

Sparse filter models for solving permutation indeterminacy in convolutive blind source separation

Prasad Sudhakar, Rémi Gribonval

► **To cite this version:**

Prasad Sudhakar, Rémi Gribonval. Sparse filter models for solving permutation indeterminacy in convolutive blind source separation. Rémi Gribonval. SPARS'09 - Signal Processing with Adaptive Sparse Structured Representations, Apr 2009, Saint Malo, France. 2009. <inria-00369554>

HAL Id: inria-00369554

<https://hal.inria.fr/inria-00369554>

Submitted on 20 Mar 2009

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Sparse filter models for solving permutation indeterminacy in convolutive blind source separation

Prasad Sudhakar and Rémi Gribonval
 METISS Team, Centre de recherche INRIA-Rennes Bretagne Atlantique
 Campus de Beaulieu-35042, Rennes CEDEX, France
 Email: firstname.lastname@irisa.fr
<http://www.irisa.fr/metiss>

Abstract—Frequency-domain methods for estimating mixing filters in convolutive blind source separation (BSS) suffer from permutation and scaling indeterminacies in sub-bands. Solving these indeterminacies are critical to such BSS systems. In this paper, we propose to use sparse filter models to tackle the permutation problem. It will be shown that the ℓ_1 -norm of the filter matrix increases with permutations and with this motivation, an algorithm is then presented which aims to solve the permutations in the absence of any scaling. Experimental evidence to show the behaviour of ℓ_1 -norm of the filter matrix to sub-band permutations is presented. Then, the performance of our proposed algorithm is presented, both in noiseless and noisy cases.

Index Terms—convolutive BSS, permutation ambiguity, sparsity, ℓ_1 -minimization

I. INTRODUCTION

The problem of source separation arises in various contexts such as speech enhancement/recognition, biomedical signal processing, wireless telecommunication, etc. Mathematical models of BSS with different levels of complexity have been proposed by the signal processing community. The most difficult but close to reality model is of convolutive mixtures. The underlying model of having M mixtures $x_m(t), m = 1 \dots M$ from N source signals $s_n(t), n = 1 \dots N$, given a discrete time index t , is given by

$$x_m(t) = \sum_{n=1}^N \sum_{k=0}^{K-1} a_{mnk} s_n(t-k) + v_m(t) \quad (1)$$

with $v_m(t)$ the noise term. In the matrix notation it can be written as

$$\mathbf{x}(t) = \sum_{k=0}^{K-1} \mathbf{A}_k \mathbf{s}(t-k) + \mathbf{v}(t) \quad (2)$$

where $\mathbf{x}(t), \mathbf{v}(t)$ are $m \times 1$ vectors, \mathbf{A}_k is an $M \times N$ matrix which contains the filter coefficients at k^{th} index. The notation $\mathbf{A}_{mn}(t) = a_{mnt}$ will also be used for each mixing filter, which is of length K . The ultimate objective of a BSS system is to recover back the original source signals $s_n, n = 1 \dots N$ given only the mixtures $x_m(t), m = 1 \dots M$.

A standard approach to separate sources is to first estimate the mixing matrix $\mathbf{A}(t)$ and then recover the sources $s_n(t)$ [1]. Fig. 1 shows the block diagram of such a system. This

paper focusses on the second block of the system: estimation of the mixing matrix.

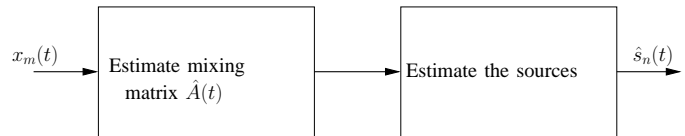


Fig. 1. Block diagram of a typical BSS system

II. PERMUTATION PROBLEM DESCRIPTION

Several methods have been proposed by the signal processing community to estimate the mixing matrices in convolutive BSS. Pedersen et. al. present an excellent survey of the existing methods [1]. Broadly, the techniques can be classified into time-domain and frequency-domain techniques. Both approaches have their own advantages and disadvantages. They are summarized in table 3 of [1].

The context of our problem arises in the frequency-domain convolutive BSS approach. A survey of these techniques is provided in [2]. The principle behind the frequency-domain techniques is that the convolutive mixture case is transformed (under the narrowband assumption) into complex-valued instantaneous mixture case for each frequency bin:

$$\mathbf{x}(f, t) = \mathbf{A}(f) \mathbf{s}(f, t) + \mathbf{v}(f, t) \quad (3)$$

where $f = 1 \dots F$ are the sub-band frequencies.

Several algorithms to estimate mixing matrix in case of instantaneous mixtures have been developed by the source separation community [3], [4]. Frequency domain convolutive BSS systems use one of these standard algorithms on each frequency bin to provide an estimate $\hat{\mathbf{A}}(f)$ of $\mathbf{A}(f)$. However, as the estimation is done on each frequency bin independently, this approach suffers from permutation and scaling indeterminacies in each sub-band f . Specifically, the estimated $\hat{\mathbf{A}}(f)$ is related to the true filter matrix $\mathbf{A}(f)$, for each f in the following form

$$\hat{\mathbf{A}}(f) = \mathbf{A}(f) \mathbf{\Lambda}(f) \mathbf{P}(f) \quad (4)$$

where $\mathbf{P}(f)$ is the frequency-dependent permutation matrix, $\mathbf{\Lambda}(f)$ is a diagonal matrix containing the arbitrary scaling factors.

The frequency-domain methods have to invariably solve the permutation and scaling indeterminacy to eventually estimate $\mathbf{A}(t)$ up to a unique global permutation and scaling $\hat{\mathbf{A}}(t) = \mathbf{A}(t)\mathbf{\Lambda}\mathbf{P}$.

A. Existing approaches to solve the described problem

There are two main kinds of methods to solve the permutation indeterminacy in the sub-bands of the estimated mixing filters [1].

The first set of techniques use consistency measures across the frequency sub-bands of the filters to recover the correct permutations, such as inter-frequency smoothness, etc. This category also includes the beamforming approach to identify the direction of arrival of sources and then adjust the permutations [5].

The second set of techniques use the consistency of the spectrum of the recovered signals to achieve the same. The consistency across the spectrum of the recovered signals is applicable for only those signals which have strong correlation across sub-bands, such as speech [6].

There has also been some effort to combine the above mentioned approaches for better performance [7]. Based on the different definitions of consistency, methods to correct permutations have been proposed in the literature. A categorization of these methods based on the definition of consistency has been presented in [1] (Table 4).

III. PROPOSED APPROACH: SPARSE FILTER MODELS

In our work, we propose to use a special type of consistency that can be assumed on the mixing filters: sparsity. That is, the number S of non-negligible coefficients in each filter $\mathbf{A}_{mn}(t)$ is significantly less than its length K . The motivation behind this approach is that the acoustic room impulse responses (RIR) tend to have a few reflection paths relative to its duration. Hence, one can approximate the acoustic RIRs by sparse filters.

The idea behind our approach is that when the mixing filters are sparse, the permutations in the sub-bands disturb the sparse structure of the filters. That is, each filter in the reconstructed filter matrix $\hat{\mathbf{A}}(t)$ after sub-band permutations will no longer be sparse. So, one can hope to correct the permutations by bringing back the sparse structure in each of the filters.

Ideally, the sparsity of the filter matrix is measured by its ℓ_0 -norm, defined as

$$\|\mathbf{A}(t)\|_0 = \sum_{m=1}^M \sum_{n=1}^N \sum_{t=1}^K |\mathbf{A}_{mn}(t)|^0. \quad (5)$$

Lesser the norm, sparser is the filter matrix. However, the ℓ_1 -norm of the filter matrix defined by

$$\|\mathbf{A}(t)\|_1 = \sum_{m=1}^M \sum_{n=1}^N \sum_{t=1}^K |\mathbf{A}_{mn}(t)| \quad (6)$$

is also a sparsity promoting norm. It has been shown by Donoho that working with ℓ_1 -norm is as effective as working with ℓ_0 -norm while looking for sparse solutions to linear

systems [8]. Furthermore ℓ_1 -norm being convex, has certain computational advantages over the former and is robust to noise.

In this paper, we concentrate primarily only on the permutation problem and will assume that the estimated filter sub-bands are free from any scaling indeterminacy. That is

$$\hat{\mathbf{A}}(f) = \mathbf{A}(f)\mathbf{P}(f), f = 1 \dots F. \quad (7)$$

Hence, when $\hat{\mathbf{A}}(t)$ is reconstructed using $\hat{\mathbf{A}}(f)$, given by Eqn. (7) one can expect an increase in $\|\hat{\mathbf{A}}(t)\|_1$. In Sec. IV we experimentally show that this is indeed the case.

Then in Sec. V we show that if there is no scaling indeterminacy, a simple algorithm which minimizes the ℓ_1 -norm of the filter matrix solves the permutation indeterminacy, even under noisy conditions.

IV. SOURCE PERMUTATIONS AND ℓ_1 -NORM

This section presents some preliminary experimental study on how the ℓ_1 -norm of the filter matrix $\mathbf{A}(t)$ is affected by permutation in the sub-bands.

For experimental purposes, 50 different filter matrices were synthetically created with $N = 5$ sources and $M = 3$ channels, and each filter having a length of $K = 1024$ with $S = 10$ non-zero coefficients. The non-zero coefficients were i.i.d. Gaussian with mean zero and variance one. The locations of non-zero coefficients were selected uniformly at random. Then for each such instance of filter matrix $\mathbf{A}(t)$ the discrete Fourier transform (DFT) $\hat{\mathbf{A}}_{mn}(f)$ was computed for each filter $\mathbf{A}_{mn}(t)$, to obtain the frequency-domain representation $\hat{\mathbf{A}}(f)$ of $\mathbf{A}(t)$. These filters were used in the following two kinds of experiments.

A. Random source permutations

In a practical scenario, each sub-band can have a random permutation. So, for each filter matrix in the frequency domain, the sources were permuted randomly in an increasing of number of sub-bands (chosen at random), and their ℓ_1 -norms were computed. The positive and negative frequency sub-bands were permuted identically.

For one such experimental instance, Fig. 2 shows the variation of ℓ_1 -norm against increasing number of randomly permuted sub-bands. The circle shows the norm of the true filter matrix. Each star represents the norm after randomly permuting the sub-bands at random locations. Note the gradual increase in the norm as the number of sub-bands being permuted increases. Similar experiments were conducted with combinations of $M = 3, N = 4$ and $M = 2, N = 3$ and $S = 10, 15, 20$, leading to similar observations.

B. Sensitivity of ℓ_1 -norm to permutations

In order to show that even a *single* permutation in only one sub-band can increase the norm, only two sources, chosen at random were permuted in increasing number of sub-bands.

For one such instance, Fig. 3 shows a plot of the variation in ℓ_1 -norm with the number of sub-bands permuted. The circle in the plot shows the ℓ_1 -norm of the true filter matrix. Each star shows the norm after permuting the sources 2 and 3.

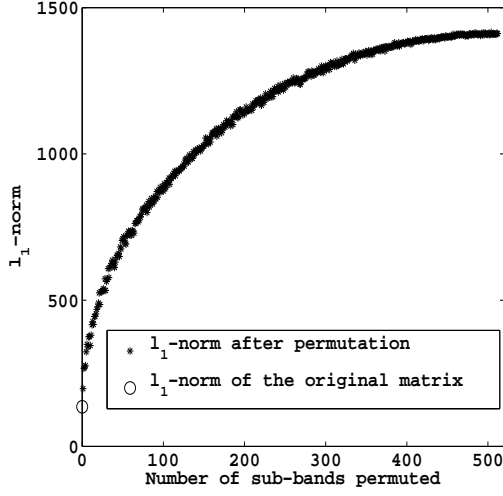


Fig. 2. The variation in the ℓ_1 -norm of the filter matrix against the number of sub-bands permuted when all the sources are permuted randomly

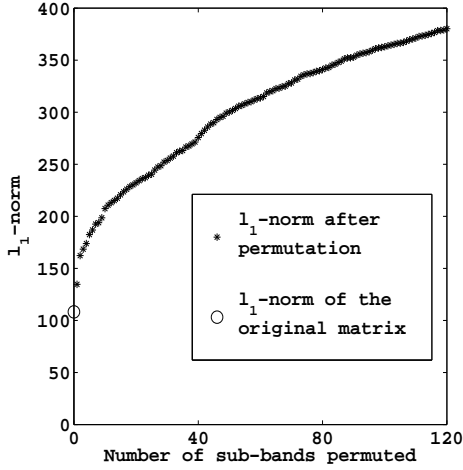


Fig. 3. The variation in the ℓ_1 -norm of the filter matrix against the number of sub-bands permuted when only sources 2 and 3 are permuted.

V. PROPOSED ALGORITHM

With the inspiration from the previous section, we present an algorithm in this section to solve the permutation indeterminacy.

Assumption: The estimate $\hat{\mathbf{A}}(f)$ of the sub-band coefficients $\mathbf{A}(f)$ are provided by some other independent technique and $\hat{\mathbf{A}}(f) = \mathbf{A}(f)\mathbf{P}(f)$, $f = 1 \dots F$.

Figure 4 shows the block diagram of our approach. The output from the first block is assumed to be available, and the algorithm presented in this section is what makes up the second block.

The absence of scaling indeterminacy while estimating $\hat{\mathbf{A}}(f)$ is not totally a realistic assumption. But we feel that this assumption aids the understanding of connections between the

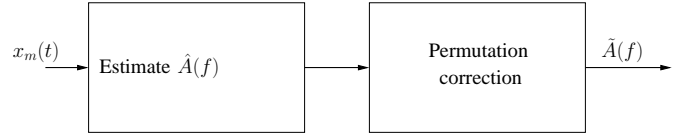


Fig. 4. Block diagram of our approach to solve permutation indeterminacy

time-domain sparsity of the filters and controlled disturbances in the frequency-domain. This understanding will be later useful in handling the arbitrary scaling of sub-band coefficients.

A. Description

We denote the set of all the possible source permutations by \mathcal{P} ($|\mathcal{P}| = N!$) and the inverse discrete Fourier transform by $IDFT$. At each sub-band (starting from the first sub-band), every possible permutation $\mathbf{P} \in \mathcal{P}$ is explored, keeping the rest of the sub-bands fixed, and only that permutation is retained which minimizes the ℓ_1 -norm. This ensures that ℓ_1 -norm of the filter matrix is lowered to the minimum possible extent by aligning that particular sub-band. At the end of one such iteration through all the sub-bands, the norm of the filter matrix would have significantly reduced. However, as the sub-bands were locally examined, the resulting norm may not be the minimum. Hence, the entire process of locally aligning the sub-bands is iterated several times until the difference in the norms between iterations falls below a preset threshold θ . Algorithm 1 contains the pseudo-code of the description.

Input: $\hat{\mathbf{A}}, \theta$: The estimated sub-band coefficients of $\mathbf{A}(t)$ and a threshold

Output: $\tilde{\mathbf{A}}$: The sub-band coefficient matrix after solving for the permutations

- (1) Initialize $\tilde{\mathbf{A}} \leftarrow \hat{\mathbf{A}}$;
- (2) Update all sub-bands;

foreach $f = 1 : F$ **do**
 $old\tilde{\mathbf{A}} \leftarrow \tilde{\mathbf{A}}$;
foreach $\mathbf{P} \in \mathcal{P}$ **do**
 $\tilde{\mathbf{A}}(f) \leftarrow \hat{\mathbf{A}}(f)\mathbf{P}$;
 $val(\mathbf{P}) \leftarrow \|IDFT(\tilde{\mathbf{A}}(f'))\|_1, f' = 1 \dots F$;
end
 $\mathbf{P}(f) \leftarrow \arg \min_{\mathbf{P} \in \mathcal{P}} val(\mathbf{P})$;
 $\tilde{\mathbf{A}}(f) \leftarrow \hat{\mathbf{A}}(f)\mathbf{P}(f)$;
end
- (3) Test if the algorithm should stop;
if $\|\hat{\mathbf{A}}(t)\|_1 \geq \|old\hat{\mathbf{A}}(t)\|_1 - \theta$ **then** Output $\tilde{\mathbf{A}}$ **else**
 Go to step (2)

Algorithm 1: Algorithm to solve the permutation indeterminacy by minimizing the ℓ_1 -norm of the time domain filter matrix

B. Objective of the algorithm

The aim of the algorithm is to obtain the sub-band matrix $\tilde{\mathbf{A}}(f)$ which would have the minimum ℓ_1 -norm in its corresponding time domain representation $\tilde{\mathbf{A}}(t)$. However, currently

we do not have analytical proof about the convergence of the algorithm to the global minimum but only ample experimental evidence that minimizing the ℓ_1 -norm corrects the permutations. It should also be noted that once the algorithm stops, the sources will be globally permuted.

C. Complexity

A brute force approach to solve the ideal ℓ_1 -minimization problem would need $N!^K K \log K$ operations. In our approach, each outer iteration needs to inspect $N! \times K$ permutations and at each step, an inverse FFT having a complexity of $K \log K$ has to be performed. Hence, the complexity in each outer iteration is $N!K^2 \log K$. This is still expensive because it grows in factorial with the number of sources, but it is tractable for small problem sizes.

VI. RESULTS

In this section we present an illustration of the algorithm presented above.

A. The no noise, no scaling case

Firstly, we consider the case where the sub-band coefficients are assumed to be estimated without noise and scaling ambiguity. 20 filter matrices with the specifications $N = 3, M = 2, K = 1024$ and $S = 10$ were created, transformed using DFT and the sub-bands permuted in a similar way as explained in Sec. IV-A. These were the input to the algorithm. The value of θ was set to 0.0001 in all the experiments.

The output was transformed back to time domain to compute the reconstruction error. In all the experiments, the output filters were identical to the true filters up to a global permutation and within a numerical precision in Matlab.

B. Effect of noise

The estimation of $\hat{\mathbf{A}}(f)$ by an actual BSS algorithm invariably involves some level of noise (as well as scaling, which we do not deal with here). Hence, the permutation solving algorithm needs to be robust to certain level of estimation noise. Experiments were conducted by permuting the sub-bands and adding noise to the coefficients:

$$\hat{\mathbf{A}}(f) = \mathbf{A}(f)\mathbf{P}(f) + \mathbf{N}(f) \quad (8)$$

where $\mathbf{N}(f)$ is i.i.d. complex Gaussian with mean zero and variance σ^2 .

For illustration, Fig. 5 shows an instance of the reconstructed filter matrix using Algorithm 1 with the input corrupted by additive complex Gaussian noise with $\sigma^2 = 0.2$. Each filter had 10 significant coefficients which have been faithfully recovered, along with some amount of noise.

For quantifying the effect of noise, an input-output SNR analysis was performed. The input SNR defined in Eqn. (9) was varied between -10 dB and 40 dB in steps of 5 dB and the corresponding output SNR defined in Eqn. (10) was computed. The problem size chosen for the SNR analysis was $N = 3, M = 2, K = 1024$. 20 independent experiments were

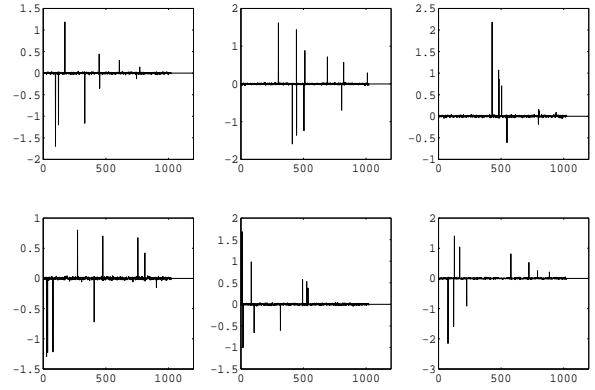


Fig. 5. Reconstructed filter matrix with additive noise in the sub-bands

conducted for each input SNR level, and the output SNRs were averaged to obtain each point.

$$SNR_{in} = 20 \log_{10} \left(\frac{\|\mathbf{A}(t)\|_2}{\|\mathbf{N}(t)\|_2} \right) \quad (9)$$

$$SNR_{out} = 20 \log_{10} \left(\frac{\|\mathbf{A}(t)\|_2}{\|\mathbf{A}(t) - \hat{\mathbf{A}}(t)\|_2} \right) \quad (10)$$

The experiments were repeated for $S = 10$ and $S = 25$.

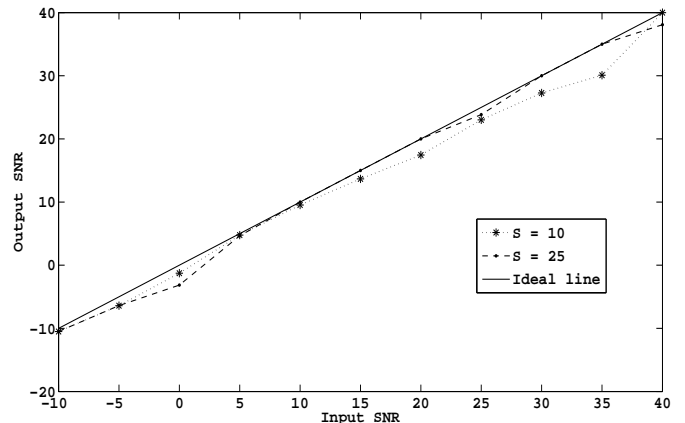


Fig. 6. Plot of output SNR versus input SNR in dB. The definitions of the SNR are given in equations (9) and (10)

Figure. 6 shows the variation of output SNR (in dB) with input SNR. Due to the absence of scaling, a perfectly reconstructed filter will have the same SNR as the input. The thick line shows the ideal relationship for reference. In the range of 5 to 10 dB input SNR, the curve for $S = 10$ coincides with the ideal line. For $S = 25$, the curve coincides with the ideal line for range of 5 to 20 dB input SNR. At other places, both the curves closely follow the ideal line suggesting perfect reconstruction in most number of experiments.

VII. CONCLUSION AND FUTURE WORK DIRECTION

Frequency-domain estimation of mixing matrices in convolutive BSS suffers from the indeterminacies of source

permutations and scaling of sub-bands. Hence, solving the permutation indeterminacy is an important aspect of such BSS systems. In this paper, it has been shown that in the absence of scaling, the ℓ_1 -norm of the filter matrix is very sensitive to permutations in sub-bands. An algorithm has been presented based on the minimization of ℓ_1 -norm of the filter matrix to solve for the permutations. Experimental results show that in the absence of scaling, the ℓ_1 -minimization principle to solve the permutations performs well even in the presence of noise.

Though the absence of scaling is not a realistic assumption, it can be a first step towards sparsity motivated permutation solving methods. Also, the complexity of the algorithm grows with the $N!$ and K^2 , which is expensive even for moderate values for the number of sources N and filter length K .

Our future work focusses on the theoretical analysis of the variation of ℓ_1 -norm with permutations. This may help to replace the combinatorial optimization step by an efficient convex optimization formulation. Further, working with real world mixing filters, and devising ℓ_1 -norm based methods to solve the sub-band scaling ambiguity are some of the interesting extensions to the presented work.

REFERENCES

- [1] M. S. Pedersen, J. Larsen, U. Kjems, and L. C. Parra, "A survey of convolutive blind source separation," *Springer Handbook on Speech Processing and Speech Communication*, 2006.
- [2] S. Makino, H. Sawada, R. Mukai, and S. Araki, "Blind source separation of convolutive mixtures of speech in frequency domain," *IEICE Trans. Fundamentals*, vol. E88-A, no. 7, pp. 1640–1655, July 2005.
- [3] A. J. Bell and T. J. Sejnowski, "An information-maximization approach to blind separation and blind deconvolution," *Neural Computation*, vol. 7, no. 6, pp. 1129–1159, 1995. [Online]. Available: <http://citeseer.ist.psu.edu/bell95informationmaximization.html>
- [4] P. Smaragdis, "Blind separation of convolved mixtures in the frequency domain," *Neural Computation*, vol. 22, pp. 21–34, 1998.
- [5] M. Z. Ikram and D. R. Morgan, "A beamforming approach to permutation alignment for multichannel frequency-domain blind source separation," in *Proc. of ICASSP*, May 2002, pp. 881–884.
- [6] N. Murata, S. Ikeda, and A. Ziehe, "An approach to blind source separation based on temporal structure of speech signals," RIKEN Brain Science Institute, Tech. Rep. BSIS Technical reports 98-2, 1998.
- [7] H. Sawada, R. Mukai, S. Araki, and S. Makino, "A robust and precise method for solving the permutation problem of frequency-domain blind source separation," *IEEE Trans. on Speech and Audio Processing*, vol. 12, no. 5, pp. 530–538, 2004.
- [8] D. Donoho, "For most large underdetermined systems of linear equations, the minimal ℓ_1 norm near-solution approximates the sparsest near-solution," *Communications on Pure and Applied Mathematics*, vol. 59, no. 7, pp. 907–934, July 2006.