

Nonnegative matrix factorization with the Itakura-Saito divergence. With application to music analysis.

Cédric Févotte, Nancy Bertin and Jean-Louis Durrieu*

Abstract

This article presents theoretical, algorithmic and experimental results about nonnegative matrix factorization (NMF) with the Itakura-Saito (IS) divergence. We describe how IS-NMF is underlain by a well-defined statistical model of superimposed Gaussian components and is equivalent to maximum likelihood estimation of variance parameters. This setting can accommodate regularization constraints on the factors through Bayesian priors. In particular, inverse-Gamma and Gamma Markov chain priors are considered in this work. Estimation can be carried out using a space-alternating generalized expectation-maximization (SAGE) algorithm; this leads to a novel type of NMF algorithm, whose convergence to a stationary point of the IS cost function is guaranteed.

We also discuss the links between the IS divergence and other cost functions used in NMF, in particular the Euclidean distance and the generalized Kullback-Leibler (KL) divergence. As such, we describe how IS-NMF can also be performed using a gradient multiplicative algorithm (a standard algorithm structure in NMF) whose convergence is observed in practice, though not proven.

Finally, we report a furnished experimental comparative study of Euclidean-NMF, KL-NMF and IS-NMF algorithms applied to the power spectrogram of a short piano sequence recorded in real conditions, with various initializations and model orders. Then we show how IS-NMF can successfully be employed for denoising and *upmix* (mono to stereo conversion) of an original piece of early jazz music. These experiments indicate that IS-NMF correctly captures the semantics of audio and is better suited to the representation of music signals than NMF with the usual Euclidean and KL costs.

Keywords: Nonnegative matrix factorization (NMF), unsupervised machine learning, Bayesian linear regression, space-alternating generalized expectation-maximization (SAGE), music transcription, single-channel source separation, audio restoration, computational auditory scene analysis (CASA).

1 Introduction

Nonnegative matrix factorization (NMF) is a now popular dimension reduction technique, employed for non-subtractive, part-based representation of nonnegative data. Given a data matrix \mathbf{V} of dimensions $F \times N$ with nonnegative entries, NMF is the problem of finding a factorization

$$\mathbf{V} \approx \mathbf{W}\mathbf{H} \tag{1}$$

where \mathbf{W} and \mathbf{H} are nonnegative matrices of dimensions $F \times K$ and $K \times N$, respectively. K is usually chosen such that $F K + K N \ll F N$, hence reducing the data dimension. Note that the factorization is in general only approximate, so that the terms “approximate nonnegative matrix factorization” or “nonnegative matrix approximation” also appear in the literature. NMF has been used for various problems in diverse fields. To cite a few, let us mention the problems of learning parts of faces and semantic features of text ([Lee and Seung, 1999](#)), polyphonic music transcription ([Smaragdakis and Brown, 2003](#)), object characterization by reflectance spectra analysis ([Berry et al., 2007](#)), portfolio diversification ([Drakakis et al., 2007](#)), as well as Scotch whiskies clustering ([Young et al., 2006](#)).

*The authors are with LTCI (CNRS - TELECOM ParisTech), 37-39, rue Dareau, 75014 Paris, France.

In the literature, the factorization (1) is usually sought after through the minimization problem

$$\min_{\mathbf{W}, \mathbf{H} \geq \mathbf{0}} D(\mathbf{V}|\mathbf{WH}) \quad (2)$$

where $D(\mathbf{V}|\mathbf{WH})$ is a cost function defined by

$$D(\mathbf{V}|\mathbf{WH}) = \sum_{f=1}^F \sum_{n=1}^N d([\mathbf{V}]_{fn} | [\mathbf{WH}]_{fn}) \quad (3)$$

and where $d(x|y)$ is a scalar cost function. Popular choices are the Euclidean distance, that we here define as

$$d_{EUC}(x|y) = \frac{1}{2}(x - y)^2 \quad (4)$$

and the (generalized) Kullback-Leibler (KL) divergence, also referred to as I-divergence, defined by

$$d_{KL}(x|y) = x \log \frac{x}{y} - x + y. \quad (5)$$

Both cost functions are positive, and take value zero if and only if $x = y$.

Lee and Seung (2001) proposed gradient descent algorithms to solve the minimization problem (2) under the latter two cost functions. Using a suitable step size, the gradient descent update rules are turned into multiplicative rules, under which the cost function is shown to be non-increasing. The simplicity of the update rules has undoubtedly contributed to the popularity of NMF, and most of the above-mentioned applications are based on Lee and Seung's algorithm for minimization of either the Euclidean distance or the KL divergence.

Nevertheless, some papers have considered NMF under other cost functions, and other algorithmic structures. In particular Cichocki and co-authors have devised several types of NMF algorithms for cost functions such as Csiszár divergences (including Amari's α -divergence) and the β -divergence in (Cichocki et al., 2006b), with several other cost functions considered in (Cichocki et al., 2006a). Also, Dhillon and Sra (2005) have described multiplicative algorithms for the wide family of Bregman divergences. The choice of the NMF cost function should be driven by the type of data to analyze, and if a good deal of literature is devoted to improving performance of algorithms given a cost function, little literature has been devoted to how to choose a cost function with respect to (wrt) a particular type of data and application.

In this paper we are specifically interested in NMF with the Itakura-Saito (IS) divergence and we are going to demonstrate its relevance to the decomposition of audio spectra. The expression of the IS divergence is given by

$$d_{IS}(x|y) = \frac{x}{y} - \log \frac{x}{y} - 1. \quad (6)$$

This divergence was obtained by Itakura and Saito (1968) from the maximum likelihood (ML) estimation of short-time speech spectra under autoregressive modeling. It was presented as "a measure of the goodness of fit between two spectra" and became popular in the speech community during the seventies. It was in particular praised for the good perceptual properties of the reconstructed signals it led to (Gray et al., 1980).

As we shall see, this divergence has other interesting properties. It is in particular scale-invariant, meaning that low energy components of \mathbf{V} bear the same relative importance as high energy ones. This is relevant to situations in which the coefficients of \mathbf{V} have a large dynamic range, such as in audio short-term spectra. The IS divergence also leads to desirable statistical interpretations of the NMF problem. Indeed, we describe how NMF can in this case be recast as maximum likelihood (ML) estimation of \mathbf{W} and \mathbf{H} in superimposed signals under simple Gaussian assumptions. Equivalently, we describe how IS-NMF can be interpreted as ML of \mathbf{W} and \mathbf{H} in multiplicative Gamma noise.

The IS divergence belongs to the class of Bregman divergences and is a limit case of the β -divergence. Thus, the gradient descent multiplicative rules given in (Dhillon and Sra, 2005) and (Cichocki et al., 2006b) – which coincide in the IS case – can be applied. If convergence of this algorithm is observed in practice, its proof is still an open problem. The statistical framework

going along with IS-NMF allows to derive a new type of minimization method, derived from space-alternating expectation-maximization (SAGE), a variant of the standard expectation-maximization (EM) algorithm. This method leads to new update rules, which do not possess a multiplicative structure. The EM setting guarantees convergence of this algorithm to a stationary point of the cost function. Moreover, the statistical framework opens doors to Bayesian approaches for NMF, allowing elaborate priors on \mathbf{W} and \mathbf{H} , for which maximum a posteriori (MAP) estimation can again be performed using SAGE. Examples of such priors, yielding regularized estimates of \mathbf{W} and \mathbf{H} , are presented in this work.

IS-NMF underlies previous works in the area of automatic music transcription and single-channel audio source separation, but never explicitly so. In particular our work builds on (Benaroya et al., 2003, 2006; Abdallah and Plumbley, 2004; Plumbley et al., 2006) and the connections between IS-NMF and these papers will be discussed.

This article is organized as follows. Section 2 addresses general properties of IS-NMF. The relation between the IS divergence and other cost functions used in NMF is discussed in Section 2.1, Section 2.2 addresses scale invariance and Section 2.3 describes the statistical interpretations of IS-NMF. Section 3 presents two IS-NMF algorithms; an existing multiplicative algorithm is described in Section 3.1 while Section 3.2 introduces a new algorithm derived from SAGE. Section 4 reports an experimental comparative study of Euclidean-NMF, KL-NMF or IS-NMF algorithms applied to the power spectrogram of a short piano sequence recorded in real conditions, with various initializations and model orders. These experiments show that IS-NMF correctly captures the semantics of the signal and is better suited to the representation of audio than NMF with the usual Euclidean and KL costs. Section 5 presents how IS-NMF can accommodate regularization constraints on \mathbf{W} and \mathbf{H} within a Bayesian framework and how SAGE can be adapted to MAP estimation. In particular, we give update rules for IS-NMF with Gamma and inverse-Gamma Markov chain priors on the rows of \mathbf{H} . In Section 6, we present audio restoration results of an original early recording of jazz music; we show how the proposed regularized IS-NMF algorithms can successfully be employed for denoising and *upmix* (mono to stereo conversion) of the original data. Finally, conclusions and perspectives of this work are given in Section 7.

2 Properties of NMF with the Itakura-Saito divergence

In this section we address the links between the IS divergence and other cost functions used for NMF, then we discuss its scale invariance property and finally describe the statistical interpretations of IS-NMF.

2.1 Relation to other divergences used in NMF

β -divergence As observed by Cichocki et al. (2006b,a), the IS divergence is a limit case of the β -divergence introduced by Eguchi and Kano (2001), that we here define as

$$d_{\beta}(x|y) \stackrel{\text{def}}{=} \begin{cases} \frac{1}{\beta(\beta-1)} (x^{\beta} + (\beta-1)y^{\beta} - \beta x y^{\beta-1}) & \beta \in \mathbb{R} \setminus \{0, 1\} \\ x(\log x - \log y) + (y - x) & \beta = 1 \\ \frac{x}{y} - \log \frac{x}{y} - 1 & \beta = 0 \end{cases} \quad (7)$$

Eguchi and Kano (2001) assume $\beta > 1$, but the definition domain can very well be extended to $\beta \in \mathbb{R}$. The β -divergence is shown to be continuous in β by using the identity $\lim_{\beta \rightarrow 0} (x^{\beta} - y^{\beta})/\beta = \log(x/y)$. It was considered in NMF by Cichocki et al. (2006b) and also coincides up to a factor $1/\beta$ with the generalized divergence of Kompass (2007) which, in the context of NMF as well, was separately constructed so as to interpolate between the KL divergence ($\beta = 1$) and the Euclidean distance ($\beta = 2$). Note that the derivative of $d_{\beta}(x|y)$ wrt y is also continuous in β , and simply writes

$$\nabla_y d_{\beta}(x|y) = y^{\beta-2} (y - x). \quad (8)$$

The derivative shows that $d_{\beta}(x|y)$, as a function of y , has a single minimum in $y = x$ and that it increases with $|y - x|$, justifying its relevance as a measure of fit. Figure 1 represents the Euclidean, KL and IS divergences for $x = 1$.

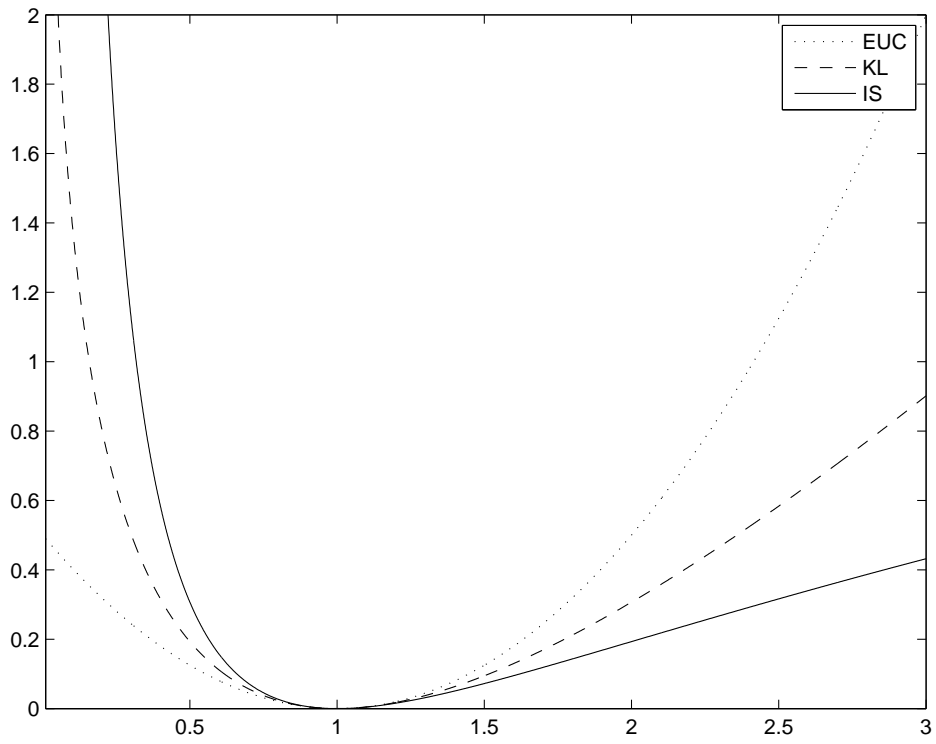


Figure 1: Euclidean, KL and IS costs $d(x|y)$ as a function of y and for $x = 1$. The Euclidean and KL divergences are convex on $(0, \infty)$. The IS divergence is convex on $(0, 2x]$ and concave on $[2x, \infty)$.

Using equation (8), the gradients of criterion $D_\beta(\mathbf{V}|\mathbf{WH})$ wrt \mathbf{W} and \mathbf{H} simply writes

$$\nabla_{\mathbf{H}} D_\beta(\mathbf{V}|\mathbf{WH}) = \mathbf{W}^T \left((\mathbf{WH})^{[\beta-2]} \cdot (\mathbf{WH} - \mathbf{V}) \right) \quad (9)$$

$$\nabla_{\mathbf{W}} D_\beta(\mathbf{V}|\mathbf{WH}) = \left((\mathbf{WH})^{[\beta-2]} \cdot (\mathbf{WH} - \mathbf{V}) \right) \mathbf{H}^T \quad (10)$$

where \cdot denotes Hadamard entrywise product and $\mathbf{A}^{[n]}$ denotes the matrix with entries $[\mathbf{A}]_{ij}^n$. The multiplicative gradient descent approach taken in (Lee and Seung, 2001; Cichocki et al., 2006b) is equivalent to updating each parameter by multiplying its value at previous iteration by the ratio of the negative and positive parts of the derivative of the criterion wrt this parameter, namely $\theta \leftarrow \theta \cdot [\nabla f(\theta)]_- / [\nabla f(\theta)]_+$, where $\nabla f(\theta) = [\nabla f(\theta)]_+ - [\nabla f(\theta)]_-$ and the summands are both nonnegative. This ensures nonnegativity of the parameter updates, provided initialization with a nonnegative value. A fixed point θ^* of the algorithm implies either $\nabla f(\theta^*) = 0$ or $\theta^* = 0$. This leads to the following updates,

$$\mathbf{H} \leftarrow \mathbf{H} \cdot \frac{\mathbf{W}^T ((\mathbf{WH})^{[\beta-2]} \cdot \mathbf{V})}{\mathbf{W}^T (\mathbf{WH})^{[\beta-1]}} \quad (11)$$

$$\mathbf{W} \leftarrow \mathbf{W} \cdot \frac{((\mathbf{WH})^{[\beta-2]} \cdot \mathbf{V}) \mathbf{H}^T}{(\mathbf{WH})^{[\beta-1]} \mathbf{H}^T} \quad (12)$$

where $\frac{\mathbf{A}}{\mathbf{B}}$ denotes the matrix $\mathbf{A} \cdot \mathbf{B}^{[-1]}$. Lee and Seung (1999) showed that $D_\beta(\mathbf{V}|\mathbf{WH})$ is non-increasing under the latter updates for $\beta = 2$ (Euclidean distance) and $\beta = 1$ (KL divergence). Kompass (2007) generalizes the proof to the case $1 \leq \beta \leq 2$. In practice, we observe that the criterion is still nonincreasing under updates (11) and (12) for $\beta < 1$ and $\beta > 2$ (and in particular for $\beta = 0$, corresponding to the IS divergence), but no proof is available. Indeed, the proof given by Kompass makes use of the convexity of $d_\beta(x|y)$ as a function of y , which is only true for $1 \leq \beta \leq 2$. In the rest of the paper, the term ‘‘EUC-NMF’’ will be used as a shorthand for ‘‘Euclidean-NMF’’.

Bregman divergences The IS divergence belongs to the class of Bregman divergences, defined as $d_\phi(x|y) = \phi(x) - \phi(y) - \nabla \phi(y)(x - y)$, where ϕ is a strictly convex function of \mathbb{R} that has a continuous derivative $\nabla \phi$. The IS divergence is obtained with $\phi(y) = -\log(y)$. Using the same approach as in previous paragraph, Dhillon and Sra (2005) derive the following update rules for minimization of $D_\phi(\mathbf{V}|\mathbf{WH})$,

$$\mathbf{H} \leftarrow \mathbf{H} \cdot \frac{\mathbf{W}^T (\nabla^2 \phi(\mathbf{WH}) \cdot \mathbf{V})}{\mathbf{W}^T (\nabla^2 \phi(\mathbf{WH}) \cdot \mathbf{WH})} \quad (13)$$

$$\mathbf{W} \leftarrow \mathbf{W} \cdot \frac{(\nabla^2 \phi(\mathbf{WH}) \cdot \mathbf{V}) \mathbf{H}^T}{(\nabla^2 \phi(\mathbf{WH}) \cdot \mathbf{WH}) \mathbf{H}^T} \quad (14)$$

Again, the authors observed in practice continual descent of $D_\phi(\mathbf{V}|\mathbf{WH})$ under these rules, but a proof of convergence is yet to be found. Note that equations (11) and (12) coincide with equations (13) and (14) for the IS divergence.

2.2 Scale invariance

The following property holds for any value of β ,

$$d_\beta(\gamma x | \gamma y) = \gamma^\beta d_\beta(x|y). \quad (15)$$

It implies that the IS divergence is scale-invariant (i.e, $d_{IS}(\gamma x | \gamma y) = d_{IS}(x|y)$), and is the only one of the β -divergence family to possess this property. The scale invariance means that same relative weight is given to small and large coefficients of \mathbf{V} in cost function (3), in the sense that a bad fit of the factorization for a low-power coefficient $[\mathbf{V}]_{fn}$ will cost as much as a bad fit for higher power coefficient $[\mathbf{V}]_{f'n'}$. On the opposite, factorizations obtained with $\beta > 0$ (such as with the Euclidean distance or the KL divergence) will rely more heavily on the largest coefficients and less precision is to be expected in the estimation of the low-power components.

The scale invariance of the IS divergence is relevant to decomposition of audio spectra, which typically exhibit exponential power decrease along frequency f and also usually comprise low-power transient components such as note attacks together with higher power components such as tonal parts of sustained notes. The results of the decomposition of a piano spectrogram presented in Section 4 confirm these expectations by showing that IS-NMF extracts components corresponding to very low residual noise and hammer hits on the strings with great accuracy. These components are either ignored or severely degraded when using Euclidean or KL divergences.

2.3 Statistical interpretations

We now turn to statistical interpretations of IS-NMF which will eventually lead to a new EM-based algorithm, described in Section 3.

2.3.1 Notations

The entries of matrices \mathbf{V} , \mathbf{W} and \mathbf{H} are denoted v_{fn} , w_{fk} and h_{kn} respectively. Lower case bold letters will in general denote columns, such that $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_K]$, while lower case plain letters with a single index denote rows, such that $\mathbf{H} = [h_1^T, \dots, h_K^T]^T$. We also define the matrix $\hat{\mathbf{V}} = \mathbf{WH}$, whose entries are denoted \hat{v}_{fn} . Where these conventions clash, the intended meaning should be clear from the context.

2.3.2 Sum of Gaussian components

Theorem 1 (IS-NMF as ML estimation in sum of Gaussian components). Consider the generative model defined by, $\forall n = 1, \dots, N$

$$\mathbf{x}_n = \sum_{k=1}^K \mathbf{c}_{k,n} \quad (16)$$

where \mathbf{x}_n and $\mathbf{c}_{k,n}$ belong to $\mathbb{C}^{F \times 1}$ and

$$\mathbf{c}_{k,n} \sim \mathcal{N}_c(0, h_{kn} \text{diag}(\mathbf{w}_k)), \quad (17)$$

where $\mathcal{N}_c(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the proper multivariate complex Gaussian distribution and where the components $\mathbf{c}_{1,n}, \dots, \mathbf{c}_{K,n}$ are mutually independent and individually independently distributed. Define \mathbf{V} as the matrix with entries $v_{fn} = |x_{fn}|^2$. Then, maximum likelihood estimation of \mathbf{W} and \mathbf{H} from $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ is equivalent to NMF of \mathbf{V} into $\mathbf{V} \approx \mathbf{WH}$, where the Itakura-Saito divergence is used.

Proof. Under assumptions of Theorem 1 and using the expression of $\mathcal{N}_c(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ given in Appendix A, the minus log likelihood function $C_{ML,1}(\mathbf{W}, \mathbf{H}) \stackrel{\text{def}}{=} -\log p(\mathbf{X}|\mathbf{W}, \mathbf{H})$ simply factorizes as

$$C_{ML,1}(\mathbf{W}, \mathbf{H}) = -\sum_{n=1}^N \sum_{f=1}^F \log \mathcal{N}_c \left(x_{fn} | 0, \sum_k w_{fk} h_{kn} \right) \quad (18)$$

$$= NF \log \pi + \sum_{n=1}^N \sum_{f=1}^F \log \left(\sum_k w_{fk} h_{kn} \right) + \frac{|x_{fn}|^2}{(\sum_k w_{fk} h_{kn})} \quad (19)$$

$$\stackrel{c}{=} \sum_{n=1}^N \sum_{f=1}^F d_{IS} \left(|x_{fn}|^2 \mid \sum_k w_{fk} h_{kn} \right) \quad (20)$$

where $\stackrel{c}{=}$ denotes equality up to constant terms. The minimization of $C_{ML,1}(\mathbf{W}, \mathbf{H})$ wrt \mathbf{W} and \mathbf{H} thus amounts to the NMF $\mathbf{V} \approx \mathbf{WH}$ with the IS divergence. Note that Theorem 1 holds also for real-valued Gaussian components. In that case $C_{ML,1}(\mathbf{W}, \mathbf{H})$ equals $D_{IS}(\mathbf{V}|\mathbf{WH})$ up to a constant and a factor 1/2. \square

The generative model (16) was introduced by Benaroya et al. (2003, 2006) for single-channel audio source separation. In that context, $\mathbf{x}_n = [x_{1n}, \dots, x_{fn}, \dots, x_{Fn}]^T$ is the Short-Time Fourier Transform (STFT) of an audio signal x , where $n = 1, \dots, N$ is a frame index and $f = 1, \dots, F$ is a frequency index. The signal x is assumed to be the sum of two sources $x = s_1 + s_2$ and the STFTs of the sources are modeled as $\mathbf{s}_{1,n} = \sum_{k=1}^{K_1} \mathbf{c}_{k,n}$ and $\mathbf{s}_{2,n} = \sum_{k=K_1+1}^{K_1+K_2} \mathbf{c}_{k,n}$, with $K_1 + K_2 = K$. This means that each source STFT is modeled as a sum of *elementary components* each characterized by a Power Spectral Density (PSD) \mathbf{w}_k modulated in time by frame-dependent activation coefficients h_{kn} . The PSDs characterizing each source are learnt on training data, before the mixture spectrogram $|\mathbf{X}|^{[2]}$ is decomposed onto the known dictionary $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_{K_1}, \mathbf{w}_{K_1+1}, \dots, \mathbf{w}_{K_1+K_2}]$. However, in these papers, the PSDs and the activation coefficients are estimated separately using somewhat ad-hoc strategies (the PSDs are learnt with vector quantization) and the equivalence between ML estimation and IS-NMF is not fully exploited.

Complex Gaussian modelling of STFT frames of audio signals has been widely used in signal processing and has proven to be a satisfying model for many applications, in particular for audio denoising, see, e.g., (Cohen and Gannot, 2007) for a review. But while denoising settings typically assume one observation frame \mathbf{x}_n to be the sum of a source frame and a noise frame, IS-NMF extends in essence this modelling by assuming that one observation frame is the sum of several Gaussian frames with different covariances.

The generative model (16) may also be viewed as a generalization of well-known models of composite signals. For example, inference in superimposed components with Gaussian structure can be tracked back to (Feder and Weinstein, 1988). In the latter paper however, the components are assumed stationary and solely modeled by their PSD \mathbf{w}_k , which is in turn parametrized by a set of parameters of interest $\boldsymbol{\theta}_k$, to be estimated. One extension brought in equation (16) is the addition of the amplitude parameters \mathbf{H} . This however has the inconvenience of making the total number of parameters $FK + KN$ dependent of N , with the consequence of losing asymptotical optimality properties of ML estimation. But note that it is precisely the addition of the amplitude parameters in the model that allows \mathbf{W} to be treated as a set of possibly identifiable parameters. Indeed, if h_{kn} is set to 1 for all k and n the variance of \mathbf{x}_n becomes $\sum_k \mathbf{w}_k$ for all n , i.e, is equal to the sum of the parameters. This would obviously make each PSD \mathbf{w}_k not uniquely identifiable.

Very interestingly, the equivalence between IS-NMF and ML inference in sum of Gaussian components provides means of reconstructing the components $\mathbf{c}_{k,n}$ with a sense of statistical optimality, which contrasts with NMF using other divergences where methods of reconstructing components from the factorization \mathbf{WH} are somewhat ad-hoc (see below). Indeed, given \mathbf{W} and \mathbf{H} , minimum mean square error (MMSE) estimates can be obtained through Wiener filtering, such that

$$\hat{c}_{k,fn} = \frac{w_{fk} h_{kn}}{\sum_{l=1}^K w_{fl} h_{ln}} x_{fn}. \quad (21)$$

The Wiener gains summing up to 1 for a fixed entry (f, n) the decomposition is conservative, i.e,

$$\mathbf{x}_n = \sum_{k=1}^K \hat{\mathbf{c}}_{k,n}. \quad (22)$$

Note that a consequence of Wiener reconstruction is that the phase of all components $\hat{c}_{k,fn}$ is equal to the phase of x_{fn} .

Most works in audio have considered the NMF of magnitude spectra $|\mathbf{X}|$ instead of power spectra $|\mathbf{X}|^{[2]}$, see, e.g., (Smaragdis and Brown, 2003; Smaragdis, 2007; Virtanen, 2007; Bertin et al., 2007). In that case, it can be noted (see, e.g., Virtanen et al. (2008)), that KL-NMF is related to the ML problem of estimating \mathbf{W} and \mathbf{H} in the model structure

$$|\mathbf{x}_n| = \sum_{k=1}^K |\mathbf{c}_{k,n}| \quad (23)$$

under Poissonian assumptions, i.e., $|c_{k,fn}| \sim \mathcal{P}(w_{fk}h_{kn})$, where $\mathcal{P}(\lambda)$ is the Poisson distribution, defined in Appendix A. Indeed, the sum of Poisson random variables being Poissonian itself (with the shape parameters summing up as well), one obtains $|x_{fn}| \sim \mathcal{P}(\sum_{k=1}^K w_{fk}h_{kn})$. Then, it can easily be seen that the likelihood $-\log p(\mathbf{X}|\mathbf{W}, \mathbf{H})$ is equal up to a constant to $D_{KL}(|\mathbf{X}| \parallel \mathbf{W}\mathbf{H})$. Here, \mathbf{W} is homogeneous to a magnitude spectrum and not to a power spectrum. After factorization, component estimates are typically formed using the phase of the observations (Virtanen, 2007), such that

$$\hat{c}_{k,fn} = w_{fk} h_{kn} \arg(x_{fn}), \quad (24)$$

where $\arg(x)$ denotes the phase of complex scalar x . This approach is worth a few comments. First, the Poisson distribution is formerly only defined for integers, which impairs statistical interpretation of KL-NMF on non-countable data such as audio spectra (but one could assume an appropriate data scaling and a very fine quantization to work around this).¹ Second, this approach enforces nonnegativity in a somehow arbitrary way by taking the absolute value of data \mathbf{X} . In contrast, with the Gaussian modelling, nonnegativity arises naturally through the variance fitting problem equivalence. Similarly, the reconstruction method enforces the components to have same phase as observation coefficients, while this is only a consequence of Wiener filtering in the Gaussian modelling framework. Last, the component reconstruction method is not statistically-grounded and is not conservative, i.e. $\mathbf{x}_n \approx \sum_{k=1}^K \hat{\mathbf{c}}_{k,n}$. Note that Wiener reconstruction is used with KL-NMF of the magnitude spectrum $|\mathbf{X}|$ by Smaragdis (2007), where it is presented as spectral filtering and its conservativity is pointed out.

2.3.3 Multiplicative noise

Theorem 2 (IS-NMF as ML estimation in Gamma multiplicative noise). Consider the generative model

$$\mathbf{V} = (\mathbf{W}\mathbf{H}) \cdot \mathbf{E} \quad (25)$$

where \mathbf{E} is multiplicative independent and identically-distributed (i.i.d.) Gamma noise with mean 1. Then, maximum likelihood estimation of \mathbf{W} and \mathbf{H} is equivalent to NMF of \mathbf{V} into $\mathbf{V} \approx \mathbf{W}\mathbf{H}$, where the Itakura-Saito divergence is used.

Proof. Let us note $\{e_{fn}\}$ the entries of \mathbf{E} . We have $v_{fn} = \hat{v}_{fn} e_{fn}$, with $p(e_{fn}) = \mathcal{G}(e_{fn}|\alpha, \beta)$, and where $\mathcal{G}(x|\alpha, \beta)$ is the Gamma probability density function (pdf) defined in Appendix A. Under the i.i.d. noise assumption, the minus log likelihood $C_{ML,2}(\mathbf{W}, \mathbf{H}) \stackrel{\text{def}}{=} -\log p(\mathbf{V}|\mathbf{W}, \mathbf{H})$ writes

$$C_{ML,2}(\mathbf{W}, \mathbf{H}) = -\sum_{f,n} \log p(v_{fn}|\hat{v}_{fn}) \quad (26)$$

$$= -\sum_{f,n} \log \mathcal{G}(v_{fn}/\hat{v}_{fn}|\alpha, \beta) / \hat{v}_{fn} \quad (27)$$

$$\stackrel{c}{=} \beta \sum_{f,n} \frac{v_{fn}}{\hat{v}_{fn}} - \frac{\alpha}{\beta} \log \frac{v_{fn}}{\hat{v}_{fn}} - 1 \quad (28)$$

The ratio α/β is simply the mean of the Gamma distribution. When it is equal to 1, we obtain that $C_{ML,2}(\boldsymbol{\theta})$ is equal to $D_{IS}(\mathbf{V}|\hat{\mathbf{V}}) = D_{IS}(\mathbf{V}|\mathbf{W}\mathbf{H})$ up to a positive factor and a constant. \square

The multiplicative noise equivalence explains the scale invariance of the IS divergence, because the noise acts as a scale factor on \hat{v}_{fn} . On the opposite, EUC-NMF is equivalent to ML likelihood estimation of \mathbf{W} and \mathbf{H} in additive i.i.d. Gaussian noise. The influence of additive noise is greater on coefficients of $\hat{\mathbf{V}}$ with small amplitude (i.e, low SNR) than on the largest ones. As to KL-NMF, it neither corresponds to multiplicative nor additive noise but actually corresponds to ML

¹Actually, KL-NMF has interesting parallels with inference in probabilistic latent variable models of histogram data, see (Shashanka et al., 2008a).

Algorithm 1 IS-NMF/MU

Input : nonnegative matrix \mathbf{V}

Output : nonnegative matrices \mathbf{W} and \mathbf{H} such that $\mathbf{V} \approx \mathbf{WH}$

Initialize \mathbf{W} and \mathbf{H} with nonnegative values

for $i = 1 : n_{iter}$ **do**

$$\mathbf{H} \leftarrow \mathbf{H} \cdot \frac{\mathbf{W}^T ((\mathbf{WH})^{[-2]}) \mathbf{V}}{\mathbf{W}^T (\mathbf{WH})^{[-1]}}$$

$$\mathbf{W} \leftarrow \mathbf{W} \cdot \frac{((\mathbf{WH})^{[-2]}) \mathbf{V} \mathbf{H}^T}{(\mathbf{WH})^{[-1]} \mathbf{H}^T}$$

Normalize \mathbf{W} and \mathbf{H}

end for

estimation in Poisson noise.² To summarize, we have

$$\text{EUC-NMF: } p(v_{fn}|\hat{v}_{fn}) = \mathcal{N}(v_{fn}|\hat{v}_{fn}, \sigma^2), \quad (29)$$

$$\text{KL-NMF: } p(v_{fn}|\hat{v}_{fn}) = \mathcal{P}(v_{fn}|\hat{v}_{fn}), \quad (30)$$

$$\text{IS-NMF: } p(v_{fn}|\hat{v}_{fn}) = \frac{1}{\hat{v}_{fn}} \mathcal{G}\left(\frac{v_{fn}}{\hat{v}_{fn}}|\alpha, \alpha\right), \quad (31)$$

and in all cases, $E\{v_{fn}|\hat{v}_{fn}\} = \hat{v}_{fn}$.

Theorem 2 reports in essence how Abdallah and Plumbley (2004) derive a “statistically motivated error measure”, which happens to be the IS divergence, in the very similar context of nonnegative sparse coding (see also developments in Plumbley et al. (2006)). Pointing out the scale invariance of this measure, this work leads Virtanen (2007) to consider the IS divergence (but again without referring to it as such) for NMF in the context of single-channel source separation, but the algorithm is applied to the magnitude spectra instead of the power spectra, losing statistical coherence, and the sources are reconstructed through equation (24) instead of Wiener filtering.

3 Algorithms for NMF with the Itakura-Saito divergence

In this section we describe two algorithms for IS-NMF. The first one has a multiplicative structure and is only a special case of the derivations of Section 2.1. The second one is of a novel type, EM-based, and is derived from the statistical presentation of IS-NMF as given in Theorem 1.

3.1 Multiplicative gradient descent algorithm

A multiplicative gradient descent IS-NMF algorithm is obtained by either setting $\beta = 0$ in (11) and (12) or setting $\phi(y) = -\log(y)$ in (13) and (14). The resulting update rules coincide and lead to Algorithm 1. These update rules were also obtained by Abdallah and Plumbley (2004), prior to (Dhillon and Sra, 2005; Cichocki et al., 2006b). In the following, we refer to this algorithm as “IS-NMF/MU”. This algorithm includes a normalization step at every iteration, which eliminates trivial scale indeterminacies leaving the cost function unchanged. We impose $\|\mathbf{w}_k\|_2 = 1$ and scale h_k accordingly. Again, we emphasize that continual descent of the cost function is observed in practice with this algorithm, but that a proof of convergence is yet to be found.

3.2 SAGE algorithm

We now describe an EM-based algorithm for estimation of the parameters $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{H}\}$, derived from the statistical formalism introduced in Theorem 1. The additive structure of the generative model (16) allows to update the parameters describing each component $\mathbf{C}_k \stackrel{\text{def}}{=} [\mathbf{c}_{k,1}, \dots, \mathbf{c}_{k,N}]$ separately, using SAGE (Fessler and Hero, 1994). SAGE is an extension of EM for data models with particular structures, including data generated by superimposed components. It is known

²KL-NMF is wrongly presented as ML in additive Poisson noise in numerous publications.

to converge faster in iterations than standard EM, though one iteration of SAGE is usually more computationally demanding than EM as it usually requires to update the sufficient statistics “more often”. Let us consider a partition of the parameter space $\boldsymbol{\theta} = \bigcup_{k=1}^K \boldsymbol{\theta}_k$ with

$$\boldsymbol{\theta}_k = \{\mathbf{w}_k, h_k\}, \quad (32)$$

where we recall that \mathbf{w}_k is the k^{th} column of \mathbf{W} and h_k is the k^{th} row of \mathbf{H} . The SAGE algorithm involves choosing for each subset of parameters $\boldsymbol{\theta}_k$ a *hidden-data space* which is complete for this particular subset. Here, the hidden-data space for $\boldsymbol{\theta}_k$ is simply chosen to be $\mathbf{C}_k \stackrel{\text{def}}{=} [\mathbf{c}_{k,1}, \dots, \mathbf{c}_{k,N}]$. An EM-like functional is then built for *each* subset $\boldsymbol{\theta}_k$ as the conditional expectation of the minus log likelihood of \mathbf{C}_k , which writes

$$Q_k^{ML}(\boldsymbol{\theta}_k | \boldsymbol{\theta}') \stackrel{\text{def}}{=} - \int_{\mathbf{C}_k} \log p(\mathbf{C}_k | \boldsymbol{\theta}_k) p(\mathbf{C}_k | \mathbf{X}, \boldsymbol{\theta}') d\mathbf{C}_k. \quad (33)$$

One iteration i of the SAGE algorithm then consists of computing (E-step) and minimizing (M-step) $Q_k^{ML}(\boldsymbol{\theta}_k | \boldsymbol{\theta}')$ for $k = 1, \dots, K$. Note that $\boldsymbol{\theta}'$ always contains the most up-to-date parameter values, and not only the values at iteration $i - 1$ like in standard EM. This leads to the above-mentioned increase in computational burden, which is mild in our case.

The derivations of the SAGE algorithm for IS-NMF are detailed in Appendix B. However, for a fixed k , the E-step merely consists of computing the posterior power \mathbf{V}_k of component \mathbf{C}_k , defined by $[\mathbf{V}_k]_{fn} = v_{k,fn} = |\mu_{k,fn}^{post}|^2 + \lambda_{k,fn}^{post}$, where $\mu_{k,fn}^{post}$ and $\lambda_{k,fn}^{post}$ are the posterior mean and variance of $c_{k,fn}$, given by

$$\mu_{k,fn}^{post} = \frac{w_{fk} h_{kn}}{\sum_l w_{fl} h_{ln}} x_{fn}, \quad (34)$$

$$\lambda_{k,fn}^{post} = \frac{w_{fk} h_{kn}}{\sum_l w_{fl} h_{ln}} \sum_{l \neq k} w_{fl} h_{ln}. \quad (35)$$

The M-step is then shown to amount to the following one-component NMF problem

$$\min_{\mathbf{w}_k, h_k \geq 0} D_{IS}(\mathbf{V}'_k | \mathbf{w}_k h_k) \quad (36)$$

where \mathbf{V}'_k denotes \mathbf{V}_k as computed from $\boldsymbol{\theta}'$. Interestingly, in the one-component case, the gradients simplify to

$$\nabla_{h_{kn}} Q_k^{ML}(\mathbf{w}_k, h_k | \boldsymbol{\theta}') = \frac{F}{h_{kn}} - \frac{1}{h_{kn}^2} \sum_{f=1}^F \frac{v'_{k,fn}}{w_{fk}}, \quad (37)$$

$$\nabla_{w_{fk}} Q_k^{ML}(\mathbf{w}_k, h_k | \boldsymbol{\theta}') = \frac{N}{w_{fk}} - \frac{1}{w_{fk}^2} \sum_{n=1}^N \frac{v'_{k,fn}}{h_{kn}}. \quad (38)$$

The gradients are easily zeroed, leading to the following updates

$$h_{kn}^{(i+1)} = \frac{1}{F} \sum_f \frac{v'_{k,fn}}{w_{fk}^{(i)}}, \quad (39)$$

$$w_{fk}^{(i+1)} = \frac{1}{N} \sum_n \frac{v'_{k,fn}}{h_{kn}^{(i+1)}}, \quad (40)$$

which guarantees $Q_k^{ML}(\mathbf{w}_k^{(i+1)}, h_k^{(i+1)} | \boldsymbol{\theta}') \leq Q_k^{ML}(\mathbf{w}_k^{(i)}, h_k^{(i)} | \boldsymbol{\theta}')$. This can also be written in matrix form, as shown in Algorithm 2, which summarizes the SAGE algorithm for IS-NMF. In the following, we will refer to this algorithm as “IS-NMF/EM”.

IS-NMF/EM and IS-NMF/MU have same complexity $\mathcal{O}(12 F K N)$ per iteration, but can lead to different run times, as shown in the results below. Indeed, in our Matlab implementation,

Algorithm 2 IS-NMF/EM

Input : nonnegative matrix \mathbf{V}
Output : nonnegative matrices \mathbf{W} and \mathbf{H} such that $\mathbf{V} \approx \mathbf{WH}$
Initialize \mathbf{W} and \mathbf{H} with nonnegative values
for $i = 1 : n_{iter}$ **do**
 for $k = 1 : K$ **do**
 Compute $\mathbf{G}_k = \frac{\mathbf{w}_k h_k}{\mathbf{WH}}$ % Wiener gain
 Compute $\mathbf{V}_k = \mathbf{G}_k^{[2]} \cdot \mathbf{V} + (1 - \mathbf{G}_k) \cdot (\mathbf{w}_k h_k)$ % Posterior power of \mathbf{C}_k
 $h_k \leftarrow \frac{1}{F} (\mathbf{w}_k^{[-1]})^T \mathbf{V}_k$ % Update row k of \mathbf{H}
 $\mathbf{w}_k \leftarrow \frac{1}{N} \mathbf{V}_k (h_k^{[-1]})^T$ % Update column k of \mathbf{W}
 Normalize \mathbf{w}_k and h_k
 end for
end for

% Note that \mathbf{WH} needs to be computed only once, at initialization, and be
subsequently updated as $\mathbf{WH} = \mathbf{w}_k^{old} h_k^{old} + \mathbf{w}_k^{new} h_k^{new}$.

the operations in IS-NMF/MU can be efficiently vectorized using matrix entrywise multiplication, while IS-NMF/EM requires looping over the components, which is more time consuming.

The convergence of IS-NMF/EM to a stationary point of $D_{IS}(\mathbf{V}|\mathbf{WH})$ is granted by property of SAGE. However, it can only converge to a point in the interior domain of the parameter space, i.e, \mathbf{W} and \mathbf{H} cannot take entries equal to zero. This is seen in equation (36): if either w_{fk} or h_{kn} is zero, then the cost $d_{IS}(v'_{k,fn}|w_{fk}h_{kn})$ becomes infinite. This is not a feature shared by IS-NMF/MU, which does not a priori exclude zero coefficients in \mathbf{W} and \mathbf{H} (but excludes $\hat{v}_{fn} = 0$, which would lead to a division by zero). However, because zero coefficients are invariant under multiplicative updates (see Section 2.1), if IS-NMF/MU attains a fixed point solution with zero entries, then it cannot be determined if the limit point is a stationary point. Yet, if the limit point does not take zero entries (i.e, belongs to the interior of the parameter space) then it is a stationary point, which may or may not be a local minimum. This is stressed by [Berry et al. \(2007\)](#) for EUC-NMF but holds for IS-NMF/MU as well.

Note that SAGE has been used in the context of single-channel source separation by [Ozerov et al. \(2007\)](#) for inference on a model somehow related to the IS-NMF model (16). Indeed, these authors address voice/music separation using a generative model of the form $\mathbf{x}_n = \mathbf{c}_{V,n} + \mathbf{c}_{M,n}$ where the first component represents voice while the second one represents music. Then, each component is given a Gaussian mixture model (GMM). The GMM parameters for voice are learnt from training data, while the music parameters are adapted to data. Though related, the GMM and NMF models are quite different in essence. The first one expresses the signal as a sum of two components that can each take different states. The second one expresses the signal as a sum of K components, each representative of one object. It cannot be claimed that one model is better than the other, but that they rather address different characteristics. It is anticipated that the two models can be used jointly within the SAGE framework, for example, by modelling voice $\mathbf{c}_{V,n}$ with a GMM (i.e, a specific component with many states) and music $\mathbf{c}_{M,n}$ with a NMF model (i.e, a composite signal with many components).

4 Analysis of a short piano excerpt

In this section we report an experimental comparative study of the above-mentioned NMF algorithms applied to the spectrogram of a short monophonic piano sequence. In a first step we compare the results of multiplicative Euclidean, KL and IS NMF algorithms for several values of K , before we more specifically compare the multiplicative and EM-based algorithms for IS-NMF in a second step.

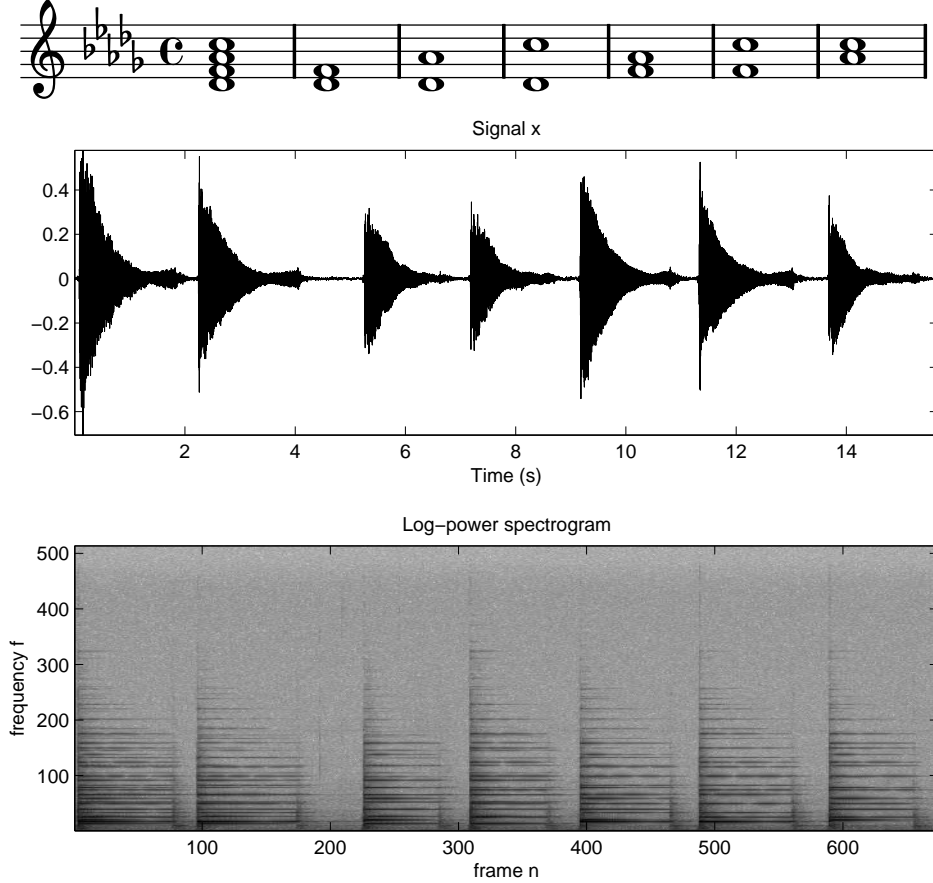


Figure 2: Three representations of data; (top): original score, (middle): time-domain recorded signal x , (bottom): log-power spectrogram $\log |\mathbf{X}| \cdot [2]$. The four notes read D_4^b (pitch 61), F_4 (pitch 65), A_4^b (pitch 68) and C_5 (pitch 72). They all together form a D^b major seventh chord. In the recorded interpretation the third chord is slightly out of tempo.

4.1 Experimental setup

A real piano sequence, played from score given in figure 2 on a Yamaha DisKlavier MX100A upright piano, was recorded in a small size room by a Schoeps omnidirectional microphone, placed about 15 cm (6 in) above the opened body of the piano. The sequence is composed of 4 notes, played all at once in the first measure and then played by pairs in all possible combinations in the subsequent measures. The 15.6 seconds long recorded signal was downsampled to $\nu_s = 22050$ Hz, yielding $T = 339501$ samples. A STFT \mathbf{X} of x was computed using a sinebell analysis window of length $L = 1024$ (46 ms) with 50 % overlap between two frames, leading to $N = 674$ frames and $F = 513$ frequency bins. The time-domain signal x and its log-power spectrogram are represented on figure 2.

IS-NMF/MU, IS-NMF/EM and the multiplicative gradient descent NMF algorithms with Euclidean and KL costs were implemented in Matlab and run on data $\mathbf{V} = |\mathbf{X}| \cdot [2]$. Note that in the following the terms “EUC-NMF” and “KL-NMF” will implicitly refer to the multiplicative implementation of these NMF techniques. All algorithms were run for several values of the number of components, more specifically for $K = 1, \dots, 10$. For each value of K , 10 runs of each algorithm were produced from 10 random initializations of \mathbf{W} and \mathbf{H} , chosen, in Matlab notations, as $\mathbf{W} = \text{abs}(\text{randn}(F, K)) + \text{ones}(F, K)$ and $\mathbf{H} = \text{abs}(\text{randn}(K, N)) + \text{ones}(K, N)$. The algorithms were run for $n_{iter} = 5000$ iterations.

K	1	2	3	4	5	10	$\mathcal{O}(\cdot)$
EUC-NMF	17	18	20	24	27	37	$4 FKN + 2 K^2(F + N)$
KL-NMF	90	90	92	100	107	117	$8 FKN$
IS-NMF/MU	127	127	129	135	138	149	$12 FKN$
IS-NMF/EM	81	110	142	171	204	376	$12 FKN$

Table 1: Run times in seconds of 1000 iterations of the NMF algorithms applied to the piano data, implemented in Matlab on a 2.16 GHz Intel Core 2 Duo iMac with 2 GB RAM. The run times include the computation of the cost function at each iteration (for possible convergence monitoring). The last column shows the algorithm complexities per iteration, expressed in number of flops (addition, soustraction, multiplication, division). The complexity of EUC-NMF assumes $K < F, N$.

4.2 Pitch estimation

In the following results, it will be observed that some of the basis elements (columns of \mathbf{W}) have a pitched structure, characteristic of individual musical notes. If pitch estimation is not the objective *per se* of the following study, it is informative to check if correct pitch values can be inferred from the factorization. As such, a fundamental frequency (or pitch) estimator is applied using the method described in (Vincent et al., 2007). It consists in computing dot products of \mathbf{w}_k with a set of J frequency combs and retaining the pitch number corresponding to the largest dot product. Each comb is a cosine function with period f_j , scaled and shifted to the amplitude interval $[0 \ 1]$, that takes its maximum value 1 at bins multiple of f_j . The set of fundamental frequency bins $f_j = \frac{\nu_j}{\nu_s} L$ is indexed on the MIDI logarithmic scale, i.e, such that

$$\nu_j = 440 \times 2^{\frac{p_j - 69}{12}}. \quad (41)$$

The piano note range usually goes from $p_{min} = 21$, i.e, note A_0 with fundamental frequency $f_{min} = 27.5$ Hz, to $p_{max} = 108$, i.e, note C_8 with frequency $f_{max} = 4186$ Hz. Two adjacent keys are separated by a semitone ($\Delta p = 1$). The MIDI pitch number of the notes pictured on figure 2 are 61 (D_4^b), 65 (F_4), 68 (A_4^b) and 72 (C_5), and were chosen arbitrarily. In our implementation of the pitch estimator, the MIDI range was sampled from 20.6 to 108.4 with step 0.2. In the following, an arbitrary pitch value of 0 will be given to unpitched basis elements; the classification of pitched and unpitched elements was done manually by looking at the basis elements and listening to the component reconstructions.

4.3 Results and discussion

Convergence behavior and algorithm complexities Run times of 1000 iterations of each of the four algorithms are shown in Table 1, together with the algorithm complexities. Figure 3 shows for each algorithm and for every value of K the final cost values of the 10 runs, after the 5000 algorithm iterations. A first observation is that the minimum and maximum cost values differs, for $K > 4$ in the Euclidean case, $K > 3$ in the KL case and $K > 2$ in the IS case. This either means that the algorithms have failed to converge after 5000 iterations in some cases, or suggests the presence of local minima. Figure 4 displays for all 4 algorithms the evolution of the cost functions along the 5000 iterations for all of the 10 runs, in the particular case $K = 6$.

Evolution of the factorizations with order K In this paragraph we examine in details the underlying semantics of the factorizations obtained with all three cost functions. We here only address the comparison of factorizations obtained from the three multiplicative algorithms. IS-NMF/EM and IS-NMF/MU will be more specifically compared in the next paragraph. Otherwise stated, the factorizations studied below are those obtained from the run yielding the minimum cost value among the 10 runs. Figures 5 to 8 display the columns of \mathbf{W} and corresponding rows of \mathbf{H} . The columns of \mathbf{W} are represented against frequency bin f on the left (in \log_{10} amplitude scale) while the rows of \mathbf{H} are represented against frame index n on the right (in linear amplitude

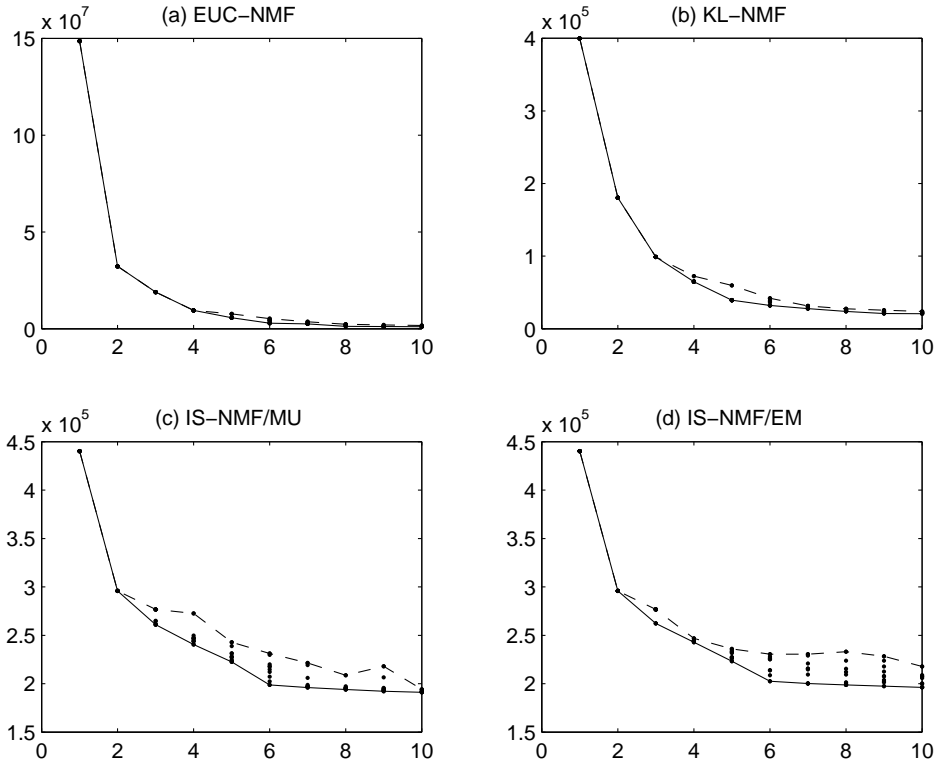


Figure 3: Cost values after 5000 iterations, obtained from 10 different random initializations. (a): Euclidean distance, (b): KL divergence, (c): IS divergence (using IS-NMF/MU), (d): IS divergence (using IS-NMF/EM). On each plot, the solid line connects all minimum cost values while the dashed line connects all maximum cost values.

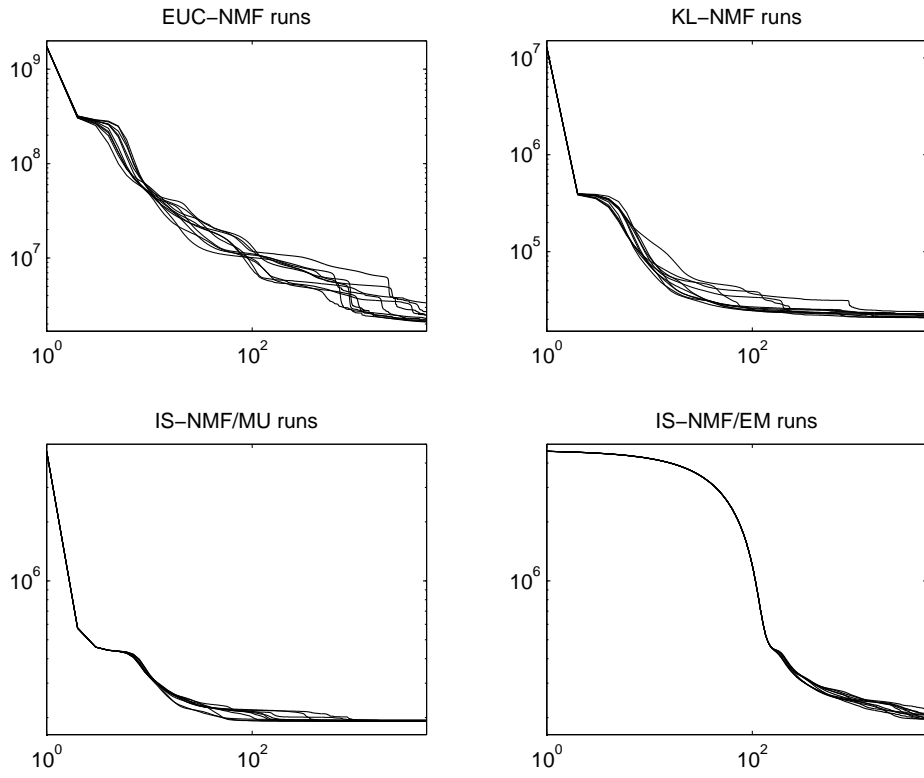


Figure 4: Evolution in log-log scale of the cost functions along the 5000 iterations of all 10 runs of the 4 algorithms, in the specific case $K = 6$.

scale). Pitched components are displayed first (top to down, in ascending order of estimated pitch value), followed by the unpitched components. For sake of conciseness only part of the results are reproduced in this article but we emphasize that the factorizations obtained with all four algorithms for $K = 4, 5, 6$ are available online at (Companion Web Site), together with sound reconstructions of the individual components. Component STFTs $\hat{\mathbf{C}}_k$ were computed by applying the Wiener filter (21) to \mathbf{X} using the factors \mathbf{W} and \mathbf{H} obtained with all 3 cost functions. Time-domain components c_k were then reconstructed by inverting the STFTs using an adequate overlap-add procedure with dual synthesis window. By conservativity of Wiener reconstruction and linearity of the inverse-STFT, the time-domain decomposition is also conservative, i.e, such that

$$x = \sum_{k=1}^K c_k. \quad (42)$$

Common sense suggests that choosing as many components as notes forms a sensible guess for the value of K , so as to obtain a *meaningful* factorization of $|\mathbf{X}|^{[2]}$ where each component would be expected to represent one and only one note. The factorizations obtained with all three costs for $K = 4$ prove that this is not the case. Euclidean and KL-NMF rather successfully extracts notes 65 and 68 into separate components (second and third), but notes 61 and 72 are melted into the first component while a fourth component seems to capture transient events corresponding to the note attacks (sound of the hammer hitting the string) and the sound produced by the release of the sustain pedal. The first two components obtained with IS-NMF have a similar interpretation to those given by EUC-NMF and KL-NMF. However the two other components differ in nature: the third component comprise note 68 *and* transients, while the fourth component is akin to residual noise. It is interesting to notice how this last component, though of much lower energy than the others components (in the order of 1 compared to 10^4 for the others) bears equal importance in the decomposition. This is undoubtedly a consequence of the scale invariance property of the IS divergence discussed in Section 2.2.

A fully separated factorization (at least as intended) is obtained for $K = 5$ with KL-NMF, as displayed on figure 5. This results in four components each made up of a single note and a fifth component containing sound events corresponding to note attacks and pedal releases. However these latter events are not well localized in time, and suffer from an unnatural tremolo effect (oscillating variations in amplitudes), as can be heard from the reconstructed sound files. Surprisingly, the decomposition obtained with EUC-NMF by setting $K = 5$, results in splitting the second component of the $K = 4$ decomposition in two components with estimated pitches 65 and 65.4, instead of actually demixing the third component which comprised notes 61 and 72. As for IS-NMF, the first component now groups notes 61 and 68, the second and third components respectively capture notes 65 and 72, the fourth component is still akin to residual noise, while the fifth component perfectly renders the attacks and releases.

Full separation of the individual notes is finally obtained with Euclidean and IS costs for $K = 6$, as shown on figures 6 and 7. KL-NMF produces an extra component (with pitch estimate 81) which is not clearly interpretable, and is in particular not akin to residual noise as could have been hoped for. The decomposition obtained with the IS cost describes as follows. The four first components correspond to individual notes whose pitch estimate match exactly the pitches of the notes played. The visual aspect of the PSDs is much better than the basis elements learnt from EUC-NMF and KL-NMF. The fifth component captures the hammer hits and pedal releases with great accuracy and the sixth component is akin to residual noise.

When the decomposition is carried on beyond $K = 6$, it is observed that EUC-NMF and KL-NMF split existing components into several subcomponents (such as components capturing sustained and decaying parts of one note) with pitch in the neighborhood of the note fundamental frequency. On the opposite, IS-NMF/MU spends the extra components in fine-tuning the representation of the low energy components, i.e, residual noise and transient events (as such, the hammer hits and pedal releases eventually get split in two distincts components). As such, for $K = 10$, the pitch estimates reads EUC-NMF: [61 64.8 64.8 65 65 65.8 68 68.4 72.2 0], KL-NMF:

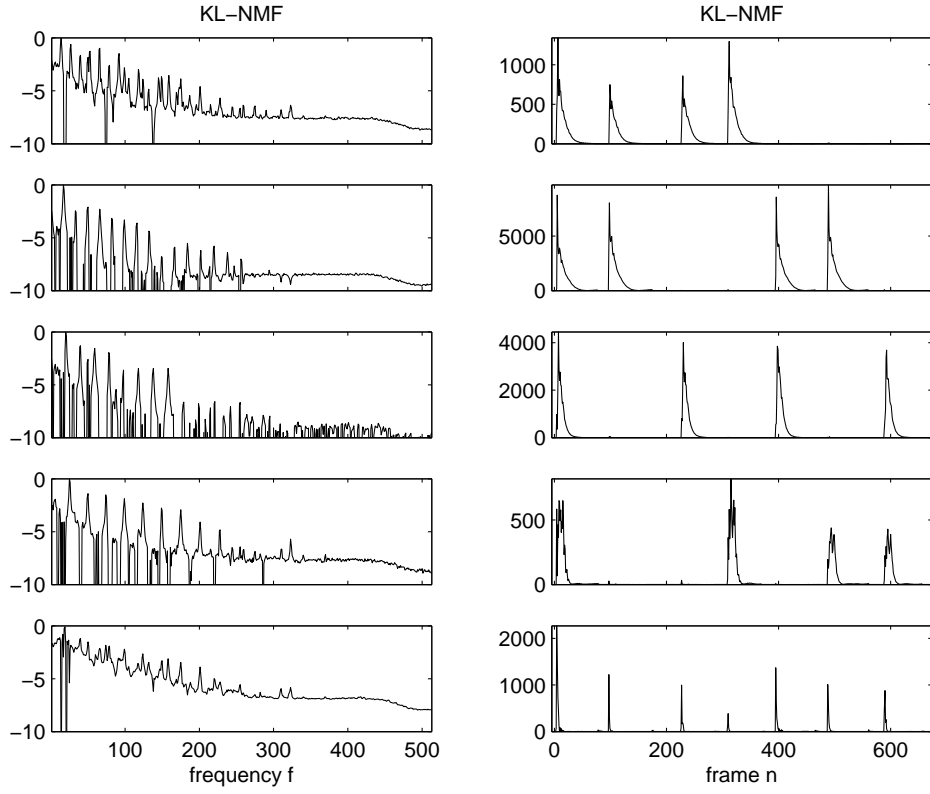


Figure 5: KL-NMF with $K = 5$. Pitch estimates: $[61 \ 65 \ 68 \ 72.2 \ 0]$. Left : columns of \mathbf{W} (\log_{10} scale). Right : rows of \mathbf{H} .

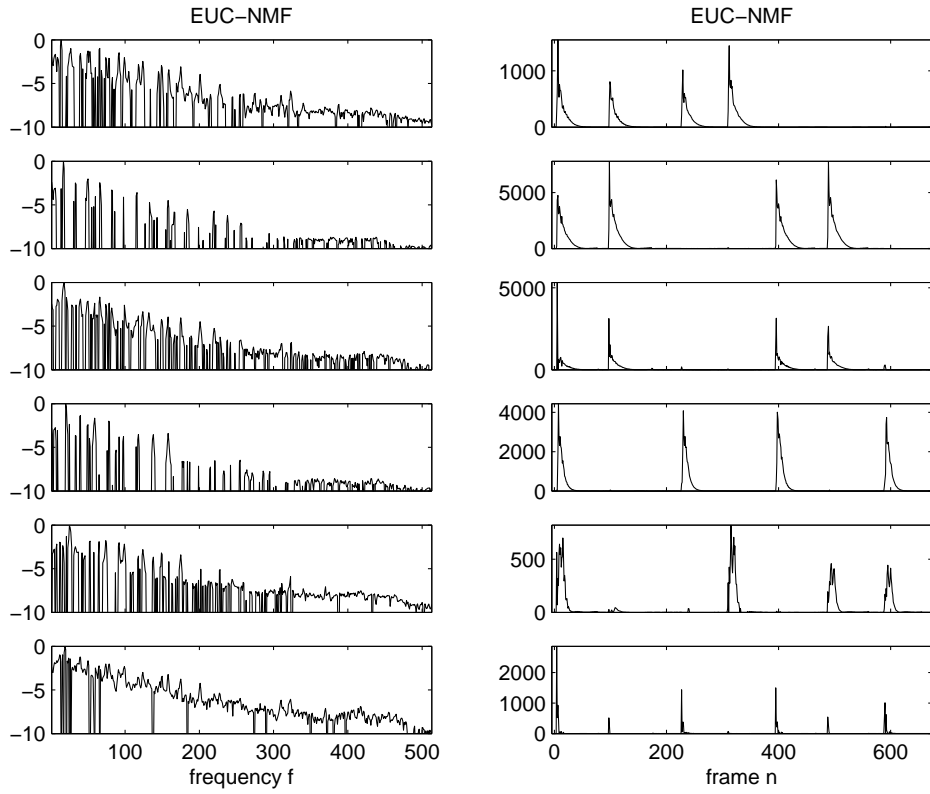


Figure 6: EUC-NMF with $K = 6$. Pitch estimates: [61 65 65.4 68 72 0]. Left : columns of \mathbf{W} (\log_{10} scale). Right : rows of \mathbf{H} .

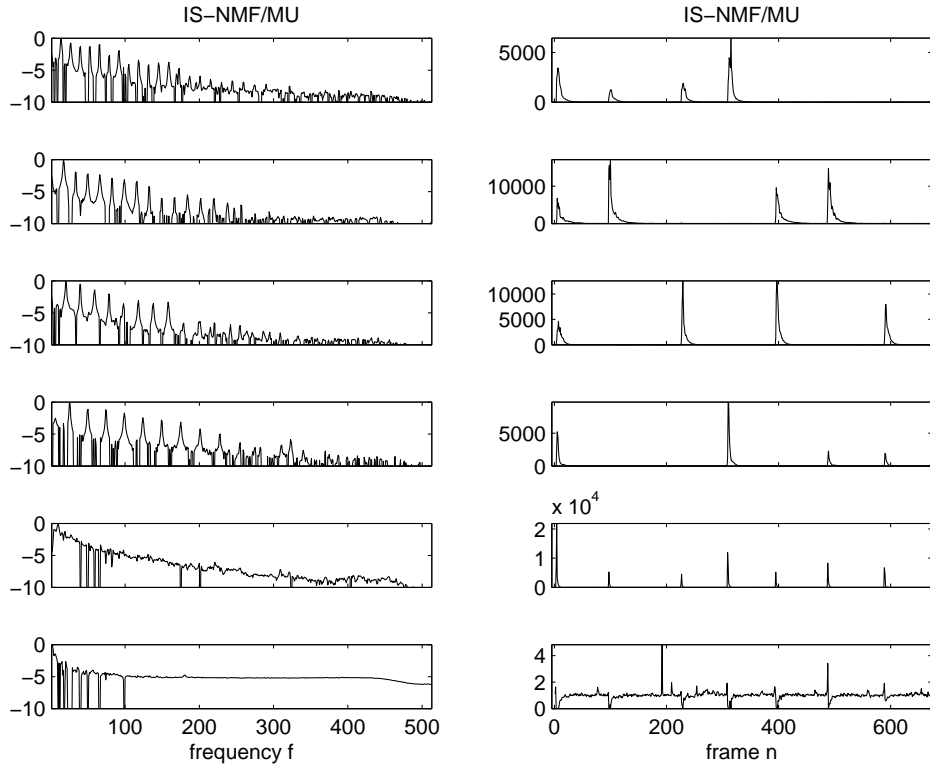


Figure 7: IS-NMF/MU with $K = 6$. Pitch estimates: [61 65 68 72 0 0]. Left : columns of \mathbf{W} (\log_{10} scale). Right : rows of \mathbf{H} .

[61 61 65 65 66 68 72 80.2 0 0], IS-NMF/MU: [61 61 65 68 72 0 0 0 0]. If note 61 is indeed split in 2 components with IS-NMF/MU, one of the two components is actually inaudible.

The message we want to bring out from this experimental study is the following. The nature of the decomposition obtained with IS-NMF, and its progression as K increases, is in accord with an *object-based* representation of music, close to our own comprehension of sound. Entities with well-defined semantics emerge from the decomposition (individual notes, hammer hits, pedal releases, residual noise) while the decompositions obtained from the Euclidean and KL costs are less interpretable from this perspective. We need to mention that these conclusions do not always hold when the factorization is not the one yielding the lowest cost values from the 10 runs. As such, we also examined the factorizations with highest cost values (with all three cost functions) and we found out that they did not reveal the same semantics, which was in turn not always easily interpretable. The upside however is that lowest IS cost values correspond to the most desirable factorizations indeed, so that IS-NMF “makes sense”.

Comparison of multiplicative and EM-based IS-NMF Algorithms IS-NMF/MU and IS-NMF/EM are designed to address the same task of minimizing the cost $D_{IS}(\mathbf{V}|\mathbf{WH})$, so that the achieved factorization should be identical in nature provided they complete this task. As such, the progression of the factorization provided by IS-NMF/EM is similar to the one observed for IS-NMF/MU and described in the previous paragraph. However, the resulting factorizations are not exactly equivalent, because IS-NMF/EM does not inherently allow zeros in the factors (see Section 3.2). This feature can be desirable for \mathbf{W} as the presence of sharp notches in the spectrum may not be physically realistic for audio, but can be considered a drawback as far as \mathbf{H} is concerned. Indeed, the rows of \mathbf{H} being akin to activation coefficients, when a sound object k is not present in frame n , then h_{kn} should be strictly zero. These remarks probably explain the factorization obtained from IS-NMF/EM with $K = 6$, displayed on figure 8. The notches present in the PSDs learnt with IS-NMF/MU, as seen on figure 7, have disappeared from the PSDs on figure 8, which exhibit better regularity. Unfortunately, IS-NMF/EM does not fully separate out the note attacks in the fifth component, like IS-NMF/MU does. Indeed, parts of the attacks appear in the second component, and the rest appears in the fifth component, which also contains the pedal releases. This is possibly explained by the a priori high sparsity of a transients component, which can be handled by IS-NMF/MU but not IS-NMF/EM (because it does not allow zero values in \mathbf{H}). Note that increasing the number of components K or the number of algorithm iterations n_{iter} does not solve this specific issue.

Regarding compared convergence of the algorithms, IS-NMF/MU decreases the cost function much faster in the initial iterations and, with this data set, attains lower final cost values than IS-NMF/EM, as shown on figure 3 or figure 4 for $K = 6$. As already mentioned, though the two algorithms have the same complexity, the run time per iteration of IS-NMF/MU is smaller than IS-NMF/EM for $K > 3$, see Table 1.

5 Regularized IS-NMF

We now describe how the statistical setting going along with IS-NMF can be exploited to incorporate regularization constraints/prior information in the factors estimates.

5.1 Bayesian setting

We consider a Bayesian setting where \mathbf{W} and \mathbf{H} are given (independent) prior distributions $p(\mathbf{W})$ and $p(\mathbf{H})$. We are looking for a joint MAP estimate of \mathbf{W} and \mathbf{H} through minimization of criterion

$$C_{MAP}(\mathbf{W}, \mathbf{H}) \stackrel{\text{def}}{=} -\log p(\mathbf{W}, \mathbf{H}|\mathbf{X}) \quad (43)$$

$$\stackrel{c}{=} D_{IS}(\mathbf{V}|\mathbf{WH}) - \log p(\mathbf{W}) - \log p(\mathbf{H}) \quad (44)$$

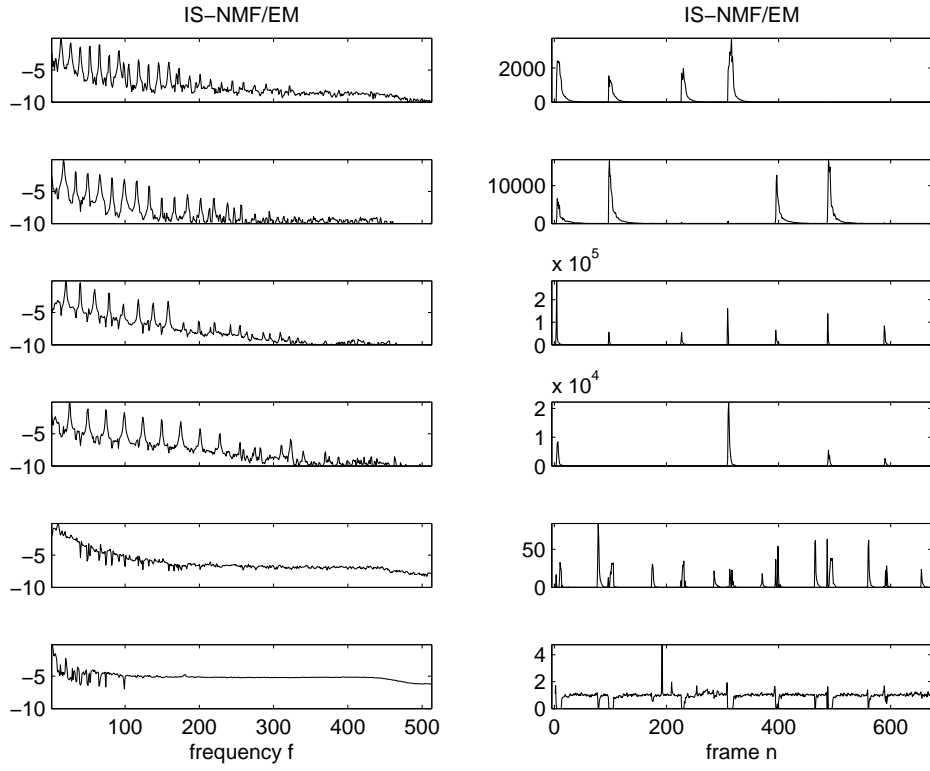


Figure 8: IS-NMF/EM with $K = 6$. Pitch estimates: [61 65 68 72 0 0]. Left : columns of \mathbf{W} (\log_{10} scale). Right : rows of \mathbf{H} .

When independent priors of the form $p(\mathbf{W}) = \prod_k p(\mathbf{w}_k)$ and $p(\mathbf{H}) = \prod_k p(h_k)$ are used, then the SAGE algorithm presented in Section 3.2 can again be used for MAP estimation. In that case, the functionals to be minimized for each component k write

$$Q_k^{MAP}(\boldsymbol{\theta}_k|\boldsymbol{\theta}') \stackrel{\text{def}}{=} - \int_{\mathbf{C}_k} \log p(\boldsymbol{\theta}_k|\mathbf{C}_k) p(\mathbf{C}_k|\mathbf{X}, \boldsymbol{\theta}') d\mathbf{C}_k \quad (45)$$

$$\stackrel{c}{=} Q_k^{ML}(\mathbf{w}_k, h_k|\boldsymbol{\theta}') - \log p(\mathbf{w}_k) - \log p(h_k) \quad (46)$$

Thus, the E-step still amounts to computing $Q_k^{ML}(\mathbf{w}_k, h_k|\boldsymbol{\theta}')$, as done in Section 3.2, and only the M-step is changed by the regularization constraints $-\log p(\mathbf{w}_k)$ and $-\log p(h_k)$ which now need to be taken into account.

Next we more specifically consider Markov chain priors favoring smoothness over the rows of \mathbf{H} . In the following results no prior structure will be assumed for \mathbf{W} (i.e. \mathbf{W} is estimated through ML). However, we stress that the methodology presented for the rows of \mathbf{H} can equivalently be transposed to the columns of \mathbf{W} , that prior structures can be imposed on both \mathbf{W} and \mathbf{H} and that these structures need not to belong to the same class of models. Note also that since the components are treated separately, they can each be given a different type of model (for example some components could be assigned a GMM, as discussed at the end of Section 3.2).

We assume the following prior structure for h_k ,

$$p(h_k) = \prod_{n=2}^N p(h_{kn}|h_{k(n-1)}) p(h_{k1}), \quad (47)$$

where $p(h_{kn}|h_{k(n-1)})$ is a pdf with mode $h_{k(n-1)}$. The motivation behind this prior is to constrain h_{kn} not to differ significantly from its value at entry $n-1$, hence favoring smoothness of the estimate. Possible pdf choices are, for $n = 2, \dots, N$,

$$p(h_{kn}|h_{k(n-1)}) = \mathcal{IG}(h_{kn}|\alpha, (\alpha+1)h_{k(n-1)}) \quad (48)$$

and

$$p(h_{kn}|h_{k(n-1)}) = \mathcal{G}(h_{kn}|\alpha, (\alpha-1)/h_{k(n-1)}) \quad (49)$$

where $\mathcal{G}(x|\alpha, \beta)$ is the previously introduced Gamma pdf, with mode $(\alpha-1)/\beta$ (for $\alpha \geq 1$) and $\mathcal{IG}(x|\alpha, \beta)$ is the inverse-Gamma pdf (see Appendix A), with mode $\beta/(\alpha+1)$. Both priors are constructed so that their mode is obtained for $h_{kn} = h_{k(n-1)}$. α is a shape parameter that controls the sharpness of the prior around its mode. A high value of α will increase sharpness and will thus accentuate smoothness of h_k while a low value of α will render the prior more diffuse and thus less constraining. The two priors become actually very similar for large values of α , as shown on figure 9. In the following, h_{k1} is assigned the scale-invariant Jeffreys noninformative prior $p(h_{k1}) \propto 1/h_{k1}$.

5.2 New updates

Under prior structure (47), the derivative of $Q_k^{MAP}(\mathbf{w}_k, h_k|\boldsymbol{\theta}')$ wrt h_{kn} writes, $\forall n = 2, \dots, N-1$,

$$\begin{aligned} \nabla_{h_{kn}} Q_k^{MAP}(\mathbf{w}_k, h_k|\boldsymbol{\theta}') = \\ \nabla_{h_{kn}} Q_k^{ML}(\mathbf{w}_k, h_k|\boldsymbol{\theta}') - \nabla_{h_{kn}} \log p(h_{k(n+1)}|h_{kn}) - \nabla_{h_{kn}} \log p(h_{kn}|h_{k(n-1)}) \end{aligned} \quad (50)$$

This is shown to be equal to

$$\nabla_{h_{kn}} Q_k^{MAP}(\mathbf{w}_k, h_k|\boldsymbol{\theta}') = \frac{1}{h_{kn}^2} (p_2 h_{kn}^2 + p_1 h_{kn} + p_0) \quad (51)$$

where the values of p_0 , p_1 and p_2 are specific to the type of prior employed (Gamma or inverse-Gamma chains), as given in Table 2. Updating h_{kn} then simply amounts to solving an order 2 polynomial. The polynomial has only one nonnegative root, given by

$$h_{kn} = \frac{\sqrt{p_1^2 - 4p_2 p_0} - p_1}{2p_2}. \quad (52)$$

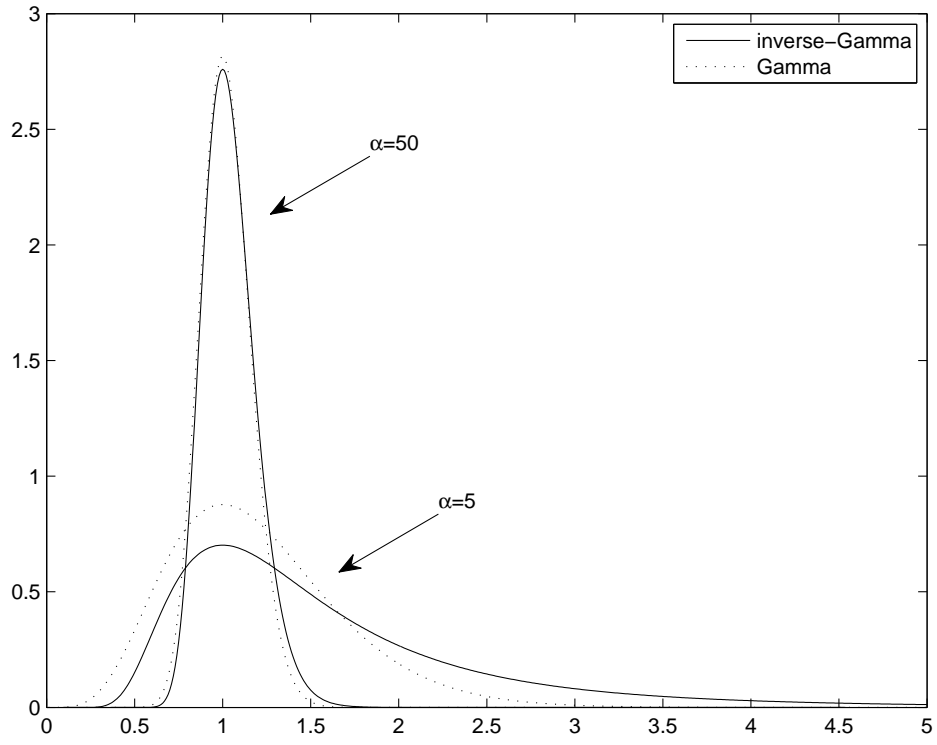


Figure 9: Prior pdfs $\mathcal{IG}(h_{kn}|\alpha-1, \alpha h_{k(n-1)})$ (solid line) and $\mathcal{G}(h_{kn}|\alpha+1, \alpha/h_{k(n-1)})$ (dashed line), for $h_{k(n-1)} = 1$ and for $\alpha = \{5, 50\}$.

inverse-Gamma Markov chain			
	p_2	p_1	p_0
h_{k1}	$(\alpha + 1)/h_{k2}$	$F - \alpha + 1$	$-F \hat{h}_{k1}^{ML}$
h_{kn}	$(\alpha + 1)/h_{k(n+1)}$	$F + 1$	$-F \hat{h}_{kn}^{ML} - (\alpha + 1) h_{k(n-1)}$
h_{kN}	0	$F + \alpha + 1$	$-F \hat{h}_{kN}^{ML} - (\alpha + 1) h_{k(N-1)}$
Gamma Markov chain			
	p_2	p_1	p_0
h_{k1}	0	$F + \alpha + 1$	$-F \hat{h}_{k1}^{ML} - (\alpha - 1) h_{k2}$
h_{kn}	$(\alpha - 1)/h_{k(n-1)}$	$F + 1$	$-F \hat{h}_{kn}^{ML} - (\alpha - 1) h_{k(n+1)}$
h_{kN}	$(\alpha - 1)/h_{k(N-1)}$	$F - \alpha + 1$	$-F \hat{h}_{kN}^{ML}$

Table 2: Coefficients of the order 2 polynomial to solve in order to update h_{kn} in Bayesian IS-NMF with a Markov chain prior. \hat{h}_{kn}^{ML} denotes the ML update, given by equation (39).

The coefficients h_{k1} and h_{kN} at the borders of the Markov chain require specific updates, but they also only require solving polynomials of order either 2 or 1, with coefficients given in Table 2 as well.

Note that the difference between the updates with the Gamma and inverse-Gamma chains prior mainly amounts to interchanging the positions of $h_{k(n-1)}$ and $h_{k(n+1)}$ in p_0 and p_2 . Interestingly, it can be noticed that using a backward Gamma chain prior $p(h_k) = \prod_{n=1}^{N-1} p(h_{kn}|h_{k(n+1)}) p(h_{kN})$ with shape parameter α is actually equivalent (in terms of MAP updates) to using a forward inverse-Gamma chain prior as in equation (47) with shape parameter $\alpha - 2$. Respectively, using a backward inverse-Gamma chain prior with shape parameter α is equivalent to using a forward Gamma chain prior with shape parameter $\alpha + 2$.

Note that [Virtanen et al. \(2008\)](#) have recently considered Gamma chains for regularization of KL-NMF. The modelling proposed in their work is however different than ours. Their Gamma chain prior is constructed in a hierarchical setting, i.e, by introducing extra auxiliary variables, so as to ensure conjugacy of the priors with the Poisson observation model. Estimation of the factors is then carried out with the standard gradient descent multiplicative approach and single-channel source separation results are presented from the factorization of the magnitude spectrogram $|\mathbf{X}|$ with component reconstruction (24). Regularized NMF algorithms for the Euclidean and KL costs with norm-2 constraints on $h_{kn} - h_{k(n-1)}$ have also been considered by [Chen et al. \(2006\)](#) and [Virtanen \(2007\)](#). Finally, we also wish to mention that [Shashanka et al. \(2008b\)](#) have recently derived in a Bayesian setting a regularized version of KL-NMF with sparsity constraints.

6 Learning the semantics of music with IS-NMF

The aim of the experimental study proposed in Section 4 was to analyze the results of several NMF algorithms on a short, simple and well-defined musical sequence, with respect to the cost function, initialization and model order. We now present results of NMF on a long polyphonic recording. Our goal is to examine how much of the semantics can NMF learn from the signal, with a fixed number of components and a fixed random initialization. This is not easily assessed numerically in the most general context, but quantitative evaluations could be performed on specific tasks in simulation settings. Such tasks could include music transcription, like in ([Abdallah and Plumbley, 2004](#)), single-channel source separation, like in [Benaroya et al. \(2003, 2006\)](#) or content-based music retrieval based on NMF features.

Rather than choosing and addressing one of these specific tasks, we here propose to use NMF in a real-case audio restoration scenario, where the purpose is to denoise and upmix original monophonic material (one channel) to stereo (two channels). This task is very close to single-channel source separation, with the difference that we are here not aiming at perfectly separating each of the sources, but rather isolating subsets of coherent components that can be given different

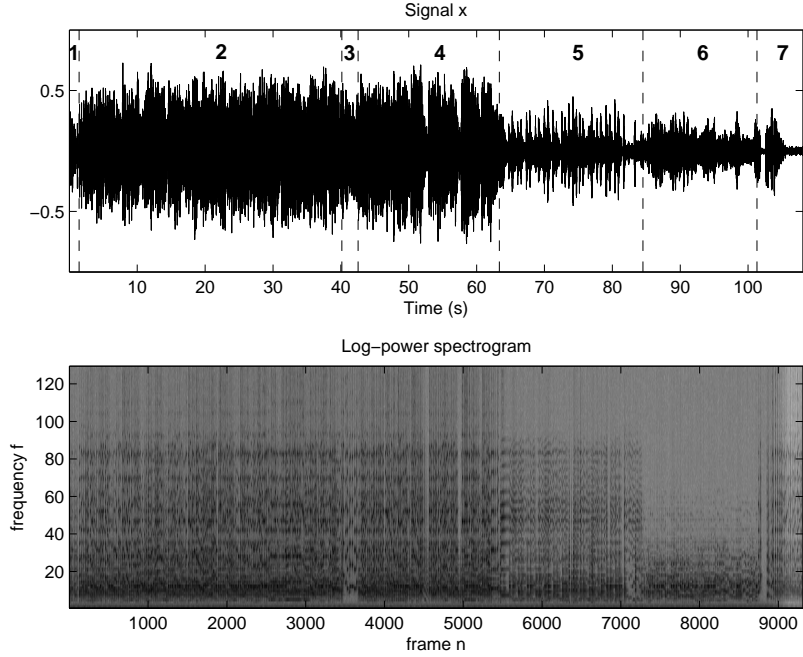


Figure 10: Original Louis Armstrong data; (top): time-domain recorded signal x , (bottom): log-power spectrogram. The vertical dashed lines on top plot identify successive phases in the music piece, that we annotated manually: (2,4,7) all instruments, (1) clarinet only, (3) trumpet solo, (5) clarinet and piano, (6) piano solo.

directions of arrival in the stereo remaster so as to render a sensation of spatial diversity. We will show in particular that the addition of smoothness constraints on the rows of \mathbf{H} lead to more pleasant component reconstructions, and better brings out the pitched structure of some of the learnt PSDs.

6.1 Experimental setup

We address the decomposition of a 108 seconds-long music excerpt from *My Heart (Will Always Lead Me Back To You)* recorded by Louis Armstrong and His Hot Five in the twenties. The band features (to our best hearing) a trumpet, a clarinet, a trombone, a piano and a double bass. The data is original unprocessed mono material containing substantial noise. The signal was downsampled to $\nu_s = 11025$ kHz, yielding $T = 1191735$ samples. The STFT \mathbf{X} of x was computed using a sinebell analysis window of length $L = 256$ (23 ms) with 50 % overlap between two frames, leading to $N = 9312$ frames and $F = 129$ frequency bins. The time-domain signal x and its log-power spectrogram are represented on figure 10.

We applied EUC-NMF, KL-NMF, IS-NMF/MU and IS-NMF/EM to $\mathbf{V} = |\mathbf{X}|^{[2]}$, as well as a regularized version of IS-NMF, as described in Section 5. We used the inverse-Gamma Markov chain prior (48) with α arbitrarily set to 10. We will refer to this algorithm as “IS-NMF/IG”. Among many trials, this value of α provided a good trade-off between smoothness of the component reconstructions and adequacy to data. Experiments with the Gamma Markov chain prior (48) did not lead to significant differences in the results and are not reported here.

The number of components K was arbitrarily set to 10. All five algorithms were run for $n_{iter} = 5000$ iterations and were initialized with same random values. For comparison, we have also applied KL-NMF to the magnitude spectrogram $|\mathbf{X}|$ with component reconstruction (24), as this can be considered state of the art methodology for NMF-based single-channel audio source

separation (Virtanen, 2007).

6.2 Results and discussion

For sake of conciseness we here only display the decomposition obtained with IS-NMF/IG, see figure 11, because it leads to the best results as far as our audio restoration task is concerned, but we stress that all decompositions and component reconstructions obtained from all NMF algorithms are available online at (Companion Web Site). Figure 11 displays the estimated basis functions \mathbf{W} in log-scale on the left, and represents on the right the time-domain signal components reconstructed from Wiener filtering.

Figure 12 displays the evolution of the IS cost along the 5000 iterations with IS-NMF/MU, IS-NMF/EM and IS-NMF/IG. In this case, IS-NMF/EM achieves a lower cost than IS-NMF/MU. The run times of 1000 iterations of the algorithms were respectively, EUC-NMF: 1.9 min, KL-NMF: 6.8 min, IS-NMF/MU: 8.7 min, IS-NMF/EM: 23.2 min and IS-NMF/IG: 32.2 min.

The comparison of the decompositions obtained with the three cost functions (Euclidean, KL and IS), through visual inspection of \mathbf{W} and listening of the components c_k , shows again that the IS divergence leads to the most interpretable results. In particular, some of the columns of matrix \mathbf{W} produced by all three IS-NMF algorithms have a clear pitched structure, which indicates that some notes have been extracted. Furthermore, one of the components captures the hiss noise from the recording. Discarding this component from the reconstruction of x yields satisfying denoising (this is particularly noticeable during the piano solo, where the input SNR is low). Very surprisingly, most of the rhythmic accompaniment (piano and double bass) is isolated in a single component (component 1 of IS-NMF/MU, component 2 of IS-NMF/EM and IS-NMF/IG), though its spectral content is clearly evolving in time. A similar effect happens with IS-NMF/IG and the trombone, which is mostly contained by component 7.

While we do not have a definite explanation for this, we believe that this is a consequence of Wiener reconstruction. Indeed the Wiener component reconstruction is only seen as a set of K masking filters applied to x_{fn} , so that it does not constrain the spectrum of component k to be exactly \mathbf{w}_k (up to amplitude h_{kn}), like the reconstruction method (24) does. So if one assumes that the NMF model (16) adequately captures some of the sound entities present in the mix (in our case that would be the preponderant notes or chords and the noise), then the other entities are bound to be relegated in remaining components, by conservativity of the decomposition $x = \sum_{k=1}^K c_k$.

As anticipated, the addition of frame-persistency constraints with IS-NMF/IG impacts the learnt basis \mathbf{W} . In particular, some of the components exhibit a more pronounced pitched structure. But more importantly, the regularization yields more pleasant sound reconstructions, this is particularly noticeable when listening to the accompaniment component obtained from IS-NMF/MU (component 1) or IS-NMF/EM (component 2) on the one side and from IS-NMF/IG (component 2) on the other side. Note also that in every case the sound quality of Wiener reconstructions is far better than state of the art KL-NMF of $|\mathbf{X}|$ and ad-hoc reconstruction (24).

To conclude this study, we provide online a restored version of the original recording, produced from the IS-NMF/IG decomposition. This is to our best knowledge the first use of NMF in a real-case audio restoration scenario. The restoration includes denoising (by discarding component 9, which is regarded as noise) and upmixing. A stereo mix is produced by dispatching parts of each component to left and right channels, hence simulating directions of arrival. As such, we manually created a mix where the components are arranged from 54° left to 54° right, such that the wind instruments (trumpet, clarinet, trombone) are placed left and the stringed instruments (piano, double bass) are placed right. While this stereo mix does render a sensation of spatialization we emphasize that its quality could undoubtedly be improved with appropriate sound engineering skills.

The originality of our restoration approach lays in 1) the joint noise removal and upmix (as opposed to a suboptimal sequential approach) and 2) the genuine content-based remastering, as opposed to standard techniques based, e.g, on phase delays and/or equalization.

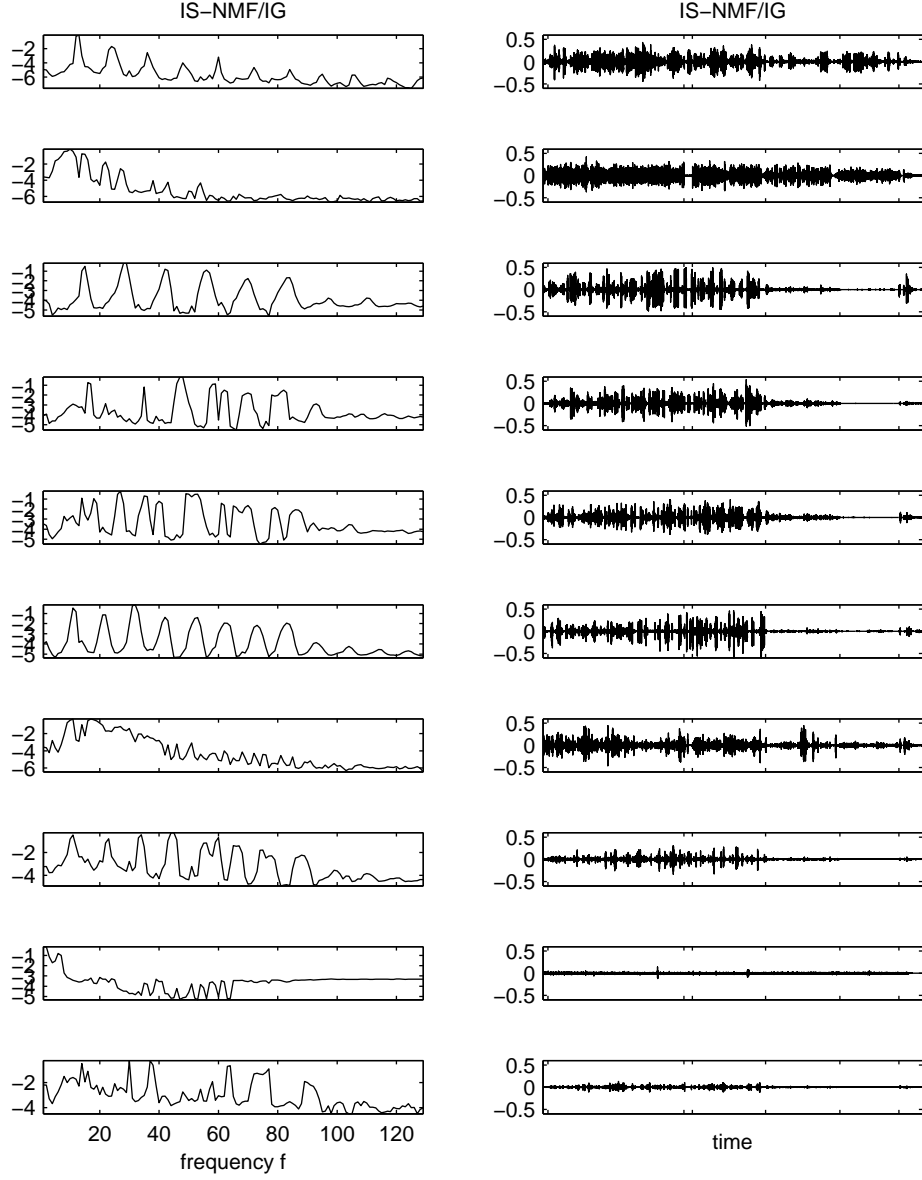


Figure 11: Decomposition of Louis Armstrong music data with IS-NMF/IG. Left : columns of \mathbf{W} (\log_{10} scale). Right : reconstructed components c_k ; the \underline{x} -axis ticks correspond to the temporal segmentation border lines displayed with signal x on figure 10. Component 2 captures most of the accompaniment, component 7 most of the trombone and component 9 most of the hiss noise. Summing up the other components leads to extracting the trumpet and clarinet, together with some piano notes.

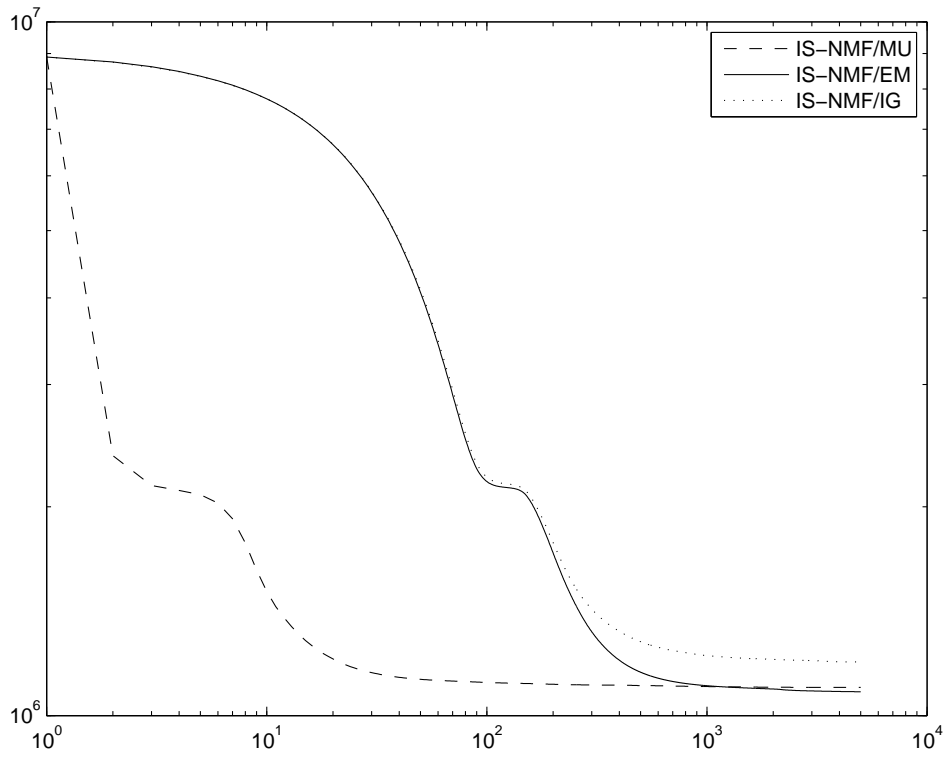


Figure 12: Evolution in log-log scale of the IS cost function along the 5000 iterations of IS-NMF/MU, IS-NMF/EM and IS-NMF/IG, initialized with same random values, with $K = 10$.

7 Conclusions

We have presented modelling and algorithmic aspects of NMF with the Itakura-Saito divergence. On the modelling side, we wish to bring out the following three features of IS-NMF that have been demonstrated in this paper;

- 1) IS-NMF is underlain by a statistical model of superimposed Gaussian components,
- 2) this model is relevant to the representation of audio signals,
- 3) this model can accommodate regularization constraints through Bayesian approaches.

On the algorithmic side, we have proposed a novel type of NMF algorithm, IS-NMF/EM, derived from SAGE, a variant of the EM algorithm. The convergence of this algorithm to a stationary point of the cost function $D_{IS}(\mathbf{V}|\mathbf{WH})$ is guaranteed by property of EM. This new algorithm was compared to an existing algorithm, IS-NMF/MU, whose convergence is not proved, though observed in practice. This article also reports an experimental comparative study of the standard EUC-NMF and KL-NMF algorithms, together with the two described IS-NMF algorithms, applied to a given data set (a short piano sequence), with various random initializations and model orders. Such a furnished experimental study was to our best knowledge not yet available. This article also reports a proof of concept of the use of IS-NMF for audio restoration, with a real example. Finally, we also believe to have shed light on the statistical implications of NMF with all of three cost functions.

We have shown how smoothness constraints on \mathbf{W} and \mathbf{H} can easily be handled in a Bayesian setting with IS-NMF. As such, we have shown how Markov chains prior structures can improve both the auditory quality of the component reconstructions and the interpretability of the basis elements. The Bayesian setting opens doors to even more elaborate prior structures that can better fit the specificities of data. For music signals we believe that two promising lines of research lay in 1) the use of switching state models for the rows of \mathbf{H} that explicitly model the possibility for h_{kn} to be strictly zero with a certain prior probability (and time persistency could be favored by modelling the state sequence with a discrete Markov chain), and 2) the use of models that explicitly take into account the pitched structure of some of the columns of \mathbf{W} , and where the fundamental frequency could act as a model parameter. These models fit into the problem of object-based representation of sound, which is an active area of research in the music information retrieval and auditory scene analysis communities.

In Section 4 we have compared the factorization results of a short piano power spectrogram, obtained from three costs functions, given a common algorithmic structure, i.e. standard multiplicative updates. The experiments illustrate the slow convergence of this type of algorithm, which has been already pointed out in other papers, e.g. (Cichocki et al., 2006a; Berry et al., 2007; Lin, 2007). If the proposed IS-NMF/EM does not improve on this issue, its strength is however to offer enough flexibility to accommodate Bayesian approaches. Now that we believe to have made our point that the IS cost is well suited to the factorization of audio power spectrograms (i.e. independently of the type of algorithm used), further work will address the development of faster IS-NMF algorithms. Following developments for other cost functions, we intend to investigate projected gradient techniques (Lin, 2007), exponentiated gradient descent and generalizations (Cichocki et al., 2006a), quasi-Newton second-order methods (Zdunek and Cichocki, 2007), as well as multilayered approaches (Cichocki and Zdunek, 2006).

Key issues that still need to be resolved in NMF concern identifiability and order selection. A related issue is the investigation into the presence of local minima in the cost functions, and ways to avoid them. In that matter, Markov chain Monte Carlo (MCMC) sampling techniques could be used as a diagnostic tool to better understand the topography of the criteria to minimize. While it is not clear whether these techniques can be applied to EUC-NMF or KL-NMF, they can readily be applied to IS-NMF, using its underlain Gaussian composite structure the same way IS-NMF/EM does. As to the avoidance of local minima, techniques inheriting from simulated annealing could be applied with IS-NMF, either in MCMC or EM inference.

Regarding order selection, usual criteria such as the Bayesian information criterion (BIC) or Akaike's criterion (see, e.g., [Stoica and Selén \(2004\)](#)) cannot be directly applied to IS-NMF, because the number of parameters ($F K + K N$) is not constant wrt the number of observations N . This feature breaks the validity of the assumptions in which these criteria have been designed. As such, a final promising line of research concerns the design of methods characterizing $p(\mathbf{V}|\mathbf{W})$ instead of $p(\mathbf{V}|\mathbf{W}, \mathbf{H})$, treating \mathbf{H} as a latent variable, like in independent component analysis ([MacKay, 1996](#); [Lewicki and Sejnowski, 2000](#)). Besides allowing for model order selection, such approaches would lead to more reliable estimation of the basis \mathbf{W} .

A Standard distributions

Proper complex Gaussian

$$\mathcal{N}_c(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = |\pi \boldsymbol{\Sigma}|^{-1} \exp(-(\mathbf{x} - \boldsymbol{\mu})^H \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}))$$

Poisson

$$\mathcal{P}(x|\lambda) = \exp(-\lambda) \frac{\lambda^x}{x!}$$

Gamma

$$\mathcal{G}(u|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} u^{\alpha-1} \exp(-\beta u), u \geq 0$$

inverse-Gamma

$$\mathcal{IG}(u|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} u^{-(\alpha+1)} \exp(-\frac{\beta}{u}), u \geq 0$$

The inverse-Gamma distribution is the distribution of $1/X$ when X is Gamma distributed.

B Derivations of the SAGE algorithm

In this appendix we detail the derivations leading to Algorithm 2. The functions involved in the definition of $Q_k^{ML}(\boldsymbol{\theta}_k|\boldsymbol{\theta}')$, given by equation (33), can be derived as follows. The hidden-data minus log likelihood writes

$$-\log p(\mathbf{C}_k|\boldsymbol{\theta}_k) = -\sum_{n=1}^N \sum_{f=1}^F \log \mathcal{N}_c(c_{k,fn}|0, h_{kn} w_{fk}) \quad (53)$$

$$\stackrel{c}{=} \sum_{n=1}^N \sum_{f=1}^F \log(w_{fk} h_{kn}) + \frac{|c_{k,fn}|^2}{w_{fk} h_{kn}}. \quad (54)$$

Then, the hidden-data posterior is obtained through Wiener filtering, yielding

$$p(\mathbf{C}_k|\mathbf{X}, \boldsymbol{\theta}) = \prod_{n=1}^N \prod_{f=1}^F \mathcal{N}_c(c_{k,fn}|\mu_{k,fn}^{post}, \lambda_{k,fn}^{post}), \quad (55)$$

with $\mu_{k,fn}^{post}$ and $\lambda_{k,fn}^{post}$ given by equations (34) and (35). The E-step is performed by taking the expectation of (54) wrt the hidden-data posterior, leading to

$$Q_k^{ML}(\boldsymbol{\theta}_k|\boldsymbol{\theta}') \stackrel{c}{=} \sum_{n=1}^N \sum_{f=1}^F \log(w_{fk} h_{kn}) + \frac{|\mu_{k,fn}^{post'}|^2 + \lambda_{k,fn}^{post'}}{w_{fk} h_{kn}} \quad (56)$$

$$\stackrel{c}{=} \sum_{n=1}^N \sum_{f=1}^F d_{IS}(|\mu_{k,fn}^{post'}|^2 + \lambda_{k,fn}^{post'} | w_{fk} h_{kn}). \quad (57)$$

The M-step thus amounts to minimizing $D_{IS}(\mathbf{V}'_k | \mathbf{w}_k h_k)$ wrt to $\mathbf{w}_k \geq 0$ and $h_k \geq 0$, as stated in Section 3.2.

Acknowledgement

The authors acknowledge Roland Badeau, Olivier Cappé, Jean-François Cardoso, Maurice Charbit and Alexey Ozerov for discussions related to this work. Many thanks to Gaël Richard for comments and suggestions about this manuscript and to Simon Godsill for helping us with the Louis Armstrong original data.

References

- S. A. Abdallah and M. D. Plumbley. Polyphonic transcription by nonnegative sparse coding of power spectra. In *5th International Symposium Music Information Retrieval (ISMIR'04)*, pages 318–325, Barcelona, Spain, Oct. 2004.
- L. Benaroya, R. Gribonval, and F. Bimbot. Non negative sparse representation for Wiener based source separation with a single sensor. In *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP'03)*, pages 613–616, Hong Kong, 2003.
- L. Benaroya, R. Blouet, C Févotte, and I. Cohen. Single sensor source separation using multiple-window STFT representation. In *Proc. International Workshop on Acoustic Echo and Noise Control (IWAENC'06)*, Paris, France, Sep. 2006.
- M. W. Berry, M. Browne, A. N. Langville, V. P. Pauca, and R. J. Plemmons. Algorithms and applications for approximate nonnegative matrix factorization. *Computational Statistics & Data Analysis*, 52(1):155–173, Sep. 2007.
- N. Bertin, R. Badeau, and G. Richard. Blind signal decompositions for automatic transcription of polyphonic music: NMF and K-SVD on the benchmark. In *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP'07)*, Honolulu, Hawaii, USA, 2007.
- Z. Chen, A. Cichocki, and T. M. Rutkowski. Constrained non-negative matrix factorization method for EEG analysis in early detection of Alzheimer’s disease. In *Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP'06)*, Toulouse, France, May 2006.
- A. Cichocki and R. Zdunek. Multilayer nonnegative matrix factorization. *Electronics Letters*, 42(16):947–948, 2006.
- A. Cichocki, S.-I. Amari, R. Zdunek, R. Kompass, G. Hori, and Z. He. Extended SMART algorithms for non-negative matrix factorization. In *Proc. International Conference on Artificial Intelligence and Soft Computing (ICAISC'06)*, pages 548–562, Zakopane, Poland, June 2006a.
- A. Cichocki, R. Zdunek, and S. Amari. Csiszar’s divergences for non-negative matrix factorization: Family of new algorithms. In *6th International Conference on Independent Component Analysis and Blind Signal Separation (ICA'06)*, pages 32–39, Charleston SC, USA, Mar. 2006b.
- I. Cohen and S. Gannot. Spectral enhancement methods. In M. M. Sondhi J. Benesty and Y. Huang, editors, *Springer Handbook of Speech Processing*, chapter 45. Springer, 2007.
- Companion Web Site. <http://www.tsi.enst.fr/~fevotte/Samples/is-nmf>.
- I. S. Dhillon and S. Sra. Generalized nonnegative matrix approximations with Bregman divergences. *Advances in Neural Information Processing Systems*, 19, 2005.
- K. Drakakis, S. Rickard, R. de Frein, and A. Cichocki. Analysis of financial data using non-negative matrix factorization. *International Journal of Mathematical Sciences*, 6(2), June 2007.
- S. Eguchi and Y. Kano. Robustifying maximum likelihood estimation. Technical report, Institute of Statistical Mathematics, June 2001. Research Memo. 802. http://www.ism.ac.jp/~eguchi/pdf/Robustify_MLE.pdf.
- M. Feder and E. Weinstein. Parameter estimation of superimposed signals using the EM algorithm. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 36(4):477–489, Apr. 1988.
- J. A. Fessler and A. O. Hero. Space-alternating generalized expectation-maximization algorithm. *IEEE Transactions on Signal Processing*, 42(10):2664–2677, Oct. 1994.
- R. M. Gray, A. Buzo, A. H. Gray, and Y. Matsuyama. Distortion measures for speech processing. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 28(4):367–376, Aug. 1980.
- F. Itakura and S. Saito. Analysis synthesis telephony based on the maximum likelihood method. In *Proc 6th International Congress on Acoustics*, pages C–17 – C–20, Tokyo, Japan, Aug. 1968.

- R. Kompass. A generalized divergence measure for nonnegative matrix factorization. *Neural Computation*, 19(3):780–791, 2007.
- D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *Advances in Neural and Information Processing Systems 13*, pages 556–562, 2001.
- D. D. Lee and H. S. Seung. Learning the parts of objects with nonnegative matrix factorization. *Nature*, 401:788–791, 1999.
- M. S. Lewicki and T. J. Sejnowski. Learning overcomplete representations. *Neural Computation*, 12:337–365, 2000.
- C.-J. Lin. Projected gradient methods for nonnegative matrix factorization. *Neural Computation*, 19:2756–2779, 2007.
- D. MacKay. Maximum likelihood and covariant algorithms for independent component analysis. <http://www.inference.phy.cam.ac.uk/mackay/ica.pdf>, 1996. Unpublished.
- A. Ozerov, P. Philippe, F. Bimbot, and R. Gribonval. Adaptation of Bayesian models for single-channel source separation and its application to voice/music separation in popular songs. *IEEE Transactions on Audio, Speech, and Language Processing*, 15(5):1564–1578, Jul. 2007.
- M. D. Plumbley, S. A. Abdallah, T. Blumensath, and M. E. Davies. Sparse representations of polyphonic music. *Signal Processing*, 86(3):417–431, Mar. 2006.
- M. Shashanka, B. Raj, and P. Smaragdis. Probabilistic latent variable models as nonnegative factorizations. *Computational Intelligence and Neuroscience*, 2008(Article ID 947438):8 pages, 2008a. doi:10.1155/2008/947438.
- M. Shashanka, B. Raj, and P. Smaragdis. Sparse overcomplete latent variable decomposition of counts data. In J.C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 1313–1320. MIT Press, Cambridge, MA, 2008b.
- P. Smaragdis. Convolutional speech bases and their application to speech separation. *IEEE Transactions on Audio, Speech, and Language Processing*, 15(1):1–12, Jan. 2007.
- P. Smaragdis and J. C. Brown. Non-negative matrix factorization for polyphonic music transcription. In *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, Oct. 2003.
- P. Stoica and Y. Selén. Model-order selection: a review of information criterion rules. *IEEE Signal Processing Magazine*, 21(4):36–47, Jul. 2004.
- E. Vincent, N. Bertin, and R. Badeau. Two nonnegative matrix factorization methods for polyphonic pitch transcription. In *Proc. Music Information Retrieval Evaluation eXchange (MIREX)*, 2007.
- T. Virtanen. Monaural sound source separation by non-negative matrix factorization with temporal continuity and sparseness criteria. *IEEE Transactions on Audio, Speech and Language Processing*, 15(3):1066–1074, Mar. 2007.
- T. Virtanen, A. T. Cemgil, and S. Godsill. Bayesian extensions to non-negative matrix factorisation for audio signal modelling. In *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP'08)*, pages 1825–1828, Las Vegas, Nevada, USA, Apr. 2008.
- S. S. Young, P. Fogel, and D. Hawkins. Clustering Scotch whiskies using non-negative matrix factorization. *Joint Newsletter for the Section on Physical and Engineering Sciences and the Quality and Productivity Section of the American Statistical Association*, 14(1):11–13, June 2006.
- R. Zdunek and A. Cichocki. Nonnegative matrix factorization with constrained second-order optimization. *Signal Processing*, 87(8):1904–1916, Aug. 2007.