

GNMF with Newton-Based Methods

Rafał Zdunek¹, Anh-Huy Phan², and Andrzej Cichocki^{2,3,4}

¹ Department of Electronics, Wrocław University of Technology,
Wybrzeże Wyspiańskiego 27, 50-370 Wrocław, Poland

`rafal.zdunek@pwr.wroc.pl`

² Laboratory for Advanced Brain Signal Processing
RIKEN BSI, Wako-shi, Japan

³ Warsaw University of Technology, Poland

⁴ Systems Research Institute, Polish Academy of Science (PAN), Poland

Abstract. Several variants of Nonnegative Matrix Factorization (NMF) have been proposed for supervised classification of various objects. Graph regularized NMF (GNMF) incorporates the information on the data geometric structure to the training process, which considerably improves the classification results. However, the multiplicative algorithms used for updating the underlying factors may result in a slow convergence of the training process. To tackle this problem, we propose to use the Spectral Projected Gradient (SPG) method that is based on quasi-Newton methods. The results are presented for image classification problems.

Keywords: NMF, Graph-regularized NMF, SPG, Image classification.

1 Introduction

Nonnegative Matrix Factorization (NMF) [1] decomposes a nonnegative matrix into lower-rank factor matrices that have nonnegative entries and usually some physical meaning. When NMF is applied to the matrix of training samples, we obtain sparse nonnegative feature vectors and coefficients of their nonnegative combination. The vectors of the coefficients lie in a low-dimensional latent component space. Hence, NMF is often regarded as a dimensionality reduction technique, and it has been widely applied for classification of various objects [2–6].

As reported in [7], the factor matrices obtained with NMF are generally non-unique. Several attempts have been done to additionally constrain them to satisfy a certain degree of sparsity, smoothness, uncorrelatedness, or orthogonality [2]. Cai *et al.* [8, 9] noticed that the projection from the high-dimensional observation space to the low-dimensional space should preserve the data geometrical structure. That is, any training samples forming one class should, after being projected, belong to the same class in the latent component space. Thus, they proposed Graph regularized NMF (GNMF) [8] that constrains one of the factor matrices with the information on the data geometric structure encoded in the nearest-neighbor graph of the training samples. This constraint was imposed to NMF by a specifically designed regularization term in the objective function

that was then minimized with the standard multiplicative algorithm [2]. Guan *et al.* [10] considerably accelerated the convergence of GNMF by using additive gradient descent updates.

In this paper, we propose to improve the convergence rate of GNMF updates even more, by applying another Newton-based methods that provide the estimates according to the Karush-Kuhn-Tucker (KKT) optimality conditions. First, we formulate the Quadratic Programming (QP) problems for minimizing the penalized objective function. The QP problems can be efficiently solved with many numerical algorithms. To tackle large-scale classification problems, we suggest to use the modified Spectral Projected Gradient (SPG) method that belongs to the class of quasi-Newton methods. Moreover, we also propose to control the penalty parameters iteratively by some schedule included in the alternating update scheme.

The paper is organized in the following way. The next section discusses the Graph-regularized NMF. Section 3 is concerned with the optimization algorithms. The numerical experiments for image classification problems are presented in Section 4. Finally, the conclusions are drawn in Section 5.

2 Graph-Regularized NMF

Let $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_T] \in \mathbb{R}_+^{I \times T}$, where $\mathbf{y}_t \in \mathbb{R}_+^I$ is the t -th training sample. Applying NMF to \mathbf{Y} , we get $\mathbf{Y} \cong \mathbf{A}\mathbf{X}$, where the columns of the matrix $\mathbf{A} \in \mathbb{R}_+^{I \times J}$ represent the feature vectors, and the columns of the matrix $\mathbf{X} \in \mathbb{R}_+^{J \times T}$ are encoding vectors.

In several variants of NMF, the objective function can be expressed by the quadratic function:

$$\Psi(\mathbf{A}, \mathbf{X}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + \frac{\alpha_X}{2} \text{tr}(\mathbf{X}\mathbf{L}_X\mathbf{X}^T) + \frac{\alpha_A}{2} \text{tr}(\mathbf{A}^T\mathbf{L}_A\mathbf{A}), \quad (1)$$

where $\mathbf{L}_X \in \mathbb{R}^{T \times T}$ and $\mathbf{L}_A \in \mathbb{R}^{I \times I}$ are symmetric weighting matrices. In supervised classification, the matrix \mathbf{L}_X contains the information on assignments of the training samples to their classes. In DNMF [4], it is determined by the matrix of inner- and outer-class scattering. In GNMF [8], \mathbf{L}_X is the graph Laplacian matrix that represents a data geometrical structure in the observation space. It takes form: $\mathbf{L}_X = \mathbf{D} - \mathbf{W}$, where $\mathbf{W} = [w_{nm}] \in \mathbb{R}_+^{T \times T}$ contains the entries that determine the edges in the nearest neighbor graph of the observed points, and $\mathbf{D} = \text{diag} \left(\sum_{m=1}^T w_{nm} \right) \in \mathbb{R}_+^{T \times T}$. The edges can be determined by the hard connections:

$$w_{nm} = \begin{cases} 1, & \text{if } \mathbf{y}_n \in \mathcal{N}_p(\mathbf{y}_m), \text{ or } \mathbf{y}_m \in \mathcal{N}_p(\mathbf{y}_n), \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

where $\mathcal{N}_p(\mathbf{y}_t)$ is the p nearest neighbor of the sample \mathbf{y}_t . We can also use the Heat kernel weighting:

$$w_{nm} = \begin{cases} \exp \left\{ -\frac{\|\mathbf{y}_n - \mathbf{y}_m\|_2^2}{\sigma} \right\}, & \text{if } \mathbf{y}_n \in \mathcal{N}_p(\mathbf{y}_m), \text{ or } \mathbf{y}_m \in \mathcal{N}_p(\mathbf{y}_n), \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

or the cosine measure:

$$w_{nm} = \begin{cases} \mathbf{y}_n^T \mathbf{y}_m, & \text{if } \mathbf{y}_n \in \mathcal{N}_p(\mathbf{y}_m), \text{ or } \mathbf{y}_m \in \mathcal{N}_p(\mathbf{y}_n), \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

The matrix \mathbf{L}_A in (1) can enforce the smoothness in the feature vectors (the column vectors in \mathbf{A}) or other modality. We assumed the simplest approach to the smoothness by setting $\mathbf{L}_A = \mathbf{I}_I$, where $\mathbf{I}_I \in \mathbb{R}_+^{I \times I}$ is an identity matrix.

3 Algorithm

Since the matrix \mathbf{L}_X in (1) is a symmetric and positive definite, the regularization term $\text{tr}(\mathbf{X}\mathbf{L}_X\mathbf{X}^T)$ can be reformulated as follows:

$$\Psi_r(\mathbf{X}) = \text{tr}(\mathbf{X}\mathbf{L}_X\mathbf{X}^T) = \|\mathbf{X}\mathbf{L}_X^{\frac{1}{2}}\|_F^2 = \|(\mathbf{L}_X^{\frac{1}{2}} \otimes \mathbf{I}_J)\mathbf{x}\|_2^2 = \mathbf{x}^T(\mathbf{L}_X \otimes \mathbf{I}_J)\mathbf{x}, \quad (5)$$

where $\mathbf{x} = \text{vec}(\mathbf{X}) \in \mathbb{R}^{JT}$ is a vectorized form of \mathbf{X} , and \otimes stands for the Kronecker product.

Considering the function (5), the minimization problem: $\min_{\mathbf{X}} \Psi(\mathbf{A}, \mathbf{X})$, s.t. $\mathbf{X} \geq \mathbf{0}$ can be expressed in terms of the Quadratic Programming (QP) problem: $\min_{\mathbf{x}} \frac{1}{2}\mathbf{x}^T \mathbf{Q}_X \mathbf{x} + \mathbf{c}_X^T \mathbf{x}$, s.t. $\mathbf{x} \geq 0$, where $\mathbf{Q}_X = \mathbf{I}_T \otimes \mathbf{A}^T \mathbf{A} + \alpha_X \mathbf{L}_X \otimes \mathbf{I}_J \in \mathbb{R}^{JT \times JT}$ and $\mathbf{c}_X = -\text{vec}(\mathbf{A}^T \mathbf{Y}) \in \mathbb{R}^{JT}$.

Similarly, the matrix \mathbf{A} can be also computed by formulating the QP problem: $\min_{\mathbf{a}} \frac{1}{2}\mathbf{a}^T \mathbf{Q}_A \mathbf{a} + \mathbf{c}_A^T \mathbf{a}$, s.t. $\mathbf{a} \geq 0$, where $\mathbf{a} = \text{vec}(\mathbf{A}^T) \in \mathbb{R}^{IJ}$, $\mathbf{Q}_A = (\mathbf{X}\mathbf{X}^T + \alpha_A \mathbf{I}_J) \otimes \mathbf{I}_I \in \mathbb{R}^{IJ \times IJ}$ and $\mathbf{c}_A = -\text{vec}(\mathbf{Y}\mathbf{X}^T) \in \mathbb{R}^{IJ}$.

Since the function (1) is quadratic with respect to both arguments \mathbf{A} and \mathbf{X} (but not jointly), the matrices \mathbf{Q}_A and \mathbf{Q}_X are equivalent to the Hessian matrices for \mathbf{A} and \mathbf{X} , respectively. When $\alpha_A > 0$, the matrix \mathbf{Q}_A is positive definite. Under the assumption of positive definiteness of the matrix \mathbf{L}_X , the matrix \mathbf{Q}_X is also positive definite. Hence, both QP problems are strictly convex. To solve such problems, we can use many numerical algorithms such as the Active-Set (AS), Interior-Point (IP), and Spectral Projected Gradient (SPG) [11]. These algorithm are based on the Newton or quasi-Newton updates.

Note that the matrix \mathbf{Q}_A has a block-diagonal structure, and hence the updates of \mathbf{A} might be considerably accelerated by transforming the nonnegative least-squares problem: $\min_{\mathbf{A} \geq 0} \frac{1}{2}\|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + \frac{\alpha_A}{2}\|\mathbf{A}\|_F^2$ to the normal equations $\mathbf{X}\mathbf{X}^T \mathbf{A}^T = \mathbf{X}\mathbf{Y}^T$ subject to the nonnegativity constraints $\mathbf{A} \geq \mathbf{0}$. Then, the solution can be efficiently searched with the FC-NNLS algorithm that was proposed by Benthem and Keenan [12], and then adapted to NMF problems in [13].

The updates for \mathbf{X} cannot be accelerated in the similar way, however, there is still a possibility of applying some quasi-Newton method without formulating the Hessian \mathbf{Q}_X . Note that the matrix \mathbf{Q}_X is very large when the number of training samples is large, and it is rather a dense matrix due to the matrix \mathbf{L}_X . One of these possibilities is to use the SPG method [14] that combines the standard gradient projection scheme with the nonmonotonic Barzilai-Borwein

(BB) method [11]. It is used for minimization of convex functions subject to box-constraints.

In the SPG method, the descent direction $\mathbf{p}_t^{(k)}$ for updating the vector \mathbf{x}_t in the k -th iteration is defined as follows:

$$\mathbf{p}_t^{(k)} = \left[\mathbf{x}_t^{(k)} - (\alpha_t^{(k)})^{-1} \nabla_{\mathbf{x}_t} \Psi(\mathbf{A}, \mathbf{x}_t^{(k)}) \right]_+ - \mathbf{x}_t^{(k)}, \quad (6)$$

for $\alpha_t^{(k)} > 0$ selected in such a way that the matrix $\alpha_t^{(k)} \mathbf{I}_J$ approximates the Hessian matrix.

In [2], this method was adopted to parallel processing of all column vectors in \mathbf{X} . Using this approach, we have the update rule:

$$\mathbf{X}^{(k+1)} = \mathbf{X}^{(k)} + \mathbf{P}^{(k)} \mathbf{Z}^{(k)}, \quad (7)$$

where $\mathbf{Z}^{(k)} = \text{diag}\{\boldsymbol{\eta}^{(k)}\}$. The column vectors of $\mathbf{P}^{(k)} \in \mathbb{R}^{J \times T}$ and the entries of the vector $\boldsymbol{\eta}^{(k)} \in \mathbb{R}_+^T$ are descent directions and steplengths for updating the vectors $\{\mathbf{x}_t\}$, respectively. According to (6), the matrix $\mathbf{P}^{(k)}$ has the form:

$$\mathbf{P}^{(k)} = \left[\mathbf{X}^{(k)} - \mathbf{G}_X^{(k)} \mathbf{D}^{(k)} \right]_+ - \mathbf{X}^{(k)}, \quad (8)$$

where $\mathbf{G}_X^{(k)} = \nabla_{\mathbf{X}} \Psi(\mathbf{A}, \mathbf{X}^{(k)}) \in \mathbb{R}^{J \times T}$ and $\mathbf{D}^{(k)} = \text{diag}\{(\alpha_t^{(k)})^{-1}\} \in \mathbb{R}^{T \times T}$.

The coefficients $\{\alpha_t^{(k)}\}$ can be obtained from the secant equation that is given by $\mathbf{S}^{(k)} \text{diag}\{\alpha_t^{(k+1)}\} = \mathbf{W}^{(k)}$, where $\mathbf{S}^{(k)} = \mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}$ and $\mathbf{W}^{(k)} = \nabla_{\mathbf{X}} \Psi(\mathbf{A}, \mathbf{X}^{(k+1)}) - \nabla_{\mathbf{X}} \Psi(\mathbf{A}, \mathbf{X}^{(k)})$. For the minimization of the objective function (1) with respect to \mathbf{X} , the matrix $\mathbf{W}^{(k)}$ takes the form: $\mathbf{W}^{(k)} = \mathbf{A}^T \mathbf{A} \mathbf{S}^{(k)} + \alpha_X \mathbf{S}^{(k)} \mathbf{L}_X$. From (7) we have: $\mathbf{S}^{(k)} = \mathbf{P}^{(k)} \mathbf{Z}^{(k)}$. In consequence, the secant equation leads to:

$$\begin{aligned} \alpha^{(k+1)} &= \frac{\text{diag}\left\{(\mathbf{S}^{(k)})^T \mathbf{W}^{(k)}\right\}}{\text{diag}\left\{(\mathbf{S}^{(k)})^T \mathbf{S}^{(k)}\right\}} = \frac{\text{diag}\left\{(\mathbf{S}^{(k)})^T \mathbf{A}^T \mathbf{A} \mathbf{S}^{(k)} + \alpha_X (\mathbf{S}^{(k)})^T \mathbf{S}^{(k)} \mathbf{L}_X\right\}}{\text{diag}\left\{(\mathbf{S}^{(k)})^T \mathbf{S}^{(k)}\right\}} \\ &= \frac{\text{diag}\left\{(\mathbf{P}^{(k)})^T \mathbf{A}^T \mathbf{A} \mathbf{P}^{(k)} + \alpha_X (\mathbf{P}^{(k)})^T \mathbf{P}^{(k)} \mathbf{Z}^{(k)} \mathbf{L}_X (\mathbf{Z}^{(k)})^{-1}\right\}}{\text{diag}\left\{(\mathbf{P}^{(k)})^T \mathbf{P}^{(k)}\right\}} \\ &= \frac{\mathbf{1}_J^T \left[\mathbf{P}^{(k)} \circledast \left(\mathbf{A}^T \mathbf{A} \mathbf{P}^{(k)} + \alpha_X \mathbf{P}^{(k)} \mathbf{Z}^{(k)} \mathbf{L}_X (\mathbf{Z}^{(k)})^{-1} \right) \right]}{\mathbf{1}_J^T \left[\mathbf{P}^{(k)} \circledast \mathbf{P}^{(k)} \right]}, \end{aligned} \quad (9)$$

where \circledast stands for the Hadamard product, and the operation $\text{diag}\{\mathbf{M}\}$ creates a vector containing the main diagonal entries of a matrix \mathbf{M} . Note that the matrix $\mathbf{Z}^{(k)}$ is diagonal, so the product $\mathbf{Z}^{(k)} \mathbf{L}_X (\mathbf{Z}^{(k)})^{-1}$ can be readily calculated.

The steplengths can be estimated by solving the minimization problem:

$$\boldsymbol{\eta}_*^{(k)} = \arg \min_{\boldsymbol{\eta}^{(k)}} \Psi \left(\mathbf{A}, \mathbf{X}^{(k)} + \mathbf{P}^{(k)} \text{diag}\{\boldsymbol{\eta}^{(k)}\} \right). \quad (10)$$

If $\alpha_X = 0$, the problem (10) can be expressed in a closed-form. Otherwise, iterative updates must be used, e.g. the Armijo rule [11].

The final form of the modified SPG algorithm is given by Algorithm 1. The final form of the NMF algorithm used in the training process is given by Algorithm 2.

Algorithm 1. SPG algorithm

Input : $\mathbf{Y} \in \mathbb{R}_+^{I \times T}$, $\mathbf{A} \in \mathbb{R}_+^{I \times J}$, $\mathbf{X}^{(0)} \in \mathbb{R}_+^{J \times T}$ - initial guess, k_{max} - number of iterations for SPG updates, $\alpha_{min} > 0$, $\alpha_{max} > 0$, $\forall t : \bar{\alpha}_t^{(0)} = \frac{1}{2}\alpha_{max}$,

Output: $\hat{\mathbf{X}}$ - estimated factor matrices,

- 1 **for** $k = 0, 1, \dots, k_{max}$ **do**
- 2 $\mathbf{G}_X^{(k)} = \nabla_{\mathbf{X}} \Psi(\mathbf{A}, \mathbf{X}^{(k)}) = \mathbf{A}^T(\mathbf{A}\mathbf{X}^{(k)} - \mathbf{Y}) + \alpha_X \mathbf{X}^{(k)} \mathbf{L}_X$; // Gradient
- 3 $\mathbf{P}^{(k)} = \left[\mathbf{X}^{(k)} - \mathbf{G}_X^{(k)} \text{diag}\{(\bar{\alpha}_t^{(k)})^{-1}\} \right]_+$ - $\mathbf{X}^{(k)}$; // Descent direction
- 4 $\bar{\eta}^{(k)} = \max\{0, \min\{1, \eta^{(k)}\}\}$; // where $\eta^{(k)}$ is estimated with (10)
- 5 $\mathbf{X}^{(k+1)} = \mathbf{X}^{(k)} + \mathbf{P}^{(k)} \text{diag}\{\bar{\eta}^{(k)}\}$;
- 6 $\bar{\alpha}^{(k+1)} = \max\{\alpha_{min}, \min\{\alpha_{max}, \alpha^{(k+1)}\}\}$; // where $\alpha^{(k+1)}$ is set to (9)

In the training process, we obtain the nonnegative matrices \mathbf{A} and \mathbf{X} . The column vectors of \mathbf{X} contain the discriminant information. To classify the test sample $\tilde{\mathbf{y}}$, first we need to project it onto the subspace spanned by the column vectors of the matrix \mathbf{A} . As a result, we obtain $\tilde{\mathbf{x}} \in \mathbb{R}_+^J$. This step can be carried out with the SPG, assuming $\alpha_X = 0$. Then, the following problem is solved: $t_* = \arg \min_{1 \leq t \leq T} \|\tilde{\mathbf{x}} - \mathbf{x}_t\|_2$, which gives the index t_* of the class to which the sample $\tilde{\mathbf{y}}$ is classified.

4 Experiments

The experiments are carried out for classification of facial images taken from the ORL database¹. It contains 400 frontal facial images of 40 people (10 pictures per person). We selected 8 training images randomly from each class, and the remaining 2 images are used for testing.

We test the following NMF algorithms: MUE (standard multiplicative Lee-Seung algorithm for the Euclidean distance) [1], GNMF [8], MD-NMF [10], standard projected ALS [2], LPG (Lin's Projected Gradient) [15], IP (Interior-Point NMF) [16], regularized FC-NNLS [13], SPG-NMF (Algorithm 2). For the SPG algorithm, we found the optimal parameters: $\alpha_X = 10^{-5}$, $\bar{\alpha} = 10^{-12}$, $\alpha_0 = 0.01$, and $k_{max} = \min\{k, 50\}$, where k is the alternating step in Algorithm 2. The matrix \mathbf{L}_X is determined using the hard connection criterion given by (2). The iterative process is terminated after 50 alternating steps.

¹ <http://people.cs.uchicago.edu/~dinoj/vis/orl/>

Algorithm 2. SPG-NMF Algorithm

Input : $\mathbf{Y} \in \mathbb{R}^{I \times T}$, J - lower rank, α_0 - initial regularization parameter,
Output: Factor matrices: $\mathbf{A} \in \mathbb{R}_+^{I \times J}$ and $\mathbf{X} \in \mathbb{R}_+^{J \times T}$

- 1 **Initialize**: \mathbf{A} and \mathbf{X} with nonnegative random numbers;
- 2 Replace negative entries (if any) in \mathbf{Y} with zero-value, $k = 0$;
- 3 **repeat**
- 4 $\alpha_A^{(k)} = \max \{ \bar{\alpha}, 2^{-k} \alpha_0 \}$; // Regularization parameter schedule
- 5 $\mathbf{X}^{(k+1)} = \text{SPG}(\mathbf{Y}, \mathbf{A}^{(k)}, \mathbf{X}^{(k)}, \alpha_X)$;
- 6 $\bar{d}_j^{(k+1)} = \sum_{t=1}^T x_{jt}^{(k+1)}$,
 $\mathbf{X}^{(k+1)} \leftarrow \text{diag} \left\{ \left(\bar{d}_j^{(k+1)} \right)^{-1} \right\} \mathbf{X}^{(k+1)}$, $\mathbf{A}^{(k)} \leftarrow \mathbf{A}^{(k)} \text{diag} \left\{ \bar{d}_j^{(k+1)} \right\}$;
- 7 $\bar{\mathbf{A}}^{(k+1)} = \text{FCNNLS}(\mathbf{Y}^T, (\mathbf{X}^{(k+1)})^T, (\mathbf{A}^{(k)})^T, \alpha_A^{(k)})$;
- 8 $\mathbf{A}^{(k+1)} = (\bar{\mathbf{A}}^{(k+1)})^T$;
- 9 $\bar{a}_{ij}^{(k+1)} = \sum_{i=1}^I a_{ij}^{(k+1)}$,
 $\mathbf{X}^{(k+1)} \leftarrow \text{diag} \left\{ \bar{a}_{ij}^{(k+1)} \right\} \mathbf{X}^{(k+1)}$, $\mathbf{A}^{(k+1)} \leftarrow \mathbf{A}^{(k+1)} \text{diag} \left\{ \left(\bar{a}_{ij}^{(k+1)} \right)^{-1} \right\}$;
- 10 $k \leftarrow k + 1$;
- 11 **until** Stop criterion is satisfied;

The NMF algorithms are initialized with uniformly distributed random matrices, and tested for various values of the related parameters. Fig. 1 presents the mean recognition rate versus the number of components (parameter J) obtained with different NMF algorithms.

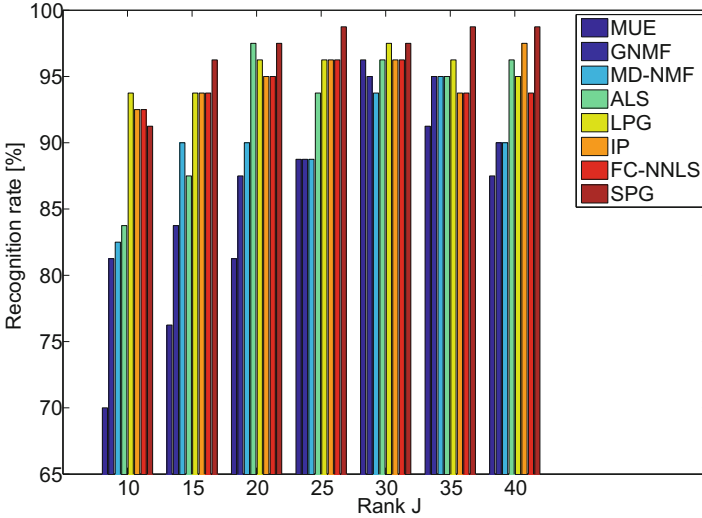


Fig. 1. Recognition rate obtained using various NMF algorithms versus the number of components J

The normalized residual errors versus the number of iterations for the selected NMF algorithms are plotted in Fig. 2.

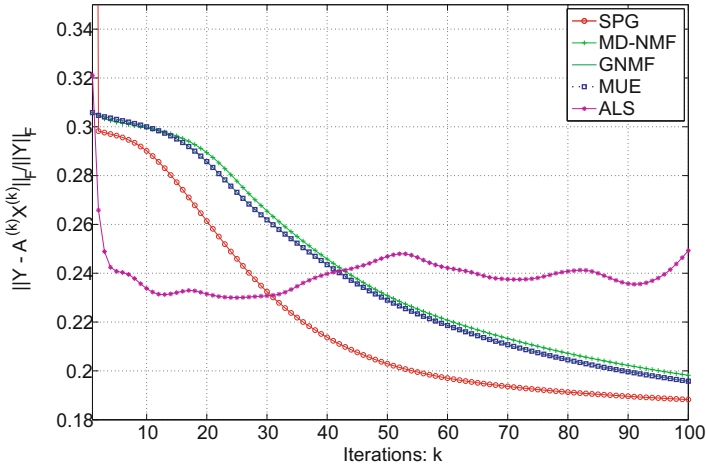


Fig. 2. Normalized residual errors versus alternating iterations

5 Conclusions

The results presented in Fig. 1 demonstrate that the SPG-NMF algorithm outperforms the other tested algorithms in terms of the recognition rate for $J > 10$. Usually an increase in the factorization rank leads to a higher recognition rate. The experiments also confirm that the NMF algorithms based on Newton-like methods (SPG, IP, FC-NNLS and ALS) converge faster than the multiplicative algorithms. This can be observed in Fig. 2 where the SPG-NMF algorithm demonstrates a better convergence behavior than the others. Initially the projected ALS algorithm converges faster but it does not guarantee a monotonic convergence. As observed in Fig. 2 the residual error of the SPG-NMF decreases monotonically with alternating steps. This behavior is also justified by the fact that both SPG and FC-NNLS algorithms converge to the solution optimal according to the KKT conditions. Moreover, the convergence of the SPG-NMF is faster than for the multiplicative algorithms since the SPG is a quasi-Newton method, i.e., the gradient direction is scaled using the information on the Hessian approximation.

Summing up, the experiments showed that the proposed algorithm works very efficiently for the facial classification problem. The usefulness of the proposed algorithm in other applications of NMF will be analyzed in the further research.

References

1. Lee, D.D., Seung, H.S.: Learning the parts of objects by non-negative matrix factorization. *Nature* 401, 788–791 (1999)

2. Cichocki, A., Zdunek, R., Phan, A.H., Amari, S.I.: *Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation*. Wiley and Sons (2009)
3. Qin, L., Zheng, Q., Jiang, S., Huang, Q., Gao, W.: Unsupervised texture classification: Automatically discover and classify texture patterns. *Image and Vision Computing* 26(5), 647–656 (2008)
4. Zafeiriou, S., Tefas, A., Buciu, I., Pitas, I.: Exploiting discriminant information in nonnegative matrix factorization with application to frontal face verification. *IEEE Transactions on Neural Networks* 17(3), 683–695 (2006)
5. Guillamet, D., Vitria, J.: Classifying faces with nonnegative matrix factorization. In: *Proc. 5th Catalan Conference for Artificial Intelligence, Castello de la Plana, Spain*, pp. 24–31 (2002)
6. Benetos, E., Kotti, M., Kotropoulos, C.: Musical instrument classification using non-negative matrix factorization algorithms and subset feature selection. In: *Proc. of 2006 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2006), Toulouse, France* (2006)
7. Donoho, D., Stodden, V.: When does non-negative matrix factorization give a correct decomposition into parts? In: *Thrun, S., Saul, L., Schölkopf, B. (eds.) Advances in Neural Information Processing Systems (NIPS), vol. 16*. MIT Press, Cambridge (2004)
8. Cai, D., He, X., Wu, X., Han, J.: Nonnegative matrix factorization on manifold. In: *Proc. 8-th IEEE International Conference on Data Mining (ICDM), pp. 63–72* (2008)
9. Cai, D., He, X., Han, J., Huang, T.: Graph regularized nonnegative matrix factorization for data representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33(8), 1548–1560 (2011)
10. Guan, N., Tao, D., Luo, Z., Yuan, B.: Manifold regularized discriminative nonnegative matrix factorization with fast gradient descent. *IEEE Transactions on Image Processing* 20(7), 2030–2048 (2011)
11. Nocedal, J., Wright, S.J.: *Numerical Optimization*. Springer Series in Operations Research. Springer, New York (1999)
12. Benthem, M.H.V., Keenan, M.R.: Fast algorithm for the solution of large-scale non-negativity-constrained least squares problems. *Journal of Chemometrics* 18, 441–450 (2004)
13. Kim, H., Park, H.: Non-negative matrix factorization based on alternating non-negativity constrained least squares and active set method. *SIAM Journal in Matrix Analysis and Applications* 30(2), 713–730 (2008)
14. Birgin, E.G., Martnez, J.M., Raydan, M.: Nonmonotone spectral projected gradient methods on convex sets. *SIAM Journal on Control and Optimization* 10, 1196–1211 (2000)
15. Lin, C.J.: Projected gradient methods for non-negative matrix factorization. *Neural Computation* 19(10), 2756–2779 (2007)
16. Zdunek, R.: Spectral signal unmixing with interior-point nonnegative matrix factorization. In: *Villa, A.E.P., Duch, W., Érdi, P., Masulli, F., Palm, G. (eds.) ICANN 2012, Part I. LNCS, vol. 7552, pp. 65–72*. Springer, Heidelberg (2012)