

UNDERDETERMINED SOURCE SEPARATION FOR COLORED SOURCES

Stefan Winter^{†,1}, Walter Kellermann[‡], Hiroshi Sawada[†], Shoji Makino[†]

[†]NTT Communication Science Laboratories,
NTT Corporation
2-4 Hikaridai Seika-cho Soraku-gun Kyoto
619-0237 Japan

[‡]Department of Multimedia Communication and Signal
Processing, University Erlangen-Nuremberg
Cauerstr. 7, 91058 Erlangen
Germany

ABSTRACT

This contribution focuses on the source separation stage as important part of underdetermined blind source separation (BSS). So far nearly all approaches for underdetermined BSS assume independently, identically distributed (i.i.d.) sources. They completely ignore the redundancy that is in the temporal structure of colored sources like speech signals. Instead, we propose multivariate models based on the multivariate Student's t or multivariate Gaussian distribution and investigate their potential for underdetermined BSS. We provide a simple yet effective filter based on the sources' autocorrelations for recovering the sources as basis for further advances in underdetermined BSS. The challenge is estimating the filter coefficients blindly. The experimental results support the idea that source separation for underdetermined BSS can be reduced to the separation of their autocorrelations.

1. INTRODUCTION

Blind source separation (BSS) describes techniques that aim at separating P signals if only Q mixed versions of the original signals are available. The need for BSS arises for example if the signals of simultaneous conversations are captured by several microphones. Most BSS approaches assume that there are at least as many microphones as source signals ($Q \geq P$), which is called (over-) determined BSS.

Instead, we consider underdetermined BSS, where we have less microphones than source signals ($Q < P$). Only few approaches have been proposed so far [2, 5, 17, 18], and the separation quality in terms of interference suppression and signal distortion is still not as good as with (over-) determined BSS. This is in particular true if wideband signals like speech signals are involved. The difficulty is that in contrast to (over-) determined BSS the solution of underdetermined BSS goes beyond system identification. Even if the mixing system is fully identified, additional effort is required to separate the mixtures.

So far, nearly all approaches to solve the latter problem assume independently, identically distributed (i.i.d.) sources [2, 5, 17, 18]. While this assumption may serve well as a first order approximation, it completely ignores the redundancy that is in the temporal structure of colored sources like speech signals. Therefore, we propose multivariate models, that take the temporal correlation explicitly into account and investigate their potential. We consider here only linear, instantaneous mixtures of speech signals in the time domain. This paper concentrates on the separation and assumes that the mixing matrix is known.

After formulating the problem analytically in Sec. 2, we provide statistical models for Bayesian inference in Sec. 3. In Sec. 4 we propose a closed-form and numerical approach for minimum mean square error (MMSE) source estimation based on the posterior distribution of the sources. Section 5 elaborates more on autocorrelation estimation, which turns out to be essential for high-quality underdetermined source separation. Section 6 presents experimental results which are discussed in Section 7.

2. PROBLEM FORMULATION

With $s_p(t) \in \mathbb{R}$ denoting the p -th source signal ($1 \leq p \leq P$) and $\tilde{\mathbf{A}} \in \mathbb{R}^{Q \times P}$ the mixing matrix, we obtain mixed signals $x_q(t) \in \mathbb{R}$ ($1 \leq q \leq Q, Q < P$) by

$$\begin{bmatrix} x_1(t) \\ \vdots \\ x_Q(t) \end{bmatrix} = \tilde{\mathbf{A}} \begin{bmatrix} s_1(t) \\ \vdots \\ s_P(t) \end{bmatrix} + \mathbf{n}(t). \quad (1)$$

$\mathbf{n}(t) \in \mathbb{R}^P$ denotes noise added to the sensors. Let

$$\mathbf{S}_p = [s_p(0) \quad \cdots \quad s_p(T-1)]^T, \quad 1 \leq p \leq P \quad (2)$$

denote a frame of length T of the p -th original speech signal with autocorrelation

$$r_p(\tau) = E\{s_p(t)s_p(t+\tau)\}. \quad (3)$$

$(\cdot)^T$ denotes the transpose. We summarize the P frames of the different source signals by the vector

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_1 \\ \vdots \\ \mathbf{S}_P \end{bmatrix} \in \mathbb{R}^{P \cdot T} \quad (4)$$

For each source \mathbf{S}_p the frame-dependent autocorrelation matrix $\mathbf{R}_p \in \mathbb{R}^{T \times T}$ is given by a symmetric Toeplitz matrix. Its first row is defined by

$$[r_p(0) \quad \cdots \quad r_p(T-1)] \quad (5)$$

We summarize the P autocorrelation matrices by

$$\mathbf{R} = \begin{bmatrix} \mathbf{R}_1 & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \mathbf{R}_P \end{bmatrix} \in \mathbb{R}^{P \cdot T \times P \cdot T} \quad (6)$$

Extending the scalar elements \tilde{A}_{qp} of $\tilde{\mathbf{A}}$ to diagonal matrices

$$\mathbf{A}_{qp} = \begin{bmatrix} \tilde{A}_{qp} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \tilde{A}_{qp} \end{bmatrix} \in \mathbb{R}^{T \times T}, \quad (7)$$

we define an extended mixing matrix

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1P} \\ \vdots & \ddots & \vdots \\ \mathbf{A}_{Q1} & \cdots & \mathbf{A}_{QP} \end{bmatrix} \in \mathbb{R}^{Q \cdot T \times P \cdot T} \quad (8)$$

Similar to the source signals we define

$$\mathbf{X}_q = [x_q(0) \quad \cdots \quad x_q(T-1)]^T, \quad 1 \leq q \leq Q \quad (9)$$

¹The author is on leave from the Department of Multimedia Communication and Signal Processing, University Erlangen-Nuremberg

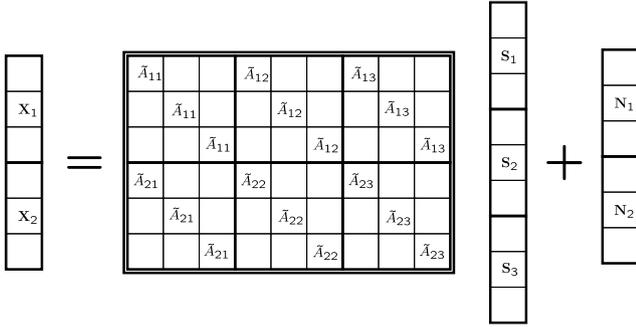


Figure 1: Mixing process for underdetermined BSS with $Q = 2$, $P = 3$

and summarize the \mathbf{X}_q by

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_Q \end{bmatrix} \in \mathbb{R}^{Q \cdot T} \quad (10)$$

This results in the compact description of the mixing process for one frame by

$$\mathbf{X} = \mathbf{A}\mathbf{S} + \mathbf{N} \quad (11)$$

as illustrated in Fig. 1. $\mathbf{N} \in \mathbb{R}^{Q \cdot T}$ is derived from $\mathbf{n}(t)$ in a similar way as \mathbf{X} is derived from $x(t)$.

The final goal in BSS is the estimation of signals \mathbf{Y} that resemble the original signals \mathbf{S} as closely as possible. Since only the mixed signals are available, this is at best possible up to arbitrary permutation and scaling. With (over-) determined BSS it is sufficient to estimate the mixing matrix \mathbf{A} or its inverse. Since the mixing matrix is square and therefore invertible (assuming that \mathbf{A} is non-singular), the inverse can be used to separate the signals. In contrast to (over-) determined BSS, here the estimation of the mixing matrix \mathbf{A} is not sufficient. Even if \mathbf{A} is available, estimating the original signals from the mixtures poses a problem on its own, since \mathbf{A} cannot be simply inverted. In the following we concentrate on estimating the original signals and assume that the mixing matrix is known or can otherwise be estimated [2, 14, 15, 16].

3. MODELS

The most general approach for estimating the unknown source signals is based on Bayesian inference. It yields the posterior distribution $p(\boldsymbol{\theta}|\mathbf{X})$ of the desired parameters $\boldsymbol{\theta}$ (here: source signals) by accounting for their prior distribution $p(\boldsymbol{\theta})$ and the likelihood $p(\mathbf{X}|\boldsymbol{\theta})$ of the observed data \mathbf{X} (here: mixed signals). From Bayes' rule we obtain [12]

$$p(\boldsymbol{\theta}|\mathbf{X}) \propto p(\mathbf{X}|\boldsymbol{\theta}) \cdot p(\boldsymbol{\theta}). \quad (12)$$

So far, nearly all approaches to underdetermined BSS assume independent and identical prior distributions. In contrast, we propose in this section two multivariate priors tailored to speech signals. We also define the likelihood derived from the commonly assumed Gaussian sensor noise model.

In both speech models we make the common assumption that the speech signals are mutually independent.

$$p(\mathbf{S}) = \prod_p p(\mathbf{S}_p) \quad (13)$$

We further assume that the speech signal is approximately stationary within a frame of appropriately chosen length T .

In our first speech model, we model each frame \mathbf{S}_p by a zero mean, multivariate Student's t distribution

$$p(\mathbf{S}_p) = \mathcal{ST}(\mathbf{S}_p|\alpha_p, \boldsymbol{\Sigma}_p) \quad (14)$$

with degrees of freedom α_p and scale $\boldsymbol{\Sigma}_p$. Modeling the speech signal by a multivariate Student's t distribution has several advantages:

- It is a good approximation of real speech signals. This does not just follow from experimental results but is also based on results with spherically invariant random processes (SIRPs) [3] and the univariate Student's t distribution [5].
- The multivariate Student's t distribution can be expressed as a multivariate Gaussian mixture model (MGM) with an infinite number of mixtures [1] as follows

$$\mathcal{ST}(\mathbf{S}_p|\alpha_p, \boldsymbol{\Sigma}_p) = \int \mathcal{N}(\mathbf{S}_p|\mathbf{0}, \mathbf{R}_p) \cdot \mathcal{IW}(\mathbf{R}_p|\alpha_p, \boldsymbol{\Sigma}_p) d\mathbf{R}_p. \quad (15)$$

\mathcal{N} and \mathcal{IW} denote a multivariate Gaussian and inverse Wishart distribution, respectively. According to this relation, the multivariate Student's t distribution is the marginal distribution obtained from the joint distribution of the speech signal and its autocorrelation. This leads to a hierarchical Bayesian model with a multivariate Gaussian distribution for the speech signal \mathbf{S}_p

$$p(\mathbf{S}_p|\mathbf{R}_p) = \mathcal{N}(\mathbf{S}_p|\mathbf{0}, \mathbf{R}_p) \quad (16)$$

and an inverse Wishart distribution for the autocorrelation matrix \mathbf{R}_p

$$p(\mathbf{R}_p|\alpha_p, \boldsymbol{\Sigma}_p) = \mathcal{IW}(\mathbf{R}_p|\alpha_p, \boldsymbol{\Sigma}_p). \quad (17)$$

The expected value of \mathbf{R}_p is [8]

$$E\{\mathbf{R}_p\} = \frac{1}{\alpha_p - T - 1} \boldsymbol{\Sigma}_p^{-1}. \quad (18)$$

The hierarchical model will be used in Sec. 4 to obtain conditional posterior PDFs.

- The inverse Wishart distribution is the conjugate prior of a multivariate Gaussian likelihood with respect to its correlation matrix [7]. This means that the corresponding posterior distribution is again an inverse Wishart distribution, which is also exploited in Sec. 4.

The degrees of freedom α_p and the scale $\boldsymbol{\Sigma}_p$ in (17) can either be considered deterministic or random. We summarize all parameters in the parameter vector $\boldsymbol{\theta}$ and obtain the joint prior distribution as

$$p_{\boldsymbol{\theta}}(\boldsymbol{\theta}) = p(\mathbf{S}, \mathbf{R}, \mathbf{A}, \sigma^2, \boldsymbol{\alpha}, \boldsymbol{\Sigma}) \quad (19)$$

$$= p(\mathbf{S}|\mathbf{R})p(\mathbf{R}|\boldsymbol{\alpha}, \boldsymbol{\Sigma})p(\mathbf{A})p(\sigma^2)p(\boldsymbol{\alpha})p(\boldsymbol{\Sigma}) \quad (20)$$

with $\boldsymbol{\alpha} = [\alpha_1 \ \dots \ \alpha_P]^T$. The relation between the different variables is illustrated in Fig. 2.

As a special case of the first speech model we propose an analytically feasible, multivariate Gaussian prior as second speech model

$$p(\mathbf{S}_p) = \lim_{\alpha_p \rightarrow \infty} \mathcal{ST}(\mathbf{S}_p|\alpha_p, \boldsymbol{\Sigma}_p) = \mathcal{N}(\mathbf{S}_p|\mathbf{0}, \mathbf{R}_p) \quad (21)$$

with zero mean vector and covariance matrix \mathbf{R}_p .

For deriving the likelihood we assume additive white Gaussian noise with identical variance σ^2 across all sensors. Based on the mixing process (11), this assumption results in the Gaussian likelihood

$$p(\mathbf{X}|\mathbf{A}, \mathbf{S}, \sigma^2) = \mathcal{N}(\mathbf{X}|\mathbf{A}\mathbf{S}, \sigma^2). \quad (22)$$

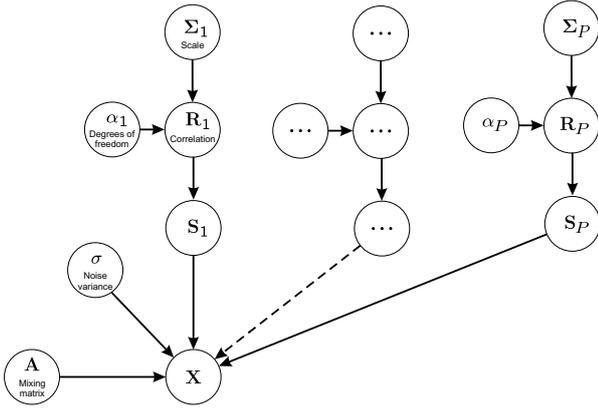


Figure 2: Graphical model

4. BAYESIAN INFERENCE

The general goal in Bayesian inference is to estimate the posterior PDF of unknown parameters based on the available data. Except for few special cases like the Gaussian speech model (21) with Gaussian likelihood (22), the posterior cannot be derived in closed-form. Instead, it requires numerical methods based for example on Monte Carlo techniques [4]. Based on our multivariate models in the previous section, we derive in the following MMSE estimates of the sources in closed-form and by numerical approximation. Further details are provided in the appendix.

4.1 Multivariate Student's t prior

In order to estimate the sources based on the multivariate Student's t prior (14), we employ a Gibbs sampler, which is a Monte Carlo technique. In contrast to our multivariate approach, the Gibbs sampler was applied in [5] to an i.i.d. speech model based on a univariate Student's t distribution.

4.1.1 Basic idea

The Gibbs sampler approximates the posterior PDF by sampling iteratively from appropriately chosen conditional PDFs. In fact, the conditional PDFs form a first-order Markov chain, whose stationary PDF is the desired posterior PDF [4, 8].

Let $\theta = \{\theta_1, \dots, \theta_K\}$ denote the unknown parameters with prior PDF $p_\theta(\theta)$ and \mathbf{X} the available data with likelihood $p_{\mathbf{X}|\theta}(\mathbf{X}|\theta)$. We further define θ_{-i} as summarizing all parameters without parameter θ_i . For the Gibbs sampler we need the conditional PDFs $p(\theta_i|\theta_{-i}, \mathbf{X})$, which can be shown to be [5]

$$p(\theta_i|\theta_{-i}, \mathbf{X}) \sim p_\theta(\theta) p_{\mathbf{X}|\theta}(\mathbf{X}|\theta) \quad (23)$$

Once they are derived, we can sample iteratively from each conditional PDF using the previous samples as summarized in Algorithm 1. As soon as a sufficient number of samples is drawn, they ap-

Algorithm 1: General Gibbs sampler

```

Initialize  $\theta$  (e.g. randomly) ;
for  $k=1..K$  do
  for  $i$  do
     $\theta_i^{(k+1)} \sim p(\theta_i|\theta_{-i}^{(k,k+1)}, \mathbf{x})$ ;
  end
end
end

```

proximate the posterior PDF. They can further be used to calculate point estimates \mathbf{Y} (MAP, MMSE, ...) of the original signals \mathbf{S} . In order to avoid transient effects a certain number of initial samples should be discarded [9, 13].

4.1.2 Implementation

We now derive the conditional PDFs for our specific problem based on the multivariate Student's t speech model in Sec. 3.

For the conditional PDF of the speech signal we obtain

$$p(\mathbf{S}|\theta_{-S}, \mathbf{X}) \propto p(\mathbf{S}|\mathbf{R}) \cdot p(\mathbf{X}|\mathbf{A}, \mathbf{S}, \sigma^2) \quad (24)$$

$$\propto \mathcal{N}(\mathbf{S}|\mathbf{0}, \mathbf{R}) \cdot \mathcal{N}(\mathbf{X}|\mathbf{A}\mathbf{S}, \sigma^2) \quad (25)$$

$$\propto \mathcal{N}(\mathbf{S}|\boldsymbol{\mu}_S, \boldsymbol{\Sigma}_S) \quad (26)$$

with

$$\boldsymbol{\Sigma}_S = \left(\mathbf{R}^{-1} + \frac{1}{\sigma^2} \mathbf{A}^T \mathbf{A} \right)^{-1} \quad (27)$$

and

$$\boldsymbol{\mu}_S = \left(\sigma^2 \mathbf{R}^{-1} + \mathbf{A}^T \mathbf{A} \right)^{-1} \mathbf{A}^T \mathbf{X}. \quad (28)$$

For the conditional PDF of the autocorrelation we obtain

$$p(\mathbf{R}|\theta_{-R}, \mathbf{X}) \propto p(\mathbf{S}|\mathbf{R}) \cdot p(\mathbf{R}|\boldsymbol{\alpha}, \boldsymbol{\Sigma}) \quad (29)$$

$$\propto \prod_p \mathcal{N}(\mathbf{S}_p|\mathbf{0}, \mathbf{R}_p) \cdot \mathcal{IW}(\mathbf{R}_p|\alpha_p, \boldsymbol{\Sigma}_p)$$

$$\propto \prod_p \mathcal{IW}\left(\mathbf{R}_p|\alpha_p + 1, \left(\mathbf{S}_p \mathbf{S}_p^T + \boldsymbol{\Sigma}_p^{-1}\right)^{-1}\right)$$

To successfully implement this Gibbs sampler, appropriate choices for the degrees of freedom α_p and the scale $\boldsymbol{\Sigma}_p$ are necessary. It is common in Bayesian inference to manually tune the degrees of freedom. Usually there is also sufficient prior data available to train the scale (e.g. speaker adaptation for speech recognition). Here we face the problem, that we do not have sufficient prior data. We consider the following options:

- We could consider the scale as random variable as well and choose an uninformative prior (e.g. Jeffrey's prior, uniform prior, ...). So far experiments with uninformative priors did not yet lead to reasonable results. Therefore, we don't consider this option here.
- Knowing that the choice of the scale $\boldsymbol{\Sigma}_p$ directly influences the autocorrelation by (18), we can set it to a specific value, determined for example by some preprocessing (e.g. ℓ_1 -norm minimization [17]) as detailed in the next section.

4.2 Multivariate Gaussian prior

Plugging the Gaussian prior (21) and the likelihood (22) in (12) we obtain for the Gaussian posterior similar to (24)

$$p(\mathbf{S}|\mathbf{X}) \propto p(\mathbf{S}|\mathbf{R}) \cdot p(\mathbf{X}|\mathbf{A}, \mathbf{S}, \sigma^2) \quad (30)$$

$$\propto \mathcal{N}(\mathbf{S}|\boldsymbol{\mu}_S, \boldsymbol{\Sigma}_S).$$

with $\boldsymbol{\Sigma}_S$ and $\boldsymbol{\mu}_S$ given by (27) and (28), respectively.

Since the MMSE estimate of a normally distributed random variable is its mean, the resulting estimate of the sources is

$$\mathbf{Y} = \underbrace{\left(\sigma^2 \mathbf{R}^{-1} + \mathbf{A}^T \mathbf{A} \right)^{-1} \mathbf{A}^T \mathbf{X}}_{=\mathbf{W}}. \quad (31)$$

The unmixing filter \mathbf{W} is a generalized inverse modified by the autocorrelation matrix. The influence of the autocorrelation matrix \mathbf{R} is determined by the noise variance σ^2 . Without modification by the autocorrelation matrix (31) could not be determined, since $\mathbf{A}^T \mathbf{A}$ is a singular matrix.

5. AUTOCORRELATION

For the proposed speech models the autocorrelation plays a crucial role. With the multivariate Gaussian model, the estimate is used directly. With the multivariate Student's t model, the estimate can be used to set the scale Σ_p , if it is treated as a deterministic value. It then determines the expected value of the autocorrelation. Therefore, for high quality separation in terms of interference suppression and signal distortion [6] we need a good estimate of the autocorrelation. In this section we will discuss possible options for estimating the autocorrelation. In general, we can distinguish between non-parametric and parametric approaches.

5.1 Non-parametric estimation

Non-parametric estimation has the advantage that no assumption about the structure is necessary. For our purpose with only a limited number T of samples of the signal in question available, we choose a biased estimation technique [10]. It guarantees that the estimated autocorrelation is positive semi-definite. For a signal frame $s(t)$, $0 \leq t \leq T-1$ the biased autocorrelation estimate is given by

$$r(\tau) = E\{s(t)s(t+\tau)\} \quad (32)$$

$$\approx \frac{1}{T} \sum_{t=0}^{T-\tau-1} s(t)s(t+\tau) \quad (33)$$

5.2 Parametric estimation

Non-parametric estimation of the autocorrelation suffers from an insufficient amount of data compared to the number of parameters (if we define all time lags of the autocorrelation as parameters). This results in high variance and low resolution [11]. In general, it is desirable to reduce the number of parameters by an appropriate parametric model that exploits the inherent structure.

Based on the harmonicity of speech we propose approximating the biased autocorrelation by M harmonic, decaying cosine functions.

$$r_p(\tau) \approx \left(\sum_{m=1}^M \gamma_{m,p} \cos(p \cdot f_{0,p} \cdot \tau) \right) \exp(-\beta_p \tau) \quad (34)$$

Here, $\gamma_{m,p}$ denotes the gain of each cosine, $f_{0,p}$ is the fundamental frequency and β_p is the decay factor. This reduces the number of parameters to $M+2$ for each source.

5.2.1 Least square estimation

If we have access to the autocorrelation that we want to approximate, we can use a least square approach to estimate the parameters. We apply this method to the following scenarios.

- To explore the possibilities of the parametric approach, we can use the autocorrelations of the original signals to estimate the parameters.
- Other parameter estimation methods might need to be initialized. Then we can use initial estimates of the sources (e.g. ℓ_1 -norm minimization [17]) and use their autocorrelations.
- Linear mixtures of mutually independent signals have the property, that their autocorrelations are a linear mixture of the autocorrelation of the original signals.

$$r_{xq}(\tau) = E\{x_q(t)x_q(t-\tau)\} \quad (35)$$

$$= E\left\{ \left(\sum_p \tilde{A}_{qp} s_p(t) \right) \cdot \left(\sum_p \tilde{A}_{qp} s_p(t-\tau) \right) \right\}$$

$$= E\left\{ \sum_p \tilde{A}_{qp}^2 s_p(t)s_p(t-\tau) \right\} \quad (36)$$

$$= \sum_p \tilde{A}_{qp}^2 E\{s_p(t)s_p(t-\tau)\} \quad (37)$$

$$= \sum_p \tilde{A}_{qp}^2 r_{s_p}(\tau) \quad (38)$$

If we assume, that the fundamental frequency $f_{0,p}$ is different for each speaker, we can estimate the parameters from the mixtures for each source in parallel.

In each scenario a non-parametric estimation is the basis for the parameter estimation. As a general least square cost function for all scenarios we propose

$$\mathcal{J}(\gamma, \beta, \mathbf{f}_0) = \sum_q \sum_{\tau} (r_q(\tau) -$$

$$\sum_p \tilde{A}_{q,p}^2 \left(\left(\sum_{m=1}^M \gamma_{m,p} \cos(m \cdot f_{0,p} \cdot \tau) \right) \exp(-\beta_p \tau) \right)^2$$

6. EXPERIMENTAL RESULTS

We performed experiments with three speech signals (two male, one female) of 1.62 seconds (8kHz sampling rate). We generated two instantaneous mixtures with the mixing matrix

$$\tilde{\mathbf{A}} = \begin{bmatrix} 1 & 1 & 1 \\ 1.3 & -0.9 & 0.8 \end{bmatrix} \quad (40)$$

and assumed that the mixing matrix is available. The signals were processed in frames of $T=256$ samples and shifted by 64 samples. No noise was added to the mixtures, but the noise variance was set to $\sigma^2 = 10^{-4}$. This led to a good influence of the autocorrelation matrix \mathbf{R} . For approximations of the autocorrelation we set $M=2$. The sources had equal variance of 0.0037. We compared the following approaches:

- L1 Separating the mixtures by ℓ_1 -norm minimization in the frequency domain, which is based on a univariate speech model [17].
- L1G Using the results from L1 to estimate the autocorrelation. This estimate was used to further enhance the separated signals by the Gibbs sampler (24)-(29) ($K=1000$ iterations).
- OG Estimating the signals with the Gibbs sampler (24)-(29) ($K=10$ iterations), whereby the autocorrelation was estimated from the original signals (33).
- OI With the autocorrelation estimated from the original signals (33), the signals were obtained with the modified generalized inverse (31).
- OAI Estimating the signals with the modified generalized inverse (31), whereby each autocorrelation was approximated by two harmonic, exponentially weighted cosines based on the original signals (34).
- OSI Estimating the signals with the modified generalized inverse (31), whereby the autocorrelation matrix was a diagonal matrix filled with the first lag of the autocorrelation of the original signals (33).
- MAI Estimating the signals with the modified generalized inverse (31), whereby the autocorrelation was approximated directly from the mixtures (34).

For the evaluation of the separation results we used the signal-to-distortion ratio (SDR), signal-to-interference ratio (SIR) and signal-to-artifact ratio (SAR) as defined in [6]. The averaged results of the different approaches are shown in Table 1.

7. CONCLUSION

The experimental results in Table 1 together with subjective evaluation suggest that knowing the autocorrelation together with the mixing matrix is sufficient to perform high quality underdetermined source separation in noiseless environments (OG and OI compared to L1G, OSI and MAI). In other words, the problem of underdetermined source separation can be reduced to the estimation of the underlying autocorrelations once the mixing matrix is available. The

Method	SDR	SIR	SAR
LI	9.28	14.84	11.79
LIG	5.29	17.65	7.25
OG	12.3	35.61	12.56
OI	18.21	34.33	18.51
OAI	10.42	20.75	11.40
OSI	11.11	18.08	12.89
MAI	3.45	13.63	5.83

Table 1: Experimental results

proposed parameterization of the autocorrelation did not yet lead to satisfying results (OAI, MAI). Therefore, either better autocorrelation models or improved non-parametric estimation methods are necessary.

The experimental results obtained by the MMSE approach OI resemble the results of the Gibbs sampler OG. This suggests that the exact prior distribution of the sources plays only a minor role as long as the autocorrelation is taken into account. Therefore, the comparatively easy and fast solution given by (31) might provide a basis for algorithms that yield better autocorrelation estimates and eventually underdetermined source separation.

Possible options for improving the autocorrelation estimate include linear prediction based models and independent component analysis (ICA) techniques. Linear prediction has already proven in speech coding, that it allows very compact yet high quality descriptions of speech signals. ICA exploits the statistical independence between the sources and might be helpful for estimating the involved parameters.

A. CONDITIONAL PDFS FOR GIBBS SAMPLER

- Conditional PDF of speech signal

$$\begin{aligned}
p(\mathbf{S}|\boldsymbol{\theta}_{-S}, \mathbf{X}) &\propto p(\mathbf{S}|\mathbf{R}) \cdot p(\mathbf{X}|\mathbf{A}, \mathbf{S}, \sigma^2) \\
&\propto \mathcal{N}(\mathbf{S}|\mathbf{R}) \cdot \mathcal{N}(\mathbf{X}|\mathbf{A}\mathbf{S}, \sigma^2) \\
&\propto \exp\left(-\frac{1}{2}\mathbf{S}^T \mathbf{R}^{-1} \mathbf{S}\right) \cdot \\
&\quad \exp\left(-\frac{1}{2}(\mathbf{X} - \mathbf{A}\mathbf{S})^T \frac{\mathbf{I}}{\sigma^2} (\mathbf{X} - \mathbf{A}\mathbf{S})\right) \\
&\propto \exp\left(-\frac{1}{2}\left(\underbrace{\mathbf{S}^T \left(\mathbf{R}^{-1} + \frac{1}{\sigma^2} \mathbf{A}^T \mathbf{A}\right) \mathbf{S}}_{:=\mathbf{S}^T \boldsymbol{\Sigma}_S^{-1} \mathbf{S}} - \right.\right. \\
&\quad \left.\left. \underbrace{-\mathbf{S}^T \frac{1}{\sigma^2} \mathbf{A}^T \mathbf{X}}_{:=\boldsymbol{\Sigma}_S^{-1} \boldsymbol{\mu}_S} - \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{A} \mathbf{S} + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X}\right)\right) \\
&= \mathcal{N}(\mathbf{S}|\boldsymbol{\mu}_S, \boldsymbol{\Sigma}_S)
\end{aligned} \tag{41}$$

with $\boldsymbol{\Sigma}_S$ and $\boldsymbol{\mu}_S$ given by (27) and (28), respectively.

- Conditional PDF of autocorrelation

$$\begin{aligned}
p(\mathbf{R}|\boldsymbol{\theta}_{-R}, \mathbf{X}) &\propto p(\mathbf{S}|\mathbf{R}) \cdot p(\mathbf{R}|\boldsymbol{\alpha}, \boldsymbol{\Sigma}) \\
&\propto \prod_p \mathcal{N}(\mathbf{S}_p|\mathbf{0}, \mathbf{R}_p) \cdot \prod_p \mathcal{IW}(\mathbf{R}_p|\alpha_p, \boldsymbol{\Sigma}_p) \\
&\propto \prod_p \mathcal{N}(\mathbf{S}_p|\mathbf{0}, \mathbf{R}_p) \cdot \mathcal{IW}(\mathbf{R}_p|\alpha_p, \boldsymbol{\Sigma}_p) \\
&\propto \prod_p |\mathbf{R}_p|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}\mathbf{S}_p^T \mathbf{R}_p^{-1} \mathbf{S}_p\right).
\end{aligned} \tag{42}$$

$$\begin{aligned}
&|\mathbf{R}_p|^{-\frac{1}{2}(\alpha_p+T+1)} \exp\left(-\frac{1}{2}\text{tr}\left(\mathbf{R}_p^{-1} \boldsymbol{\Sigma}_p^{-1}\right)\right) \\
&\propto \prod_p |\mathbf{R}_p|^{-\frac{1}{2}(\alpha_p+T+2)} \cdot \\
&\quad \exp\left(-\frac{1}{2}\text{tr}\left(\mathbf{R}_p^{-1} \mathbf{S}_p \mathbf{S}_p^T\right) - \frac{1}{2}\text{tr}\left(\mathbf{R}_p^{-1} \boldsymbol{\Sigma}_p^{-1}\right)\right) \\
&= \prod_p \mathcal{IW}\left(\mathbf{R}_p|\alpha_p+1, \left(\mathbf{S}_p \mathbf{S}_p^T + \boldsymbol{\Sigma}_p^{-1}\right)^{-1}\right)
\end{aligned}$$

References

- [1] J.M. Bernardo and A.F.M. Smith. *Bayesian theory*. Wiley series in probability and statistics. Wiley, 2000.
- [2] A. Blin, S. Araki, and S. Makino. Underdetermined blind separation of convolutive mixtures of speech using time-frequency mask and mixing matrix estimation. *IEICE Trans. Fundamentals*, E88-A(7):1693–1700, 2005.
- [3] H. Brehm and W. Stammerl. Description and generation of spherically invariant speech-model signals. *Signal Processing*, 12:119–141, 1987.
- [4] A. Doucet and X. Wang. Monte Carlo methods for signal processing: a review in the statistical signal processing context. *IEEE Signal Processing Magazine*, 22(6):152–170, Nov 2005.
- [5] C. Févotte and S. J. Godsill. A bayesian approach for blind separation of sparse sources. Technical report, Cambridge University Engineering Dept., January 2005.
- [6] C. Févotte, R. Gribonval, and E. Vincent. BSS_EVAL toolbox user guide – Revision 2.0. Technical Report 1706, IRISA, April 2005.
- [7] C. Fraley and A.E. Raftery. Bayesian regularization for normal mixture estimation and model-based clustering. Technical Report 486, University of Washington, Statistics, August 2005.
- [8] A. Gelman, J.B. Carlin, H.S. Stern, and D.B. Rubin. *Bayesian data analysis*. Chapman & Hall, 1995.
- [9] D.J.C. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 7 edition, Aug 2004.
- [10] D.G. Manolakis, V.K. Ingle, and S.M. Kogon. *Statistical and adaptive signal processing*. McGraw-Hill, 2000.
- [11] J.M. Mendel. Tutorial on higher-order statistics (spectra) in signal processing and system theory: theoretical results and some applications. *Proceedings of the IEEE*, 79(3):278–305, March 1991.
- [12] A. Papoulis and S.U. Pillai. *Probability, random variables, and stochastic processes*. McGraw-Hill, 4th edition, 2002.
- [13] D.B. Rowe. *Multivariate Bayesian statistics: models for source separation and signal unmixing*. Chapman & Hall/CRC, 2003.
- [14] F.J. Theis. *Mathematics in independent component analysis*. PhD thesis, University of Regensburg, 2002.
- [15] L. Vielva, I. Santamaria, C. Pantaleon, J. Ibanez, and D. Erdogmus. Estimation of the mixing matrix for underdetermined blind source separation using spectral estimation techniques. In *Proc. EUSIPCO 2002*, volume 1, pages 557–560, Sep 2002.
- [16] S. Winter, H. Sawada, S. Araki, and S. Makino. Overcomplete BSS for convolutive mixtures based on hierarchical clustering. In *Proc. ICA 2004*, pages 652–660, Sept. 2004.
- [17] S. Winter, H. Sawada, and S. Makino. On real and complex valued L1-norm minimization for overcomplete blind source separation. In *2005 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 86–89, New Paltz, NY, USA, 2005.
- [18] M. Zibulevsky and B.A. Pearlmutter. Blind source separation by sparse decomposition. *Neural Computations*, 13(4):863–882, 2001.