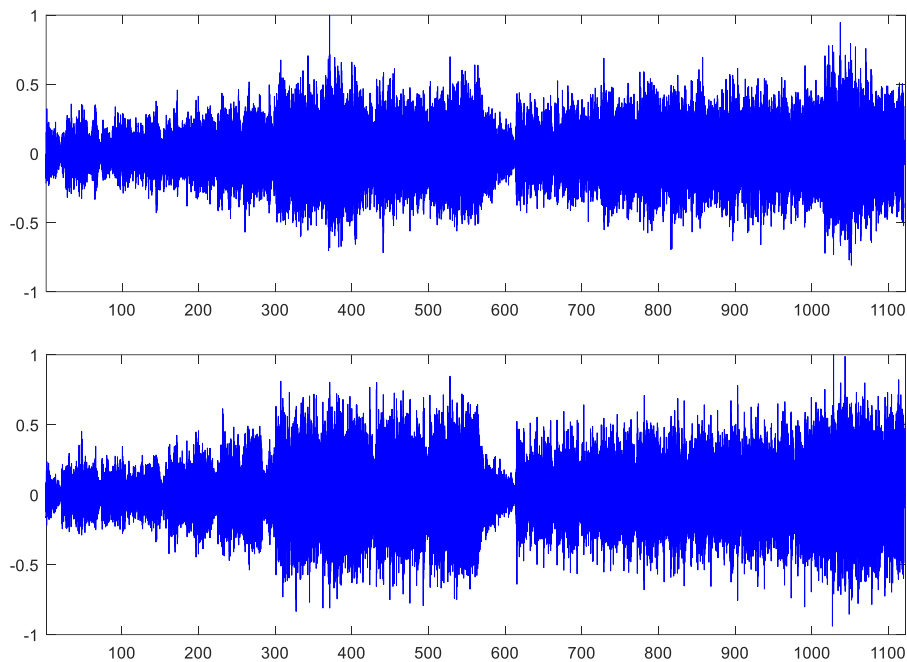
The original male music is:



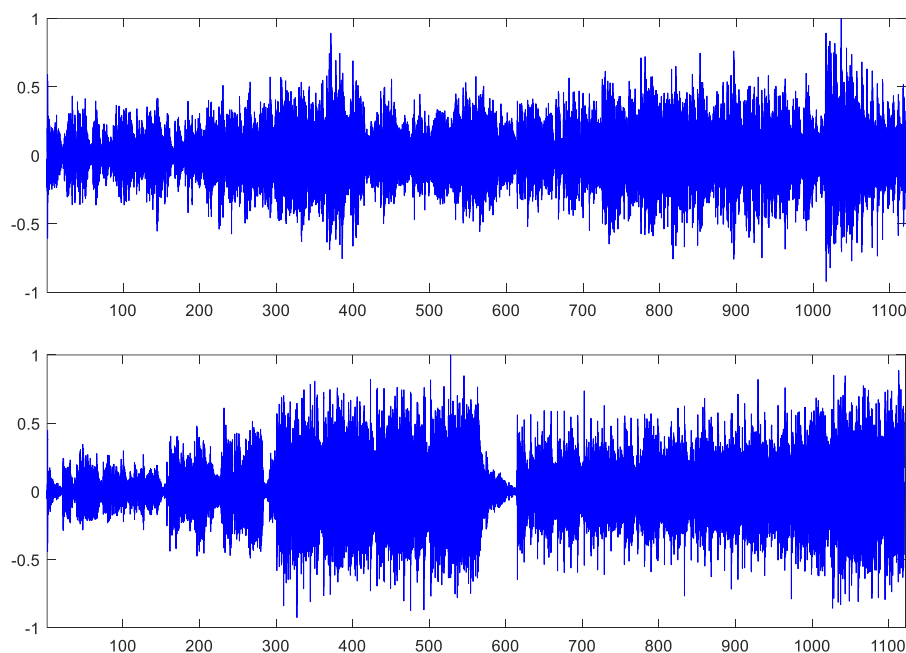
The original female music is:



The original female music is: the mixed signal received by the two microphones is:



After separation by IVA method, the two-channel signals are:



By comparing the signals received by the two microphones with the original two-channel signals, it can be seen that the mixed signals received by the microphones are indistinct from the music waveform. By comparing the separated signal with the original two-channel signal, it can be found that the separated two-channel signal well corresponds to the original audio signal. Therefore, independent vector analysis method can achieve good results in music signal separation.

4. Conclusion

Due to the fact that independent vector analysis can utilize the dependency among the frequencies, it is superior to the independent component analysis. Since the rhythm of music, there may be a strong correlation between different frequency points of the same sound source, so IVA method can be used

for signal separation. This paper discussed the application of IVA method for music sound separation. Simulation results show that the music sound can indeed be separated by IVA method.

References

- [1] Hyvärinen, Aapo, J. Hurri, and P. O. Hoyer. (2014) Independent Component Analysis. *IEEE Transactions on Neural Networks* 15.2: 529–529.
- [2] Comon, Pierre. (1994) Independent component analysis, a new concept? *Signal Processing* 36.3:287-314.
- [3] Kim, Taesu, Torbjørn Eltoft, and T. W. Lee. (2006) Independent Vector Analysis: An Extension of ICA to Multivariate Components. In: *International Conference on Independent Component Analysis and Signal Separation* Springer. Berlin, Heidelberg.
- [4] Ono, Nobutaka. (2011) Stable and fast update rules for independent vector analysis based on auxiliary function technique. In: *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics, WASPAA 2011, New Paltz, NY, USA*.
- [5] Belouchrani, A., et al. (1997) A blind source separation technique using second-order statistics. *IEEE Transactions on Signal Processing* 45.2:434-444.
- [6] Adali, Tulay, Levin-Schwartz, Yuri, and Calhoun, Vince D. (2015) Multimodal Data Fusion Using Source Separation: Two Effective Models Based on ICA and IVA and Their Properties. *Proceedings of the IEEE* 103.9:1478.
- [7] Himawan, Ivan, I. Mccowan, and M. Lincoln. (2007) Microphone Array Beamforming Approach to Blind Speech Separation. *Machine Learning for Multimodal Interaction, 4th International Workshop, MLMI 2007, Brno, Czech Republic, June 28-30, 2007, Revised Selected Papers* DBLP.