

Sparse NMF – half-baked or well done?

Le Roux, J.; Weninger, F.J.; Hershey, J.R.

TR2015-023 March 2015

Abstract

Non-negative matrix factorization (NMF) has been a popular method for modeling audio signals, in particular for single-channel source separation. An important factor in the success of NMF-based algorithms is the "quality" of the basis functions that are obtained from training data. In order to model rich signals such as speech or wide ranges of non-stationary noises, NMF typically requires using a large number of basis functions. However, without additional constraints, using a large number of bases leads to trivial solutions where the bases can indiscriminately model any signal. Two main approaches have been considered to cope with this issue: introducing sparsity on the activation coefficients, or skipping training altogether and randomly selecting basis functions as a subset of the training data ("exemplar-based NMF"). Surprisingly, the sparsity route is widely regarded as leading to similar or worse results than the simple and extremely efficient (no training!) exemplar-based approach. Only a small fraction of researchers have realized that sparse NMF works well if implemented correctly. However, to our knowledge, no thorough comparison has been presented in the literature, and many researchers in the field may remain unaware of this fact. We review exemplar-based NMF as well as two versions of sparse NMF, a simplistic ad hoc one and a principled one, giving a detailed derivation of the update equations for the latter in the general case of beta divergences, and we perform a thorough comparison of the three methods on a speech separation task using the 2nd CHiME Speech Separation and Recognition Challenge dataset. Results show that, contrary to a popular belief in the community, learning basis functions using NMF with sparsity, if done the right way, leads to significant gains in source-to-distortion ratio with respect to both exemplar-based NMF and the ad hoc implementation of sparse NMF.

Mitsubishi Electric Research Laboratories

This work may not be copied or reproduced in whole or in part for any commercial purpose. Permission to copy in whole or in part without payment of fee is granted for nonprofit educational and research purposes provided that all such whole or partial copies include the following: a notice that such copying is by permission of Mitsubishi Electric Research Laboratories, Inc.; an acknowledgment of the authors and individual contributions to the work; and all applicable portions of the copyright notice. Copying, reproduction, or republishing for any other purpose shall require a license with payment of fee to Mitsubishi Electric Research Laboratories, Inc. All rights reserved.

Sparse NMF – half-baked or well done?

Jonathan Le Roux

Mitsubishi Electric Research Labs (MERL)
Cambridge, MA, USA
leroux@merl.com

Felix Weninger

TUM
Munich, Germany
felix@weninger.de

John R. Hershey

Mitsubishi Electric Research Labs (MERL)
Cambridge, MA, USA
hershey@merl.com

Abstract

Non-negative matrix factorization (NMF) has been a popular method for modeling audio signals, in particular for single-channel source separation. An important factor in the success of NMF-based algorithms is the “quality” of the basis functions that are obtained from training data. In order to model rich signals such as speech or wide ranges of non-stationary noises, NMF typically requires using a large number of basis functions. However, without additional constraints, using a large number of bases leads to trivial solutions where the bases can indiscriminately model any signal. Two main approaches have been considered to cope with this issue: introducing sparsity on the activation coefficients, or skipping training altogether and randomly selecting basis functions as a subset of the training data (“exemplar-based NMF”). Surprisingly, the sparsity route is widely regarded as leading to similar or worse results than the simple and extremely efficient (no training!) exemplar-based approach. Only a small fraction of researchers have realized that sparse NMF works well if implemented correctly. However, to our knowledge, no thorough comparison has been presented in the literature, and many researchers in the field may remain unaware of this fact. We review exemplar-based NMF as well as two versions of sparse NMF, a simplistic ad hoc one and a principled one, giving a detailed derivation of the update equations for the latter in the general case of beta divergences, and we perform a thorough comparison of the three methods on a speech separation task using the 2nd CHiME Speech Separation and Recognition Challenge dataset. Results show that, contrary to a popular belief in the community, learning basis functions using NMF with sparsity, if done the right way, leads to significant gains in source-to-distortion ratio with respect to both exemplar-based NMF and the ad hoc implementation of sparse NMF.

1 Contributions of this report

- Experimental comparison of exemplar-based NMF (ENMF), sparse NMF with basis renormalization in the objective function (SNMF), sparse NMF with basis renormalization after each update (NMF+S) on a supervised audio source separation task: **new**
- Detailed derivation of multiplicative update equations for SNMF with beta divergence from a general perspective of gradients with unit-norm constraints: **discussion on gradients with unit-norm constraints adapted from our previous work [1]; some elements similar to the derivation for the convolutive NMF case in [2]**

- Proof that the tangent gradient and the natural gradient are equivalent for the unit L^2 -norm constraint: **new**
- An efficient MATLAB[®] implementation of SNMF with β -divergence: **new**

2 Non-negative matrix factorization and sparsity

Non-negative matrix factorization (NMF) is a popular approach to analyzing non-negative data. It has been successfully used over the past 15 years in a surprisingly wide range of applications (e.g., [3,4] and references therein). In many cases, it is used as a dimension reduction technique, where the goal is to factor some non-negative data matrix $\mathbf{M} \in \mathbb{R}^{F \times T}$ into the product of two non-negative matrices $\mathbf{W} \in \mathbb{R}^{F \times R}$ and $\mathbf{H} \in \mathbb{R}^{R \times T}$, where the inner dimension R of the product is much smaller than either dimensions F and T of the original matrix. Without any sparsity penalty on either or both of the factors \mathbf{W} and \mathbf{H} , applying NMF with R of the order of, or larger than F or T would lead to meaningless results, as there are infinitely many different factorizations that can exactly lead to \mathbf{M} . In the context of audio signal processing for example, NMF is typically applied to the magnitude or power spectrogram \mathbf{M} of a corpus of sounds. The matrix \mathbf{W} can then be interpreted as a set of spectral basis functions (also referred to as a dictionary), and the matrix \mathbf{H} as the collection of their activations at each time frame. However, whatever the size of the corpus, vanilla NMF will be unable to learn a meaningful dictionary with more elements than there are time frequency bins. To alleviate this issue, one can introduce sparsity penalties on the factors, for example on the activations \mathbf{H} in the context of audio. This ensures that, even though there are many elements in the dictionary, only a small number will be active at the same time. This enables the model to learn a meaningful representation of the data with many more elements than what vanilla NMF would allow. Another approach to solving this issue of representing rich datasets is to obtain the dictionary not by learning, but by sampling the data, i.e., by using randomly selected data samples as the dictionary elements. Both approaches have been explored. We present here an overview of these approaches, and thoroughly compare their performance on a supervised audio source separation task.

In unsupervised application scenarios, both \mathbf{W} and \mathbf{H} are optimized at test time so that their product fits the data; in supervised application scenarios, \mathbf{W} is typically learned at training time, again by optimizing \mathbf{W} and $\mathbf{H}^{\text{training}}$ so that their product fits some training data [2, 5], while at test time, only \mathbf{H} is optimized given \mathbf{W} . In any case, both types of scenario involve solving the following problem:

$$\overline{\mathbf{W}}, \overline{\mathbf{H}} = \arg \min_{\mathbf{W}, \mathbf{H}} D(\mathbf{M} | \mathbf{W}\mathbf{H}) + \mu |\mathbf{H}|_1, \quad (1)$$

where D is a cost function that is minimized when $\mathbf{M} = \mathbf{W}\mathbf{H}$. We abuse notations on D and assume that the cost function computed on matrices is equal to the sum of the element-wise costs over all elements. Here we use the β -divergence, D_β , which for $\beta = 0$ yields the Itakura-Saito (IS) distance, for $\beta = 1$ yields the generalized Kullback-Leibler (KL) divergence, and for $\beta = 2$ yields the Euclidean distance. The β -divergence is defined in general as:

$$D_\beta(x|y) \stackrel{\text{def}}{=} \begin{cases} \frac{1}{\beta(\beta-1)}(x^\beta - y^\beta - \beta y^{\beta-1}(x-y)) & \text{if } \beta \in \mathbb{R} \setminus \{0, 1\} \\ x(\log x - \log y) + (y-x) & \text{if } \beta = 1 \\ \frac{x}{y} - \log \frac{x}{y} - 1 & \text{if } \beta = 0 \end{cases}. \quad (2)$$

The L^1 sparsity constraint with weight μ is added to favor solutions where few basis vectors are active at a time. To avoid scaling indeterminacies, some normalization constraint is typically assumed on \mathbf{W} , e.g., $\|\mathbf{W}\|_2 = 1$ where $\|\cdot\|_2$ denotes the L^2 norm.

If the NMF bases \mathbf{W} are held fixed, such as at test time in supervised scenarios, the optimal activations $\hat{\mathbf{H}}$ are estimated such that

$$\hat{\mathbf{H}} = \arg \min_{\mathbf{H}} D(\mathbf{M} | \mathbf{W}\mathbf{H}) + \mu |\mathbf{H}|_1. \quad (3)$$

A convenient algorithm [6] for minimizing (3) that preserves non-negativity of \mathbf{H} by multiplicative updates is given by iterating

$$\mathbf{H} \leftarrow \mathbf{H} \otimes \frac{\mathbf{W}^\top (\mathbf{M} \otimes \mathbf{\Lambda}^{\beta-2})}{\mathbf{W}^\top \mathbf{\Lambda}^{\beta-1} + \mu} \quad (4)$$

until convergence, with $\Lambda := \mathbf{W}\mathbf{H}$ and exponents applied element-wise.

It remains to determine how to obtain the NMF bases \mathbf{W} . In NMF without a sparsity cost on \mathbf{H} , i.e. $\mu = 0$, an optimization scheme that alternates between updating \mathbf{H} given \mathbf{W} as in (4):

$$\mathbf{H} \leftarrow \mathbf{H} \otimes \frac{\mathbf{W}^\top (\mathbf{M} \otimes \Lambda^{\beta-2})}{\mathbf{W}^\top \Lambda^{\beta-1}}, \quad (5)$$

and updating \mathbf{W} given \mathbf{H} in a similar way:

$$\mathbf{W} \leftarrow \mathbf{W} \otimes \frac{(\Lambda^{\beta-2} \otimes \mathbf{M}) \mathbf{H}^\top}{\Lambda^{\beta-1} \mathbf{H}^\top}, \quad (6)$$

where Λ is kept updated to $\mathbf{W}\mathbf{H}$, can be shown to be non-decreasing with respect to the objective function [6]. However, optimizing (1) including the sparsity cost requires special care concerning the normalization constraint on \mathbf{W} , as we shall see below. Another option that has been explored is to avoid training altogether and simply select the NMF bases \mathbf{W} randomly from training samples. We shall now investigate these options in more details.

3 Obtaining NMF Bases

3.1 Learning the bases: Sparse NMF (SNMF) or NMF with sparsity (NMF+S)

Since the L^1 sparsity constraint on \mathbf{H} in (1) is not scale-invariant, it can be trivially minimized by scaling of the factors. As mentioned above, this scaling indeterminacy can be avoided by imposing some sort of normalization on the bases \mathbf{W} . Two approaches for enforcing this normalization have been predominant in the literature.

A tempting but ad hoc practice [7] is to perform standard NMF optimization with sparsity in the \mathbf{H} update, using (4) then (6), and rescaling \mathbf{W} and \mathbf{H} such that \mathbf{W} has unit norm after each iteration. We here denote this approach NMF+S, for NMF with sparsity.

Another approach [2, 8] is to directly reformulate the objective function including a column-wise normalized version of \mathbf{W} [9], leading to an approach which we refer to as sparse NMF (SNMF):

$$\widetilde{\mathbf{W}}, \widetilde{\mathbf{H}} = \arg \min_{\mathbf{W}, \mathbf{H}} D_\beta(\mathbf{S} \mid \widetilde{\mathbf{W}}\mathbf{H}) + \mu \|\mathbf{H}\|_1, \quad (7)$$

where $\widetilde{\mathbf{W}} = \left[\frac{\mathbf{w}_1}{\|\mathbf{w}_1\|} \cdots \frac{\mathbf{w}_R}{\|\mathbf{w}_R\|} \right]$ is the column-wise normalized version of \mathbf{W} . The update for \mathbf{H} given \mathbf{W} is essentially the same as before, except that it now involves $\widetilde{\mathbf{W}}$ instead of \mathbf{W} :

$$\mathbf{H} \leftarrow \mathbf{H} \otimes \frac{\widetilde{\mathbf{W}}^\top (\mathbf{M} \otimes \Lambda^{\beta-2})}{\widetilde{\mathbf{W}}^\top \Lambda^{\beta-1} + \mu},$$

now with $\Lambda := \widetilde{\mathbf{W}}\mathbf{H}$. However, the update for \mathbf{W} given \mathbf{H} will be different from (6), which was obtained without normalization. We derive this update equation in Section 5.

Importantly, as pointed out by [9], SNMF is not equivalent to NMF+S. Indeed, the a posteriori rescaling of \mathbf{W} and \mathbf{H} in NMF+S changes the value of the sparsity cost function, which is not the case for (7) as \mathbf{H} is not rescaled. One can in fact see that the NMF+S optimization scheme may actually lead to an increase in the cost function, as shown in Figure 1. In the course of optimization with SNMF, the evolution of the KL divergence part and the sparsity part of the objective function shows that a trade-off between the two is reached, while the total cost function decreases monotonously; with NMF+S, an increase of the sparsity part and of the total cost function can be observed after a few iterations. In practice, the iterations are stopped when the total objective function starts increasing.

Multiplicative update algorithms to optimize (7) with the Euclidean distance and KL divergence have appeared in the literature. We only found the case of an arbitrary $\beta \geq 0$ presented in [2] for the case of convolutive NMF. We give here a detailed derivation for the case of NMF, using a similar perspective to [2], but including a more general account on optimization of objective function with normalized parameters.

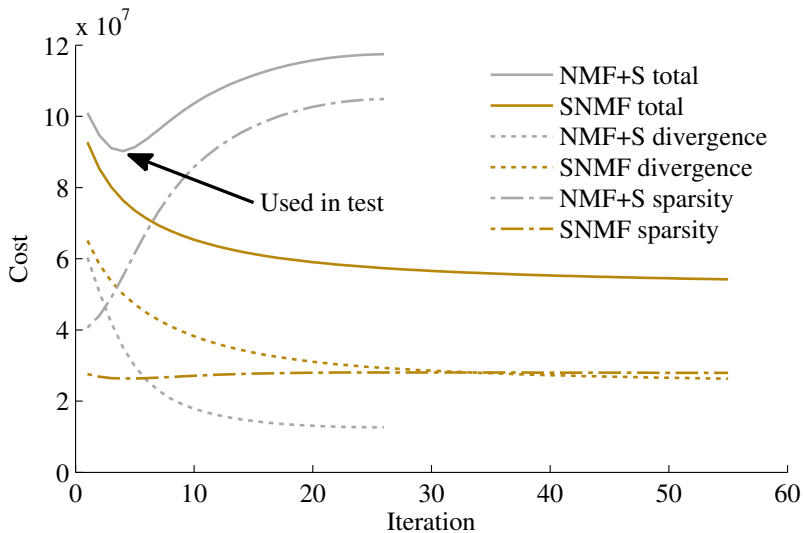


Figure 1: Evolution of the cost function for SNMF and NMF+S with the KL divergence ($\beta = 1$), for $\mu = 5$ and $K = 1000$, while training on a database of speech power spectra as described in Section 7.

3.2 Sampling the bases: Exemplar-based NMF (ENMF)

Recently, exemplar-based approaches, where every basis function corresponds to an observation in the training data, have become popular for large-scale factorizations of audio signals [10, 11]. It has been argued that they lead to similar or better performance compared with NMF with sparsity, with the extra advantage that they do not require training [12] (see also [13] for a similar study on complex NMF [14]). It is however hard to find thorough comparisons of ENMF with NMF methods where the basis functions are trained using sparsity constraints. We will show experimentally on a large supervised audio source separation task that, even though it may have been true that ENMF performed better when compared with the ad-hoc NMF+S, it is actually not the case anymore if compared with SNMF, the proper implementation of sparsity in NMF.

4 Optimization of objective functions with normalized parameters

We explained above how in NMF+S the update of \mathbf{W} is performed in two steps: first a multiplicative update without normalization constraint, followed by normalization to ensure that the bases \mathbf{W} respect the unit-norm constraint. A big issue with this approach is that it requires a corresponding rescaling of \mathbf{H} , and this rescaling may increase the contribution of the sparsity term to the objective function, potentially resulting in an increase of the objective function compared to the value before normalization of \mathbf{W} . Even without sparsity term, while the normalization of \mathbf{W} and rescaling of \mathbf{H} do not change the value of the objective function, alleviating the scaling indeterminacy between \mathbf{W} and \mathbf{H} may still help the optimization process. Updating \mathbf{W} while taking into account as much as possible the unit-norm normalization on the bases is thus important with or without a sparsity penalty.

A heuristic approach to deriving multiplicative update equations for NMF relies on computing the gradient of the objective function and splitting it into the difference of two non-negative terms. We thus consider here the general problem of computing a notion of the gradient of an objective function with normalized parameters. Two notions of gradient appear to be natural candidates for consideration here: the natural gradient [15] of the objective function with (non-normalized) parameters on the unit-norm constraint manifold, and the classical gradient (w.r.t. the non-normalized parameters) of the objective function in which the bases appear explicitly normalized [9, 16]. It turns out that for a unit-norm normalization with the L^2 -norm, these two notions of gradient are actually equivalent.

4.1 Gradient with normalized parameters

In general, let $\mathcal{I}(w)$ be an objective function with parameter $w \in \mathbb{R}^n$, and let $\tilde{\mathcal{I}}(w) = \mathcal{I}\left(\frac{w}{\|w\|}\right)$, where $\|w\|$ is any vector norm, be the corresponding objective function where the parameters are explicitly normalized. In the following, we shall denote $v = \left. \frac{\partial \|u\|}{\partial u} \right|_{u=w}$ and $\tilde{w} = \frac{1}{\|w\|}w$. Then

$$\frac{\partial \tilde{\mathcal{I}}}{\partial w_i}(w) = \sum_j \frac{\partial \tilde{w}_j}{\partial w_i} \cdot \frac{\partial \mathcal{I}}{\partial w_j}(\tilde{w}), \quad (8)$$

and using

$$\frac{\partial \tilde{w}_j}{\partial w_i} = \begin{cases} \frac{1}{\|w\|} - \frac{w_i v_i}{\|w\|^2} & \text{if } j = i \\ -\frac{w_j v_i}{\|w\|^2} & \text{if } j \neq i \end{cases} \quad (9)$$

we obtain

$$\frac{\partial \tilde{\mathcal{I}}}{\partial w_i}(w) = \frac{1}{\|w\|} \frac{\partial \mathcal{I}}{\partial w_i}(\tilde{w}) - v_i \sum_j \frac{w_j}{\|w\|^2} \cdot \frac{\partial \mathcal{I}}{\partial w_j}(\tilde{w}), \quad (10)$$

and finally

$$\nabla \tilde{\mathcal{I}}(w) = \frac{1}{\|w\|} \left(\nabla \mathcal{I}(\tilde{w}) - (\tilde{w} \cdot \nabla \mathcal{I}(\tilde{w}))v \right), \quad (11)$$

or written in matrix form

$$\nabla \tilde{\mathcal{I}}(w) = \frac{1}{\|w\|} (\text{Id} - v \tilde{w}^\top) \nabla \mathcal{I}(\tilde{w}). \quad (12)$$

For the L^2 norm, we have $v = \frac{w}{\|w\|}$, and Eq. (12) simplifies to

$$\nabla \tilde{\mathcal{I}}(w) = \frac{1}{\|w\|} (\text{Id} - vv^\top) \nabla \mathcal{I}(\tilde{w}). \quad (13)$$

This is all we need to compute the update equations in Section 5, but we study further the relationship between this notion of gradient and the so-called tangent and natural gradients.

4.2 Relation to the tangent gradient

The tangent gradient, introduced formally in [17] (it had been used earlier, e.g., by Krasulina [18] or Reddy et al. [19]), is defined in general for a unit-norm constraint (with arbitrary norm) as the orthogonal projection of the gradient onto the tangent plane to the unit-norm surface at \tilde{w} . It can be computed as

$$\nabla^{(T)} \mathcal{I}(\tilde{w}) = \left(\text{Id} - \frac{v v^\top}{\|v\|_2^2} \right) \nabla \mathcal{I}(\tilde{w}), \quad (14)$$

where $v = \left. \frac{\partial \|u\|}{\partial u} \right|_{u=\tilde{w}}$. Note that even for a unit-norm constraint with any arbitrary norm $\|\cdot\|$, it is still the L^2 norm of v which appears in the expression of the tangent gradient, due to the orthogonal projection being performed in Euclidean space.

Interestingly, for the L^2 -norm constraint, we see from Eq. 13 that the gradient $\nabla \tilde{\mathcal{I}}$ is equal up to a scaling factor to the tangent gradient of \mathcal{I} taken at point \tilde{w} . Indeed, we then have $v = \tilde{w}$ and $\|v\|_2 = 1$, and thus

$$\nabla \tilde{\mathcal{I}}(w) = \frac{1}{\|w\|} \nabla^{(T)} \mathcal{I}(\tilde{w}). \quad (15)$$

4.3 Relation to the natural gradient

The natural gradient [15] of a function \mathcal{I} defined on a manifold \mathcal{S} is a modification of its standard gradient according to the local curvature of the parameter space. It can be obtained from the standard gradient $\nabla \mathcal{I}$ as

$$\nabla^{(N)} \mathcal{I}(s) = \mathbf{G}^{-1}(s) \nabla \mathcal{I}(s), \forall s \in \mathcal{S}, \quad (16)$$

where \mathcal{S} is here the constraint manifold, and \mathbf{G} is the Riemannian metric tensor of the constraint manifold at s .

In the particular case of the L^2 -norm constraint, the tangent gradient of the original objective function \mathcal{I} can be shown to be equal, on the unit-norm surface, to its natural gradient:

$$\nabla^{(N)}\mathcal{I}(\tilde{w}) = \nabla^{(T)}\mathcal{I}(\tilde{w}), \quad \forall \tilde{w} \text{ s.t. } \|\tilde{w}\|_2 = 1. \quad (17)$$

However, while this is reported as a fact in the literature [17], to our knowledge, no proof is available. We give a detailed proof in Appendix A.

Combining with Eq. 15, we obtain:

$$\nabla\tilde{\mathcal{I}}(\tilde{w}) = \nabla^{(T)}\mathcal{I}(\tilde{w}) = \nabla^{(N)}\mathcal{I}(\tilde{w}), \quad \forall \tilde{w} \text{ s.t. } \|\tilde{w}\|_2 = 1. \quad (18)$$

Altogether, starting from a point \tilde{w} on the unit-norm manifold for L^2 , the gradient of the objective function $\tilde{\mathcal{I}}$ with explicitly normalized parameters is equal to the natural gradient of the original objective function \mathcal{I} as defined on the unit-norm manifold.

5 Update equations for SNMF with β -divergence

We come back to the derivation of the update equations for (7). The update equations for \mathbf{H} are unchanged, except that they now use the normalized version $\tilde{\mathbf{W}}$ of \mathbf{W} . Let us derive multiplicative update equations for \mathbf{W} by splitting the expression for the gradient into positive and negative parts, similarly to [20].

We first compute the gradient of $D_\beta(\mathbf{M} | \mathbf{WH})$ with respect to \mathbf{W}_i :

$$\nabla_{\mathbf{W}_i} D_\beta(\mathbf{M} | \mathbf{WH}) = \left((\mathbf{WH})^{\beta-2} \otimes (\mathbf{WH} - \mathbf{M}) \right) \mathbf{H}_i^\top, \quad (19)$$

where the \otimes product and the operation of raising to an exponent are considered element-wise. To avoid clutter, we will use the notation $\mathbf{\Lambda} = \tilde{\mathbf{W}}\mathbf{H}$. From there and the result of Section 4.1, we can compute the gradient of $D_\beta(\mathbf{M} | \tilde{\mathbf{W}}\mathbf{H})$ with respect to \mathbf{W}_i :

$$\nabla_{\mathbf{W}_i} D_\beta(\mathbf{M} | \mathbf{\Lambda}) = \frac{1}{\|\tilde{\mathbf{W}}_i\|} \left(\text{Id} - \tilde{\mathbf{W}}_i \tilde{\mathbf{W}}_i^\top \right) \left(\mathbf{\Lambda}^{\beta-2} \otimes (\mathbf{\Lambda} - \mathbf{M}) \right) \mathbf{H}_i^\top. \quad (20)$$

We split the gradient into a positive part $[\nabla_{\mathbf{W}_i} D_\beta(\mathbf{M} | \mathbf{\Lambda})]_+$ and a negative part $[\nabla_{\mathbf{W}_i} D_\beta(\mathbf{M} | \mathbf{\Lambda})]_-$:

$$[\nabla_{\mathbf{W}_i} D_\beta(\mathbf{M} | \mathbf{\Lambda})]_+ = \frac{1}{\|\tilde{\mathbf{W}}_i\|} \left(\mathbf{\Lambda}^{\beta-1} + \tilde{\mathbf{W}}_i \tilde{\mathbf{W}}_i^\top \left(\mathbf{\Lambda}^{\beta-2} \otimes \mathbf{M} \right) \right) \mathbf{H}_i^\top \quad (21)$$

$$[\nabla_{\mathbf{W}_i} D_\beta(\mathbf{M} | \mathbf{\Lambda})]_- = \frac{1}{\|\tilde{\mathbf{W}}_i\|} \left(\mathbf{\Lambda}^{\beta-2} \mathbf{M} + \tilde{\mathbf{W}}_i \tilde{\mathbf{W}}_i^\top \mathbf{\Lambda}^{\beta-1} \right) \mathbf{H}_i^\top \quad (22)$$

We thus obtain the following multiplicative updates for \mathbf{W}_i :

$$\mathbf{W}_i \leftarrow \mathbf{W}_i \otimes \frac{[\nabla_{\mathbf{W}_i} D_\beta(\mathbf{M} | \mathbf{\Lambda})]_-}{[\nabla_{\mathbf{W}_i} D_\beta(\mathbf{M} | \mathbf{\Lambda})]_+} \quad (23)$$

$$= \mathbf{W}_i \otimes \frac{\left(\mathbf{\Lambda}^{\beta-2} \otimes \mathbf{M} + \tilde{\mathbf{W}}_i \tilde{\mathbf{W}}_i^\top \mathbf{\Lambda}^{\beta-1} \right) \mathbf{H}_i^\top}{\left(\mathbf{\Lambda}^{\beta-1} + \tilde{\mathbf{W}}_i \tilde{\mathbf{W}}_i^\top \left(\mathbf{\Lambda}^{\beta-2} \otimes \mathbf{M} \right) \right) \mathbf{H}_i^\top}. \quad (24)$$

Note that this update can be computed efficiently for the whole matrix using the following equivalent form:

$$\mathbf{W} \leftarrow \mathbf{W} \otimes \frac{\left(\mathbf{\Lambda}^{\beta-2} \otimes \mathbf{M} \right) \mathbf{H}^\top + \tilde{\mathbf{W}} \otimes \left(\mathbf{1}\mathbf{1}^\top \left(\tilde{\mathbf{W}} \otimes \left(\mathbf{\Lambda}^{\beta-1} \mathbf{H}^\top \right) \right) \right)}{\mathbf{\Lambda}^{\beta-1} \mathbf{H}^\top + \tilde{\mathbf{W}} \otimes \left(\mathbf{1}\mathbf{1}^\top \left(\tilde{\mathbf{W}} \otimes \left(\left(\mathbf{\Lambda}^{\beta-2} \otimes \mathbf{M} \right) \mathbf{H}^\top \right) \right) \right)}, \quad (25)$$

where $\mathbf{1}$ is a column vector with all elements equal to 1. The operation $\tilde{\mathbf{W}} \otimes (\mathbf{1}\mathbf{v}^\top)$, where \mathbf{v} is a column vector, can be carried out efficiently in Matlab using a command such as

Algorithm 1 NMF+S: NMF with Sparsity, half-baked

Inputs: \mathbf{M} , R , $\beta \geq 0$, $\mu \geq 0$
 \mathbf{W} initialized randomly or by sampling from the data, \mathbf{H} initialized randomly
 $\mathbf{H} \leftarrow \tilde{\mathbf{H}}$, where $\tilde{\mathbf{H}} = \left[\|\mathbf{w}_1\| \mathbf{h}_1; \dots; \|\mathbf{w}_R\| \mathbf{h}_R \right]$
 $\mathbf{W} \leftarrow \tilde{\mathbf{W}}$
 $\Lambda = \mathbf{W}\mathbf{H}$
repeat
 $\mathbf{H} \leftarrow \mathbf{H} \otimes \frac{\mathbf{W}^\top (\mathbf{M} \otimes \Lambda^{\beta-2})}{\mathbf{W}^\top \Lambda^{\beta-1} + \mu}$
 $\Lambda = \mathbf{W}\mathbf{H}$
 $\mathbf{W} \leftarrow \mathbf{W} \otimes \frac{(\Lambda^{\beta-2} \otimes \mathbf{M})\mathbf{H}^\top}{\Lambda^{\beta-1}\mathbf{H}^\top}$
 $\mathbf{H} \leftarrow \tilde{\mathbf{H}}$, $\mathbf{W} \leftarrow \tilde{\mathbf{W}}$
 $\Lambda = \mathbf{W}\mathbf{H}$
until convergence
return \mathbf{W} , \mathbf{H}

Algorithm 2 SNMF: Sparse NMF, well done

Inputs: \mathbf{M} , R , $\beta \geq 0$, $\mu \geq 0$
 \mathbf{W} initialized randomly or by sampling from the data, \mathbf{H} initialized randomly
 $\mathbf{W} \leftarrow \tilde{\mathbf{W}}$
 $\Lambda = \mathbf{W}\mathbf{H}$
repeat
 $\mathbf{H} \leftarrow \mathbf{H} \otimes \frac{\tilde{\mathbf{W}}^\top (\mathbf{M} \otimes \Lambda^{\beta-2})}{\tilde{\mathbf{W}}^\top \Lambda^{\beta-1} + \mu}$
 $\Lambda = \mathbf{W}\mathbf{H}$
 $\mathbf{W} \leftarrow \mathbf{W} \otimes \frac{(\Lambda^{\beta-2} \otimes \mathbf{M})\mathbf{H}^\top + \tilde{\mathbf{W}} \otimes (\mathbf{1}\mathbf{1}^\top (\tilde{\mathbf{W}} \otimes (\Lambda^{\beta-1}\mathbf{H}^\top)))}{\Lambda^{\beta-1}\mathbf{H}^\top + \tilde{\mathbf{W}} \otimes (\mathbf{1}\mathbf{1}^\top (\tilde{\mathbf{W}} \otimes ((\Lambda^{\beta-2} \otimes \mathbf{M})\mathbf{H}^\top)))}$
 $\mathbf{W} \leftarrow \tilde{\mathbf{W}}$
 $\Lambda = \mathbf{W}\mathbf{H}$
until convergence
return \mathbf{W} , \mathbf{H}

`bsxfun(@times, Wtilde, v')`. For convenience of computation, we renormalize \mathbf{W} at the end of each iteration, but we do not rescale \mathbf{H} accordingly. Because only $\tilde{\mathbf{W}}$ appears in the objective function, this does not change the value of the objective function, and it avoids having to keep both a normalized and an unnormalized version of \mathbf{W} . The algorithm would lead to exactly the same $\tilde{\mathbf{W}}$ if we did not normalize \mathbf{W} and performed the iterative updates. The training procedure with SNMF is summarized in Algorithm 2. The source code for an efficient MATLAB[®] implementation is publicly available for research-only purposes at the following address: <ftp://ftp.merl.com/pub/leroux/sparseNMF.zip>.

As a comparison, for NMF+S, the updates for \mathbf{W} are simply:

$$\mathbf{W} \leftarrow \mathbf{W} \otimes \frac{(\Lambda^{\beta-2} \otimes \mathbf{M})\mathbf{H}^\top}{\Lambda^{\beta-1}\mathbf{H}^\top}, \quad (26)$$

followed by re-normalization of \mathbf{W} and the corresponding rescaling of \mathbf{H} . The training procedure with NMF+S is summarized in Algorithm 1.

6 Source separation using NMF

NMF has been a very popular algorithm in the audio community, and in particular it is commonly used for challenging single-channel audio source separation tasks, such as speech enhancement in the presence of non-stationary noises [10, 21]. In this context, the basic idea is to represent the

features of the sources via sets of basis functions and their activation coefficients, one set per source. Mixtures of signals are then analyzed using the concatenated sets of basis functions, and each source is reconstructed using its corresponding activations and basis set. In order to accurately represent each source, the sets of basis functions need to be large, counting potentially hundreds or thousands of basis functions. While vanilla NMF can be useful to obtain low-rank approximations, it leads to trivial solutions when used with a number of basis functions comparable to or higher than the dimension of the feature space in which it operates.

In the context of audio source separation, NMF operates on a matrix of F -dimensional non-negative spectral features, usually the power or magnitude spectrogram of the mixture, $\mathbf{M} = [\mathbf{m}_1 \cdots \mathbf{m}_T]$, where T is the number of frames and $\mathbf{m}_t \in \mathbb{R}_+^F$, $t = 1, \dots, T$ are obtained by short-time Fourier analysis of the time-domain signal. With L sources, each source $l \in \{1, \dots, L\}$ is represented using a matrix containing R_l non-negative basis column vectors, $\mathbf{W}^l = \{\mathbf{w}_r^l\}_{r=1}^{R_l}$, multiplied by a matrix of activation column vectors $\mathbf{H}^l = \{\mathbf{h}_t^l\}_{t=1}^T$, for each time t . The r th row of \mathbf{H}^l contains the activations for the corresponding basis \mathbf{w}_r^l at each time t . From this, a factorization

$$\mathbf{M} \approx \sum_l \mathbf{S}^l \approx [\mathbf{W}^1 \cdots \mathbf{W}^S][\mathbf{H}^1; \cdots; \mathbf{H}^S] = \mathbf{W}\mathbf{H} \quad (27)$$

is obtained, where we use the notation $[\mathbf{a}; \mathbf{b}]$ for $[\mathbf{a}^\top \mathbf{b}^\top]^\top$. An approach related to Wiener filtering is typically used to reconstruct each source while ensuring that the source estimates sum to the mixture:

$$\hat{\mathbf{S}}^l = \frac{\mathbf{W}^l \mathbf{H}^l}{\sum_l \mathbf{W}^l \mathbf{H}^l} \otimes \mathbf{M}, \quad (28)$$

where \otimes denotes element-wise multiplication and the quotient line element-wise division. In our study, all \mathbf{W}^l are learned or obtained in advance from training data, using one of the algorithms of Section 3.1. At test time, only the activation matrices $\hat{\mathbf{H}} = [\hat{\mathbf{H}}^1; \cdots; \hat{\mathbf{H}}^S]$ are estimated, using the update equations (4), so as to find a (local) minimum of (3). This is called *supervised NMF* [22]. In the supervised case, the activations for each frame are independent from the other frames ($\mathbf{m}_t \approx \sum_l \mathbf{W}^l \mathbf{h}_t^l$). Thus, source separation can be performed on-line and with latency corresponding to the window length plus the computation time to obtain the activations for one frame [21].

Since sources often have similar characteristics in the short-term observations (such as unvoiced phonemes and broadband noise, or voiced phonemes and music), it seems beneficial to use information from multiple time frames. In our study, this is done by stacking features from multiple contiguous time frames into ‘supervectors’: the observation \mathbf{m}_t' at time t corresponds to the observations $[\mathbf{m}_{t-T_L}; \cdots; \mathbf{m}_t; \cdots; \mathbf{m}_{t+T_R}]$ where T_L and T_R are the left and right context sizes. Missing observations at the beginning and end of the data are replaced by copying the first and last observations, i.e., $\mathbf{m}_{t:t < 0} := \mathbf{m}_1$ and $\mathbf{m}_{t:t > T} := \mathbf{m}_T$. Analogously, each basis element \mathbf{w}_k^l will model a sequence of spectra, stacked into a column vector. For readability, we omit the $'$ from the matrix names and assume that all features correspond to stacked column vectors of short-time spectra.

7 Experiments and Results

The algorithms are evaluated on the corpus of the 2nd CHiME Speech Separation and Recognition Challenge, which is publicly available¹. The task is to separate speech from noisy and reverberated mixtures. The noise was recorded in a home environment with mostly non-stationary noise sources such as children, household appliances, television, radio, etc. Training, development, and test sets of noisy mixtures along with noise-free reference signals are created from the Wall Street Journal (WSJ-0) corpus of read speech and a corpus of training noise recordings. The dry speech recordings are convolved with room impulse responses from the same environment where the noise corpus is recorded. The training set consists of 7 138 utterances at six signal-to-noise ratios (SNRs) from -6 to 9 dB, in steps of 3 dB. The development and test sets consist of 410 and 330 utterances at each of these SNRs, for a total of 2 460 / 1 980 utterances. Our evaluation measure for speech separation is source-to-distortion ratio (SDR) [23]. By construction of the WSJ-0 corpus, our evaluation is speaker-independent. Furthermore, the background noise in the development and test set is disjoint from the training noise, and a different room impulse response is used to convolve the dry utterances.

¹http://spandh.dcs.shef.ac.uk/chime_challenge/ – as of Feb. 2014

7.1 Feature extraction

Each feature vector (in the mixture, source, and reconstructed source spectrograms as well as the basis vectors) covers nine consecutive frames ($T_L = 8$, $T_R = 0$) obtained as short-time Fourier spectral magnitudes, using 25 ms window size, 10 ms window shift, and the square root of the Hann window. Since no information from the future is used ($T_R = 0$), the observation features (\mathbf{m}_t) can be extracted on-line. In analogy to the features in \mathbf{M} , each column of $\hat{\mathbf{S}}^l$ corresponds to a sliding window of consecutive reconstructed frames. Only the last frame in each sliding window is reconstructed, which leads to an on-line algorithm.

7.2 Practical implementation using exemplar-based and sparse NMF

We use the same number R of basis vectors for speech and noise ($R^1 = R^2 = R$). We run an experiment for $R = 100$ and $R = 1000$. The maximum number of iterations at test time is set to $Q = 25$ based on the trade-off between SDR and complexity – running NMF until convergence increased SDR only by about 0.1 dB SDR in preliminary experiments. At training time, we use up to $Q = 100$ iterations. As described above, we consider three different approaches to obtain NMF basis functions. In ENMF, the set of basis functions \mathbf{W} corresponds to $R^1 + R^2 = 2R$ randomly selected spectral patches of speech and noise, each spanning $T_L + 1 = 9$ frames, from the isolated CHiME speech and background noise training sets. For SNMF and NMF+S, basis training is performed separately for the speech and the noise: setting \mathbf{S}^1 to the spectrograms of the concatenated noise-free CHiME training set and \mathbf{S}^2 to those of the corresponding background noise in the multi-condition training set yields bases \mathbf{W}^l , $l = 1, 2$. Due to space complexity, we use only 10 % of the training utterances. We initialize \mathbf{W} using the exemplar bases. We found that this provided fast convergence of the objective especially for large values of μ . The update procedure for NMF+S follows Algorithm 1, while that for SNMF follows Algorithm 2. Note that, in SNMF, even though normalization is explicitly taken into account in the objective function, we still re-normalize \mathbf{W} at the end of every iteration so that the only difference with NMF+S is limited to the update equations for \mathbf{W} .

7.3 Results on the CHiME development and test set

We compare the three approaches for three different settings of the β -divergence: $\beta = 0$ (IS divergence), $\beta = 1$ (KL divergence), and $\beta = 2$ (squared Euclidean distance). For each divergence, we investigate two settings of the basis size R , $R \in \{100, 1000\}$, and multiple settings of the sparsity weight μ in (3), $\mu \in \{0, 0.1, 0.2, 0.5, 1, 2, 5, 10, 20, 50, 100\}$. At test time, all settings are explored using the bases trained using the corresponding setting for SNMF and NMF+S, and the exemplar bases for ENMF.

Figure 2 compares the average SDR obtained on the CHiME development set by using ENMF, NMF+S and SNMF, for various sparsity parameters as well as basis sizes R , for the IS distance ($\beta = 0$). Figure 3 and Figure 4 show the same comparison for the KL divergence ($\beta = 1$) and the Euclidean distance ($\beta = 2$), respectively. It can be seen that SNMF outperforms ENMF for all sparsity values for the IS distance, and for higher sparsity values ($\mu \geq 0.5$) for the KL divergence and the Euclidean distance. One can also see that SNMF consistently and significantly outperforms its ad hoc counterpart NMF+S for all divergence functions, number of basis functions, and sparsity weights, except for very high sparsity weights ($\mu \geq 50$) for the Euclidean distance; note however that the performance for such sparsity weights is much lower than the optimal performance obtained by SNMF, which is reached consistently for sparsity weights around 5 to 10 in our experiments. It is also interesting to note that the performance of NMF+S is generally either comparable to or worse than that of the exemplar-based basis selection method ENMF, which has the extra advantage that it does not require any training. This is arguably one of the reasons why sampling the basis functions instead of training them has been so popular in NMF-based works.

Overall, the best performance on the development set is obtained by SNMF, for the KL-divergence with $R = 1000$ and $\mu = 5$, with an SDR of 9.20 dB. The best performance for ENMF is obtained for the KL divergence with $R = 1000$ and $\mu = 5$, with an SDR of 7.49 dB, while that of NMF+S is obtained for the Itakura-Saito distance with $R = 1000$ and $\mu = 20$, with an SDR of 7.69 dB.

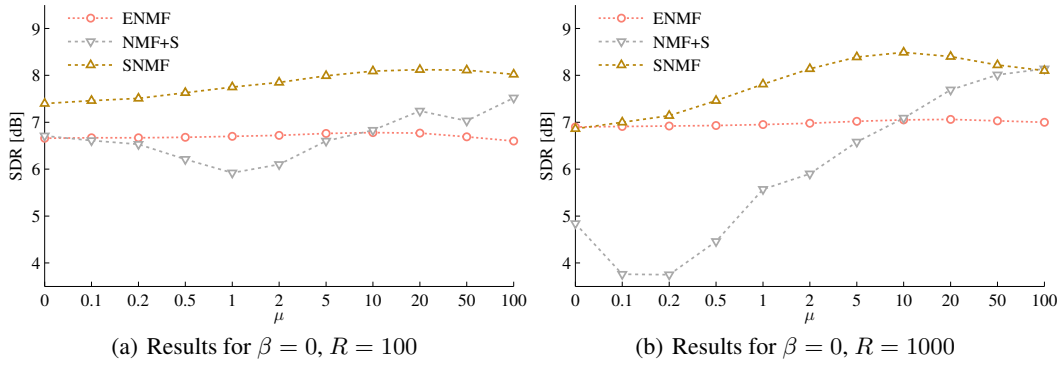


Figure 2: Average SDR obtained with the Itakura-Saito divergence ($\beta = 0$) for various sparsity weights μ on the CHiME Challenge (WSJ-0) development set.

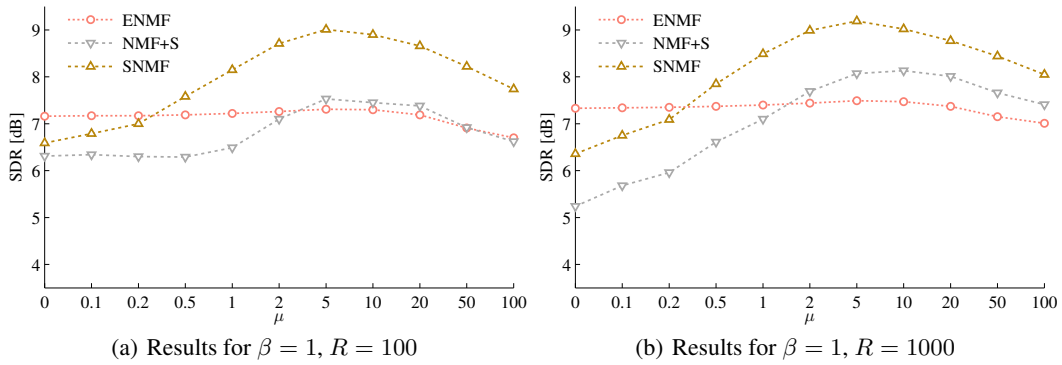


Figure 3: Average SDR obtained with the Kullback-Leibler divergence ($\beta = 1$) for various sparsity weights μ on the CHiME Challenge (WSJ-0) development set.

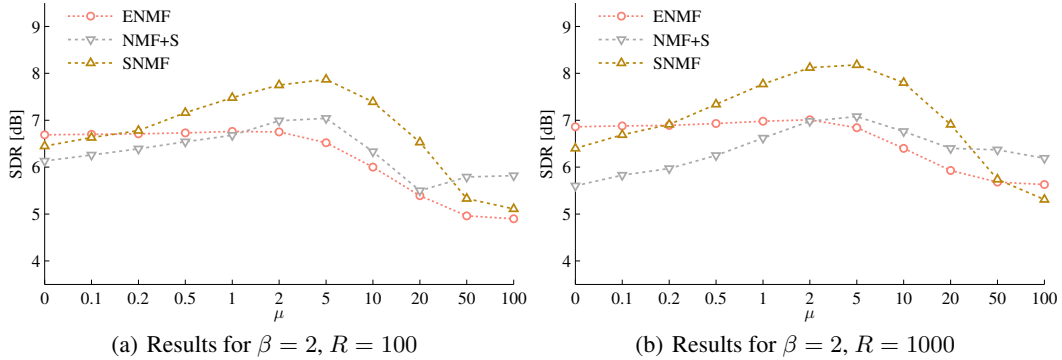


Figure 4: Average SDR obtained with the Euclidean distance ($\beta = 2$) for various sparsity weights μ on the CHiME Challenge (WSJ-0) development set.

Table 1 shows the results on the CHiME test set by ENMF, NMF+S and SNMF for the KL divergence, using $\mu = 5$ as tuned on the development set, for $R = 1000$. The results mirror those obtained on the development set.

7.4 Influence of initialization

We compare here various procedures to initialize the basis functions and activations before training. The default procedure, referred to as eW-rH, is the one used in the above experiments, where the basis functions \mathbf{W} are obtained as the exemplar bases, and \mathbf{H} is randomly initialized with uniform

SDR [dB]	Input SNR [dB]						Avg.
	-6	-3	0	3	6	9	
Noisy	-2.27	-0.58	1.66	3.40	5.20	6.60	2.34
ENMF	3.01	5.58	7.60	9.58	11.79	13.67	8.54
NMF+S	3.60	5.98	7.58	9.19	10.95	12.16	8.24
SNMF	5.48	7.53	9.19	10.88	12.89	14.61	10.10

Table 1: Source separation performance on the CHiME Challenge (WSJ-0) test set using KL-divergence, $\mu = 5$ and $R = 1000$.

SDR [dB]	Input SNR [dB]						Avg.
	-6	-3	0	3	6	9	
NMF+S rW-rH	-0.24	2.63	5.57	7.48	9.54	11.91	6.15
NMF+S eW-rH	2.91	5.40	7.72	9.12	10.80	12.47	8.07
NMF+S eW-oH	2.48	5.00	7.46	8.99	10.80	12.71	7.91
SNMF rW-rH	4.66	6.82	8.78	10.18	11.86	13.69	9.33
SNMF eW-rH	4.37	6.59	8.66	10.06	11.81	13.67	9.19
SNMF eW-oH	3.63	5.98	8.21	9.66	11.48	13.43	8.73

Table 2: Influence of initialization of \mathbf{W} and \mathbf{H} during training on source separation performance, measured on the CHiME Challenge (WSJ-0) development set using KL-divergence, $\mu = 5$ and $R = 1000$. “rW-rH”: initialize both \mathbf{W} and \mathbf{H} randomly; “eW-rH”: initialize \mathbf{W} using ENMF and \mathbf{H} randomly; “eW-oH”: initialize \mathbf{W} using ENMF, then initialize \mathbf{H} by optimizing it with fixed \mathbf{W} .

distribution on the open interval $(0, 1)$. We consider two other initialization procedures: in one, referred to as rW-rH, both \mathbf{W} and \mathbf{H} are randomly initialized; in the other, referred to as eW-oH, we start with the same initialization as eW-rH, then optimize \mathbf{H} with fixed \mathbf{W} until convergence, after which both \mathbf{W} and \mathbf{H} are optimized as usual. The intent here is to look for the best initialization strategy. We ran these experiments in a single setting, using the KL divergence, with $\mu = 5$ and $K = 1000$. Results are reported in Table 2.

Interestingly, optimizing \mathbf{H} (eW-oH) led to the worst results for SNMF, while completely random initialization (rW-rH) led to the best results, only outperforming our default setup (eW-rH) by a small margin. This hints at the tendency of the algorithm to get stuck into local minima, which is expected with multiplicative update type algorithms for NMF. On the other hand, these results show that SNMF can robustly estimate good-performing basis functions from a random initialization. Such is not the case for NMF+S, where rW-rH led to the worst results: one explanation is that the algorithm being wrong in the first place, it benefits from being initialized with the somewhat relevant basis functions obtained with ENMF.

8 Conclusion

We presented a thorough comparison of three methods for obtaining NMF basis functions from training data: ENMF, which relies on sampling from the data; NMF+S, which attempts to perform unsupervised NMF with sparsity using a simple ad-hoc implementation; and SNMF, which actually optimizes a sparse NMF objective function. We showed that, on a large and challenging speech separation task, SNMF does lead to better basis functions than ENMF, while NMF+S leads to poorer results than either. Since NMF+S is so temptingly simple to implement, it may have convinced some researchers that using samples was better than learning the bases. Exemplar-based methods may still be useful in situations where training complexity is an issue. However, when the best performance is desired, SNMF should be preferred.

Appendix A Proof of equality of the tangent gradient and the natural gradient for the L^2 -norm constraint

The natural gradient [15] is a modification of the standard gradient according to the local curvature of the parameter space. It can be obtained from the standard gradient as

$$\nabla^{(N)}\mathcal{I}(s) = \mathbf{G}^{-1}(s)\nabla\mathcal{I}(s), \forall s \in \mathcal{S}, \quad (29)$$

where \mathcal{S} is a manifold, and \mathbf{G} is the Riemannian metric tensor of the manifold at s . The metric tensor at point s defines the dot product on the tangent space at s . For \mathbb{R}^n with Cartesian coordinates, the metric tensor is simply the identity matrix, and the dot product between vectors is defined as usual. Note that the definition of the manifold and \mathbf{G} encompasses the coordinate system in which s is expressed. For example, $\mathbb{R}^2 \setminus \{(0,0)\}$ with Cartesian coordinates is a different manifold from $\mathbb{R}^2 \setminus \{(0,0)\}$ with polar coordinates. In order to obtain the natural gradient as expressed in another coordinate system, one needs to perform a change of variables. This will be our strategy here: compute the metric tensor using a spherical coordinate system, in which it is easier to define the unit-norm manifold, and use a change of variables to express the natural gradient in the Cartesian coordinate system. From here on, we assume that \mathcal{S} is the unit L^2 -norm constraint manifold in \mathbb{R}^n , or in other words the Cartesian n -sphere of radius 1.

First, we derive the metric tensor of the n -sphere. We denote by f the mapping from the spherical coordinate system in n dimensions to the Cartesian coordinates, $f : (r, \phi_1, \dots, \phi_{n-1}) \mapsto (x_1, \dots, x_n)$ such that $f_i(r, \phi_1, \dots, \phi_{n-1}) = x_i$ with

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_i \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} = \begin{bmatrix} r \cos \phi_1 \\ r \sin \phi_1 \cos \phi_2 \\ \vdots \\ r \sin \phi_1 \cdots \sin \phi_{i-1} \cos \phi_i \\ \vdots \\ r \sin \phi_1 \cdots \sin \phi_{n-2} \cos \phi_{n-1} \\ r \sin \phi_1 \cdots \sin \phi_{n-2} \sin \phi_{n-1} \end{bmatrix} \quad (30)$$

Then, the metric tensor in the spherical coordinate system $\mathbf{G}^{(r,\phi)}$ can be obtained from that in the Cartesian coordinate system $\mathbf{G}^{(\mathbf{x})}$ using a (reverse) change of coordinates:

$$\mathbf{G}^{(r,\phi)} = (D_{r,\phi}f)^\top \mathbf{G}^{(\mathbf{x})} D_{r,\phi}f, \quad (31)$$

where $D_{r,\phi}f$ is the Jacobian matrix of f . For convenience, we introduce the notation:

$$D_{r,\phi}f = [\nabla_r f \mid D_\phi f] = [\nabla_r \mathbf{x} \mid D_\phi \mathbf{x}]. \quad (32)$$

As $\mathbf{G}^{(\mathbf{x})}$ is equal to Id, we obtain:

$$\mathbf{G}^{(r,\phi)} = (D_{r,\phi}f)^\top D_{r,\phi}f. \quad (33)$$

Let us explicitly compute $D_{r,\phi}f$. First, it is easy to see that $\nabla_r f = \frac{1}{r}\mathbf{x}$, and

$$\nabla_{\phi_i} f = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ -r \sin \phi_1 \cdots \sin \phi_{i-1} \sin \phi_i \\ r \sin \phi_1 \cdots \sin \phi_{i-1} \cos \phi_i \cos \phi_{i+1} \\ \vdots \\ r \sin \phi_1 \cdots \sin \phi_{i-1} \cos \phi_i \sin \phi_{i+1} \cdots \sin \phi_{n-1} \cos \phi_{n-1} \\ r \sin \phi_1 \cdots \sin \phi_{i-1} \cos \phi_i \sin \phi_{i+1} \cdots \sin \phi_{n-1} \sin \phi_{n-1} \end{bmatrix} \leftarrow i\text{-th row},$$

where the cases $i = 1$ and $i = n - 1$ can be easily inferred from the above.

We can see that

$$G_{rr}^{(r,\phi)} = 1, \quad (34)$$

$$G_{\phi_i\phi_i}^{(r,\phi)} = r^2 \prod_{l=1}^{i-1} \sin^2 \phi_l, \quad (35)$$

$$G_{r\phi_i}^{(r,\phi)} = 0, \quad (36)$$

$$G_{\phi_i\phi_j}^{(r,\phi)} = 0, \quad i \neq j. \quad (37)$$

The computation of $G_{rr}^{(r,\phi)}$ is straightforward. For $G_{\phi_i\phi_i}^{(r,\phi)}$, the result can be obtain by noticing that

$$G_{\phi_i\phi_i}^{(r,\phi)} = r^2 \left(\prod_{l=1}^{i-1} \sin^2 \phi_l \right) (\sin^2 \phi_i + \cos^2 \phi_i \cdot \eta_{i+1}) \quad (38)$$

where η_l is defined by backward recurrence using

$$\eta_l = (\cos^2 \phi_l + \sin^2 \phi_l \cdot \eta_{l+1}), \text{ and } \eta_{n-1} = \cos^2 \phi_{n-1} + \sin^2 \phi_{n-1} = 1. \quad (39)$$

One can easily see that $\eta_l = 1, \forall l \in i+1, \dots, n-2$ as well, and thus that the factor $(\sin^2 \phi_i + \cos^2 \phi_i \cdot \eta_{i+1})$ simplifies to 1. We now compute $G_{r\phi_i}^{(r,\phi)}$:

$$G_{r\phi_i}^{(r,\phi)} = r \left(\prod_{l=1}^{i-1} \sin^2 \phi_l \right) (-\cos \phi_i \sin \phi_i + \sin \phi_i \cos \phi_i \cdot \eta_{i+1}). \quad (40)$$

As seen above, $\eta_{i+1} = 1$, and the right-hand term thus cancels out. Similarly, we compute $G_{\phi_i\phi_j}^{(r,\phi)}$ with $i < j$ (the case of $j > i$ is symmetric):

$$G_{\phi_i\phi_j}^{(r,\phi)} = r^2 \left(\prod_{\substack{l=1 \\ l \neq i}}^{j-1} \sin^2 \phi_l \right) (-\cos \phi_i \cos \phi_j)(\sin \phi_i \sin \phi_j) + (\cos \phi_i \sin \phi_j)(\sin \phi_i \cos \phi_j) \eta_{j+1}, \quad (41)$$

which again cancels out. Note that the above result is only given for the sake of completeness, and we actually do not use the value of $G_{\phi_i\phi_j}^{(r,\phi)}$ for $i \neq j$ in the remainder.

Altogether, we can write

$$\mathbf{G}^{(r,\phi)} = \left[\begin{array}{c|ccc} 1 & 0 & \cdots & 0 \\ \hline 0 & & & \\ \vdots & & \mathbf{G}^{(\phi)} & \\ 0 & & & \end{array} \right] \quad (42)$$

The quantity we are interested in is the inverse of the metric tensor of the unit-norm manifold, $(\mathbf{G}^{(\phi)})^{-1}$, expressed in Cartesian coordinates, i.e., $(D_\phi f)(\mathbf{G}^{(\phi)})^{-1}(D_\phi f)^\top$. But noticing that

$$(D_{r,\phi} f)(\mathbf{G}^{(r,\phi)})^{-1}(D_{r,\phi} f)^\top = \text{Id}, \quad (43)$$

we obtain

$$\text{Id} = \left[\begin{array}{c|c} \vdots & \\ \hline \frac{1}{r} \mathbf{x} & D_\phi f \\ \vdots & \end{array} \right] \left[\begin{array}{c|ccc} 1 & 0 & \cdots & 0 \\ \hline 0 & & & \\ \vdots & & \mathbf{G}^{(\phi)} & \\ 0 & & & \end{array} \right]^{-1} \left[\begin{array}{c|c} \vdots & \\ \hline \frac{1}{r} \mathbf{x} & D_\phi f \\ \vdots & \end{array} \right]^\top \quad (44)$$

$$= \left[\begin{array}{c|c} \vdots & \\ \hline \frac{1}{r} \mathbf{x} & (D_\phi f)(\mathbf{G}^{(\phi)})^{-1} \\ \vdots & \end{array} \right] \left[\begin{array}{c|c} \vdots & \\ \hline \frac{1}{r} \mathbf{x} & D_\phi f \\ \vdots & \end{array} \right]^\top \quad (45)$$

$$= \frac{1}{r^2} \mathbf{x} \mathbf{x}^\top + (D_\phi f)(\mathbf{G}^{(\phi)})^{-1}(D_\phi f)^\top \quad (46)$$

which leads to:

$$(D_{\phi}f)(\mathbf{G}^{(\phi)})^{-1}(D_{\phi}f)^{\top} = \text{Id} - \frac{1}{r^2}\mathbf{x}\mathbf{x}^{\top}. \quad (47)$$

For a point \mathbf{x} on the unit-norm surface, we recognize the projection operator used to compute the tangent gradient in (14), and we thus obtain:

$$\nabla^{(N)}\mathcal{I}(s) = (D_{\phi}f)(\mathbf{G}^{(\phi)})^{-1}(D_{\phi}f)^{\top}\nabla\mathcal{I}(s) = \nabla^{(T)}\mathcal{I}(s), \forall s \in \mathcal{S}. \quad (48)$$

Appendix B Source code for SNMF

B.1 example.m

```
1 % Example script for Sparse NMF with beta-divergence distortion function
2 % and L1 penalty on the activations.
3 %
4 % If you use this code, please cite:
5 % J. Le Roux, J. R. Hershey, F. Weninger,
6 % "Sparse NMF - half-baked or well done?,"
7 % MERL Technical Report, TR2015-023, March 2015
8 % @TechRep{LeRoux2015mar,
9 %   author = {{Le Roux}, J. and Hershey, J. R. and Weninger, F.},
10 %   title = {Sparse {NMF} - half-baked or well done?},
11 %   institution = {Mitsubishi Electric Research Labs (MERL)},
12 %   number = {TR2015-023},
13 %   address = {Cambridge, MA, USA},
14 %   month = mar,
15 %   year = 2015
16 % }
17 %
18 %*****
19 % Copyright (C) 2015 Mitsubishi Electric Research Labs (Jonathan Le Roux,
20 %   Felix Weninger, John R. Hershey)
21 % Apache 2.0 (http://www.apache.org/licenses/LICENSE-2.0)
22 %*****
23
24 % You need to provide a non-negative matrix v to be factorized.
25
26 params = struct;
27
28 % Objective function
29 params.cf = 'kl'; % 'is', 'kl', 'ed'; takes precedence over setting the beta value
30 % alternately define: params.beta = 1;
31 params.sparsity = 5;
32
33 % Stopping criteria
34 params.max_iter = 100;
35 params.conv_eps = 1e-3;
36 % Display evolution of objective function
37 params.diplay = 0;
38
39 % Random seed: any value over than 0 sets the seed to that value
40 params.random_seed = 1;
41
42 % Optional initial values for W
43 %params.init_w
44 % Number of components: if init_w is set and r larger than the number of
45 % basis functions in init_w, the extra columns are randomly generated
46 params.r = 500;
47 % Optional initial values for H: if not set, randomly generated
48 %params.init_h
49
50 % List of dimensions to update: if not set, update all dimensions.
51 %params.w_update_ind = true(r,1); % set to false(r,1) for supervised NMF
52 %params.h_update_ind = true(r,1);
53
54 [w, h, objective] = sparse_nmf(v, params);
```

B.2 sparse_nmf.m

```
1 function [w, h, objective] = sparse_nmf(v, params)
2
3 % SPARSE_NMF Sparse NMF with beta-divergence reconstruction error,
4 % L1 sparsity constraint, optimization in normalized basis vector space.
5 %
6 % [w, h, objective] = sparse_nmf(v, params)
7 %
8 % Inputs:
9 % v: matrix to be factorized
10 % params: optional parameters
11 %   beta:   beta-divergence parameter (default: 1, i.e., KL-divergence)
12 %   cf:     cost function type (default: 'kl'; overrides beta setting)
13 %           'is': Itakura-Saito divergence
14 %           'kl': Kullback-Leibler divergence
15 %           'kl': Euclidean distance
16 %   sparsity: weight for the L1 sparsity penalty (default: 0)
17 %   max_iter: maximum number of iterations (default: 100)
18 %   conv_eps: threshold for early stopping (default: 0,
19 %           i.e., no early stopping)
20 %   display: display evolution of objective function (default: 0)
21 %   random_seed: set the random seed to the given value
22 %               (default: 1; if equal to 0, seed is not set)
23 %   init_w:   initial setting for W (default: random;
24 %           either init_w or r have to be set)
25 %   r:       # basis functions (default: based on init_w's size;
26 %           either init_w or r have to be set)
27 %   init_h:   initial setting for H (default: random)
28 %   w_update_ind: set of dimensions to be updated (default: all)
29 %   h_update_ind: set of dimensions to be updated (default: all)
30 %
31 % Outputs:
32 % w: matrix of basis functions
33 % h: matrix of activations
34 % objective: objective function values throughout the iterations
35 %
36 %
37 %
38 % References:
39 % J. Eggert and E. Korner, "Sparse coding and NMF," 2004
40 % P. D. O'Grady and B. A. Pearlmutter, "Discovering Speech Phones
41 % Using Convolutional Non-negative Matrix Factorisation
42 % with a Sparseness Constraint," 2008
43 % J. Le Roux, J. R. Hershey, F. Weninger, "Sparse NMF - half-baked or well
44 % done?," 2015
45 %
46 % This implementation follows the derivations in:
47 % J. Le Roux, J. R. Hershey, F. Weninger,
48 % "Sparse NMF - half-baked or well done?,"
49 % MERL Technical Report, TR2015-023, March 2015
50 %
51 % If you use this code, please cite:
52 % J. Le Roux, J. R. Hershey, F. Weninger,
53 % "Sparse NMF - half-baked or well done?,"
54 % MERL Technical Report, TR2015-023, March 2015
55 % @TechRep{LeRoux2015mar,
56 %   author = {Le Roux, J. and Hershey, J. R. and Weninger, F.},
57 %   title = {Sparse {NMF} - half-baked or well done?},
58 %   institution = {Mitsubishi Electric Research Labs (MERL)},
59 %   number = {TR2015-023},
60 %   address = {Cambridge, MA, USA},
61 %   month = mar,
62 %   year = 2015
63 % }
64 %
65 %*****
66 % Copyright (C) 2015 Mitsubishi Electric Research Labs (Jonathan Le Roux,
67 % Felix Weninger, John R. Hershey)
68 % Apache 2.0 (http://www.apache.org/licenses/LICENSE-2.0)
69 %*****
70
71 m = size(v, 1);
72 n = size(v, 2);
73
74 if ~exist('params', 'var')
75     params = struct;
76 end
77
```

```

78 if ~isfield(params, 'max_iter')
79     params.max_iter = 100;
80 end
81
82 if ~isfield(params, 'random_seed')
83     params.random_seed = 1;
84 end
85
86 if ~isfield(params, 'sparsity')
87     params.sparsity = 0;
88 end
89
90 if ~isfield(params, 'conv_eps')
91     params.conv_eps = 0;
92 end
93
94 if ~isfield(params, 'cf')
95     params.cf = 'kl';
96 end
97
98 switch params.cf
99     case 'is'
100         params.beta = 0;
101     case 'kl'
102         params.beta = 1;
103     case 'ed'
104         params.beta = 2;
105     otherwise
106         if ~isfield(params, 'beta')
107             params.beta = 1;
108         end
109     end
110
111 if params.random_seed > 0
112     rand('seed', params.random_seed);
113 end
114
115 if ~isfield(params, 'init_w')
116     if ~isfield(params, 'r')
117         error('Number of components or initialization must be given')
118     end
119     r = params.r;
120     w = rand(m, r);
121 else
122     ri = size(params.init_w, 2);
123     w(:, 1:ri) = params.init_w;
124     if isfield(params, 'r') && ri < params.r
125         w(:, (ri + 1) : params.r) = rand(m, params.r - ri);
126         r = params.r;
127     else
128         r = ri;
129     end
130 end
131
132 if ~isfield(params, 'init_h')
133     h = rand(r, n);
134 elseif ischar(params.init_h) && strcmp(params.init_h, 'ones')
135     fprintf('sup_nmf: Initializing H with ones.\n');
136     h = ones(r, n);
137 else
138     h = params.init_h;
139 end
140
141 if ~isfield(params, 'w_update_ind')
142     params.w_update_ind = true(r, 1);
143 end
144
145 if ~isfield(params, 'h_update_ind')
146     params.h_update_ind = true(r, 1);
147 end
148
149 % sparsity per matrix entry
150 if length(params.sparsity) == 1
151     params.sparsity = ones(r, n) * params.sparsity;
152 elseif size(params.sparsity, 2) == 1
153     params.sparsity = repmat(params.sparsity, 1, n);
154 end
155
156 % Normalize the columns of W and rescale H accordingly

```

```

157 wn = sqrt(sum(w.^2));
158 w = bsxfun(@rdivide,w,wn);
159 h = bsxfun(@times, h,wn');
160
161 if ~isfield(params, 'display')
162     params.display = 0;
163 end
164
165 flr = 1e-9;
166 lambda = max(w * h, flr);
167 last_cost = Inf;
168
169 objective = struct;
170 objective.div = zeros(1,params.max_iter);
171 objective.cost = zeros(1,params.max_iter);
172
173 div_beta = params.beta;
174 h_ind = params.h_update_ind;
175 w_ind = params.w_update_ind;
176 update_h = sum(h_ind);
177 update_w = sum(w_ind);
178
179 fprintf(1,'Performing sparse NMF with beta-divergence, beta=%1f\n',div_beta);
180
181 tic
182 for it = 1:params.max_iter
183
184     % H updates
185     if update_h > 0
186         switch div_beta
187             case 1
188                 dph = bsxfun(@plus, sum(w(:, h_ind))', params.sparsity);
189                 dph = max(dph, flr);
190                 dmh = w(:, h_ind)' * (v ./ lambda);
191                 h(h_ind, :) = bsxfun(@rdivide, h(h_ind, :) .* dmh, dph);
192             case 2
193                 dph = w(:, h_ind)' * lambda + params.sparsity;
194                 dph = max(dph, flr);
195                 dmh = w(:, h_ind)' * v;
196                 h(h_ind, :) = h(h_ind, :) .* dmh ./ dph;
197             otherwise
198                 dph = w(:, h_ind)' * lambda.^(div_beta - 1) + params.sparsity;
199                 dph = max(dph, flr);
200                 dmh = w(:, h_ind)' * (v .* lambda.^(div_beta - 2));
201                 h(h_ind, :) = h(h_ind, :) .* dmh ./ dph;
202         end
203         lambda = max(w * h, flr);
204     end
205
206     % W updates
207     if update_w > 0
208         switch div_beta
209             case 1
210                 dpw = bsxfun(@plus, sum(h(w_ind, :), 2)', ...
211                     bsxfun(@times, ...
212                         sum((v ./ lambda) * h(w_ind, :)' .* w(:, w_ind)), w(:, w_ind)));
213                 dpw = max(dpw, flr);
214                 dmw = v ./ lambda * h(w_ind, :)' ...
215                     + bsxfun(@times, ...
216                         sum(bsxfun(@times, sum(h(w_ind, :), 2)', w(:, w_ind))), w(:, w_ind)));
217                 w(:, w_ind) = w(:,w_ind) .* dmw ./ dpw;
218             case 2
219                 dpw = lambda * h(w_ind, :)' ...
220                     + bsxfun(@times, sum(v * h(w_ind, :)' .* w(:, w_ind)), w(:, w_ind));
221                 dpw = max(dpw, flr);
222                 dmw = v * h(w_ind, :)' + ...
223                     bsxfun(@times, sum(lambda * h(w_ind, :)' .* w(:, w_ind)), w(:, w_ind));
224                 w(:, w_ind) = w(:,w_ind) .* dmw ./ dpw;
225             otherwise
226                 dpw = lambda.^(div_beta - 1) * h(w_ind, :)' ...
227                     + bsxfun(@times, ...
228                         sum((v .* lambda.^(div_beta - 2)) * h(w_ind, :)' .* w(:, w_ind)), ...
229                             w(:, w_ind));
230                 dpw = max(dpw, flr);
231                 dmw = (v .* lambda.^(div_beta - 2)) * h(w_ind, :)' ...
232                     + bsxfun(@times, ...
233                         sum(lambda.^(div_beta - 1) * h(w_ind, :)' .* w(:, w_ind)), w(:, w_ind));
234                 w(:, w_ind) = w(:,w_ind) .* dmw ./ dpw;
235

```

```

236     end
237     % Normalize the columns of W
238     w = bsxfun(@rdivide,w,sqrt(sum(w.^2)));
239     lambda = max(w * h, flr);
240 end
241
242
243 % Compute the objective function
244 switch div_beta
245     case 1
246         div = sum(sum(v .* log(v ./ lambda) - v + lambda));
247     case 2
248         div = sum(sum((v - lambda) .^ 2));
249     case 0
250         div = sum(sum(v ./ lambda - log ( v ./ lambda) - 1));
251     otherwise
252         div = sum(sum(v.^div_beta + (div_beta - 1)*lambda.^div_beta ...
253             - div_beta * v .* lambda.^(div_beta - 1))) / (div_beta * (div_beta - 1));
254 end
255 cost = div + sum(sum(params.sparsity .* h));
256
257 objective.div(it) = div;
258 objective.cost(it) = cost;
259
260 if params.display ~= 0
261     fprintf('iteration %d div = %.3e cost = %.3e\n', it, div, cost);
262 end
263
264 % Convergence check
265 if it > 1 && params.conv_eps > 0
266     e = abs(cost - last_cost) / last_cost;
267     if (e < params.conv_eps)
268         disp('Convergence reached, aborting iteration')
269         objective.div = objective.div(1:it);
270         objective.cost = objective.cost(1:it);
271         break
272     end
273 end
274 last_cost = cost;
275 end
276 toc
277
278 end

```

References

- [1] J. Le Roux, “Exploiting regularities in natural acoustical scenes for monaural audio signal estimation, decomposition, restoration and modification,” Ph.D. dissertation, The University of Tokyo & Université Paris VI – Pierre et Marie Curie, Mar. 2009.
- [2] P. D. Ogrady and B. A. Pearlmutter, “Discovering speech phones using convolutive non-negative matrix factorisation with a sparseness constraint,” *Neurocomputing*, vol. 72, no. 1, pp. 88–101, 2008.
- [3] A. Cichocki, R. Zdunek, A. H. Phan, and S.-i. Amari, *Nonnegative matrix and tensor factorizations: applications to exploratory multi-way data analysis and blind source separation*. John Wiley & Sons, 2009.
- [4] C. Févotte and J. Idier, “Algorithms for nonnegative matrix factorization with the beta-divergence,” *Neural Computation*, vol. 23, no. 9, pp. 2421–2456, 2011.
- [5] P. Smaragdis, “Convolutive speech bases and their application to supervised speech separation,” *IEEE Transactions on Audio, Speech and Language Processing*, vol. 15, no. 1, pp. 1–14, 2007.
- [6] C. Févotte, N. Bertin, and J.-L. Durrieu, “Nonnegative matrix factorization with the Itakura-Saito divergence: with application to music analysis,” *Neural Computation*, vol. 21, no. 3, pp. 793–830, March 2009.
- [7] W. Liu, N. Zheng, and X. Lu, “Non-negative matrix factorization for visual coding,” in *Proc. ICASSP*, vol. 3, 2003.
- [8] M. N. Schmidt and R. K. Olsson, “Single-channel speech separation using sparse non-negative matrix factorization,” in *Proc. Interspeech*, 2006.
- [9] J. Eggert and E. Körner, “Sparse coding and NMF,” in *Proc. of Neural Networks*, vol. 4, Dalian, China, 2004, pp. 2529–2533.
- [10] B. Raj, T. Virtanen, S. Chaudhuri, and R. Singh, “Non-negative matrix factorization based compensation of music for automatic speech recognition,” in *Proc. of Interspeech*, Makuhari, Japan, 2010, pp. 717–720.
- [11] J. T. Geiger, F. Weninger, A. Hurmalainen, J. F. Gemmeke, M. Wöllmer, B. Schuller, G. Rigoll, and T. Virtanen, “The TUM+TUT+KUL approach to the CHiME Challenge 2013: Multi-stream ASR exploiting BLSTM networks and sparse NMF,” in *Proc. of 2nd CHiME Workshop held in conjunction with ICASSP 2013*. Vancouver, Canada: IEEE, 2013, pp. 25–30.
- [12] P. Smaragdis, M. Shashanka, and B. Raj, “A sparse non-parametric approach for single channel separation of known sounds,” in *Proc. NIPS*, 2009, pp. 1705–1713.
- [13] B. King and L. Atlas, “Single-channel source separation using simplified-training complex matrix factorization,” in *Proc. ICASSP*, 2010, pp. 4206–4209.
- [14] H. Kameoka, N. Ono, K. Kashino, and S. Sagayama, “Complex NMF: A new sparse representation for acoustic signals,” in *Proc. ICASSP*, 2009, pp. 3437–3440.
- [15] S.-I. Amari, “Natural gradient works efficiently in learning,” *Neural Computation*, vol. 10, no. 2, pp. 251–276, 1998.
- [16] M. Mørup, M. N. Schmidt, and L. K. Hansen, “Shift invariant sparse coding of image and music data,” Technical University of Denmark, Tech. Rep. IMM2008-04659, 2008.
- [17] S. C. Douglas, S.-I. Amari, and S.-Y. Kung, “Gradient adaptation with unit-norm constraints,” Southern Methodist University, Tech. Rep. EE-99-003, 1999.
- [18] T. P. Krasulina, “Method of stochastic approximation in the determination of the largest eigenvalue of the mathematical expectation of random matrices,” *Automat. Remote Control*, vol. 2, pp. 215–221, 1970.
- [19] V. U. Reddy, B. Egardt, and T. Kailath, “Least squares type algorithm for adaptive implementation of Pisarenko’s harmonic retrieval method,” *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 30, no. 3, pp. 399–405, Jun. 1982.
- [20] D. D. Lee and H. S. Seung, “Algorithms for non-negative matrix factorization,” in *Proc. of NIPS*, Vancouver, Canada, 2001, pp. 556–562.

- [21] C. Joder, F. Weninger, F. Eyben, D. Virette, and B. Schuller, “Real-time speech separation by semi-supervised nonnegative matrix factorization,” in *Proc. of LVA/ICA*, ser. Lecture Notes in Computer Science, F. J. Theis, A. Cichocki, A. Yeredor, and M. Zibulevsky, Eds., vol. 7191. Tel Aviv, Israel: Springer, March 2012, pp. 322–329, special Session Real-world constraints and opportunities in audio source separation.
- [22] P. Smaragdis, B. Raj, and M. Shashanka, “Supervised and semi-supervised separation of sounds from single-channel mixtures,” in *Proc. of ICA*. Berlin, Heidelberg: Springer-Verlag, 2007, pp. 414–421.
- [23] E. Vincent, R. Gribonval, and C. Févotte, “Performance measurement in blind audio source separation,” *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 14, no. 4, pp. 1462–1469, Jul. 2006.