

Convex nonnegative matrix factorization with missing data

Ronan Hamon, Valentin Emiya, Cédric Févotte

► **To cite this version:**

Ronan Hamon, Valentin Emiya, Cédric Févotte. Convex nonnegative matrix factorization with missing data. IEEE International Workshop on Machine Learning for Signal Processing, Sep 2016, Vietri sul Mare, Salerno, Italy. hal-01346492

HAL Id: hal-01346492

<https://hal-amu.archives-ouvertes.fr/hal-01346492>

Submitted on 7 Oct 2016

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

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

CONVEX NONNEGATIVE MATRIX FACTORIZATION WITH MISSING DATA

Ronan Hamon, Valentin Emiya*

Cédric Févotte

Aix Marseille Univ, CNRS, LIF, Marseille, France

CNRS & IRIT, Toulouse, France

ABSTRACT

Convex nonnegative matrix factorization (CNMF) is a variant of nonnegative matrix factorization (NMF) in which the components are a convex combination of atoms of a known dictionary. In this contribution, we propose to extend CNMF to the case where the data matrix and the dictionary have missing entries. After a formulation of the problem in this context of missing data, we propose a majorization-minimization algorithm for the solving of the optimization problem incurred. Experimental results with synthetic data and audio spectrograms highlight an improvement of the performance of reconstruction with respect to standard NMF. The performance gap is particularly significant when the task of reconstruction becomes arduous, e.g. when the ratio of missing data is high, the noise is steep, or the complexity of data is high.

Index Terms— matrix factorization, nonnegativity, low-rankness, matrix completion, spectrogram inpainting

1. INTRODUCTION

Convex NMF (CNMF) [1] is a special case of nonnegative matrix factorization (NMF) [2], in which the matrix of components is constrained to be a linear combination of atoms of a known dictionary. The term “convex” refers to the constraint of the linear combination, where the combination coefficients forming each component are nonnegative and sum to 1. Compared to the fully unsupervised NMF setting, the use of known atoms is a source of supervision that may guide learning based on this additional data: in particular, an interesting case of CNMF consists in auto-encoding the data themselves, by defining the atoms as the data matrix. CNMF has been of interest in a number of contexts, such as clustering, data analysis, face recognition, or music transcription [1, 3]. It is also related to the *self-expressive* dictionary-based representation proposed in [4].

An issue that has not yet been addressed is when the data matrix has missing coefficients. Such an extension of CNMF is worth being considered, as it opens the way to data-reconstruction settings with nonnegative low-rank constraints, which covers several relevant applications. One

example concerns the field of image or audio inpainting [5, 6, 7, 8], where CNMF may improve the current reconstruction techniques. In inpainting of audio spectrograms for example, setting up the dictionary to be a comprehensive collection of notes from a specific instrument may guide the factorization toward a realistic and meaningful decomposition, increasing the quality of the reconstruction of the missing data. In this contribution, we also consider the case where the dictionary may have missing coefficients itself.

The paper is organized as follows. Section 2 formulates CNMF in the presence of missing entries in the data matrix and in the dictionary. Section 3 describes the proposed majorization-minimization (MM) algorithm. Sections 4 and 5 report experimental results with synthetic data and audio spectrograms.

2. CONVEX NONNEGATIVE MATRIX FACTORIZATION WITH MISSING DATA

2.1. Notations and definitions

For any integer N , the integer set $\{1, 2, \dots, N\}$ is denoted by $[N]$. The coefficients of a matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ are denoted by either a_{mn} or $[\mathbf{A}]_{mn}$. The element-wise matrix product, matrix division and matrix power are denoted by $\mathbf{A} \cdot \mathbf{B}$, $\frac{\mathbf{A}}{\mathbf{B}}$ and \mathbf{A}^γ , respectively where \mathbf{A} and \mathbf{B} are matrices with same dimensions and γ is a scalar. $\mathbf{0}$ and $\mathbf{1}$ denote vectors or matrices composed of zeros and ones, respectively, with dimensions that can be deduced from the context. The element-wise negation of a binary matrix \mathbf{M} is denoted by $\bar{\mathbf{M}} \triangleq \mathbf{1} - \mathbf{M}$.

2.2. NMF and Convex NMF

NMF consists in approximating a data matrix $\mathbf{V} \in \mathbb{R}_+^{F \times N}$ as the product \mathbf{WH} of two nonnegative matrices $\mathbf{W} \in \mathbb{R}_+^{F \times K}$ and $\mathbf{H} \in \mathbb{R}_+^{K \times N}$. Often, $K < \min(F, N)$, such that \mathbf{WH} is a low-rank approximation of \mathbf{V} . Every sample \mathbf{v}_n , the n -th column of \mathbf{V} , is thus decomposed as a linear combination of K elementary *components* or *patterns* $\mathbf{w}_1, \dots, \mathbf{w}_K \in \mathbb{R}_+^F$, the columns of \mathbf{W} . The coefficients of the linear combination are given by the n -th column \mathbf{h}_n of \mathbf{H} .

In [9] and [10], algorithms have been proposed for the unsupervised estimation of \mathbf{W} and \mathbf{H} from \mathbf{V} , by minimization

*This work was supported by ANR JCJC program MAD (ANR-14-CE27-0002).

of the cost function $D_\beta(\mathbf{V}|\mathbf{WH}) = \sum_{fn} d_\beta(v_{fn} | [\mathbf{WH}]_{fn})$, where $d_\beta(x|y)$ is the β -divergence defined as:

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

When ill-defined, we set by convention $d_\beta(0|0) = 0$.

CNMF is a variant of NMF in which $\mathbf{W} = \mathbf{SL}$. $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_P] \in \mathbb{R}_+^{F \times P}$ is a nonnegative matrix of *atoms* and $\mathbf{L} = [\mathbf{l}_1, \dots, \mathbf{l}_K] \in \mathbb{R}_+^{P \times K}$ is the so-called *labeling* matrix. Each dictionary element \mathbf{w}_k is thus equal to $\mathbf{S}\mathbf{l}_k$, with usually $P \gg K$, and the data is in the end decomposed as $\mathbf{V} = \mathbf{SLH}$. The scale indeterminacy between \mathbf{L} and \mathbf{H} may be lifted by imposing $\|\mathbf{l}_k\|_1 = 1$, in which case \mathbf{w}_k is precisely a convex combination of the elements of the subspace \mathbf{S} . CNMF can be related to the so-called archetypal analysis [11], but without considering any nonnegativity constraint.

The use of known examples in \mathbf{S} can then be seen as a source of supervision that guides learning. A special case of CNMF is obtained by setting $\mathbf{S} = \mathbf{V}$, thus auto-encoding the data as \mathbf{VLH} . This particular case is studied in depth in [1]. In this paper, we consider the general case for \mathbf{S} , with or without missing data.

2.3. Convex NMF with missing data

We assume that some coefficients in \mathbf{V} and \mathbf{S} may be missing. Let $\mathcal{V} \subset [F] \times [N]$ be a set of pairs of indices that locates the observed coefficients in \mathbf{V} : $(f, n) \in \mathcal{V}$ iff v_{fn} is known. Similarly, let $\mathcal{S} \subset [F] \times [P]$ be a set of pairs of indices that locates the observed coefficients in \mathbf{S} . The use of sets \mathcal{V} and \mathcal{S} may be reformulated equivalently by defining masking matrices $\mathbf{M}_\mathcal{V} \in \{0, 1\}^{F \times N}$ and $\mathbf{M}_\mathcal{S} \in \{0, 1\}^{F \times P}$ from \mathcal{V} and \mathcal{S} as

$$[\mathbf{M}_\mathcal{V}]_{fn} \triangleq \begin{cases} 1 & \text{if } (f, n) \in \mathcal{V} \\ 0 & \text{otherwise} \end{cases} \quad \forall (f, n) \in [F] \times [N] \quad (2)$$

$$[\mathbf{M}_\mathcal{S}]_{fp} \triangleq \begin{cases} 1 & \text{if } (f, p) \in \mathcal{S} \\ 0 & \text{otherwise} \end{cases} \quad \forall (f, p) \in [F] \times [P] \quad (3)$$

A major goal in this paper is to estimate \mathbf{L} , \mathbf{H} and the missing entries in \mathbf{S} , given the partially observed data matrix \mathbf{V} . Denoting by \mathbf{S}° the set of observed/known dictionary matrix coefficients, our aim is to minimize the objective function

$$C(\mathbf{S}, \mathbf{L}, \mathbf{H}) \triangleq D_\beta(\mathbf{M}_\mathcal{V} \cdot \mathbf{V} | \mathbf{M}_\mathcal{V} \cdot \mathbf{SLH}) \quad (4)$$

subject to $\mathbf{S} \in \mathbb{R}_+^{F \times P}$, $\mathbf{L} \in \mathbb{R}_+^{P \times K}$, $\mathbf{H} \in \mathbb{R}_+^{K \times N}$, and $\mathbf{M}_\mathcal{S} \cdot \mathbf{S} = \mathbf{M}_\mathcal{S} \cdot \mathbf{S}^\circ$. The particular case where the dictionary equals the data matrix itself is obtained by setting $(\mathbf{M}_\mathcal{S}, \mathbf{S}^\circ) \triangleq (\mathbf{M}_\mathcal{V}, \mathbf{V})$.

Algorithm 1 CNMF with missing data

Require: $\mathbf{V}, \mathbf{S}^\circ, \mathbf{M}_\mathcal{V}, \mathbf{M}_\mathcal{S}, \beta$

Initialize $\mathbf{S}, \mathbf{L}, \mathbf{H}$ with random nonnegative values

loop

Update \mathbf{S} :

$$\mathbf{S} \leftarrow \mathbf{M}_\mathcal{S} \cdot \mathbf{S}^\circ + \quad (5)$$

$$\bar{\mathbf{M}}_\mathcal{S} \cdot \mathbf{S} \leftarrow \left(\frac{(\mathbf{M}_\mathcal{V} \cdot (\mathbf{SLH})^{(\beta-2)} \cdot \mathbf{V}) (\mathbf{LH})^T}{(\mathbf{M}_\mathcal{V} \cdot (\mathbf{SLH})^{(\beta-1)}) (\mathbf{LH})^T} \right)^{\cdot \gamma(\beta)}$$

Update \mathbf{L} :

$$\mathbf{L} \leftarrow \mathbf{L} \cdot \left(\frac{\mathbf{S}^T (\mathbf{M}_\mathcal{V} \cdot (\mathbf{SLH})^{(\beta-2)} \cdot \mathbf{V}) \mathbf{H}^T}{\mathbf{S}^T (\mathbf{M}_\mathcal{V} \cdot (\mathbf{SLH})^{(\beta-1)}) \mathbf{H}^T} \right)^{\cdot \gamma(\beta)} \quad (6)$$

Update \mathbf{H} :

$$\mathbf{H} \leftarrow \mathbf{H} \cdot \left(\frac{(\mathbf{SL})^T (\mathbf{M}_\mathcal{V} \cdot (\mathbf{SLH})^{(\beta-2)} \cdot \mathbf{V})}{(\mathbf{SL})^T (\mathbf{M}_\mathcal{V} \cdot (\mathbf{SLH})^{(\beta-1)})} \right)^{\cdot \gamma(\beta)} \quad (7)$$

Rescale \mathbf{L} and \mathbf{H} :

$$\forall k \in [K], \mathbf{h}_k \leftarrow \|\mathbf{l}_k\|_1 \times \mathbf{h}_k \quad (8)$$

$$\mathbf{l}_k \leftarrow \frac{\mathbf{l}_k}{\|\mathbf{l}_k\|_1} \quad (9)$$

end loop

return $\mathbf{S}, \mathbf{L}, \mathbf{H}$

3. PROPOSED ALGORITHM

3.1. General description of the algorithm

Algorithm 1 extends the algorithm proposed in [9] for complete CNMF with the β -divergence to the case of missing entries in \mathbf{V} or \mathbf{S} . The algorithm is a block-coordinate descent procedure in which each block is one the three matrix factors. The updates of each block/factor is obtained via majorization-minimization (MM), a classic procedure that consists in iteratively minimizing a tight upper bound (called auxiliary function) of the objective function. In the present setting, the MM procedure leads to multiplicative updates, characteristic of many NMF algorithms, that automatically preserve nonnegativity given positive initialization.

3.2. Detailed updates

We consider the optimization of $C(\mathbf{S}, \mathbf{L}, \mathbf{H})$ with respect to each of its three arguments individually, using MM. Current updates are denoted with a tilde, i.e., $\tilde{\mathbf{S}}, \tilde{\mathbf{L}}$ and $\tilde{\mathbf{H}}$. We start by recalling the definition of an auxiliary function:

Definition 1 (Auxiliary function). The mapping $G(\mathbf{A} | \tilde{\mathbf{A}})$:

$\mathbb{R}_+^{I \times J} \times \mathbb{R}_+^{I \times J} \mapsto \mathbb{R}_+$ is an auxiliary function to $C(\mathbf{A})$ iff

$$\begin{cases} \forall \mathbf{A} \in \mathbb{R}_+^{I \times J}, & C(\mathbf{A}) = G(\mathbf{A}|\mathbf{A}) \\ \forall \mathbf{A}, \tilde{\mathbf{A}} \in \mathbb{R}_+^{I \times J}, & C(\mathbf{A}) \leq G(\mathbf{A}|\tilde{\mathbf{A}}). \end{cases} \quad (10)$$

The iterative minimization of $G(\mathbf{A}|\mathbf{A})$ with respect to \mathbf{A} , with replacement of $\tilde{\mathbf{A}}$ at every iteration, monotonically decreases the objective $C(\mathbf{A})$ until convergence. As explained in detail in [9], the β -divergence may be decomposed into the sum of a convex term $\check{d}_\beta(\cdot|\cdot)$, a concave term $\hat{d}_\beta(\cdot|\cdot)$ and a constant term cst . The first two terms can be majorized using routine Jensen and tangent inequalities, respectively, leading to tractable updates. The auxiliary functions used to derive Algorithm 1 are given by the three following propositions¹ and the monotonicity of the algorithm follows by construction.

Proposition 1 (Auxiliary function for \mathbf{S}).

Let $\tilde{\mathbf{S}} \in \mathbb{R}_+^{F \times P}$ be such that $\forall(f, n) \in [F] \times [N], \tilde{v}_{fn} > 0$, and $\forall(f, p) \in [F] \times [P], \tilde{s}_{fp} > 0$, where $\tilde{\mathbf{V}} \triangleq \tilde{\mathbf{S}}\mathbf{L}\mathbf{H}$. Then the function

$$G_S(\mathbf{S}|\tilde{\mathbf{S}}) \triangleq \sum_{fn} [\mathbf{M}_V]_{fn} \left[\check{G}_{fn}(\mathbf{S}|\tilde{\mathbf{S}}) + \hat{G}_{fn}(\mathbf{S}|\tilde{\mathbf{S}}) \right] + cst$$

$$\text{where } \check{G}_{fn}(\mathbf{S}|\tilde{\mathbf{S}}) \triangleq \sum_p \frac{[\mathbf{LH}]_{pn} \tilde{s}_{fp}}{\tilde{v}_{fn}} \check{d}_\beta \left(v_{fn} | \tilde{v}_{fn} \frac{s_{fp}}{\tilde{s}_{fp}} \right)$$

$$\begin{aligned} \hat{G}_{fn}(\mathbf{S}|\tilde{\mathbf{S}}) &\triangleq \hat{d}_\beta(v_{fn}|\tilde{v}_{fn}) \\ &+ \hat{d}'_\beta(v_{fn}|\tilde{v}_{fn}) \sum_{fp} [\mathbf{LH}]_{pn} (s_{fp} - \tilde{s}_{fp}) \end{aligned}$$

is an auxiliary function to $C(\mathbf{S}, \mathbf{L}, \mathbf{H})$ with respect to \mathbf{S} and its minimum is given by equation (5). The auxiliary function decouples with respect to the individual coefficients of \mathbf{S} and as such, the constraint $\mathbf{M}_S \cdot \mathbf{S} = \mathbf{M}_S \cdot \mathbf{S}^\circ$ is directly imposed by only updating the coefficients of \mathbf{S} with indices in $\tilde{\mathcal{S}}$.

Proposition 2 (Auxiliary function for \mathbf{L}).

Let $\tilde{\mathbf{L}} \in \mathbb{R}_+^{P \times K}$ be such that $\forall(f, n) \in [F] \times [N], \tilde{v}_{fn} > 0$ and $\forall(p, k) \in [P] \times [K], \tilde{l}_{pk} > 0$, where $\tilde{\mathbf{V}} \triangleq \mathbf{S}\tilde{\mathbf{L}}\mathbf{H}$. Then the function

$$G_L(\mathbf{L}|\tilde{\mathbf{L}}) \triangleq \sum_{fn} [\mathbf{M}_V]_{fn} \left[\check{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) + \hat{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) \right] + cst$$

$$\text{where } \check{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) \triangleq \sum_{pk} \frac{s_{fp} \tilde{l}_{pk} h_{kn}}{\tilde{v}_{fn}} \check{d}_\beta \left(v_{fn} | \tilde{v}_{fn} \frac{l_{pk}}{\tilde{l}_{pk}} \right)$$

$$\begin{aligned} \hat{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) &\triangleq \hat{d}_\beta(v_{fn}|\tilde{v}_{fn}) \\ &+ \hat{d}'_\beta(v_{fn}|\tilde{v}_{fn}) \sum_{pk} s_{fp} h_{kn} (l_{pk} - \tilde{l}_{pk}) \end{aligned}$$

¹The proof of these propositions are available in the extended version at <https://hal-amu.archives-ouvertes.fr/hal-01346492>.

is an auxiliary function to $C(\mathbf{S}, \mathbf{L}, \mathbf{H})$ with respect to \mathbf{L} and its minimum subject to $\mathbf{M}_S \cdot \mathbf{S} = \mathbf{M}_S \cdot \mathbf{S}^\circ$ for $\mathbf{M}_S \in \{0, 1\}^{F \times P}$ and $\mathbf{S}^\circ \in \mathbb{R}_+^{F \times P}$ is given by equation (6).

Proposition 3 (Auxiliary function for \mathbf{H}).

Let us define $\mathbf{W} \triangleq \mathbf{S}\mathbf{L}$ and let $\tilde{\mathbf{H}} \in \mathbb{R}_+^{K \times N}$ be such that $\forall(f, n) \in [F] \times [N], \tilde{v}_{fn} > 0$ and $\forall(k, n) \in [K] \times [N], \tilde{h}_{kn} > 0$, where $\tilde{\mathbf{V}} \triangleq \mathbf{W}\tilde{\mathbf{H}}$. Then the function

$$G_H(\mathbf{H}|\tilde{\mathbf{H}}) \triangleq \sum_{fn} [\mathbf{M}_V]_{fn} \left[\check{G}_{fn}(\mathbf{H}|\tilde{\mathbf{H}}) + \hat{G}_{fn}(\mathbf{H}|\tilde{\mathbf{H}}) \right] + cst$$

$$\text{where } \check{G}_{fn}(\mathbf{H}|\tilde{\mathbf{H}}) \triangleq \sum_k \frac{w_{fk} \tilde{h}_{kn}}{\tilde{v}_{fn}} \check{d}_\beta \left(v_{fn} | \tilde{v}_{fn} \frac{h_{kn}}{\tilde{h}_{kn}} \right)$$

$$\begin{aligned} \hat{G}_{fn}(\mathbf{H}|\tilde{\mathbf{H}}) &\triangleq \hat{d}_\beta(v_{fn}|\tilde{v}_{fn}) \\ &+ \hat{d}'_\beta(v_{fn}|\tilde{v}_{fn}) \sum_k w_{fk} (h_{kn} - \tilde{h}_{kn}) \end{aligned}$$

is an auxiliary function to $C(\mathbf{S}, \mathbf{L}, \mathbf{H})$ with respect to \mathbf{H} and its minimum is given by equation (7).

4. EXPERIMENT ON SYNTHETIC DATA

4.1. Experimental setting

The objective of this experiment is to analyze the performance of CNMF for reconstructing missing data, by comparing it with the regular NMF. We consider a data matrix \mathbf{V}^* of rank K^* synthesized under the CNMF model $\mathbf{V}^* = \mathbf{S}^* \mathbf{L}^* \mathbf{H}^*$, where the matrix of atoms \mathbf{S}^* and the ground truth factors \mathbf{L}^* and \mathbf{H}^* are generated as the absolute values of Gaussian noise. It is worth noting that \mathbf{V}^* is also consistent with a NMF model by defining $\mathbf{W}^* = \mathbf{S}^* \mathbf{L}^*$. A perturbed data matrix \mathbf{V} is obtained by considering a multiplicative noise, obtained using a Gamma distribution with mean 1 and variance $\frac{1}{\alpha}$. Hence the parameter α controls the importance of the perturbation. The mask \mathbf{M}_V of known elements in \mathbf{V} is derived by considering missing coefficients randomly and uniformly distributed over the matrix, such that the ratio of missing values is equal to σ_V . Generation of data is repeated 3 times, as well as the generation of the masks. Results are averaged over these repetitions.

From a matrix \mathbf{V} with missing entries, NMF and CNMF with missing values are applied using K components. Only the case where $\beta = 2$ has been considered in this experiment. In both algorithms, the convergence is reached when the relative difference of the cost function between two iterations is below 10^{-5} . 3 repetitions are performed using different random initialization, and the best instance (i.e., the instance which minimizes the cost function) is retained. The reconstructed data matrix is obtained as $\tilde{\mathbf{V}} = \mathbf{S}\mathbf{L}\mathbf{H}$.

The reconstruction error is obtained by computing the β -divergence between the noiseless data matrix \mathbf{V}^* , and the reconstructed matrix $\tilde{\mathbf{V}}$; the error is computed on and averaged

along the missing coefficients only, as

$$e_{\text{test}} = \frac{1}{\sum_{ij} [\bar{\mathbf{M}}_{\mathcal{V}}]_{ij}} d_{\beta}(\bar{\mathbf{M}}_{\mathcal{V}} \cdot \mathbf{V}^*, \bar{\mathbf{M}}_{\mathcal{V}} \tilde{\mathbf{V}}) \quad (11)$$

where $\bar{\mathbf{M}}_{\mathcal{V}}$ is the mask of unknown elements in \mathcal{V} . In the case of CNMF, we consider two choices for \mathbf{S} : the data matrix \mathbf{V} with missing values, and the ground truth matrix of atoms \mathbf{S}^* , considered here without missing values.

The following parameters are fixed: $F = 100$, $N = 250$, $P = 50$, $K^* = 10$, $\beta = 2$. We particularly investigate the influence of four factors: the number of estimated components $K \in [2, 14]$; the ratio of missing data $\sigma_{\mathcal{V}} \in [0.1, 0.9]$ in \mathbf{V} , i.e., of zeros in $\bar{\mathbf{M}}_{\mathcal{V}}$; the choice of the matrix $\mathbf{S} \in \{\mathbf{S}^*, \mathbf{V}\}$ for the CNMF; the noise level $\alpha \in [10, 5000]$ in \mathbf{V} (α is inversely proportional to the variance of the noise).

4.2. Results

We first focus on the influence of the number of estimated components k for the case where the true dictionary is fully known. Figure 1 displays the test error with respect to the number of estimated components K , for two levels of noise. Performance of reconstruction obtained by NMF and CNMF with $\mathbf{S} = \mathbf{S}^*$ are plotted, for different values of ratio of missing values in \mathcal{V} .

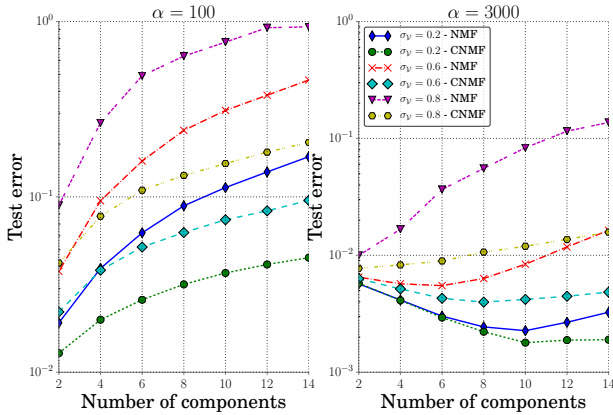


Fig. 1. Test error vs. number of components. Two levels of noise are displayed: high noise level ($\alpha = 100$, left) and low noise level ($\alpha = 3000$, right).

These results show that the noise level has a high influence on the best number of estimated components. As expected, a high noise requires a strong regularization, obtained here by selecting a low value of K . On the contrary, when the noise is low, the best choice of K is closer to the true value $K^* = 10$. Similarly, when the number of missing data is low ($\sigma_{\mathcal{V}} = 0.2$), one should set K to $K^* = 10$, either for CNMF or for NMF. When it gets higher, the optimal K gets lower in order to limit the effect of overfitting. In this case, the best number of components drops down to $K = 2$, either

for CNMF or NMF, the former still performing better than the latter. These first results outline the difference between NMF and CNMF, emphasized in the next figures.

This comparison is augmented by considering the case where $\mathbf{S} = \mathbf{V}$, with respect to the ratio $\sigma_{\mathcal{V}}$ of missing values in \mathbf{V} . Figure 2 displays the performance of NMF and CNMF as a function of $\sigma_{\mathcal{V}}$, for two noise levels. CNMF with the true atoms ($\mathbf{S} = \mathbf{S}^*$) gives the best results on the full range missing data ratio. When there are very few missing data and a low noise, NMF performs almost as well as CNMF. However, the NMF error increases much faster than the CNMF error as the number of missing data grows, or as the noise in data becomes important. This higher sensitivity of NMF to missing data may be explained by overfitting since the number of free parameters in NMF is higher than in CNMF. In the case of CNMF with $\mathbf{S} = \mathbf{V}$, the model cannot fit the data as well as CNMF with $\mathbf{S} = \mathbf{S}^*$ or as NMF. Consequently, the resulting modeling error is observed when there is few missing data, and when comparing $\mathbf{S} = \mathbf{V}$ and $\mathbf{S} = \mathbf{S}^*$ on all values. However, it performs better than NMF at high values of $\sigma_{\mathcal{V}}$ since the constraint $\mathbf{S} = \mathbf{V}$ can be seen as a regularization.

We finally investigate the robustness of methods by looking at the influence of the multiplicative noise, controlled by the parameter α , on the performance. Figure 3 shows the test error for the NMF and the CNMF with $\mathbf{S} = \mathbf{S}^*$ with respect to α and for some values of $\sigma_{\mathcal{V}}$. As expected, the test error decreases according to the variance of the noise, inversely proportional to α . If a low value of α disrupts abruptly the performance of reconstruction ($\alpha < 10^3$), the test error is slightly decreasing for $\alpha > 10^3$. When the variance of the noise is close to zero ($\alpha = 5000$), the performance of NMF and CNMF are almost the same. The performance of reconstruction differs when the variance of the noise increases, as well as the ratio of missing values in \mathbf{V} .

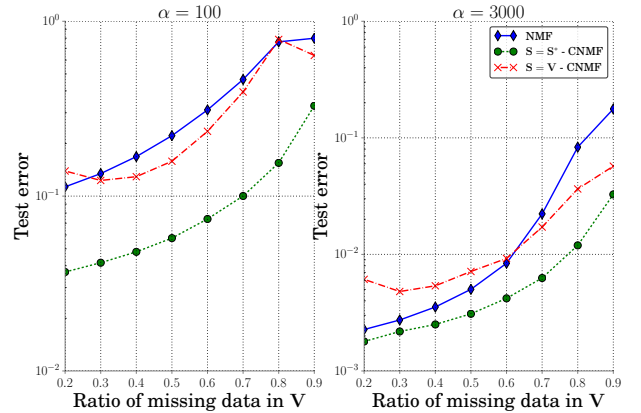


Fig. 2. Test error vs. ratio of missing data in \mathbf{V} , with $K = K^*$. Two levels of noise are displayed: high noise level ($\alpha = 100$, left) and low noise level ($\alpha = 3000$, right).

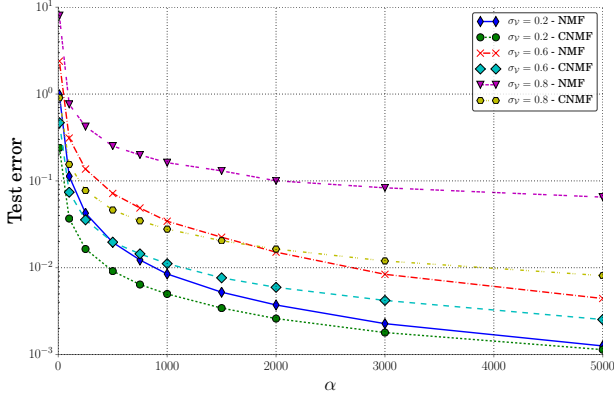


Fig. 3. Test error vs. noise level α . Each curve describes the performance of reconstruction of missing data in \mathbf{V} , with $K = K^*$, according to the method and the ratio of missing data σ_V .

5. APPLICATION TO SPECTROGRAM INPAINTING

In order to illustrate the performance of the proposed algorithm on real data, we consider spectrograms of piano music signals, which are known to be well modeled by NMF methods [12]. Indeed, NMF may provide a note-level decomposition of a full music piece, each NMF component being the estimated spectrum and activation of a single note. This approximation has proved successful and is also limited in terms of modelling error and of lack of supervision to guide the NMF algorithm. In such condition, we have designed an experiment with missing data to compare regular NMF against two CNMF variants: in the first one, we set $\mathbf{S} = \mathbf{V}$; in the second one, \mathbf{S} contains examples of all possible isolated note spectra from another piano.

5.1. Experimental setting

We consider 17 piano pieces from the MAPS dataset [13]. For each recording, the magnitude spectrogram is computed using a 46-ms sine window with 50%-overlap and 4096 frequency bins, the sampling frequency being 44.1kHz. Matrix \mathbf{V} is created from the resulting spectrogram by selecting the $F = 500$ lower-frequency part and the first five seconds, i.e., the $N = 214$ first time frames. Missing data in \mathbf{V} are artificially created by removing coefficients uniformly at random.

Three systems are compared, based on the *test error* defined as the β -divergence computed on the estimation of missing data. In all of them, the number of component K is set to the true number of notes available from the dataset annotation and we set $\beta = 1$. The first system is the regular NMF, randomly initialized. The second system is the proposed CNMF with $(\mathbf{M}_S, \mathbf{S}^o) \triangleq (\mathbf{M}_V, \mathbf{V})$. The third system is the proposed CNMF with $\mathbf{S} = \mathbf{D}$ set as a specific matrix \mathbf{D} of $P = 61$ atoms. Each atom is a single-note spectrum extracted from

the recording of another piano instrument from the MAPS dataset, from C3 to C8².

5.2. Results

Figure 4 displays the test error with respect to the ratio of missing data in \mathbf{V} , averaged over the 17 piano pieces. It clearly shows that the CNMF with the specific dictionary $\mathbf{S} = \mathbf{D}$ is much more robust to missing data than the other two systems. When less than 40% of data are missing, NMF performs slightly better; however, the NMF test error dramatically increases when more data are missing, by a factor 5.10^3 when more than 80% data are missing. This must be due to overfitting since NMF has a large number of free parameters to be estimated from very few observations when data are missing. The performance of the CNMF system with $\mathbf{S} = \mathbf{V}$ probably suffers from modelling error when very few data are missing – since the columns of \mathbf{V} may not be able to combine into K components in a convex way. In the range 50 – 70%, its performance is similar to that of NMF. Beyond this range, it seems to be less prone to overfitting than NMF, probably due to less free parameters or to a regularization effect provided by $\mathbf{S} = \mathbf{V}$.

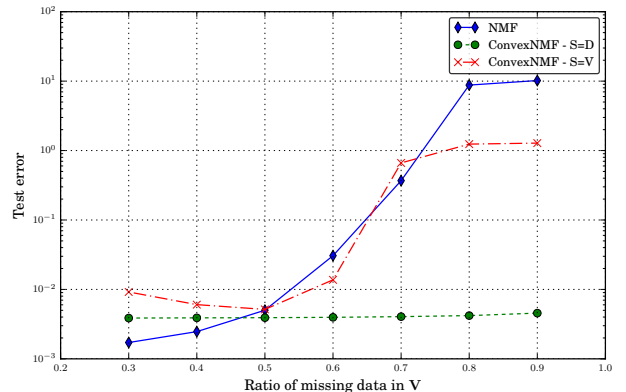


Fig. 4. Test error vs. missing data ratio in audio spectrograms.

We now investigate the influence of the “complexity” of the audio signal on the test error when the ratio of missing data is set to the high value 80%. Figure 5 displays, for each music recording, the test error with respect to the number of different pitches for all notes in the piano piece, which also equals the number of component K used by each system. CNMF with $\mathbf{S} = \mathbf{D}$ performs better than NMF whatever the number of notes and the error increases by a small factor along the represented range. NMF performs about five times worse for “easy” pieces, i.e., pieces composed of notes with about 6 different pitches and it performs about 10^4 times worse when the number of pitches is larger than 25. CNMF

²The code of the experiments is available on the webpage of the MAD project <http://mad.lif.univ-mrs.fr/>.

with $\mathbf{S} = \mathbf{V}$ performs slightly better than NMF. Since the number of components K increases equals the number of note pitches, those results confirm that NMF may highly suffer from overtraining while CNMF may not, being robust to missing data even for large values of K .

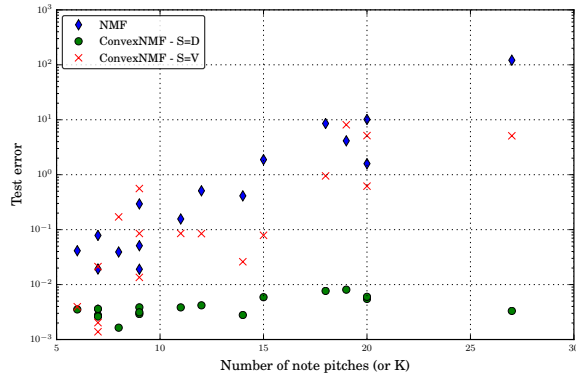


Fig. 5. Test error vs. number of different note pitches for 80% missing data (each dot represents one piece of music).

6. CONCLUSION

In this paper, we have proposed an extension of convex non-negative matrix factorization in the case where missing data, as it has been previously presented for regular NMF. The proposed method can deal with missing values both in the data matrix \mathbf{V} and in the dictionary \mathbf{S} , which is particularly useful in the case $\mathbf{S} = \mathbf{V}$ where the data is autoencoded. In this framework, an algorithm has been provided and analyzed using a Majorization-Minimization (MM) scheme to guarantee the convergence to a local minimum. A large set of experiments on synthetic data showed promising results for this variant of NMF for the task of reconstruction of missing data, and validated the value of this approach. In many situations, CNMF outperforms NMF, especially when the ratio of missing values is high and when the matrix data \mathbf{V} is noisy. This trend has been confirmed on real audio spectrograms of piano music. In particular, we have shown how the use of a generic set of isolated piano notes as atoms could dramatically enhance the robustness to missing data.

This preliminary study indicates that it is worthy of further investigation, beyond the proposed settings where missing values are uniformly distributed over the matrix. Furthermore, the influence of missing values in the dictionary has not been completely assessed, as only the case where $\mathbf{S} = \mathbf{V}$ has been taken into account. On the application side, this approach could give new insight in many problems dealing with estimation of missing data.

7. REFERENCES

- [1] C. Ding, Tao Li, and M.I. Jordan, “Convex and semi-nonnegative matrix factorizations,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 1, pp. 45–55, Jan. 2010.
- [2] D.D. Lee and H.S. Seung, “Learning the parts of objects by non-negative matrix factorization,” *Nature*, vol. 401, no. 6755, pp. 788–791, 1999.
- [3] E. Vincent, N. Bertin, and R. Badeau, “Adaptive harmonic spectral decomposition for multiple pitch estimation,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 18, no. 3, pp. 528–537, Mar. 2010.
- [4] E. Elhamifar and R. Vidal, “Sparse subspace clustering: Algorithm, theory, and applications,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 11, pp. 2765–2781, 2013.
- [5] M. Bertalmio, G. Sapiro, V. Caselles, and C. Ballester, “Image inpainting,” in *Proc of SIGGRAPH*. ACM, 2000, pp. 417–424.
- [6] P. Smaragdis, B. Raj, and M. Shashanka, “Missing data imputation for spectral audio signals,” in *Proc. of MLSP*, Grenoble, France, Sept. 2009.
- [7] J. Le Roux, H. Kameoka, N. Ono, A. De Cheveigne, and S. Sagayama, “Computational auditory induction as a missing-data model-fitting problem with bregman divergence,” *Speech Communication*, vol. 53, no. 5, pp. 658–676, 2011.
- [8] D.L. Sun and R. Mazumder, “Non-negative matrix completion for bandwidth extension: A convex optimization approach,” in *Proc. of MLSP*, Sept. 2013, pp. 1–6.
- [9] C. Févotte and J. Idier, “Algorithms for nonnegative matrix factorization with the β -divergence,” *Neural Comput.*, vol. 23, no. 9, pp. 2421–2456, 2011.
- [10] M. Nakano, H. Kameoka, J. Le Roux, Y. Kitano, N. Ono, and S. Sagayama, “Convergence-guaranteed multiplicative algorithms for nonnegative matrix factorization with β -divergence,” *Proc. of MLSP*, vol. 10, pp. 1, 2010.
- [11] A. Cutler and L. Breiman, “Archetypal analysis,” *Technometrics*, vol. 36, no. 4, pp. 338–347, 1994.
- [12] C. Févotte, N. Bertin, and J-L Durrieu, “Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis,” *Neural Comput.*, vol. 21, no. 3, pp. 793–830, 2009.
- [13] V. Emiya, R. Badeau, and B. David, “Multipitch estimation of piano sounds using a new probabilistic spectral smoothness principle,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 18, no. 6, pp. 1643–1654, 2010.

A. PROOFS

We detail the proofs of Propositions 1 and 2. The proof of Proposition 3 is straightforward using the same methodology.

We first recall preliminary elements from [9]. The β -divergence $d_\beta(x|y)$ can be decomposed as a sum of a convex term, a concave term and a constant term with respect to its second variable y as

$$d_\beta(x|y) = \check{d}_\beta(x|y) + \hat{d}_\beta(x|y) + \bar{d}_\beta(x) \quad (12)$$

This decomposition is not unique. We will use decomposition given in [9, Table 1], for which we have the following derivatives w.r.t. variable y :

$$\check{d}'_\beta(x|y) \triangleq \begin{cases} -xy^{\beta-2} & \text{if } \beta < 1 \\ y^{\beta-2}(y-x) & \text{if } 1 \leq \beta \leq 2 \\ y^{\beta-1} & \text{if } \beta > 2 \end{cases} \quad (13)$$

$$\hat{d}'_\beta(x|y) \triangleq \begin{cases} y^{\beta-1} & \text{if } \beta < 1 \\ 0 & \text{if } 1 \leq \beta \leq 2 \\ -xy^{\beta-2} & \text{if } \beta > 2 \end{cases} \quad (14)$$

A.1. Update of \mathbf{S} (Proof of Proposition 1)

We prove Proposition 1 by first constructing the auxiliary function (Proposition 4 below) and then focussing on its minimum (Proposition 5 below). Due to separability, the update of $\mathbf{S} \in \mathbb{R}_+^{F \times P}$ relies on the update of each of its columns. Hence, we only derive the update for a vector $\mathbf{s} \in \mathbb{R}_+^P$.

Definition 2 (Objective function $C_S(\mathbf{s})$). For $\mathbf{v} \in \mathbb{R}_+^N$, $\mathbf{L} \in \mathbb{R}_+^{P \times K}$, $\mathbf{H} \in \mathbb{R}_+^{K \times N}$, $\mathbf{m} \in \{0, 1\}^N$, $\mathbf{m}^o \in \{0, 1\}^P$, $\mathbf{s} \in \mathbb{R}_+^P$, let us define

$$C_S(\mathbf{s}) \triangleq \sum_n m_n \left[\check{C}_n(\mathbf{s}) + \hat{C}_n(\mathbf{s}) \right] + \bar{C} \quad (15)$$

Where

$$\check{C}_n(\mathbf{s}) \triangleq \check{d}_\beta(v_n | [\mathbf{s}^T \mathbf{LH}]_n), \hat{C}_n(\mathbf{s}) \triangleq \hat{d}_\beta(v_n | [\mathbf{s}^T \mathbf{LH}]_n) \text{ and } \bar{C} \triangleq \bar{d}_\beta(\mathbf{m}, \mathbf{v}) + \lambda \bar{d}_\beta(\mathbf{m}^o, \mathbf{s}^o) \quad (16)$$

Proposition 4 (Auxiliary function $G_S(\mathbf{s}|\tilde{\mathbf{s}})$ for $C_S(\mathbf{s})$). Let $\tilde{\mathbf{s}} \in \mathbb{R}_+^P$ be such that $\forall n, \tilde{v}_n > 0$ and $\forall p, \tilde{s}_p > 0$, where $\tilde{\mathbf{v}} \triangleq [\mathbf{s}^T \mathbf{LH}]^T$. Then the function $G_S(\mathbf{s}|\tilde{\mathbf{s}})$ defined by

$$G_S(\mathbf{s}|\tilde{\mathbf{s}}) \triangleq \sum_n m_n \left[\check{G}_n(\mathbf{s}|\tilde{\mathbf{s}}) + \hat{G}_n(\mathbf{s}|\tilde{\mathbf{s}}) \right] + \bar{C} \quad (17)$$

where

$$\check{G}_n(\mathbf{s}|\tilde{\mathbf{s}}) \triangleq \sum_p \frac{[\mathbf{LH}]_{pn} \tilde{s}_p}{\tilde{v}_n} \check{d}_\beta \left(v_n | \tilde{v}_n \frac{s_p}{\tilde{s}_p} \right) \text{ and } \hat{G}_n(\mathbf{s}|\tilde{\mathbf{s}}) \triangleq \hat{d}_\beta(v_n | \tilde{v}_n) + \hat{d}'_\beta(v_n | \tilde{v}_n) \sum_p [\mathbf{LH}]_{pn} (s_p - \tilde{s}_p). \quad (18)$$

is an auxiliary function for $C_S(\mathbf{s})$.

Proof. We trivially have $G_S(\mathbf{s}|\mathbf{s}) = C_S(\mathbf{s})$. We use the separability in n and in p in order to upper bound each convex term $\check{C}_n(\mathbf{s})$ and each concave term $\hat{C}_n(\mathbf{s})$.

Convex term $\check{C}_n(\mathbf{s})$. Let us prove that $\check{G}_n(\mathbf{s}|\tilde{\mathbf{s}}) \geq \check{C}_n(\mathbf{s})$. Let \mathcal{P} be the set of indices such that $[\mathbf{LH}]_{pn} \neq 0$ and define, for $p \in \mathcal{P}$,

$$\tilde{\lambda}_{pn} \triangleq \frac{[\mathbf{LH}]_{pn} \tilde{s}_p}{\tilde{v}_n} = \frac{[\mathbf{LH}]_{pn} \tilde{s}_p}{\sum_{p' \in \mathcal{P}} [\mathbf{LH}]_{p'n} \tilde{s}_{p'}}. \quad (19)$$

We have $\sum_{p \in \mathcal{P}} \tilde{\lambda}_{pn} = 1$ and

$$\check{G}_n(\mathbf{s}|\tilde{\mathbf{s}}) = \sum_{p \in \mathcal{P}} \tilde{\lambda}_{pn} \check{d}_\beta \left(v_n \left| \frac{[\mathbf{LH}]_{pn} s_p}{\tilde{\lambda}_{pn}} \right. \right) \geq \check{d}_\beta \left(v_n \left| \sum_{p \in \mathcal{P}} \tilde{\lambda}_{pn} \frac{[\mathbf{LH}]_{pn} s_p}{\tilde{\lambda}_{pn}} \right. \right) = \check{d}_\beta \left(v_n \left| \sum_{p=1}^P [\mathbf{LH}]_{pn} s_p \right. \right) = \check{C}_n(\mathbf{s}) \quad (20)$$

Concave term $\hat{C}_n(\mathbf{s})$. We have $\hat{G}_n(\mathbf{s}|\tilde{\mathbf{s}}) \geq \hat{C}_n(\mathbf{s})$ since $\hat{C}_n(\mathbf{s})$ is concave and $\mathbf{s} \mapsto \hat{G}_n(\mathbf{s}|\tilde{\mathbf{s}})$ is a tangent plane to $\hat{C}_n(\mathbf{s})$ in $\tilde{\mathbf{s}}$:

$$\hat{G}_n(\mathbf{s}|\tilde{\mathbf{s}}) = \hat{d}_\beta(v_n|\tilde{v}_n) + \sum_p \hat{d}'_\beta \left(v_n \left| \sum_{p'} [\mathbf{LH}]_{p'n} \tilde{s}_{p'} \right. \right) [\mathbf{LH}]_{pn} (s_p - \tilde{s}_p) = \hat{C}_n(\tilde{\mathbf{s}}) + \langle \nabla \hat{C}_n(\tilde{\mathbf{s}}), \mathbf{s} - \tilde{\mathbf{s}} \rangle \quad (21)$$

□

Proposition 5 (Minimum of $G_S(\mathbf{s}|\tilde{\mathbf{s}})$). *The minimum of $\mathbf{s} \mapsto G_S(\mathbf{s}|\tilde{\mathbf{s}})$ subject to the constraint $\mathbf{m}^o \cdot \mathbf{s} = \mathbf{m}^o \cdot \mathbf{s}^o$ is reached at \mathbf{s}^{MM} with*

$$\forall p, s_p^{MM} \triangleq \begin{cases} \tilde{s}_p \left(\frac{\sum_n m_n \tilde{v}_n^{\beta-2} v_n [\mathbf{LH}]_{pn}}{\sum_n m_n \tilde{v}_n^{\beta-1} [\mathbf{LH}]_{pn}} \right)^{\gamma(\beta)} & \text{if } m_p^o = 0 \\ s_p^o & \text{if } m_p^o = 1. \end{cases} \quad (22)$$

Proof. Since variable s_p is fixed for p such that $m_p^o = 1$, we only consider variables s_p for p such that $m_p^o = 0$. The related penalty term in $G_S(\mathbf{s}|\tilde{\mathbf{s}})$ vanishes when $m_p^o = 0$. Using (13) and (14), the minimum is obtained by cancelling the gradient

$$\nabla_{s_p} G_S(\mathbf{s}|\tilde{\mathbf{s}}) = \sum_n m_n \mathbf{LH}_{pn} \left[\check{d}'_\beta \left(v_n \left| \tilde{v}_n \frac{s_p}{\tilde{s}_p} \right. \right) + \hat{d}'_\beta(v_n|\tilde{v}_n) \right] \quad (23)$$

and by considering that the Hessian matrix is diagonal with nonnegative entries since $\check{d}_\beta(x|y)$ is convex:

$$\nabla_{s_p}^2 G_S(\mathbf{s}|\tilde{\mathbf{s}}) = \sum_n m_n \tilde{v}_n \frac{\mathbf{LH}_{pn}}{\tilde{s}_p} \check{d}''_\beta \left(v_n \left| \tilde{v}_n \frac{s_p}{\tilde{s}_p} \right. \right) \geq 0. \quad (24)$$

□

A.2. Update of \mathbf{L} (Proof of Proposition 2)

We prove Proposition 2 by first constructing the auxiliary function (Proposition 6 below) and then focussing on its minimum (Proposition 7 below). As opposed to the update of \mathbf{S} , no separability is considered here.

Definition 3 (Objective function $C_L(\mathbf{L})$). For $\mathbf{v} \in \mathbb{R}_+^F$, $\mathbf{S} \in \mathbb{R}_+^{F \times P}$, $\mathbf{H} \in \mathbb{R}_+^{K \times N}$, $\mathbf{M} \in \{0, 1\}^{F \times N}$, $\mathbf{L} \in \mathbb{R}_+^{P \times K}$, let us define

$$C_L(\mathbf{L}) \triangleq \sum_{fn} m_{fn} \left[\check{C}_{fn}(\mathbf{L}) + \hat{C}_{fn}(\mathbf{L}) \right] + \bar{C} \quad (25)$$

where

$$\check{C}_{fn}(\mathbf{L}) \triangleq \check{d}_\beta \left(v_{fn} \left| [\mathbf{SLH}]_{fn} \right. \right), \hat{C}_{fn}(\mathbf{L}) \triangleq \hat{d}_\beta \left(v_{fn} \left| [\mathbf{SLH}]_{fn} \right. \right) \text{ and } \bar{C} \triangleq \bar{d}_\beta(\mathbf{M} \cdot \mathbf{V}). \quad (26)$$

Proposition 6 (Auxiliary function $G_L(\mathbf{L}|\tilde{\mathbf{L}})$ for $C_L(\mathbf{L})$). *Let $\tilde{\mathbf{L}} \in \mathbb{R}_+^{P \times K}$ be such that $\forall f, n, \tilde{\mathbf{V}}_{fn} > 0$ and $\forall p, k, \tilde{\mathbf{L}}_{pk} > 0$, where $\tilde{\mathbf{V}} \triangleq \tilde{\mathbf{S}}\tilde{\mathbf{L}}\mathbf{H}$. Then the function $G_L(\mathbf{L}|\tilde{\mathbf{L}})$ defined by*

$$G_L(\mathbf{L}|\tilde{\mathbf{L}}) \triangleq \sum_{fn} m_{fn} \left[\check{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) + \hat{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) \right] + \bar{C} \quad (27)$$

where

$$\check{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) \triangleq \sum_{pk} \frac{s_{fp} \tilde{l}_{pk} h_{kn}}{\tilde{v}_{fn}} \check{d}_{\beta} \left(v_{fn} | \tilde{v}_{fn} \frac{l_{pk}}{\tilde{l}_{pk}} \right) \quad (28)$$

$$\hat{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) \triangleq \left[\hat{d}_{\beta}(v_{fn}|\tilde{v}_{fn}) + \hat{d}'_{\beta}(v_{fn}|\tilde{v}_{fn}) \sum_{pk} s_{fp} h_{kn} (l_{pk} - \tilde{l}_{pk}) \right] \quad (29)$$

is an auxiliary function for $C_L(\mathbf{L})$.

Proof. We trivially have $G_L(\mathbf{L}|\mathbf{L}) = C_L(\mathbf{L})$. In order to prove that $G_L(\mathbf{L}|\tilde{\mathbf{L}}) \geq C_L(\mathbf{L})$, we use the separability in f and n and we upper bound the convex terms $\check{C}_{fn}(\mathbf{L})$ and the concave terms $\hat{C}_{fn}(\mathbf{L})$.

Convex term $\check{C}_{fn}(\mathbf{L})$. Let us prove that $\check{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) \geq \check{C}_{fn}(\mathbf{L})$. Let \mathcal{P} be the set of indices such that $s_{fp} \neq 0$, \mathcal{K} be the set of indices such that $h_{kn} \neq 0$ and define, for $(p, k) \in \mathcal{P} \times \mathcal{K}$,

$$\tilde{\lambda}_{fpkn} \triangleq \frac{s_{fp} \tilde{l}_{pk} h_{kn}}{\tilde{v}_{fn}} = \frac{s_{fp} \tilde{l}_{pk} h_{kn}}{\sum_{(p', k') \in \mathcal{P} \times \mathcal{K}} s_{fp'} \tilde{l}_{p'k'} h_{k'n}}. \quad (30)$$

We have $\sum_{(p, k) \in \mathcal{P} \times \mathcal{K}} \tilde{\lambda}_{fpkn} = 1$ and

$$\check{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) = \sum_{(p, k) \in \mathcal{P} \times \mathcal{K}} \tilde{\lambda}_{fpkn} \check{d}_{\beta} \left(v_{fn} | \frac{s_{fp} l_{pk} h_{kn}}{\tilde{\lambda}_{fpkn}} \right) \quad (31)$$

$$\geq \check{d}_{\beta} \left(v_{fn} | \sum_{(p, k) \in \mathcal{P} \times \mathcal{K}} \tilde{\lambda}_{fpkn} \frac{s_{fp} l_{pk} h_{kn}}{\tilde{\lambda}_{fpkn}} \right) = \check{d}_{\beta} \left(v_{fn} | \sum_{p=1}^P \sum_{k=1}^K s_{fp} l_{pk} h_{kn} \right) = \check{C}_{fn}(\mathbf{L}) \quad (32)$$

Concave term $\hat{C}_{fn}(\mathbf{L})$. We have $\hat{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) \geq \hat{C}_{fn}(\mathbf{L})$ since $\hat{C}_{fn}(\mathbf{L})$ is concave and $\mathbf{L} \mapsto \hat{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}})$ is a tangent plane to $\hat{C}_{fn}(\mathbf{L})$ in $\tilde{\mathbf{L}}$:

$$\hat{G}_{fn}(\mathbf{L}|\tilde{\mathbf{L}}) = \hat{d}_{\beta}(v_{fn}|\tilde{v}_{fn}) + \sum_{pk} \hat{d}'_{\beta} \left(v_{fn} | \sum_{p'k'} s_{fp'} \tilde{l}_{p'k'} h_{k'n} \right) s_{fp} h_{kn} (l_{pk} - \tilde{l}_{pk}) \quad (33)$$

$$= \hat{C}_{fn}(\tilde{\mathbf{L}}) + \langle \nabla \hat{C}_{fn}(\tilde{\mathbf{L}}), \mathbf{L} - \tilde{\mathbf{L}} \rangle \quad (34)$$

□

Proposition 7 (Minimum of $G_L(\mathbf{L}|\tilde{\mathbf{L}})$). *The minimum of $\mathbf{L} \mapsto G_L(\mathbf{L}|\tilde{\mathbf{L}})$ is reached at \mathbf{L}^{MM} with*

$$\forall p, k, l_{p,k}^{MM} \triangleq \begin{cases} l_{p,k} \left(\frac{\sum_{fn} s_{fp} m_{fn} \tilde{v}_{fn}^{\beta-2} v_{fn} h_{kn}}{\sum_{fn} s_{fp} m_{fn} \tilde{v}_{fn}^{\beta-1} h_{kn}} \right)^{\gamma(\beta)} & \text{if } \mathbf{M} \neq \mathbf{0} \\ l_{p,k} & \text{otherwise.} \end{cases} \quad (35)$$

Proof. Using (13) and (14), the minimum is obtained by cancelling the gradient

$$\nabla_{l_{pk}} G_L(\mathbf{L}|\tilde{\mathbf{L}}) = \sum_{fn} m_{fn} s_{fp} h_{kn} \left[\check{d}'_{\beta} \left(v_{fn} | \tilde{v}_{fn} \frac{l_{pk}}{\tilde{l}_{pk}} \right) + \hat{d}'_{\beta}(v_{fn}|\tilde{v}_{fn}) \right] \quad (36)$$

and by considering that the Hessian matrix is diagonal with nonnegative entries since $\check{d}_{\beta}(x|y)$ is convex:

$$\nabla_{l_{pk}}^2 G_L(\mathbf{L}|\tilde{\mathbf{L}}) = \sum_{fn} m_{fn} \tilde{v}_{fn} \frac{s_{fp} h_{kn}}{\tilde{l}_{pk}} \check{d}''_{\beta} \left(v_{fn} | \tilde{v}_{fn} \frac{l_{pk}}{\tilde{l}_{pk}} \right) \geq 0. \quad (37)$$

□