

Accelerated Multiplicative Updates and Hierarchical ALS Algorithms for Nonnegative Matrix Factorization

Nicolas Gillis¹ and François Glineur²

Abstract

Nonnegative matrix factorization (NMF) is a data analysis technique used in a great variety of applications such as text mining, image processing, hyperspectral data analysis, computational biology, and clustering. In this paper, we consider two well-known algorithms designed to solve NMF problems, namely the multiplicative updates of Lee and Seung and the hierarchical alternating least squares of Cichocki et al. We propose a simple way to significantly accelerate these schemes, based on a careful analysis of the computational cost needed at each iteration, while preserving their convergence properties. This acceleration technique can also be applied to other algorithms, which we illustrate on the projected gradient method of Lin. The efficiency of the accelerated algorithms is empirically demonstrated on image and text datasets, and compares favorably with a state-of-the-art alternating nonnegative least squares algorithm.

Keywords: nonnegative matrix factorization, algorithms, multiplicative updates, hierarchical alternating least squares.

1 Introduction

Nonnegative matrix factorization (NMF) consists in approximating a nonnegative matrix M as a low-rank product of two nonnegative matrices W and H , i.e., given a matrix $M \in \mathbb{R}_+^{m \times n}$ and an integer $r < \min\{m, n\}$, find two matrices $W \in \mathbb{R}_+^{m \times r}$ and $H \in \mathbb{R}_+^{r \times n}$ such that $WH \approx M$.

With a nonnegative input data matrix M , nonnegativity constraints on the factors W and H are well-known to lead to low-rank decompositions with better interpretation in many applications such as text mining (Shahnaz et al., 2006), image processing (Lee & Seung, 1999), hyperspectral data analysis (Pauca et al., 2006), computational biology (Devarajan, 2008), and clustering (Ding et al., 2005). Unfortunately, imposing these constraints is also known to render the problem computationally difficult (Vavasis, 2009).

Since an exact low-rank representation of the input matrix does not exist in general, the quality of the approximation is measured by some criterion, typically the sum of the squares of the errors on the entries, which leads to the following minimization problem:

$$\min_{W \in \mathbb{R}^{m \times r}, H \in \mathbb{R}^{r \times n}} \|M - WH\|_F^2 \quad \text{such that} \quad W \geq 0 \text{ and } H \geq 0, \quad (\text{NMF})$$

where $\|A\|_F = (\sum_{i,j} A_{ij}^2)^{\frac{1}{2}}$ denotes the Frobenius norm of matrix A . Most NMF algorithms are iterative, and exploit the fact that (NMF) reduces to an efficiently solvable convex nonnegative least

¹University of Waterloo, Department of Combinatorics and Optimization, Waterloo, Ontario N2L 3G1, Canada. E-mail: ngillis@uwaterloo.ca. This work was carried out when the author was a Research fellow of the Fonds de la Recherche Scientifique (F.R.S.-FNRS) at Université catholique de Louvain.

²Université catholique de Louvain, CORE and ICTEAM Institute, B-1348 Louvain-la-Neuve, Belgium. E-mail: francois.glineur@uclouvain.be. This text presents research results of the Belgian Program on Interuniversity Poles of Attraction initiated by the Belgian State, Prime Minister's Office, Science Policy Programming. The scientific responsibility is assumed by the authors.

squares problem (NNLS) when one of the factors W or H is fixed. Actually, it seems that nearly all algorithms proposed for NMF adhere to the following general framework

- (0) Select initial matrices $(W^{(0)}, H^{(0)})$ (e.g., randomly). Then for $k = 0, 1, 2, \dots$, do
 - (a) Fix $H^{(k)}$ and find $W^{(k+1)} \geq 0$ such that $\|M - W^{(k+1)}H^{(k)}\|_F^2 < \|M - W^{(k)}H^{(k)}\|_F^2$.
 - (b) Fix $W^{(k+1)}$ and find $H^{(k+1)} \geq 0$ such that $\|M - W^{(k+1)}H^{(k+1)}\|_F^2 < \|M - W^{(k+1)}H^{(k)}\|_F^2$.

More precisely, at each iteration, one of the two factors is fixed and the other is updated in such a way that the objective function is reduced, which amounts to a two-block coordinate descent method. Notice that the role of matrices W and H is perfectly symmetric: if one transposes input matrix M , the new matrix M^T has to be approximated by a product $H^T W^T$, so that any formula designed to update the first factor in this product directly translates into an update for the second factor in the original problem. Formally, if the update performed in step (a) is described by $W^{(k+1)} = \text{update}(M, W^{(k)}, H^{(k)})$, an algorithm preserving symmetry will update the factor in step (b) according to $H^{(k+1)} = \text{update}(M^T, H^{(k)T}, W^{(k+1)T})^T$. In this paper, we only consider such symmetrical algorithms, and focus on the update of matrix W .

This update can be carried out in many different ways: the most natural possibility is to compute an optimal solution for the NNLS subproblem, which leads to a class of algorithms called alternating nonnegative least squares (ANLS), see, e.g., H.Kim & Park (2008). However, this computation, which can be performed with active-set-like methods (H.Kim & Park, 2008; J.Kim & Park, 2008), is relatively costly. Therefore, since an optimal solution for the NNLS problem corresponding to one factor is not required before the update of the other factor is performed, several algorithms only compute an approximate solution of the NNLS subproblem, sometimes very roughly, but with a cheaper computational cost, leading to an inexact two-block coordinate descent scheme. We now present two such procedures: the multiplicative updates of Lee and Seung and the hierarchical alternating least squares of Cichocki et al.

In their seminal papers, Lee & Seung (1999, 2001) introduce the multiplicative updates:

$$W^{(k+1)} = \text{MU}(M, W^{(k)}, H^{(k)}) = W^{(k)} \circ \frac{[MH^{(k)T}]}{[W^{(k)}H^{(k)}H^{(k)T}]},$$

where \circ (resp. $\frac{[\cdot]}{[\cdot]}$) denotes the component-wise product (resp. division) of matrices, and prove that each update monotonically decreases the Frobenius norm of the error $\|M - WH\|_F$, i.e., satisfies the description of steps (a) and (b). This technique was actually originally proposed by Daube-Witherspoon & Muehlehner (1986) to solve NNLS problems. The popularity of this algorithm came along with the popularity of NMF and many authors have studied or used this algorithm or variants to compute NMF's, see, e.g., Berry et al. (2007); Cichocki et al. (2009) and the references therein. In particular, the MATLAB[®] Statistics Toolbox implements this method.

However, MU have been observed to converge relatively slowly, especially when dealing with dense matrices M , see Han et al. (2009); Gillis & Glineur (2008) and the references therein, and many other algorithms have been subsequently introduced which perform better in most situations. For example, Cichocki et al. (2007); Cichocki & Phan (2009) and, independently, several other authors (Ho, 2008; Gillis & Glineur, 2008; Li & Zhang, 2009) proposed a technique called hierarchical alternating least squares (HALS)¹, which successively updates each column of W with an optimal and easy to compute closed-form solution. In fact, when fixing all variables but a single column $W_{:p}$ of W , the problem reduces to

$$\min_{W_{:p} \geq 0} \|M - WH\|_F^2 = \|(M - \sum_{l \neq p} W_{:l}H_{l:}) - W_{:p}H_{p:}\|_F^2 = \sum_{i=1}^m \|(M_{i:} - \sum_{l \neq p} W_{il}H_{l:}) - W_{ip}H_{p:}\|_F^2.$$

¹Ho (2008) refers to HALS as rank-one residue iteration (RRI), and Li & Zhang (2009) as FastNMF.

Because each row of W only affects the corresponding row of the product WH , this problem can be further decoupled into m independent quadratic programs in one variable W_{ip} , corresponding to the i^{th} row of M . The optimal solution W_{ip}^* of these subproblems can be easily written in closed-form

$$\begin{aligned} W_{ip}^* &= \max\left(0, \frac{(M_{i:} - \sum_{l \neq p} W_{il} H_{l:}) H_{p:}^T}{H_{p:} H_{p:}^T}\right) \\ &= \max\left(0, \frac{M_{i:} H_{p:}^T - \sum_{l \neq p} W_{il} H_{l:} H_{p:}^T}{H_{p:} H_{p:}^T}\right), \quad 1 \leq i \leq m. \end{aligned}$$

Hence HALS updates successively the columns of W , so that $W^{(k+1)} = \text{HALS}(M, W^{(k)}, H^{(k)})$ can be computed in the following way:

$$W_{:p}^{(k+1)} = \max\left(0, \frac{A_{:p} - \sum_{l=1}^{p-1} W_{:l}^{(k+1)} B_{lp} - \sum_{l=p+1}^r W_{:l}^{(k)} B_{lp}}{B_{pp}}\right),$$

successively for $p = 1, 2, \dots, r$, where $A = MH^{(k)T}$ and $B = H^{(k)}H^{(k)T}$. This amounts to approximately solving each NNLS subproblem in W with a single complete round of an exact block-coordinate descent method with r blocks of m variables corresponding to the columns of W (notice that any other ordering for the update of the columns of W is also possible).

Other approaches based on iterative methods to solve the NNLS subproblems include projected gradient descent (Lin, 2007a) or Newton-like methods (Dhillon et al., 2007; Cichocki et al., 2006); see also Cichocki et al. (2009) and the references therein.

We first analyze in Section 2 the computational cost needed to update the factors W in MU and HALS, then make several simple observations leading in Section 3 to the design of accelerated versions of these algorithms. These improvements can in principle be applied to any two-block coordinate descent NMF algorithm, as demonstrated in Section 3.4 on the projected gradient method of Lin (2007a). We mainly focus on MU, because it is by far the most popular NMF algorithm, and on HALS, because it is very efficient in practice. Section 4 studies convergence of the accelerated variants to stationary points, and shows that they preserve the properties of the original schemes. In Section 5, we experimentally demonstrate a significant acceleration in convergence on several image and text datasets, with a comparison with the state-of-the-art ANLS algorithm of J.Kim & Park (2008).

2 Analysis of the Computational Cost of Factor Updates

In order to make our analysis valid for both dense and sparse input matrices, let us introduce a parameter K denoting the number of nonzero entries in matrix M ($K = mn$ when M is dense). Factors W and H are typically stored as dense matrices throughout the execution of the algorithms. We assume that NMF achieves compression, which is often a requirement in practice. This means that storing W and H must be cheaper than storing M : roughly speaking, the number of entries in W and H must be smaller than the number of nonzero entries in M , i.e., $r(m+n) \leq K$.

Descriptions of Algorithms 1 and 2 below provide separate estimates for the number of floating point operations (flops) in each matrix product computation needed to update factor W in MU and HALS. One can check that the proposed organization of the different matrix computations (and, in particular, the ordering of the matrix products) minimizes the total computational cost (for example, starting the computation of the MU denominator WHH^T with the product WH is clearly worse than with HH^T).

MU and HALS possess almost exactly the same computational cost (the difference being a typically negligible mr flops). It is particularly interesting to observe that

Algorithm 1 MU update for $W^{(k)}$

- 1: $A = MH^{(k)T}$; $\rightarrow 2Kr$ flops
 - 2: $B = H^{(k)}H^{(k)T}$; $\rightarrow 2nr^2$ flops
 - 3: $C = W^{(k)}B$; $\rightarrow 2mr^2$ flops
 - 4: $W^{(k+1)} = W^{(k)} \circ \frac{[A]}{[C]}$; $\rightarrow 2mr$ flops
- % Total: $r(2K + 2nr + 2mr + 2m)$ flops
-

Algorithm 2 HALS update for $W^{(k)}$

- 1: $A = MH^{(k)T}$; $\rightarrow 2Kr$ flops
 - 2: $B = H^{(k)}H^{(k)T}$; $\rightarrow 2nr^2$ flops
 - 3: **for** $i = 1, 2, \dots, r$ **do**
 - 4: $C_{:k} = \sum_{l=1}^{p-1} W_{:l}^{(k+1)} B_{lk} + \sum_{l=p+1}^r W_{:l}^{(k)} B_{lk}$; $\rightarrow 2m(r-1)$ flops
 - 5: $W_{:k} = \max\left(0, \frac{A_{:k} - C_{:k}}{B_{kk}}\right)$; $\rightarrow 3m$ flops
 - 6: **end for**
- % Total: $r(2K + 2nr + 2mr + m)$ flops
-

1. Steps 1. and 2. in both algorithms are identical and do not depend on the matrix $W^{(k)}$;
2. Recalling our assumption $K \geq r(m+n)$, computation of $MH^{(k)T}$ (step 1.) is the most expensive among all steps.

Therefore, this time-consuming step should be performed sparingly, and we should take full advantage of having computed the relatively expensive $MH^{(k)T}$ and $H^{(k)}H^{(k)T}$ matrix products. This can be done by updating $W^{(k)}$ several times before the next update of $H^{(k)}$, i.e., by repeating steps 3. and 4. in MU (resp. steps 3. to 6. in HALS) several times after the computation of matrices $MH^{(k)T}$ and $H^{(k)}H^{(k)T}$. In this fashion, better solutions of the corresponding NNLS subproblems will be obtained at a relatively cheap additional cost.

The original MU and HALS algorithms do not take advantage of this fact, and alternatively update matrices W and H only once per (outer) iteration. An important question for us is now: how many times should we update W per outer iteration?, i.e., how many inner iterations of MU and HALS should we perform? This is the topic of the next section.

3 Stopping Criterion for the Inner Iterations

In this section, we discuss two different strategies for choosing the number of inner iterations: the first uses a fixed number of inner iterations determined by the flop counts, while the second is based on a dynamic stopping criterion that checks the difference between two consecutive iterates. The first approach is shown empirically to work better. We also describe a third hybrid strategy that provides a further small improvement in performance.

3.1 Fixed Number of Inner Iterations

Let us focus on the MU algorithm (a completely similar analysis holds for HALS, as both methods differ only by a negligible number of flops). Based on the flops counts, we estimate how expensive the first inner update of W would be relatively to the next ones (all performed while keeping H fixed),

which is given by the following factor ρ_W (the corresponding value for H will be denoted by ρ_H)

$$\rho_W = \frac{2Kr + 2nr^2 + 2mr^2 + 2mr}{2mr^2 + 2mr} = 1 + \frac{K + nr}{mr + m}. \quad \left(\rho_H = 1 + \frac{K + mr}{nr + n} \right).$$

Values of ρ_W and ρ_H for several datasets are given in Section 5, see Tables 1 and 2.

Notice that for $K \geq r(m + n)$, we have $\rho_W \geq 2\frac{r}{r+1}$ so that the first inner update of W is at least about twice as expensive as the subsequent ones. For a dense matrix, K is equal to mn and we actually have that $\rho_W = 1 + \frac{n(m+r)}{m(r+1)} \geq 1 + \frac{n}{r+1}$, which is typically quite large since n is often much greater than r . This means for example that, in our accelerated scheme, W could be updated about $1 + \rho_W$ times for the same computational cost as two independent updates of W in the original MU.

A simple and natural choice consists in performing inner updates of W and H a fixed number of times, depending on the values of ρ_W and ρ_H . Let us introduce a parameter $\alpha \geq 0$ such that W is updated $(1 + \alpha\rho_W)$ times before the next update of H , and H is updated $(1 + \alpha\rho_H)$ times before the next update of W . Let us also denote the corresponding algorithm MU_α (MU_0 reduces to the original MU). Therefore, performing $(1 + \alpha\rho_W)$ inner updates of W in MU_α has approximately the same computational cost as performing $(1 + \alpha)$ updates of W in MU_0 .

In order to find an appropriate value for parameter α , we have performed some preliminary tests on image and text datasets. First, let us denote $e(t)$ the Frobenius norm of the error $\|M - WH\|_F$ achieved by an algorithm within time t , and define

$$E(t) = \frac{e(t) - e_{\min}}{e(0) - e_{\min}}, \quad (3.1)$$

where $e(0)$ is the error of the initial iterate $(W^{(0)}, H^{(0)})$, and e_{\min} is the smallest error observed among all algorithms across all initializations. Quantity $E(t)$ is therefore a normalized measure of the improvement of the objective function (relative to the initial gap) with respect to time; we have $0 \leq E(t) \leq 1$ for monotonically decreasing algorithms (such as MU and HALS). The advantage of $E(t)$ over $e(t)$ is that one can meaningfully take the average over several runs involving different initializations and datasets, and display the average behavior of a given algorithm.

Figure 1 displays the average of this function $E(t)$ for dense (on the left) and sparse (on the right) matrices using the datasets described in Section 5 for five values of $\alpha = 0, 0.5, 1, 2, 4$. We observe that the original MU algorithm ($\alpha = 0$) converges significantly less rapidly than all the other tested variants (especially in the dense case). The best value for parameter α seems to be 1.

Figure 2 displays the same computational experiments for HALS². HALS with $\alpha = 0.5$ performs the best. For sparse matrices, the improvement is harder to discern (but still present); an explanation for that fact will be given at the end of Section 3.3.

3.2 Dynamic Stopping Criterion for Inner Iterations

In the previous section, a fixed number of inner iterations is performed. One could instead consider switching dynamically from one factor to the other based on an appropriate criterion. For example, it is possible to use the norm of the projected gradient as proposed by Lin (2007a). A simpler and cheaper possibility is to rely solely on the norm of the difference between two iterates. Noting $W^{(k,l)}$ the iterate after l updates of $W^{(k)}$ (while $H^{(k)}$ is being kept fixed), we stop inner iterations as soon as

$$\|W^{(k,l+1)} - W^{(k,l)}\|_F \leq \epsilon \|W^{(k,1)} - W^{(k,0)}\|_F, \quad (3.2)$$

²Because HALS involves a loop over the columns of W and rows of H , we observed that an update of HALS is noticeably slower than an update of MU when using MATLAB[®] (especially for $r \gg 1$), despite the quasi-equivalent theoretical computational cost. Therefore, to obtain fair results, we adjusted ρ_W and ρ_H by measuring directly the ratio between time spent for the first update and the next one, using the *cputime* function of MATLAB[®].

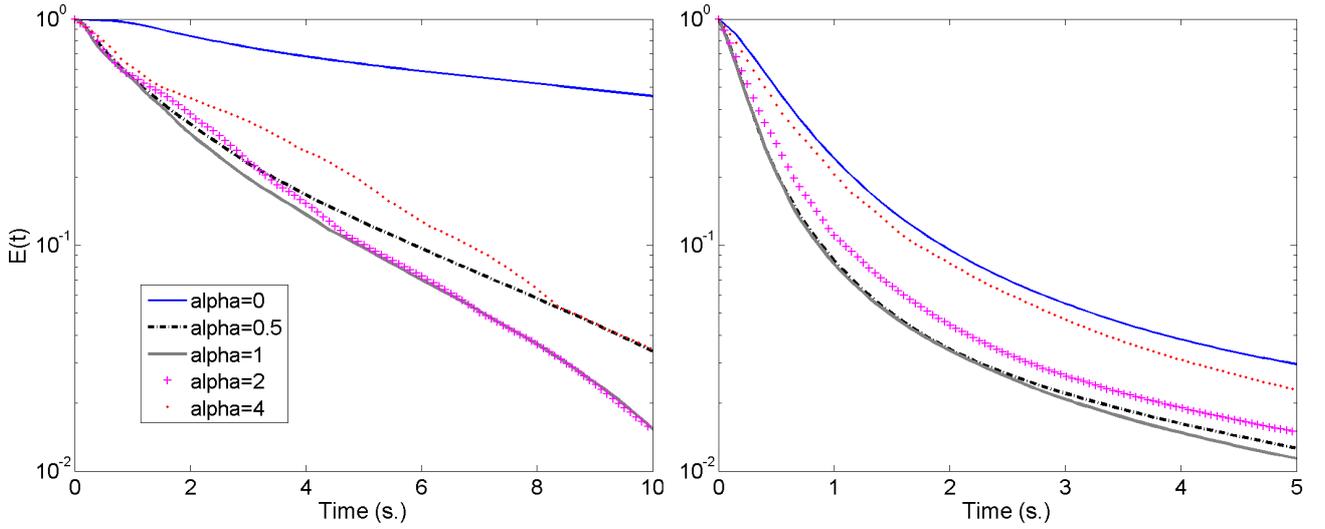


Figure 1: Average of functions $E(t)$ for MU using different values of α : (left) dense matrices, (right) sparse matrices, computed over 4 image datasets and 6 text datasets, using two different values for the rank for each dataset and 10 random initializations, see Section 5.

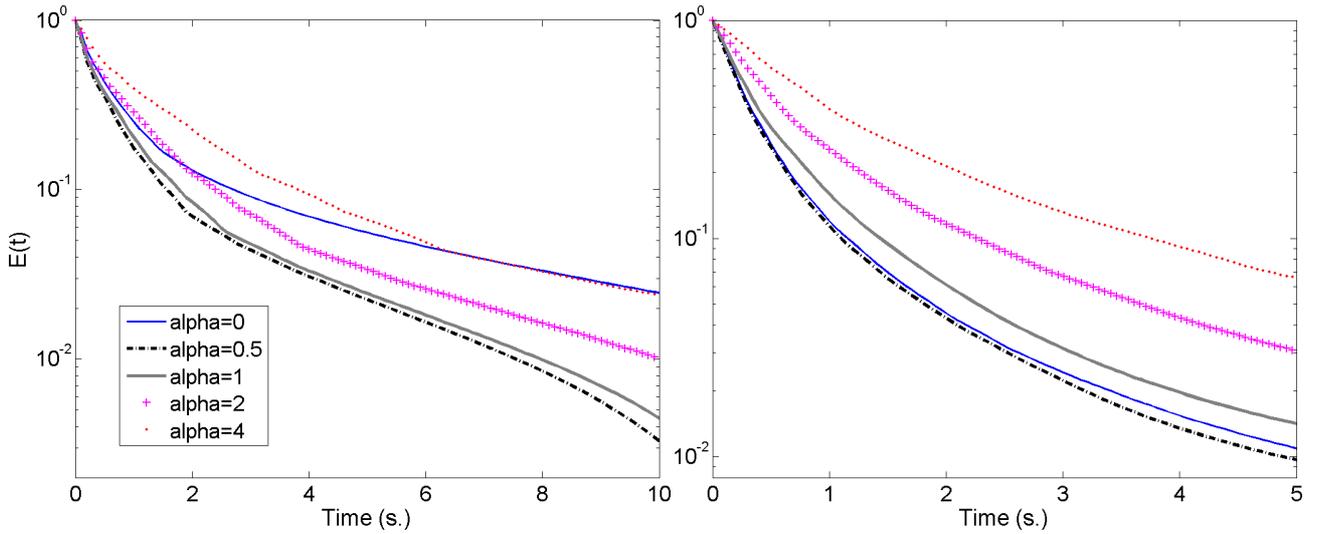


Figure 2: Average of functions $E(t)$ for HALS using different values of α : (left) dense matrices, (right) sparse matrices. Same settings as in Figure 1.

i.e., as soon as the improvement of the last update becomes negligible compared to the one obtained with the first update, but without any a priori fixed maximal number of inner iterations.

Figures 3 shows the results for MU with different values of ϵ (we also include the original MU and MU with $\alpha = 1$ presented in the previous section to serve as a comparison). Figures 4 displays the same experiment for HALS.

In both cases, we observe that the dynamic stopping criterion is not able to outperform the approach based on a fixed number of inner iterations ($\alpha = 1$ for MU, $\alpha = 0.5$ for HALS). Moreover, in the experiments for HALS with sparse matrices, it is not even able to compete with the original algorithm.

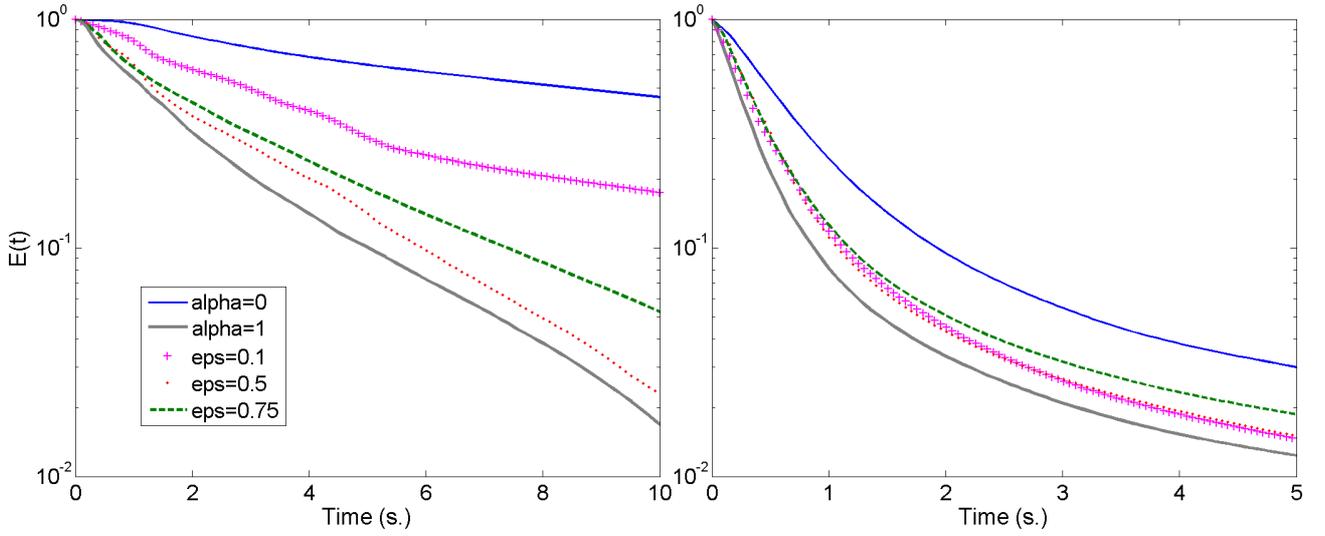


Figure 3: Average of functions $E(t)$ for MU using different values of ϵ , with $\alpha = 0$ and $\alpha = 1$ for reference (see Section 3.1): (left) dense matrices, (right) sparse matrices. Same settings as in Figure 1.

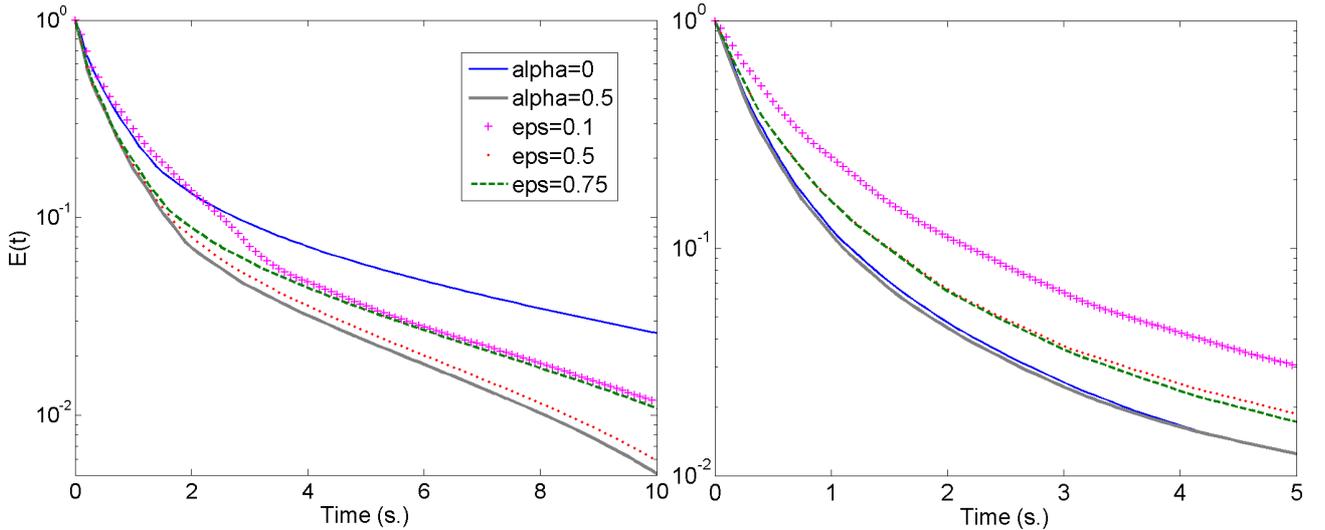


Figure 4: Average of functions $E(t)$ for HALS using different values of ϵ , with $\alpha = 0$ and $\alpha = 0.5$ for reference (see Section 3.1): (left) dense matrices, (right) sparse matrices. Same settings as in Figure 1.

3.3 A Hybrid Stopping Criterion

We have shown in the previous section that using a fixed number of inner iterations works better than a stopping criterion based solely on the difference between two iterates. However, in some circumstances, we have observed that inner iterations become ineffective before their maximal count is reached, so that it would in some cases be worth switching earlier to the other factor.

This occurs in particular when the numbers of rows m and columns n of matrix M have different orders of magnitude. For example, assume without loss of generality that $m \ll n$, so that we have $\rho_W \gg \rho_H$. Hence, on the one hand, matrix W has significantly less entries than H ($mr \ll nr$), and the corresponding NNLS subproblem features a much smaller number of variables; on the other hand, $\rho_W \gg \rho_H$ so that the above choice will lead many more updates of W performed. In other

words, many more iterations are performed on the simpler problem, which might be unreasonable. For example, for the CBCL face database (cf. Section 5) with $m = 361$, $n = 2429$ and $r = 20$, we have $\rho_H \approx 18$ and $\rho_W \approx 123$, and this large number of inner W -updates is typically not necessary to obtain an iterate close to an optimal solution of the corresponding NNLS subproblem.

Therefore, to avoid unnecessary inner iterations, we propose to combine the fixed number of inner iterations proposed in Section 3.1 with the supplementary stopping criterion described in Section 3.2. This safeguard procedure will stop the inner iterations before their maximum number $\lfloor 1 + \alpha\rho_W \rfloor$ is reached when they become ineffective (depending on parameter ϵ , see Equation (3.2)). Algorithm 3 displays the pseudocode for the corresponding accelerated MU, as well as a similar adaptation for HALS. Figures 5 and 6 displays the numerical experiments for MU and HALS respectively.

Algorithm 3 Accelerated MU and HALS

Require: Data matrix $M \in \mathbb{R}_+^{m \times n}$ and initial iterates $(W^{(0)}, H^{(0)}) \in \mathbb{R}_+^{m \times r} \times \mathbb{R}_+^{r \times n}$.

```

1: for  $k = 0, 1, 2, \dots$  do
2:   Compute  $A = MH^{(k)T}$  and  $B = H^{(k)}H^{(k)T}$ ;  $W^{(k,0)} = W^{(k)}$ ;
3:   for  $l = 1 : \lfloor 1 + \alpha\rho_W \rfloor$  do
4:     Compute  $W^{(k,l)}$  using either MU or HALS (cf. Algorithms 1 and 2);
5:     if  $\|W^{(k,l)} - W^{(k,l-1)}\|_F \leq \epsilon \|W^{(k,1)} - W^{(k,0)}\|_F$  then
6:       break;
7:     end if
8:   end for
9:    $W^{(k+1)} = W^{(k,l)}$ ;
10:  Compute  $H^{(k+1)}$  from  $H^{(k)}$  and  $W^{(k+1)}$  using a symmetrically adapted version of steps 2-9;
11: end for

```

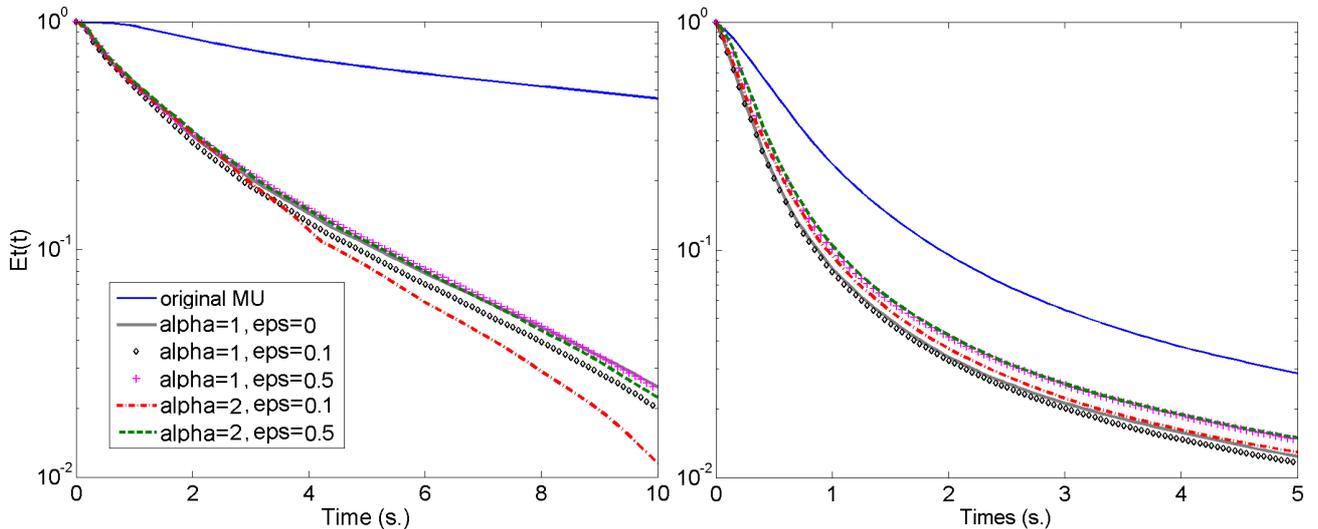


Figure 5: Average of functions $E(t)$ for MU using different values of α and ϵ : (left) dense matrices, (right) sparse matrices. Same settings as in Figure 1.

In the dense case, this safeguard procedure slightly improves performance. We also note that the best values of parameter α now seem to be higher than in the unsafeguarded case ($\alpha = 2$ versus $\alpha = 1$ for MU, and $\alpha = 1$ versus $\alpha = 0.5$ for HALS). Worse performance of those higher values of α

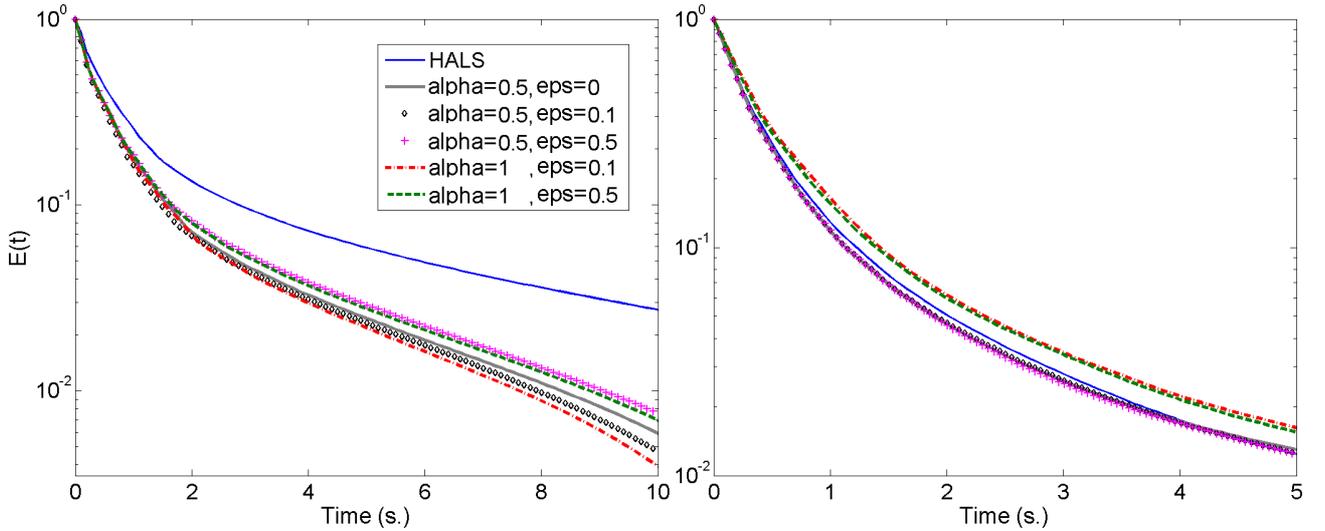


Figure 6: Average of functions $E(t)$ for HALS using different values of α and ϵ : (left) dense matrices, (right) sparse matrices. Same settings as in Figure 1.

in the unsafeguarded scheme can be explained by the fact that additional inner iterations, although sometimes useful, become too costly overall if they are not stopped when becoming ineffective.

In the sparse case, the improvement is rather limited (if not absent) and most accelerated variants provide similar performances. In particular, as already observed in Sections 3.1 and 3.2, the accelerated variant of HALS does not perform very differently from the original HALS on sparse matrices. We explain this by the fact that HALS applied on sparse matrices is extremely efficient and one inner update already decreases the objective function significantly. To illustrate this, Figure 7 shows the evolution of the relative error

$$E^k(l) = \frac{\|M - W^{(k,l)}H^{(k)}\|_F - e_{\min}^k}{\|M - W^{(k,0)}H^{(k)}\|_F - e_{\min}^k}$$

of the iterate $W^{(k,l)}$ for a sparse matrix M , where³ $e_{\min}^k = \min_{W \geq 0} \|M - WH^{(k)}\|_F$. Recall that $(W^{(k,0)}, H^{(k)})$ denotes the solution obtained after k outer iterations (starting from randomly generated matrices). For $k = 1$ (resp. $k = 20$), the relative error is reduced by a factor of more than 87% (resp. 97%) after only one inner iteration.

3.4 Application to Lin's Projected Gradient Algorithm

The accelerating procedure described in the previous sections can potentially be applied to many other NMF algorithms. To illustrate this, we have modified Lin's projected gradient algorithm (PG) (Lin, 2007a) by replacing the original dynamic stopping criterion (based on the stationarity conditions) by the hybrid strategy described in Section 3.3. It is in fact straightforward to see that our analysis is applicable in this case, since Lin's algorithm also requires the computation of HH^T and MH^T when updating W , because the gradient of the objective function in (NMF) is given by $\nabla_W \|M - WH\|_F^2 = 2WHH^T - 2MH^T$. This is also a direct confirmation that our approach can be straightforwardly applied to many more NMF algorithms than those considered in this paper.

Figure 8 displays the corresponding computational results, comparing the original PG algorithm (as available from Lin (2007a)) with its dynamic stopping criterion (based on the norm of the projected gradient) and our variants, based on a (safeguarded) fixed number of inner iterations. It demonstrates

³We have used the active-set algorithm of J.Kim & Park (2008) to compute the optimal value of the NNLS subproblem.

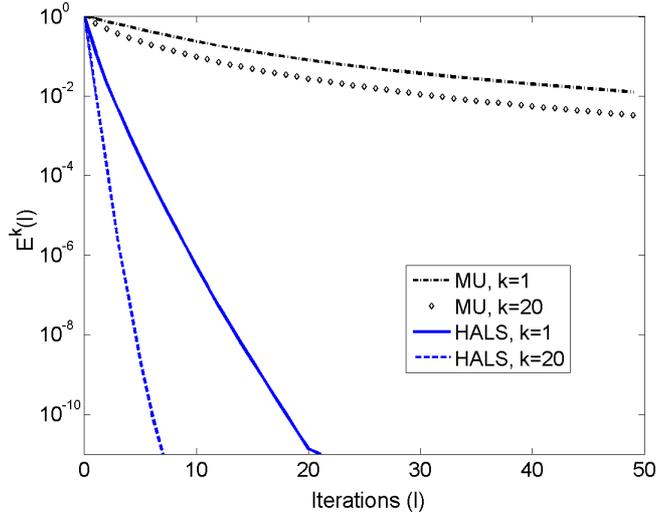


Figure 7: Evolution of the relative error $E^k(l)$ of the iterates of inner iterations in MU and HALS, solving the NNLS subproblem $\min_{W \geq 0} \|M - WH^{(k)}\|_F$ with $r = 40$ for the classic text dataset (cf. Table 2).

that our accelerated schemes perform significantly better, both in the sparse and dense cases (notice that in the sparse case, most accelerated variants perform similarly). The choice $\alpha = 0.5$ gives the best results, and the safeguard procedure does not help much; the reason being that PG converges relatively

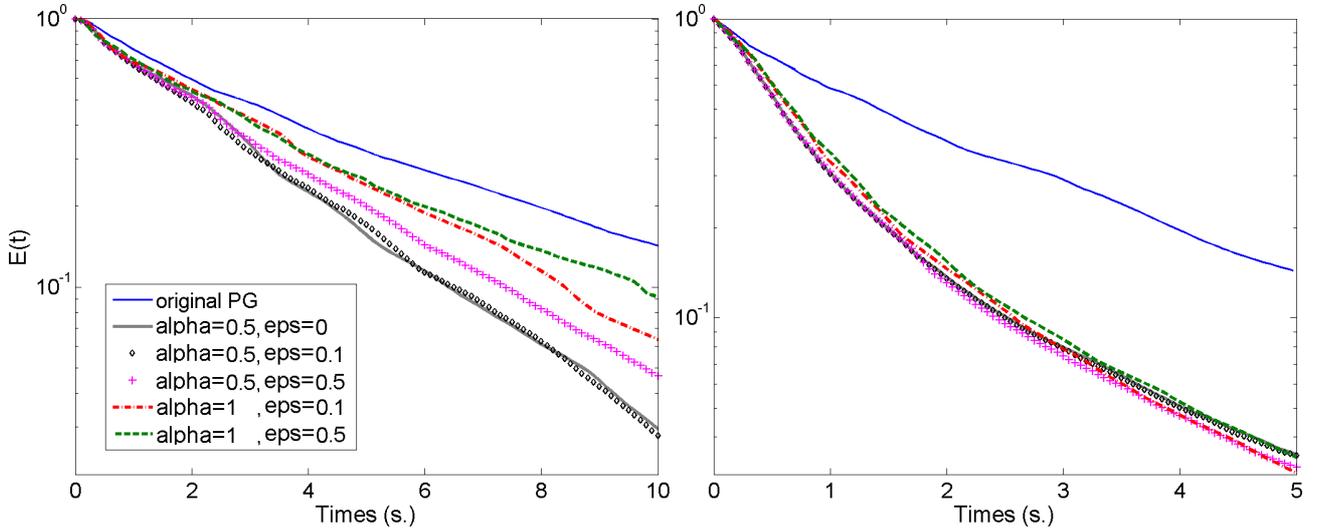


Figure 8: Average of functions $E(t)$ for the projected gradient algorithm of Lin (2007a), and its modification using a fixed number of inner iterations. Same settings as Figure 1.

slowly (we will see in Section 5 that its accelerated variant converges slower than the accelerated MU).

4 Convergence to Stationary Points

In this section, we briefly recall convergence properties of both MU and HALS, and show that they are inherited by their accelerated variants.

4.1 Multiplicative Updates

It was shown by Daube-Witherspoon & Muehlehner (1986) and later by Lee & Seung (1999) that a single multiplicative update of W (i.e., replacing W by $W \circ \frac{[MH^T]}{[WHH^T]}$ while H is kept fixed) guarantees that the objective function $\|M - WH\|_F^2$ does not increase. Since our accelerated variant simply performs several updates of W while H is unchanged (and vice versa), we immediately obtain that the objective function $\|M - WH\|_F^2$ is non-increasing under the iterations of Algorithm 3.

Unfortunately, this property does not guarantee convergence to a stationary point of (NMF), and this question on the convergence of the MU seems to be still open, see Lin (2007b). Furthermore, in practice, rounding errors might set some entries in W or H to zero, and then multiplicative updates cannot modify their values. Hence, it was observed that despite their monotonicity, MU do not necessarily converge to a stationary point, see Gonzales & Zhang (2005).

However, Lin (2007b) proposed a slight modification of MU in order to obtain the convergence to a stationary point. Roughly speaking, MU is recast as a rescaled gradient descent method and the step length is modified accordingly. Another even simpler possibility is proposed by Gillis & Glineur (2008) who proved the following theorem (see also (Gillis, 2011, §4.1) where the influence of parameter δ is discussed):

Theorem 1 (Gillis & Glineur (2008)). *For any constant $\delta > 0$, $M \geq 0$ and any⁴ $(W, H) \geq \delta$, $\|M - WH\|_F$ is nonincreasing under*

$$W \leftarrow \max\left(\delta, W \circ \frac{[MH^T]}{[WHH^T]}\right), \quad H \leftarrow \max\left(\delta, H \circ \frac{[W^T M]}{[W^T W H]}\right), \quad (4.1)$$

where the max is taken component-wise. Moreover, every limit point of the corresponding (alternated) algorithm is a stationary point of the following optimization problem

$$\min_{W \geq \delta, H \geq \delta} \|M - WH\|_F^2.$$

The proof of Theorem 1 only relies on the fact that the limit points of the updates (4.1) are fixed points (there always exists at least one limit point because the objective function is bounded below and non-increasing under updates (4.1)). Therefore, one can easily check that the proof still holds when a bounded number of inner iterations is performed, i.e., the theorem applies to our accelerated variant (cf. Algorithm 3).

It is important to realize that this is not merely a theoretical issue and that this observation can really play a crucial role in practice. To illustrate this, Figure 9 shows the evolution of the normalized objective function (cf. Equation (3.1)) using $\delta = 0$ and $\delta = 10^{-16}$ starting from the same initial matrices $W^{(0)}$ and $H^{(0)}$ randomly generated (each entry uniformly drawn between 0 and 1). We observe that, after some number of iterations, the original MU (i.e., with $\delta = 0$) get stuck while the variant with $\delta = 10^{-16}$ is still able to slightly improve W and H . Notice that this is especially critical on sparse matrices (see Figure 9, right) because many more entries of W and H are expected to be equal to zero at stationarity. For this reason, in this paper, all numerical experiments with MU use the updates from Equation (4.1) with $\delta = 10^{-16}$ (instead of the original version with $\delta = 0$).

4.2 Hierarchical Alternating Least Squares

HALS is an exact block-coordinate descent method where blocks of variables (columns of W and rows of H) are optimized in a cyclic way (first the columns of W , then the rows of H , etc.). Clearly, exact block-coordinate descent methods always guarantee the objective function to decrease. However, convergence to a stationary point requires additional assumptions. For example, Bertsekas (1999a,b) (Proposition 2.7.1) showed that if the following three conditions hold:

⁴ $(W, H) \geq \delta$ means that W and H are component-wise larger than δ .

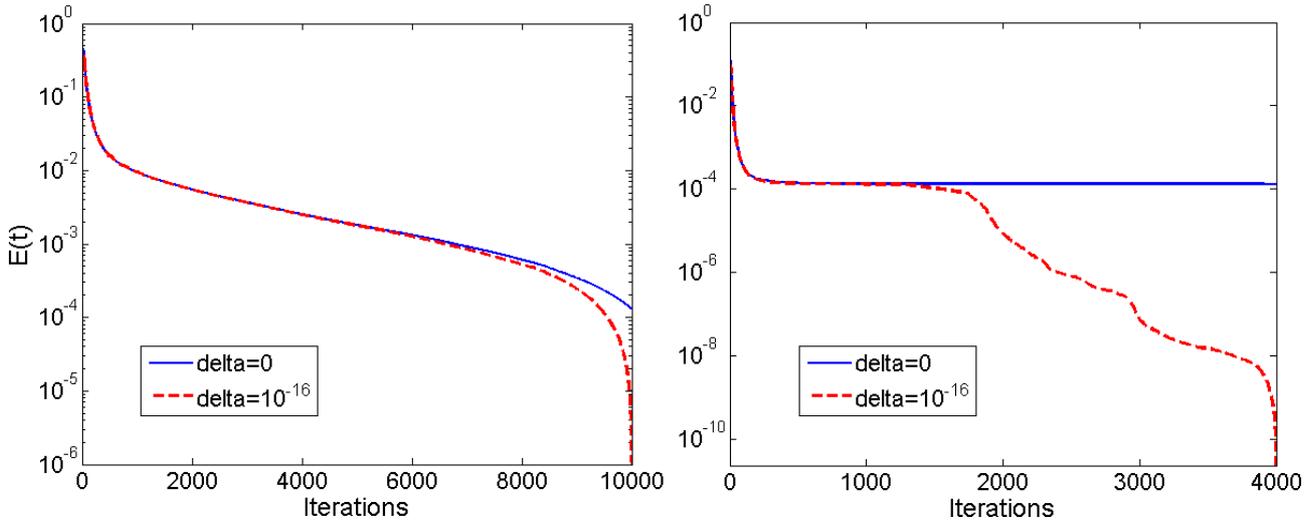


Figure 9: Functions $E(t)$ for $\delta = 0$ and $\delta = 10^{-16}$ on the (dense) ORL face dataset (cf. Table 1) and the (sparse) classic text dataset (cf. Table 2) with $r = 40$.

- each block of variables belongs to a closed convex set (which is the case here since the blocks of variables belong either to \mathbb{R}_+^m or \mathbb{R}_+^n),
- the minimum computed at each iteration for a given block of variables is uniquely attained;
- the function is monotonically nonincreasing in the interval from one iterate to the next;

then exact block-coordinate descent methods converge to a stationary point. The second and the third requirements are satisfied as long as no columns of W and no rows of H become completely equal to zero (subproblems are then strictly convex quadratic programs, whose unique optimal solutions are given by the HALS updates, see Section 1). In practice, if a column of W or a row of H becomes zero⁵, we reinitialize it to a small positive constant (we used 10^{-16}). We refer the reader to Ho (2008) and Gillis & Glineur (2008) for more details on the convergence issues related to HALS.

Because our accelerated variant of HALS is just another type of exact block-coordinate descent method (the only difference being that the variables are optimized in a different order: first several times the columns of W , then several times the rows of H , etc.), it inherits all the above properties. In fact, the statement of the above-mentioned theorem in (Bertsekas, 1999b, p.6) mentions that ‘the order of the blocks may be arbitrary as long as there is an integer K such that each block-component is iterated at least once in every group of K contiguous iterations’, which clearly holds for our accelerated algorithm with a fixed number of inner iterations and its hybrid variant⁶.

5 Numerical Experiments

In this section, we compare the following algorithms, choosing for our accelerated MU, HALS and PG schemes the hybrid stopping criterion and the best compromise for values for parameters α and ϵ according to tests performed in Section 3.

1. (MU) The multiplicative updates algorithm of Lee & Seung (2001).

⁵In practice, this typically only happens if the initial factors are not properly chosen, see Ho (2008).

⁶Note however that the accelerated algorithms based solely on the dynamic stopping criterion (Section 3.2) might not satisfy this requirement, because the number of inner iterations for each outer iteration can in principle grow indefinitely in the course of the algorithm.

2. (**A-MU**) The accelerated MU with a safeguarded fixed number of inner iterations using $\alpha = 2$ and $\epsilon = 0.1$ (cf. Algorithm 3).
3. (**HALS**) The hierarchical alternating least squares algorithm of Cichocki et al. (2007).
4. (**A-HALS**) The accelerated HALS with a safeguarded fixed number of inner iterations using $\alpha = 0.5$ and $\epsilon = 0.1$ (cf. Algorithm 3).
5. (**PG**) The projected gradient method of Lin (2007a).
6. (**A-PG**) The modified projected gradient method of Lin (2007a) using $\alpha = 0.5$ and $\epsilon = 0$ (cf. Section 3.4).
7. (**ANLS**) The alternating nonnegative least squares algorithm⁷ of J.Kim & Park (2008), which alternatively optimizes W and H exactly using a block-pivot active set method. Kim and Park showed that their method typically outperforms other tested algorithms (in particular MU and PG) on synthetic, images and text datasets.

All tests were run using MATLAB[®] 7.1 (R14), on a 3GHz Intel[®] Core[™]2 dual core processor. We present numerical results on images datasets (dense matrices, Section 5.1) and on text datasets (sparse matrices, Section 5.2). Code for all algorithms but ANLS is available at

<http://sites.google.com/site/nicolasgillis/>.

5.1 Dense Matrices - Images Datasets

Table 1 summarizes the characteristics of the different datasets.

Table 1: Image datasets.

Data	# pixels	m	n	r	$\lfloor \rho_W \rfloor$	$\lfloor \rho_H \rfloor$
ORL ¹	112 × 92	10304	400	30, 60	358, 195	13, 7
Umist ²	112 × 92	10304	575	30, 60	351, 188	19, 10
CBCL ³	19 × 19	361	2429	30, 60	12, 7	85, 47
Frey ²	28 × 20	560	1965	30, 60	19, 10	67, 36

$\lfloor x \rfloor$ denotes the largest integer smaller than x .

¹ <http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html>

² <http://www.cs.toronto.edu/~roweis/data.html>

³ <http://cbcl.mit.edu/cbcl/software-datasets/FaceData2.html>

For each dataset, we use two different values for the rank ($r = 30, 60$) and initialize the algorithms with the same 50 random factors ($W^{(0)}, H^{(0)}$) (using i.i.d. uniform random variables on $[0, 1]$)⁸. In order to assess the performance of the different algorithms, we display individually for each dataset the average over all runs of the function $E(t)$ defined in Equation (3.1), see Figure 10.

First, these results confirm what was already observed by previous works: PG performs better than MU (Lin, 2007a), ANLS performs better than MU and PG (J.Kim & Park, 2008), and HALS

⁷Code is available at <http://www.cc.gatech.edu/~hpark/>.

⁸Generating initial matrices ($W^{(0)}, H^{(0)}$) randomly typically leads to a very large initial error $e(0) = \|M - W^{(0)}H^{(0)}\|_F$. This implies that $E(t)$ will get very small after one step of any algorithm. To avoid this large initial decrease, we have scaled the initial matrices such that $\operatorname{argmin}_\alpha \|M - \alpha W^{(0)}H^{(0)}\|_F = 1$; this simply amounts to multiplying W and H by an appropriate constant, see Gillis & Glineur (2008).

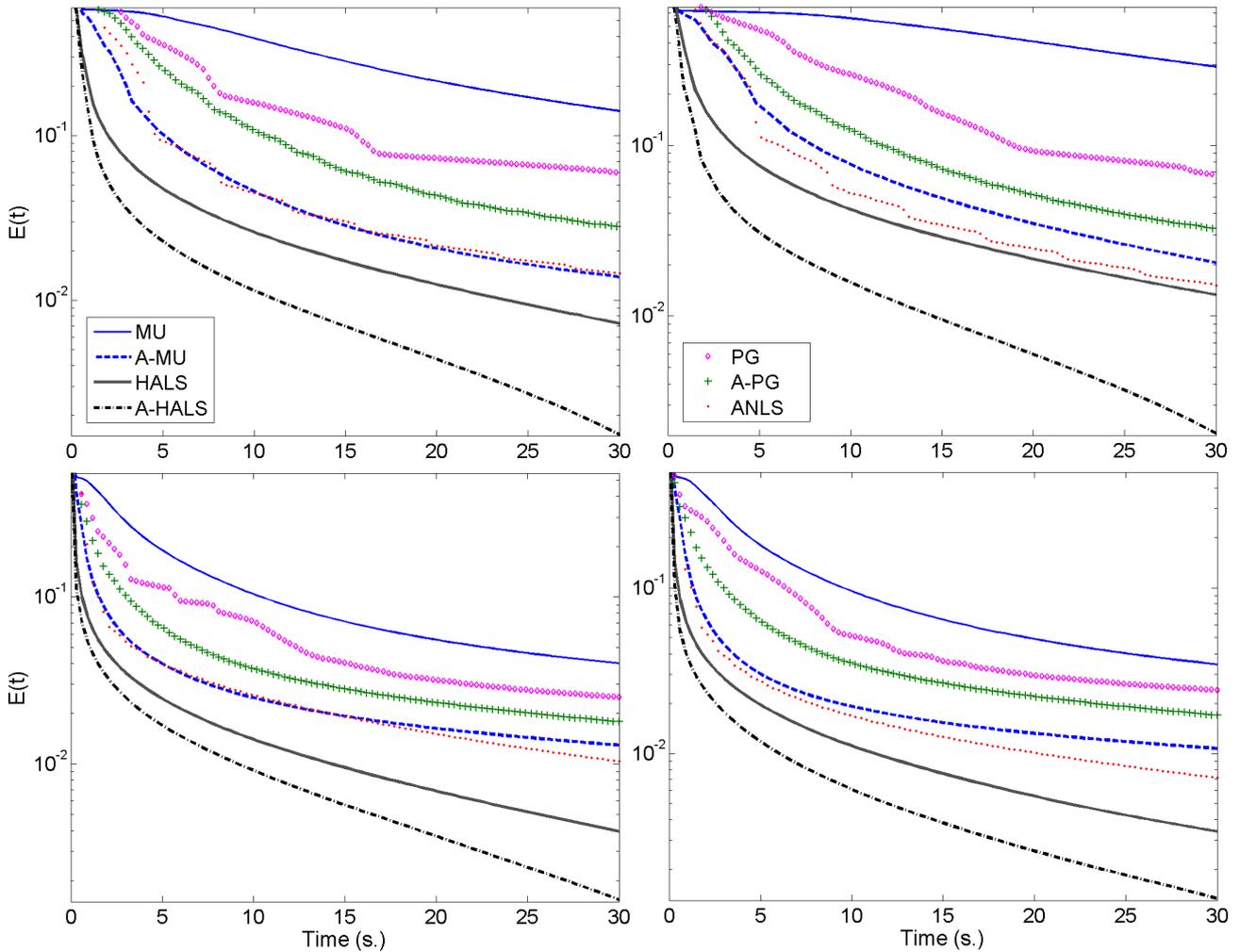


Figure 10: Average of functions $E(t)$ for different image datasets: ORL (top left), Umist (top right), CBCL (bottom left) and Frey (bottom right).

performs the best (Ho, 2008). Second, they confirm that the accelerated algorithms indeed are more efficient: A-MU (resp. A-PG) clearly outperforms MU (resp. PG) in all cases, while A-HALS is, by far, the most efficient algorithm for the tested databases. It is interesting to notice that A-MU performs better than A-PG, and only slightly worse than ANLS, often decreasing the error as fast during the first iterations.

5.2 Sparse Matrices - Text Datasets

Table 2 summarizes the characteristics of the different datasets.

The factorization rank r was set to 10 and 20. For the comparison, we used the same settings as for the dense matrices. Figure 11 displays for each dataset the evolution of the average of functions $E(t)$ over all runs. Again the accelerated algorithms are much more efficient. In particular, A-MU and A-PG converge initially much faster than ANLS, and also obtain better final solutions⁹. A-MU, HALS and A-HALS have the fastest initial convergence rates, and HALS and A-HALS generate the

⁹We also observe that ANLS no longer outperforms the original MU and PG algorithms, and only sometimes generates better solutions.

Table 2: Text mining datasets (Zhong & Ghosh, 2005) (sparsity is given in %: $100 * \#zeros / (mn)$).

Data	m	n	r	#nonzero	sparsity	$\lfloor \rho_W \rfloor$	$\lfloor \rho_H \rfloor$
classic	7094	41681	10, 20	223839	99.92	12, 9	2, 1
sports	8580	14870	10, 20	1091723	99.14	18, 11	10, 6
reviews	4069	18483	10, 20	758635	98.99	35, 22	8, 4
hitech	2301	10080	10, 20	331373	98.57	25, 16	5, 4
ohscal	11162	11465	10, 20	674365	99.47	7, 4	7, 4
la1	3204	31472	10, 20	484024	99.52	31, 21	3, 2

best solutions in all cases. Notice that A-HALS does not always obtain better final solutions than HALS (still this happens on half of the datasets), because HALS already performs remarkably well (see discussion at the end of Section 3.3). However, the initial convergence of A-HALS is in all cases at least as fast as that of HALS.

6 Conclusion

In this paper, we considered the multiplicative updates of Lee & Seung (2001) and the hierarchical alternating least squares algorithm of Cichocki et al. (2007). We introduced accelerated variants of these two schemes, based on a careful analysis of the computational cost spent at each iteration, and preserve the convergence properties of the original algorithms. The idea behind our approach is based on taking better advantage of the most expensive part of the algorithms, by repeating a (safeguarded) fixed number of times the cheaper part of the iterations. This technique can in principle be applied to most NMF algorithms; in particular, we showed how it can substantially improve the projected gradient method of Lin (2007a). We then experimentally showed that these accelerated variants, despite the relative simplicity of the modification, significantly outperform the original ones, especially on dense matrices, and compete favorably with a state-of-the-art algorithm, namely the ANLS method of J.Kim & Park (2008). A direction for future research would be to choose the number of inner iterations in a more sophisticated way, with the hope of further improving the efficiency of A-MU, A-PG and A-HALS. Finally, we observed that HALS and its accelerated version are the most efficient variants for solving NMF problems, sometimes by far.

Acknowledgments

We would like to thank the anonymous reviewers for their insightful comments, which helped to improve the paper.

References

- Berry, M.W., Browne, M., Langville, A.N., Pauca, P.V. & Plemmons, R.J. (2009). Algorithms and Applications for Approximate Nonnegative Matrix Factorization. *Computational Statistics and Data Analysis*, 52, 155–173.
- Bertsekas, D.P. (1999a). Nonlinear Programming: Second Edition. *Athena Scientific, Massachusetts*.
- Bertsekas, D.P. (1999b). Corrections for the book Nonlinear Programming: Second Edition. *Athena Scientific, Massachusetts*. Available at <http://www.athenasc.com/nlperrata.pdf>.
- Cichocki, A., Zdunek, R. & Amari, S. (2006). Non-negative Matrix Factorization with Quasi-Newton Optimization. *Lecture Notes in Artificial Intelligence, Springer, 4029*, 870–879.
- Cichocki, A., Zdunek, R. & Amari, S. (2007). Hierarchical ALS Algorithms for Nonnegative Matrix and 3D Tensor Factorization. *Lecture Notes in Computer Science, Springer, 4666*, 169–176.

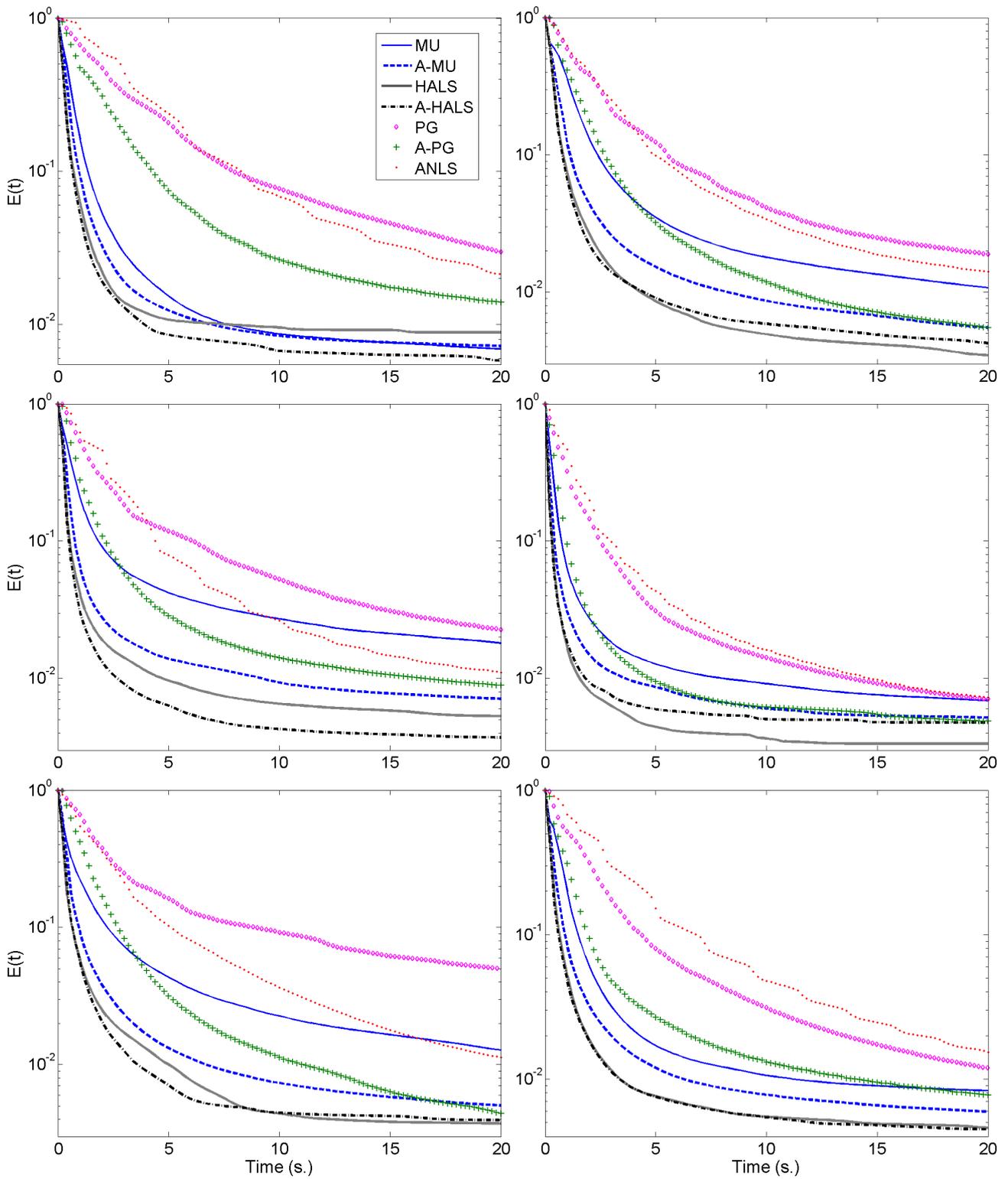


Figure 11: Average of functions $E(t)$ for text datasets: classic (top left), sports (top right), reviews (middle left), hitech (middle right), ohscal (bottom left) and la1 (bottom right).

Cichocki, A., & Phan, A.H. (2009). Fast local algorithms for large scale Nonnegative Matrix and Tensor Factorizations. *IEICE Transactions on Fundamentals of Electronics, Vol. E92-A No.3*, 708–721.

- Cichocki, A., Amari, S., Zdunek, R. & Phan, A.H. (2009). Non-negative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation. *Wiley-Blackwell*.
- Daube-Witherspoon, M.E. & Muehllehner, G. (1986). An iterative image space reconstruction algorithm suitable for volume ECT. *IEEE Trans. Med. Imaging*, 5, 61–66.
- Devarajan, K. (2008). Nonnegative Matrix Factorization: An Analytical and Interpretive Tool in Computational Biology. *PLoS Computational Biology*, 4(7), e1000029.
- Dhillon, I.S., Kim, D. & Sra, S. (2007). Fast Newton-type Methods for the Least Squares Nonnegative Matrix Approximation problem. *Proc. of SIAM Conf. on Data Mining*, 343–354.
- Ding, C., He, X. & Simon, H.D. (2005). On the Equivalence of Nonnegative Matrix Factorization and Spectral Clustering. *Proc. of SIAM Conf. on Data Mining*, 606–610.
- Gillis, N. (2011). Nonnegative Matrix Factorization: Complexity, Algorithms and Applications. *Université catholique de Louvain*, PhD Thesis.
- Gillis, N. & Glineur, F. (2008). Nonnegative Factorization and The Maximum Edge Biclique Problem. *CORE Discussion paper 2008/64*.
- Gonzales, E.F. & Zhang, Y. (2005). Accelerating the Lee-Seung algorithm for non-negative matrix factorization. *Technical report, Department of Computational and Applied Mathematics, Rice University*.
- Han, J., Han, L., Neumann, M. & U. Prasad (2009). On the rate of convergence of the image space reconstruction algorithm. *Operators and Matrices*, 3(1), 41–58.
- Ho, N.-D. (2008). Nonnegative Matrix Factorization - Algorithms and Applications. *Université catholique de Louvain*, PhD Thesis.
- Kim, H. & Park, K. (2008). Non-negative Matrix Factorization Based on Alternating Non-negativity Constrained Least Squares and Active Set Method. *SIAM J. Matrix Anal. Appl.*, 30(2), 713–730.
- Kim, J. & Park, H. (2008). Toward Faster Nonnegative Matrix Factorization: A New Algorithm and Comparisons. *Proc. of IEEE Int. Conf. on Data Mining*, 353–362.
- Lee, D.D. & Seung, H.S. (1999). Learning the Parts of Objects by Nonnegative Matrix Factorization. *Nature*, 401, 788–791.
- Dee, D.D. & Seung, H.S. (2001). Algorithms for Non-negative Matrix Factorization. *Advances in Neural Information Processing*, 13.
- Li, L. & Zhang, Y.-J. (2009). FastNMF: highly efficient monotonic fixed-point nonnegative matrix factorization algorithm with good applicability. *J. Electron. Imaging*, 18, 033004.
- Lin, C.-J. (2007a). Projected Gradient Methods for Nonnegative Matrix Factorization. *Neural Computation, MIT press*, 19, 2756–2779.
- Lin, C.-J. (2007b). On the Convergence of Multiplicative Update Algorithms for Nonnegative Matrix Factorization. *IEEE Trans. on Neural Networks*, 18(6), 1589–1596.
- Pauca, P.V., Piper, J. & Plemmons, R.J. (2006). Nonnegative matrix factorization for spectral data analysis. *Linear Algebra and its Applications*, 406(1), 29–47.
- Shahnaz, F., Berry, M.W., Langville, A.N., Pauca, V.P. & Plemmons, R.J. (2006). Document clustering using nonnegative matrix factorization. *Information Processing and Management*, 42, 373–386.
- Vavasis, S.A. (2009). On the Complexity of Nonnegative Matrix Factorization. *SIAM Journal on Optimization*, 20(3), 1364 – 1377.
- Zhong, S. & Ghosh, J. (2005). Generative model-based document clustering: a comparative study. *Knowledge and Information Systems*, 8(3), 374–384.