

An Oracle Inequality for Quasi-Bayesian Non-Negative Matrix Factorization

Pierre Alquier* & Benjamin Guedj[†]

June 26, 2018

Abstract

This is the corrected version of a paper that was published as *P. Alquier, B. Guedj, An Oracle Inequality for Quasi-Bayesian Non-negative Matrix Factorization, Mathematical Methods of Statistics, 2017, vol. 26, no. 1, pp. 55-67.*

Since then, a mistake was found in the proofs. We fixed the mistake at the price of a slightly different logarithmic term in the bound. ¹

The aim of this paper is to provide some theoretical understanding of quasi-Bayesian aggregation methods non-negative matrix factorization. We derive an oracle inequality for an aggregated estimator. This result holds for a very general class of prior distributions and shows how the prior affects the rate of convergence.

1 Introduction

Non-negative matrix factorization (NMF) is a set of algorithms in high-dimensional data analysis which aims at factorizing a large matrix M with non-negative entries. If M is an $m_1 \times m_2$ matrix, NMF consists in decomposing it as a product of two matrices of smaller dimensions: $M \simeq UV^T$ where U is $m_1 \times K$, V is $m_2 \times K$, $K \ll m_1 \wedge m_2$ and both U and V have non-negative entries. Interpreting the columns $M_{\cdot,j}$ of M as (non-negative) signals, NMF

*CREST, ENSAE, Université Paris Saclay, pierre.alquier@ensae.fr. This author gratefully acknowledges financial support from the research programme *New Challenges for New Data* from LCL and GENES, hosted by the *Fondation du Risque*, from Labex ECODEC (ANR - 11-LABEX-0047) and from Labex CEMPI (ANR-11-LABX-0007-01).

[†]Modal project-team, Inria Lille - Nord Europe research center, benjamin.guedj@inria.fr.

¹We thank Arnak Dalalyan (ENSAE) who found the mistake

amounts to decompose (exactly, or approximately) each signal as a combination of the “elementary” signals $U_{\cdot,1}, \dots, U_{\cdot,K}$:

$$M_{\cdot,j} \simeq \sum_{\ell=1}^K V_{j,\ell} U_{\cdot,\ell}. \quad (1)$$

Since the seminal paper from [Lee and Seung \(1999\)](#), NMF was successfully applied to various fields such as image processing and face classification ([Guillamet and Vitria, 2002](#)), separation of sources in audio and video processing ([Ozerov and Févotte, 2010](#)), collaborative filtering and recommender systems on the Web ([Koren et al., 2009](#)), document clustering ([Xu et al., 2003](#); [Shahnaz et al., 2006](#)), medical image processing ([Allen et al., 2014](#)) or topics extraction in texts ([Paisley et al., 2015](#)). In all these applications, it has been pointed out that NMF provides a decomposition which is usually interpretable. [Donoho and Stodden \(2003\)](#) have given a theoretical foundation to this interpretability by exhibiting conditions under which the decomposition $M \simeq UV^T$ is unique. However, let us stress that even when this is not the case, the results provided by NMF are still sensibly interpreted by practitioners.

Since a prior knowledge on the shape and/or magnitude of the signal is available in many settings, Bayesian tools have extensively been used for (general) matrix factorization ([Corander and Villani, 2004](#); [Lim and Teh, 2007](#); [Salakhutdinov and Mnih, 2008](#); [Lawrence and Urtasun, 2009](#); [Zhou et al., 2010](#)) and have been adapted for the Bayesian NMF problem ([Moussaoui et al., 2006](#); [Cemgil, 2009](#); [Févotte et al., 2009](#); [Schmidt et al., 2009](#); [Tan and Févotte, 2009](#); [Zhong and Girolami, 2009](#), among others).

The aim of this paper is to provide some theoretical analysis on the performance of an aggregation method for NMF inspired by the aforementioned Bayesian works. We propose a quasi-Bayesian estimator for NMF. By quasi-Bayesian, we mean that the construction of the estimator relies on a prior distribution π , however, it does not rely on any parametric assumptions - that is, the likelihood used to build the estimator does not have to be well-specified (it is usually referred to as a quasi-likelihood). The use of quasi-likelihoods in Bayesian estimation is advocated by [Bissiri et al. \(2016\)](#) using decision-theoretic arguments. This methodology is also popular in machine learning, and various authors developed a theoretical framework to analyze it ([Shawe-Taylor and Williamson, 1997](#); [McAllester, 1998](#); [Catoni, 2003, 2004, 2007](#), this is known as the PAC-Bayesian theory). It is also related to recent works on exponentially weighted aggregation in statistics [Dalalyan and Tsybakov \(2008\)](#); [Golubev and Ostrovski \(2014\)](#). Using these theoretical

tools, we derive an oracle inequality for our quasi-Bayesian estimator. The message of this theoretical bound is that our procedure is able to adapt to the unknown rank of M under very general assumptions for the noise.

The paper is organized as follows. Notation for the NMF framework and the definition of our quasi-Bayesian estimator are given in [Section 2](#). The oracle inequality, which is our main contribution, is given in [Section 3](#) and its proof is postponed to [Section 5](#). The computation of our estimator being completely similar to the computation of a (proper) Bayesian estimator, we end the paper with a short discussion and references to state-of-the-art computational methods for Bayesian NMF in [Section 4](#).

2 Notation

For any $p \times q$ matrix A we denote by $A_{i,j}$ its (i,j) -th entry, $A_{i,\cdot}$ its i -th row and $A_{\cdot,j}$ its j -th column. For any $p \times q$ matrix B we define

$$\langle A, B \rangle_F = \text{Tr}(AB^\top) = \sum_{i=1}^p \sum_{j=1}^q A_{i,j} B_{i,j}.$$

We define the Frobenius norm $\|A\|_F$ of A by $\|A\|_F^2 = \langle A, A \rangle_F$. Let $A_{-i,\cdot}$ denote the matrix A where the i -th column is removed. In the same way, for a vector $v \in \mathbb{R}^p$, $v_{-i} \in \mathbb{R}^{p-1}$ is the vector v with its i -th coordinate removed. Finally, let $\text{Diag}(v)$ denote the $p \times p$ diagonal matrix given by $[\text{Diag}(v)]_{i,i} = v_i$.

2.1 Model

The object of interest is an $m_1 \times m_2$ target matrix M possibly polluted with some noise \mathcal{E} . So we actually observe

$$Y = M + \mathcal{E}, \tag{2}$$

and we assume that \mathcal{E} is random with $\mathbb{E}(\mathcal{E}) = 0$. The objective is to approximate M by a matrix UV^\top where U is $m_1 \times K$, V is $m_2 \times K$ for some $K \ll m_1 \wedge m_2$, and where U , V and M all have non-negative entries. Note that, under (2), depending on the distribution of \mathcal{E} , Y might have some negative entries (the non-negativity assumption is on M rather than on Y). Our theoretical analysis only requires the following assumption on \mathcal{E} .

C1. *The entries $\mathcal{E}_{i,j}$ of \mathcal{E} are i.i.d. with $\mathbb{E}(\mathcal{E}_{i,j}) = 0$. With the notation $m(x) = \mathbb{E}[\mathcal{E}_{i,j} \mathbf{1}_{(\mathcal{E}_{i,j} \leq x)}]$ and $F(x) = \mathbb{P}(\mathcal{E}_{i,j} \leq x)$, assume that there exists a non-negative and bounded function g with $\|g\|_\infty \leq 1$ and*

$$\int_u^v m(x) dx = \int_u^v g(x) dF(x). \tag{3}$$

First, note that if (3) is satisfied for a function g with $\|g\|_\infty = \sigma^2 > 1$, we can replace (2) by the normalized model $Y/\sigma = M/\sigma + \varepsilon/\sigma$ for which C1 is satisfied. The introduction of this rather involved condition is due to the technical analysis of our estimator which is based on Theorem 2 in Section 5. Theorem 2 has first been proved by Dalalyan and Tsybakov (2007) using Stein's formula with a Gaussian noise. However, Dalalyan and Tsybakov (2008) have shown that C1 is actually sufficient to prove Theorem 2. For the sake of understanding, note that Equation 3 is fulfilled when the noise is Gaussian ($\varepsilon_{i,j} \sim \mathcal{N}(0, \sigma^2)$ with $\|g\|_\infty = \sigma^2$) or uniform ($\varepsilon_{i,j} \sim \mathcal{U}[-b, b]$ with $\|g\|_\infty = b^2/2$).

2.2 Prior

We are going to define a prior $\pi(U, V)$, where U is $m_1 \times K$ and V is $m_2 \times K$, for a fixed K . Regarding the choice of K , we prove in Section 3 that our quasi-Bayesian estimator is adaptive, in the sense that if K is chosen much larger than the actual rank of M , the prior will put very little mass on many columns of U and V , automatically shrinking them to 0. This seems to advocate for setting a large K prior to the analysis, say $K = m_1 \wedge m_2$. However, keep in mind that the algorithms discussed below have a computational cost growing with K . Anyhow, the following theoretical analysis only requires $2 \leq K \leq m_1 \wedge m_2$.

With respect to the Lebesgue measure on \mathbb{R}_+ , let us fix a density f such that

$$S_f := 1 \vee \int_0^\infty x^2 f(x) dx < +\infty.$$

For any $a, x > 0$, let

$$g_a(x) := \frac{1}{a} f\left(\frac{x}{a}\right).$$

We define the prior on U and V by

$$U_{i,\ell}, V_{i,\ell} \text{ indep. } \sim g_{\gamma_\ell}(\cdot)$$

where

$$\gamma_\ell \text{ indep. } \sim h(\cdot)$$

and h is a density on \mathbb{R}_+ . With the notation $\gamma = (\gamma_1, \dots, \gamma_K)$, define π by

$$\pi(U, V, \gamma) = \prod_{\ell=1}^K \left(\prod_{i=1}^{m_1} g_{\gamma_\ell}(U_{i,\ell}) \right) \left(\prod_{j=1}^{m_2} g_{\gamma_\ell}(V_{j,\ell}) \right) h(\gamma_\ell) \quad (4)$$

and

$$\pi(U, V) = \int_{\mathbb{R}_+^K} \pi(U, V, \gamma) d\gamma.$$

The idea behind this prior is that under h , many γ_ℓ should be small and lead to non-significant columns $U_{\cdot, \ell}$ and $V_{\cdot, \ell}$. In order to do so, we must assume that a non-negligible proportion of the mass of h is located around 0. On the other hand, a non-negligible probability must be assigned to significant values. This is the meaning of the following assumption.

C2. *There exist constants $0 < \alpha < 1$, $\beta \geq 0$ and $\delta > 0$ such that for any $0 < \varepsilon \leq \frac{1}{2\sqrt{2}S_f}$,*

$$\int_0^\varepsilon h(x) dx \geq \alpha \varepsilon^\beta \text{ and } \int_1^2 h(x) dx \geq \delta.$$

Finally, the following assumption on f is required to prove our main result.

C3. *There exist a non-increasing density \tilde{f} w.r.t. the Lebesgue measure on \mathbb{R}_+ and a constant $\mathcal{C}_f > 0$ such that for any $x > 0$,*

$$f(x) \geq \mathcal{C}_f \tilde{f}(x).$$

As shown in [Theorem 1](#), the heavier the tails of $\tilde{f}(x)$, the better the performance of Bayesian NMF.

Note that the general form of (4) encompasses as special cases almost all the priors used in the papers mentioned in the introduction. We end this subsection with classical examples of functions f and h . Regarding f :

1. Exponential prior $f(x) = \exp(-x)$ with $\tilde{f} = f$, $\mathcal{C}_f = 1$ and $S_f = 2$. This is the choice made by [Schmidt et al. \(2009\)](#). A generalization of the exponential prior is the gamma prior used in [Cemgil \(2009\)](#).
2. Truncated Gaussian prior $f(x) \propto \exp(2ax - x^2)$ with $a \in \mathbb{R}$.
3. Heavy-tailed prior $f(x) \propto \frac{1}{(1+x)^\zeta}$ with $\zeta > 1$. This choice is inspired by [Dalalyan and Tsybakov \(2008\)](#) and leads to better theoretical properties.

Regarding h :

1. The uniform distribution on $[0, 2]$ obviously satisfies [C2](#) with $\alpha = 1/2$, $\beta = 1$ and $\delta = 1/2$.
2. The inverse gamma prior $h(x) = \frac{b^a}{\Gamma(a)} \frac{1}{x^{a+1}} \exp(-\frac{b}{x})$ is classical in the literature for computational reasons (see for example [Salakhutdinov and Mnih, 2008](#); [Alquier, 2013](#)), but note that it does not satisfy [C2](#).

3. [Alquier et al. \(2014\)](#) discuss the $\Gamma(a, b)$ choice for $a, b > 0$: both gamma and inverse gamma lead to explicit conditional posteriors for γ (under a restriction on a in the second case), but the gamma distribution led to better numerical performances. When h is the density of the $\Gamma(a, b)$, [C2](#) is satisfied with $\beta = a$ and $\alpha = b^a \exp[-b/(2\sqrt{2}S_f)]/\Gamma(a + 1)$ and $\delta = \int_1^2 b^a x^{a-1} \exp(-bx) dx / \Gamma(a)$.

2.3 Quasi-posterior and estimator

We define the quasi-likelihood as

$$\widehat{L}(U, V) = \exp[-\lambda \|Y - UV^\top\|_F^2]$$

for some fixed parameter $\lambda > 0$. Note that under the assumption that $\varepsilon_{i,j} \sim \mathcal{N}(0, 1/(2\lambda))$, this would be the actual likelihood up to a multiplicative constant. As already pointed out, the use of quasi-likelihoods to define quasi-posteriors is becoming rather popular in Bayesian statistics and machine learning literatures. Here, the Frobenius norm is to be seen as a fitting criterion rather than as a ground truth. Note that other criterion were used in the literature: the Poisson likelihood ([Lee and Seung, 1999](#)), or the Itakura-Saito divergence ([Févotte et al., 2009](#)).

Definition 1. *We define the quasi-posterior as*

$$\begin{aligned} \widehat{\rho}_\lambda(U, V, \gamma) &= \frac{1}{Z} \widehat{L}(U, V) \pi(U, V, \gamma) \\ &= \frac{1}{Z} \exp[-\lambda \|Y - UV^\top\|_F^2] \pi(U, V, \gamma), \end{aligned}$$

where

$$Z := \int \exp[-\lambda \|Y - UV^\top\|_F^2] \pi(U, V, \gamma) d(U, V, \gamma)$$

is a normalization constant. The posterior mean will be denoted by

$$\widehat{M}_\lambda = \int UV^\top \widehat{\rho}_\lambda(U, V, \gamma) d(U, V, \gamma).$$

[Section 3](#) is devoted to the study the theoretical properties of \widehat{M}_λ . A short discussion on the implementation will be provided in [Section 4](#).

3 An oracle inequality

Most likely, the rank of M is unknown in practice. So, as recommended above, we usually choose K much larger than the expected order for the rank, with the hope that many columns of U and V will be shrunk to 0. The following set of matrices is introduced to formalize this idea. For any $r \in \{1, \dots, K\}$, let \mathcal{M}_r be the set of pairs of matrices (U^0, V^0) with non-negative entries such that

$$U^0 = \begin{pmatrix} U_{11}^0 & \cdots & U_{1r}^0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ U_{m_1 1}^0 & \cdots & U_{m_1 r}^0 & 0 & \cdots & 0 \end{pmatrix}, V^0 = \begin{pmatrix} V_{11}^0 & \cdots & V_{1r}^0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ V_{m_2 1}^0 & \cdots & V_{m_2 r}^0 & 0 & \cdots & 0 \end{pmatrix}.$$

We also define $\mathcal{M}_r(L)$ as the set of matrices $(U^0, V^0) \in \mathcal{M}_r$ such that, for any (i, j, ℓ) , $U_{i,\ell}^0, V_{j,\ell}^0 \leq L$.

We are now in a position to state our main theorem, in the form of the following oracle inequality.

Theorem 1. Fix $\lambda = 1/4$. Under assumptions [C1](#), [C2](#) and [C3](#),

$$\mathbb{E}(\|\widehat{M}_\lambda - M\|_F^2) \leq \inf_{1 \leq r \leq K} \inf_{(U^0, V^0) \in \mathcal{M}_r} \left\{ \|U^0 V^{0\top} - M\|_F^2 + \mathcal{R}(r, m_1, m_2, M, U^0, V^0, \beta, \alpha, \delta, K, S_f, \tilde{f}) \right\}.$$

where

$$\begin{aligned} & \mathcal{R}(r, m_1, m_2, M, U^0, V^0, \beta, \alpha, K, S_f, \tilde{f}) \\ &= 8(m_1 \vee m_2)r \log \left(\sqrt{2(m_1 \vee m_2)r} \right) \\ &+ 8(m_1 \vee m_2)r \log \left(\frac{[1 + \|U^0\|_F + \|V^0\|_F + 2\|U^0 V^{0\top} - M\|_F]^2}{\mathcal{C}_f} \right) \\ &+ 4 \sum_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} \log \left(\frac{1}{\tilde{f}(U_{i\ell}^0 + 1)} \right) + 4 \sum_{\substack{1 \leq j \leq m_2 \\ 1 \leq \ell \leq r}} \log \left(\frac{1}{\tilde{f}(V_{j\ell}^0 + 1)} \right) \\ &+ 4\beta K \log \left([1 + \|U^0\|_F + \|V^0\|_F + 2\|U^0 V^{0\top} - M\|_F]^2 \right) \\ &+ 4\beta K \log \left(2S_f \sqrt{2K(m_1 \vee m_2)} \right) \\ &+ r \left[4 \log \left(\frac{1}{\delta} \right) \right] + 4K \log \left(\frac{1}{\alpha} \right) + 4 \log(4) + 1. \end{aligned}$$

We remind the reader that the proof is given in [Section 5](#). The main message of the theorem is that \widehat{M}_λ is as close to M as would be an estimator designed with the actual knowledge of its rank (i.e., \widehat{M}_λ is adaptive to r), up to remainder terms. These terms might be difficult to read. In order to explicit the rate of convergence, we now provide a weaker version, where we assume that $M = U^0 V^{0\top}$ for some $(U^0, V^0) \in \mathcal{N}_r(L) \mathcal{N}_r(L)$; note that the estimator \widehat{M}_λ still doesn't depend on L nor on r .

Corollary 1. Fix $\lambda = 1/4$. Under assumptions [C1](#), [C2](#) and [C3](#), and when $M = U^0 V^{0\top}$ for some $(U^0, V^0) \in \mathcal{N}_r(L) \mathcal{N}_r(L)$,

$$\begin{aligned} \mathbb{E}(\|\widehat{M}_\lambda - M\|_F^2) &\leq 8(m_1 \vee m_2)r \log\left(\frac{2(m_1 \vee m_2)r(1 + 2L\sqrt{r(m_1 \vee m_2)})^2}{\mathcal{C}_f \tilde{f}(L+1)}\right) \\ &\quad + 4\beta K \log\left(2S_f \sqrt{2K(m_1 \vee m_2)}(1 + 2L\sqrt{r(m_1 \vee m_2)})^2\right) \\ &\quad + r \left[4\log\left(\frac{1}{\delta}\right)\right] + 4K \log\left(\frac{1}{\alpha}\right) + 4\log(4) + 1. \end{aligned}$$

First, note that when $L^2 = \mathcal{O}(1)$, the magnitude of the error bound is

$$(m_1 \vee m_2)r \log(m_1 m_2),$$

which is roughly the variance multiplied by the number of parameters to be estimated in any $(U^0, V^0) \in \mathcal{N}_r(L)$. Alternatively, when $M = U^0 V^{0\top}$ only for $(U^0, V^0) \in \mathcal{N}_r(L)$ for a huge L , the log term in

$$8(m_1 \vee m_2)r \log\left(\frac{(L+1)^2 m_1 m_2}{\tilde{f}(L+1)}\right)$$

becomes significant. Indeed, in the case of the truncated Gaussian prior $f(x) \propto \exp(2ax - x^2)$, the previous quantity is in

$$8(m_1 \vee m_2)rL^2 \log(Lm_1 m_2)$$

which is terrible for large L . On the contrary, with the heavy-tailed prior $f(x) \propto (1+x)^{-\zeta}$ (as in [Dalalyan and Tsybakov, 2008](#)), the leading term is

$$8(m_1 \vee m_2)r(\zeta + 2) \log(Lm_1 m_2)$$

which is way more satisfactory. Still, this prior has not received much attention from practitioners.

Remark 1. When (3) in C1 is satisfied with $\|g\|_\infty = \sigma^2 > 1$ we already remarked that it is necessary to use the normalized model $Y/\sigma = M/\sigma + \mathcal{E}/\sigma$ in order to apply Theorem 1. Going back to the original model, we get that, for $\lambda = 1/(4\sigma^2)$,

$$\mathbb{E}(\|\widehat{M}_\lambda - M\|_F^2) \leq \inf_{1 \leq r \leq K} \inf_{(U^0, V^0) \in \mathcal{M}_r} \left\{ \|U^0 V^{0\top} - M\|_F^2 + \sigma^2 \mathcal{R}(r, m_1, m_2, M, U^0, V^0, \beta, \alpha, \delta, K, S_f, \tilde{f}) \right\}.$$

4 Algorithms for Bayesian NMF

As the quasi-Bayesian estimator takes the form of a Bayesian estimator in a special model, we can obviously use tools from computational Bayesian statistics to compute it. The method of choice for computing Bayesian estimators for complex models is Monte-Carlo Markov Chain (MCMC). In the case of Bayesian matrix factorization, the Gibbs sampler was considered in the literature: see for example Salakhutdinov and Mnih (2008), Alquier et al. (2014) for the general case and Moussaoui et al. (2006), Schmidt et al. (2009) and Zhong and Girolami (2009) for NMF. The Gibbs sampler (described in its general form in Bishop, 2006, for example), is given by Algorithm 1.

Algorithm 1 Gibbs sampler.

Input Y, λ .

Initialization $U^{(0)}, V^{(0)}, \gamma^{(0)}$.

For $k = 1, \dots, N$:

For $i = 1, \dots, m_1$: draw $U_{i,\cdot}^{(k)} \sim \widehat{\rho}_\lambda(U_{i,\cdot} | V^{(k-1)}, \gamma^{(k-1)}, Y)$.

For $j = 1, \dots, m_2$: draw $V_{j,\cdot}^{(k)} \sim \widehat{\rho}_\lambda(V_{j,\cdot} | U^{(k)}, \gamma^{(k-1)}, Y)$.

For $\ell = 1, \dots, K$: draw $\gamma_\ell^{(k)} \sim \widehat{\rho}_\lambda(\gamma_\ell | U^{(k)}, V^{(k)}, Y)$.

In the aforementioned papers, there are discussions on the choices of f and h that leads to explicit forms for the conditional posteriors of $U_{i,\cdot}$, $V_{j,\cdot}$ and γ_ℓ , leading to fast algorithms. We refer the reader to these papers for detailed descriptions of the algorithm in this case, and for exhaustive simulations studies.

Optimization methods used for (non-Bayesian) NMF are much faster than the MCMC methods used for Bayesian NMF though: the original multiplicative algorithm [Lee and Seung \(1999, 2001\)](#), projected gradient descent ([Lin, 2007](#); [Guan et al., 2012](#)), second order schemes ([Kim et al., 2008](#)), linear programming ([Bittorf et al., 2012](#)), ADMM (alternative direction method of multipliers [Boyd et al., 2011](#); [Xu et al., 2012](#)), block coordinate descent [Xu and Yin \(2013\)](#) among others.

We believe that an efficient implementation of Bayesian and quasi-Bayesian methods will be based on fast optimisation methods, like Variational Bayes (VB) or Expectation-Propagation (EP) methods ([Jordan et al., 1999](#); [MacKay, 2002](#); [Bishop, 2006](#)). VB was used for Bayesian matrix factorization ([Lim and Teh, 2007](#); [Alquier et al., 2014](#)) and more recently in Bayesian NMF ([Paisley et al., 2015](#)) with promising results. Still, there is no proof that these algorithms provide valid results. To the best of our knowledge, the first attempt to study the convergence of the VB to the target distribution is studied in [Alquier et al. \(2016\)](#) for a family of problems, that do not include NMF. We believe that further investigation in this direction is necessary.

5 Proofs

This section contains the proof to the main theoretical claim of the paper ([Theorem 1](#)).

5.1 A PAC-Bayesian bound from [Dalalyan and Tsybakov \(2008\)](#)

The analysis of quasi-Bayesian estimators with PAC bounds started with [Shawe-Taylor and Williamson \(1997\)](#). McAllester improved on the initial method and introduced the name ‘‘PAC-Bayesian bounds’’ ([McAllester, 1998](#)). Catoni also improved these results to derive sharp oracle inequalities ([Catoni, 2003, 2004, 2007](#)). This methods were used in various complex models of statistical learning ([Guedj and Alquier, 2013](#); [Alquier, 2013](#); [Suzuki, 2015](#); [Mai and Alquier, 2015](#); [Guedj and Robbiano, 2015](#); [Giulini, 2015](#); [Li et al., 2016](#)). [Dalalyan and Tsybakov \(2008\)](#) proved a different PAC-Bayesian bound based on the idea of unbiased risk estimation (see [Leung and Barron, 2006](#)). We first recall its form in the context of matrix factorization.

Theorem 2. *Under [C1](#), as soon as $\lambda \leq 1/4$,*

$$\mathbb{E} \|\widehat{M}_\lambda - M\|_F^2 \leq \inf_{\rho} \left\{ \int \|UV^\top - M\|_F^2 \rho(U, V, \gamma) d(U, V, \gamma) + \frac{\mathcal{K}(\rho, \pi)}{\lambda} \right\},$$

where the infimum is taken over all probability measures ρ absolutely continuous with respect to π , and $\mathcal{K}(\mu, \nu)$ denotes the Kullback-Leibler divergence between two measures μ and ν .

We let the reader check that the proof in Dalalyan and Tsybakov (2008), stated for vectors, is still valid for matrices (also, the result Dalalyan and Tsybakov (2008) is actually stated for any σ^2 , we only use the case $\sigma^2 = 1$).

The end of the proof of Theorem 1 is organized as follows. First, we define in Section 5.2 a parametric family of probability distributions ρ :

$$\{\rho_{r,U^0,V^0,c} : c > 0, 1 \leq r \leq K, (U^0, V^0) \in \mathcal{M}_r\}.$$

We then upper bound the infimum over all ρ by the infimum over this parametric family. So, we have to calculate, or upper bound

$$\int \|UV^\top - M\|_F^2 \rho_{r,U^0,V^0,c}(U, V, \gamma) d(U, V, \gamma)$$

and

$$\mathcal{K}(\rho_{r,U^0,V^0,c}, \pi).$$

This is done in two lemmas in Section 5.3 and Section 5.4 respectively. We finally gather all the pieces together in Section 5.5, and optimize with respect to c .

5.2 A parametric family of factorizations

We define, for any $r \in \{1, \dots, K\}$ and any pair of matrices $(U^0, V^0) \in \mathcal{M}_r$, for any $0 < c \leq 1$, the density

$$\rho_{r,U^0,V^0,c}(U, V, \gamma) = \frac{\mathbf{1}_{\{\|U-U^0\|_F \leq c, \|V-V^0\|_F \leq c\}} \pi(U, V, \gamma)}{\pi(\{\|U-U^0\|_F \leq c, \|V-V^0\|_F \leq c\})}.$$

5.3 Upper bound for the integral part

Lemma 5.1. *We have*

$$\int \|UV^\top - M\|_F^2 \rho_{r,U^0,V^0,c}(U, V, \gamma) d(U, V, \gamma) \leq \|U^0 V^{0\top} - M\|_F^2 + c(1 + \|U^0\|_F + \|V^0\|_F + 2\|U^0 V^{0\top} - M\|_F)^2.$$

Proof. We have

$$\int \|UV^\top - M\|_F^2 \rho_{r,U^0,V^0,c}(U, V, \gamma) d(U, V, \gamma)$$

$$\begin{aligned}
&= \int \left(\|UV^\top - U^0V^{0\top}\|_F^2 + 2\langle UV^\top - U^0V^{0\top}, U^0V^{0\top} - M \rangle_F \right. \\
&\quad \left. + \|U^0V^{0\top} - M\|_F^2 \right) \rho_{r,U^0,V^0,c}(U, V, \gamma) d(U, V, \gamma) \\
&\leq \int \|UV^\top - U^0V^{0\top}\|_F^2 \rho_{r,U^0,V^0,c}(U, V, \gamma) d(U, V, \gamma) \\
&\quad + 2\sqrt{\int \|UV^\top - U^0V^{0\top}\|_F^2 \rho_{r,U^0,V^0,c}(U, V, \gamma) d(U, V, \gamma)} \|U^0V^{0\top} - M\|_F \\
&\quad + \|U^0V^{0\top} - M\|_F^2
\end{aligned}$$

Note that (U, V) belonging to the support of $\rho_{r,U^0,V^0,c}$ implies that

$$\begin{aligned}
\|UV^\top - U^0V^{0\top}\|_F &= \|U(V^\top - V^{0\top}) + (U - U^0)V^{0\top}\|_F \\
&\leq \|U(V^\top - V^{0\top})\|_F + \|(U - U^0)V^{0\top}\|_F \\
&\leq \|U\|_F \|V - V^0\|_F + \|U - U^0\|_F \|V^0\|_F \\
&\leq (\|U^0\|_F + c)c + c\|V^0\|_F \\
&= c(\|U^0\|_F + \|V^0\|_F + c)
\end{aligned}$$

and so

$$\begin{aligned}
&\int \|UV^\top - M\|_F^2 \rho_{r,U^0,V^0,c}(U, V, \gamma) d(U, V, \gamma) \\
&\leq c^2 (\|U^0\|_F + \|V^0\|_F + c)^2 \\
&\quad + 2c (\|U^0\|_F + \|V^0\|_F + c) \|U^0V^{0\top} - M\|_F \\
&\quad + \|U^0V^{0\top} - M\|_F^2 \\
&= c (\|U^0\|_F + \|V^0\|_F + c) [c (\|U^0\|_F + \|V^0\|_F + c) + 2\|U^0V^{0\top} - M\|_F] \\
&\quad + \|U^0V^{0\top} - M\|_F^2 \\
&\leq c (\|U^0\|_F + \|V^0\|_F + 1) [(\|U^0\|_F + \|V^0\|_F + 1) + 2\|U^0V^{0\top} - M\|_F] \\
&\quad + \|U^0V^{0\top} - M\|_F^2 \\
&\leq c [(\|U^0\|_F + \|V^0\|_F + 1) + 2\|U^0V^{0\top} - M\|_F]^2 + \|U^0V^{0\top} - M\|_F^2.
\end{aligned}$$

□

5.4 Upper bound for the Kullback-Leibler divergence

Lemma 5.2. *Under C2 and C3,*

$$\begin{aligned}
\mathcal{K}(\rho_{r,U^0,V^0,c},\pi) &\leq 2(m_1 \vee m_2)r \log\left(\frac{\sqrt{2(m_1 \vee m_2)}r}{c\mathcal{C}_f}\right) \\
&+ \sum_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} \log\left(\frac{1}{\tilde{f}(U_{i\ell}^0+1)}\right) + \sum_{\substack{1 \leq j \leq m_2 \\ 1 \leq \ell \leq r}} \log\left(\frac{1}{\tilde{f}(V_{j\ell}^0+1)}\right) \\
&+ \beta K \log\left(\frac{2S_f \sqrt{2K(m_1 \vee m_2)}}{c}\right) + K \log\left(\frac{1}{\alpha}\right) + r \log\left(\frac{1}{\delta}\right) + \log(4).
\end{aligned}$$

Proof. By definition

$$\begin{aligned}
\mathcal{K}(\rho_{r,U^0,V^0,c},\pi) &= \int \rho_{r,U^0,V^0,c}(U,V,\gamma) \log\left(\frac{\rho_{r,U^0,V^0,c}(U,V,\gamma)}{\pi(U,V,\gamma)}\right) d(U,V,\gamma) \\
&= \log\left(\frac{1}{\int \mathbf{1}_{\{\|U-U^0\|_F \leq c, \|V-V^0\|_F \leq c\}} \pi(U,V,\gamma) d(U,V,\gamma)}\right).
\end{aligned}$$

Then, note that

$$\begin{aligned}
&\int \mathbf{1}_{\{\|U-U^0\|_F \leq c, \|V-V^0\|_F \leq c\}} \pi(U,V,\gamma) d(U,V,\gamma) \\
&= \int \left(\int \mathbf{1}_{\{\|U-U^0\|_F \leq c, \|V-V^0\|_F \leq c\}} \pi(U,V|\gamma) d(U,V) \right) \pi(\gamma) d\gamma \\
&= \int \underbrace{\left(\int \mathbf{1}_{\{\|U-U^0\|_F \leq c\}} \pi(U|\gamma) dU \right)}_{=:I_1(\gamma)} \underbrace{\left(\int \mathbf{1}_{\{\|V-V^0\|_F \leq c\}} \pi(V|\gamma) dV \right)}_{=:I_2(\gamma)} \pi(\gamma) d\gamma.
\end{aligned}$$

So we have to lower bound $I_1(\gamma)$ and $I_2(\gamma)$. We deal only with $I_1(\gamma)$, as the method to lower bound $I_2(\gamma)$ is exactly the same. We define the set $E \subset \mathbb{R}^K$ as

$$E = \left\{ \gamma \in \mathbb{R}^K : \gamma_1, \dots, \gamma_r \in (1, 2] \text{ and } \gamma_{r+1}, \dots, \gamma_K \in \left(0, \frac{c}{2S_f \sqrt{2K m_1 \vee m_2}}\right] \right\}.$$

Then

$$\int I_1(\gamma) I_2(\gamma) \pi(\gamma) d\gamma \geq \int_E I_1(\gamma) I_2(\gamma) \pi(\gamma) d\gamma$$

and we focus on a lower-bound for $I_1(\gamma)$ when $\gamma \in E$.

$$I_1(\gamma) = \pi\left(\sum_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq K}} (U_{i,\ell} - U_{i,\ell}^0)^2 \leq c^2 \middle| \gamma\right)$$

$$\begin{aligned}
&= \pi \left(\sum_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} (U_{i,\ell} - U_{i,\ell}^0)^2 + \sum_{\substack{1 \leq i \leq m_1 \\ r+1 \leq \ell \leq K}} U_{i,\ell}^2 \leq c^2 \middle| \gamma \right) \\
&\geq \pi \left(\sum_{\substack{1 \leq i \leq m_1 \\ r+1 \leq \ell \leq K}} U_{i,\ell}^2 \leq \frac{c^2}{2} \middle| \gamma \right) \pi \left(\sum_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} (U_{i,\ell} - U_{i,\ell}^0)^2 \leq \frac{c^2}{2} \middle| \gamma \right) \\
&\geq \underbrace{\pi \left(\sum_{\substack{1 \leq i \leq m_1 \\ r+1 \leq \ell \leq K}} U_{i,\ell}^2 \leq \frac{c^2}{2} \middle| \gamma \right)}_{=: I_3(\gamma)} \prod_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} \pi \left((U_{i,\ell} - U_{i,\ell}^0)^2 \leq \frac{c^2}{2m_1 r} \middle| \gamma \right).
\end{aligned}$$

Now, using Markov's inequality,

$$\begin{aligned}
1 - I_3(\gamma) &= \pi \left(\sum_{\substack{1 \leq i \leq m_1 \\ r+1 \leq \ell \leq K}} U_{i,\ell}^2 \geq \frac{c^2}{2} \middle| \gamma \right) \\
&\leq 2 \frac{\mathbb{E}_\pi \left(\sum_{\substack{1 \leq i \leq m_1 \\ r+1 \leq \ell \leq K}} U_{i,\ell}^2 \middle| \gamma \right)}{c^2} \\
&= 2 \frac{\sum_{\substack{1 \leq i \leq m_1 \\ r+1 \leq \ell \leq K}} \gamma_j^2 S_f^2}{c^2} \\
&\leq \frac{1}{2},
\end{aligned}$$

and as on E , for $\ell \geq r+1$, $\gamma_j \leq c/(2S_f \sqrt{2Km_1 \vee m_2}) \leq c/(2S_f \sqrt{2Km_1})$. So

$$I_3(\gamma) \geq \frac{1}{2}.$$

Next, we remark that

$$\begin{aligned}
\pi \left((U_{i,\ell} - U_{i,\ell}^0)^2 \leq \frac{c^2}{2m_1 r} \middle| \gamma \right) &\geq \int_{U_{i,\ell}^0}^{U_{i,\ell}^0 + \frac{c}{\sqrt{2m_1 r}}} \frac{1}{\gamma_j} f \left(\frac{u}{\gamma_j} \right) du \\
&\geq \int_{U_{i,\ell}^0}^{U_{i,\ell}^0 + \frac{c}{\sqrt{2m_1 r}}} \frac{\mathcal{C}_f}{\gamma_j} \tilde{f} \left(\frac{u}{\gamma_j} \right) du.
\end{aligned}$$

Remind that $1 \leq \gamma_j \leq 2$ and \tilde{f} is non-increasing so

$$\pi \left((U_{i,\ell} - U_{i,\ell}^0)^2 \leq \frac{c^2}{2m_1 r} \middle| \gamma \right) \geq \frac{2c\mathcal{C}_f}{\sqrt{2m_1 r}} \tilde{f} \left(U_{i,\ell}^0 + \frac{c}{\sqrt{2m_1 r}} \right)$$

$$\geq \frac{2c\mathcal{C}_f}{\sqrt{2m_1r}} \tilde{f}(U_{i,\ell}^0 + 1)$$

as $c \leq 1 \leq \sqrt{m_1r}$. We plug this result and the lower-bound $I_3(\gamma) \geq 1/2$ into the expression of $I_1(\gamma)$ to get

$$I_1(\gamma) \geq \frac{1}{2} \left(\frac{2c\mathcal{C}_f}{\sqrt{2m_1r}} \right)^{m_1r} \left[\prod_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} \tilde{f}(U_{i,\ell}^0 + 1) \right].$$

Proceeding exactly in the same way,

$$I_2(\gamma) \geq \frac{1}{2} \left(\frac{2c\mathcal{C}_f}{\sqrt{2m_2r}} \right)^{m_2r} \left[\prod_{\substack{1 \leq j \leq m_1 \\ 1 \leq \ell \leq r}} \tilde{f}(V_{j,\ell}^0 + 1) \right].$$

So

$$\begin{aligned} & \int_E I_1(\gamma) I_2(\gamma) \pi(\gamma) d\gamma \\ & \geq \int_E \frac{1}{4} \left(\frac{2c\mathcal{C}_f}{\sqrt{2m_1r}} \right)^{m_1r} \left(\frac{2c\mathcal{C}_f}{\sqrt{2m_2r}} \right)^{m_2r} \left[\prod_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} \tilde{f}(U_{i,\ell}^0 + 1) \right] \left[\prod_{\substack{1 \leq j \leq m_1 \\ 1 \leq \ell \leq r}} \tilde{f}(V_{j,\ell}^0 + 1) \right] \pi(\gamma) d\gamma \\ & = \frac{1}{4} \left(\frac{2c\mathcal{C}_f}{\sqrt{2m_1r}} \right)^{m_1r} \left(\frac{2c\mathcal{C}_f}{\sqrt{2m_2r}} \right)^{m_2r} \left[\prod_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} \tilde{f}(U_{i,\ell}^0 + 1) \right] \left[\prod_{\substack{1 \leq j \leq m_1 \\ 1 \leq \ell \leq r}} \tilde{f}(V_{j,\ell}^0 + 1) \right] \int_E \pi(\gamma) d\gamma \end{aligned}$$

and

$$\begin{aligned} \int_E \pi(\gamma) d\gamma & = \left(\int_1^2 h(x) dx \right)^r \left(\int_0^{\frac{c}{2S_f \sqrt{2Km_1 \vee m_2}}} h(x) dx \right)^{K-r} \\ & \geq \delta^r \alpha^{K-r} \left(\frac{c}{2S_f \sqrt{2Km_1 \vee m_2}} \right)^{\beta(K-r)} \\ & \geq \delta^r \alpha^K \left(\frac{c}{2S_f \sqrt{2Km_1 \vee m_2}} \right)^{\beta K}, \end{aligned}$$

using [C2](#). We combine these inequalities, and we use trivia between m_1 , m_2 , $m_1 \vee m_2$ and $m_1 + m_2$ to obtain

$$\mathcal{K}(\rho_r, U^0, V^0, c, \pi) \leq 2(m_1 \vee m_2)r \log \left(\frac{\sqrt{2(m_1 \vee m_2)r}}{c\mathcal{C}_f} \right)$$

$$\begin{aligned}
& + \sum_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} \log \left(\frac{1}{\tilde{f}(U_{i\ell}^0 + 1)} \right) + \sum_{\substack{1 \leq j \leq m_2 \\ 1 \leq \ell \leq r}} \log \left(\frac{1}{\tilde{f}(V_{j\ell}^0 + 1)} \right) \\
& + \beta K \log \left(\frac{2S_f \sqrt{2K(m_1 \vee m_2)}}{c} \right) + K \log \left(\frac{1}{\alpha} \right) + r \log \left(\frac{1}{\delta} \right) + \log(4).
\end{aligned}$$

This ends the proof of the lemma. \square

5.5 Conclusion

We now plug [Lemma 5.1](#) and [Lemma 5.2](#) into [Theorem 2](#). We obtain, under [C1](#), [C2](#) and [C3](#),

$$\begin{aligned}
\mathbb{E}(\|\widehat{M}_\lambda - M\|_F^2) & \leq \inf_{1 \leq r \leq K} \inf_{(U^0, V^0) \in \mathcal{M}_r} \inf_{0 < c \leq \sqrt{Kr}} \left\{ \|U^0 V^{0\top} - M\|_F^2 \right. \\
& \quad + \frac{2(m_1 \vee m_2)r}{\lambda} \log \left(\frac{\sqrt{2(m_1 \vee m_2)r}}{c \mathcal{C}_f} \right) \\
& \quad + \frac{1}{\lambda} \sum_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} \log \left(\frac{1}{\tilde{f}(U_{i\ell}^0 + 1)} \right) + \frac{1}{\lambda} \sum_{\substack{1 \leq j \leq m_2 \\ 1 \leq \ell \leq r}} \log \left(\frac{1}{\tilde{f}(V_{j\ell}^0 + 1)} \right) \\
& \quad + \frac{\beta K}{\lambda} \log \left(\frac{2S_f \sqrt{2K(m_1 \vee m_2)}}{c} \right) + \frac{K}{\lambda} \log \left(\frac{1}{\alpha} \right) + \frac{r}{\lambda} \log \left(\frac{1}{\delta} \right) + \frac{1}{\lambda} \log(4) \\
& \quad \left. + c(1 + \|U^0\|_F + \|V^0\|_F + 2\|U^0 V^{0\top} - M\|_F)^2 \right\}.
\end{aligned}$$

Remind that we fixed $\lambda = \frac{1}{4}$. We finally choose

$$c = \frac{1}{[1 + \|U^0\|_F + \|V^0\|_F + 2\|U^0 V^{0\top} - M\|_F]^2}$$

and so the condition $c \leq 1$ is always satisfied. The inequality becomes

$$\begin{aligned}
\mathbb{E}(\|\widehat{M}_\lambda - M\|_F^2) & \leq \inf_{1 \leq r \leq K} \inf_{(U^0, V^0) \in \mathcal{M}_r} \left\{ \|U^0 V^{0\top} - M\|_F^2 \right. \\
& \quad + 8(m_1 \vee m_2)r \log \left(\sqrt{2(m_1 \vee m_2)r} \right) \\
& \quad \left. + 8(m_1 \vee m_2)r \log \left(\frac{[1 + \|U^0\|_F + \|V^0\|_F + 2\|U^0 V^{0\top} - M\|_F]^2}{\mathcal{C}_f} \right) \right\}.
\end{aligned}$$

$$\begin{aligned}
& + 4 \sum_{\substack{1 \leq i \leq m_1 \\ 1 \leq \ell \leq r}} \log \left(\frac{1}{\tilde{f}(U_{i\ell}^0 + 1)} \right) + 4 \sum_{\substack{1 \leq j \leq m_2 \\ 1 \leq \ell \leq r}} \log \left(\frac{1}{\tilde{f}(V_{j\ell}^0 + 1)} \right) \\
& + 4\beta K \log \left([1 + \|U^0\|_F + \|V^0\|_F + 2\|U^0 V^{0\top} - M\|_F]^2 \right) \\
& \quad + 4\beta K \log \left(2S_f \sqrt{2K(m_1 \vee m_2)} \right) \\
& \quad + r \left[4 \log \left(\frac{1}{\delta} \right) + 4K \log \left(\frac{1}{\alpha} \right) + 4 \log(4) + 1 \right],
\end{aligned}$$

which ends the proof.

Acknowledgements

The authors are grateful to Jialia Mei and Yohann de Castro (Université Paris-Sud) for insightful discussions and for providing many references on NMF, and to the anonymous Referee for helpful comments.

References

- G. I. Allen, L. Grose, and J. Taylor. A generalized least-square matrix decomposition. *Journal of the American Statistical Association*, 109(505): 145–159, 2014. doi: 10.1080/01621459.2013.852978. [2](#)
- P. Alquier. Bayesian methods for low-rank matrix estimation: short survey and theoretical study. In *Algorithmic Learning Theory 2013*, pages 309–323. Springer, 2013. [5](#), [10](#)
- P. Alquier, V. Cottet, N. Chopin, and J. Rousseau. Bayesian matrix completion: prior specification. Preprint arXiv:1406.1440, 2014. [6](#), [9](#), [10](#)
- P. Alquier, J. Ridgway, and N. Chopin. On the properties of variational approximations of Gibbs posteriors. *Journal of Machine Learning Research*, 17(239):1–41, 2016. [10](#)
- C. M. Bishop. *Pattern Recognition and Machine Learning*, chapter 10. Springer, 2006. [9](#), [10](#)
- P. Bissiri, C. Holmes, and S. Walker. A general framework for updating belief distributions. *Journal of the Royal Statistical Society, Series B*, 78(5), 2016. [2](#)

- V. Bittorf, B. Recht, C. Re, and J. Tropp. Factoring nonnegative matrices with linear programs. In *Advances in Neural Information Processing Systems*, pages 1214–1222, 2012. [10](#)
- S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, 2011. [10](#)
- O. Catoni. A PAC-Bayesian approach to adaptive classification. Preprint Laboratoire de Probabilités et Modèles Aléatoires, PMA-840, 2003. [2](#), [10](#)
- O. Catoni. *Statistical Learning Theory and Stochastic Optimization*. Saint-Flour Summer School on Probability Theory 2001 (Jean Picard ed.), Lecture Notes in Mathematics. Springer, 2004. [2](#), [10](#)
- O. Catoni. *PAC-Bayesian supervised classification: the thermodynamics of statistical learning*. Institute of Mathematical Statistics Lecture Notes—Monograph Series, 56. Institute of Mathematical Statistics, Beachwood, OH, 2007. [2](#), [10](#)
- A. T. Cemgil. Bayesian inference for nonnegative matrix factorisation models. *Computational Intelligence and Neuroscience*, 2009, 2009. [2](#), [5](#)
- J. Corander and M. Villani. Bayesian assessment of dimensionality in reduced rank regression. *Statistica Neerlandica*, 58:255–270, 2004. [2](#)
- A. Dalalyan and A. B. Tsybakov. Aggregation by exponential weighting, sharp PAC-Bayesian bounds and sparsity. *Machine Learning*, 72(1-2):39–61, 2008. [2](#), [4](#), [5](#), [8](#), [10](#), [11](#)
- A. S. Dalalyan and A. B. Tsybakov. Aggregation by exponential weighting and sharp oracle inequalities. In N. Bshouty and C. Gentile, editors, *Learning Theory*, volume 4539 of *Lecture Notes in Computer Science*, pages 97–111. Springer Berlin Heidelberg, 2007. [4](#)
- D. Donoho and V. Stodden. When does non-negative matrix factorization give a correct decomposition into parts? In *Advances in neural information processing systems*, 2003. [2](#)
- C. Févotte, N. Bertin, and J.-L. Durrieu. Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis. *Neural computation*, 21(3):793–830, 2009. [2](#), [6](#)

- I. Giulini. PAC-Bayesian bounds for Principal Component Analysis in Hilbert spaces. Preprint arXiv:1511.06263, 2015. [10](#)
- Y. Golubev and D. Ostrovski. Concentration inequalities for the exponential weighting method. *Mathematical Methods of Statistics*, 23(1):20–37, 2014. [2](#)
- N. Guan, D. Tao, Z. Luo, and B. Yuan. NeNMF: an optimal gradient method for nonnegative matrix factorization. *IEEE Transactions on Signal Processing*, 60(6):2882–2898, 2012. [10](#)
- B. Guedj and P. Alquier. PAC-Bayesian Estimation and Prevision in Sparse Additive Models. *Electronic Journal of Statistics*, 7:264–291, 2013. [10](#)
- B. Guedj and S. Robbiano. PAC-Bayesian High Dimensional Bipartite Ranking. Preprint arXiv:1511.02729, 2015. [10](#)
- D. Guillamet and J. Vitria. Classifying faces with nonnegative matrix factorization. In *Proc. 5th Catalan conference for artificial intelligence*, pages 24–31, 2002. [2](#)
- M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. *Machine Learning*, 37:183–233, 1999. [10](#)
- D. Kim, S. Sra, and I. S. Dhillon. Fast projection-based methods for the least squares nonnegative matrix approximation problem. *Statistical Analysis and Data Mining*, 1(1):38–51, 2008. [10](#)
- Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009. [2](#)
- N. D. Lawrence and R. Urtasun. Non-linear matrix factorization with Gaussian processes. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 601–608. ACM, 2009. [2](#)
- D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999. [2](#), [6](#), [10](#)
- D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *Advances in neural information processing systems*, pages 556–562, 2001. [10](#)

- G. Leung and A. R. Barron. Information theory and mixing least-squares regressions. *IEEE Transactions on Information Theory*, 52(8):3396–3410, 2006. [10](#)
- L. Li, B. Guedj, and S. Loustau. PAC-Bayesian online clustering. *arXiv preprint arXiv:1602.00522*, 2016. [10](#)
- Y. J. Lim and Y. W. Teh. Variational Bayesian approach to movie rating prediction. In *Proceedings of KDD Cup and Workshop*, volume 7, pages 15–21, 2007. [2](#), [10](#)
- C.-J. Lin. Projected gradient methods for nonnegative matrix factorization. *Neural computation*, 19(10):2756–2779, 2007. [10](#)
- D. J. C. MacKay. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2002. [10](#)
- T. T. Mai and P. Alquier. A Bayesian approach for matrix completion: optimal rates under general sampling distributions. *Electronic Journal of Statistics*, 9:823–841, 2015. [10](#)
- D. McAllester. Some PAC-Bayesian theorems. In *Proceedings of the Eleventh Annual Conference on Computational Learning Theory*, pages 230–234, New York, 1998. ACM. [2](#), [10](#)
- S. Moussaoui, D. Brie, A. Mohammad-Djafari, and C. Carteret. Separation of non-negative mixture of non-negative sources using a Bayesian approach and MCMC sampling. *IEEE Transactions on Signal Processing*, 54(11):4133–4145, 2006. [2](#), [9](#)
- A. Ozerov and C. Févotte. Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation. *IEEE Transactions on Audio, Speech, and Language Processing*, 18(3):550–563, 2010. [2](#)
- J. Paisley, D. Blei, and M. I. Jordan. *Bayesian nonnegative matrix factorization with stochastic variational inference*, volume Handbook of Mixed Membership Models and Their Applications, chapter 11. Chapman and Hall/CRC, 2015. [2](#), [10](#)
- R. Salakhutdinov and A. Mnih. Bayesian probabilistic matrix factorization using Markov chain Monte Carlo. In *Proceedings of the 25th international conference on Machine learning*, pages 880–887. ACM, 2008. [2](#), [5](#), [9](#)

- M. N. Schmidt, O. Winther, and L. K. Hansen. Bayesian non-negative matrix factorization. In *Independent Component Analysis and Signal Separation*, pages 540–547. Springer, 2009. [2](#), [5](#), [9](#)
- F. Shahnaz, M. W. Berry, V. P. Pauca, and R. J. Plemmons. Document clustering using nonnegative matrix factorization. *Information Processing & Management*, 42(2):373–386, 2006. [2](#)
- J. Shawe-Taylor and R. Williamson. A PAC analysis of a Bayes estimator. In *Proceedings of the Tenth Annual Conference on Computational Learning Theory*, pages 2–9, New York, 1997. ACM. [2](#), [10](#)
- T. Suzuki. Convergence rate of Bayesian tensor estimator and its minimax optimality. In *Proceedings of the 32nd International Conference on Machine Learning (Lille, 2015)*, pages 1273–1282, 2015. [10](#)
- V. Y. Tan and C. Févotte. Automatic relevance determination in nonnegative matrix factorization. In *SPARS’09-Signal Processing with Adaptive Sparse Structured Representations*, 2009. [2](#)
- W. Xu, X. Liu, and Y. Gong. Document clustering based on non-negative matrix factorization. In *Proceedings of the 26th annual international ACM SIGIR conference on research and development in information retrieval*, pages 267–273. ACM, 2003. [2](#)
- Y. Xu and W. Yin. A block coordinate descent method for regularized multi-convex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on imaging sciences*, 6(3):1758–1789, 2013. [10](#)
- Y. Xu, W. Yin, Z. Wen, and Y. Zhang. An alternating direction algorithm for matrix completion with nonnegative factors. *Frontiers of Mathematics in China*, 7(2):365–384, 2012. [10](#)
- M. Zhong and M. Girolami. Reversible jump MCMC for non-negative matrix factorization. In *International Conference on Artificial Intelligence and Statistics*, pages 663–670, 2009. [2](#), [9](#)
- M. Zhou, C. Wang, M. Chen, J. Paisley, D. Dunson, and L. Carin. Nonparametric Bayesian Matrix Completion. In *Proc. IEEE SAM*, 2010. [2](#)