# Log-based Sparse Nonnegative Matrix Factorization for Data Representation

Chong Peng[a], Yiqun Zhang[a], Yongyong Chen[b], Zhao Kang[c], Chenglizhao Chen[a,*], Qiang Cheng[d,e]

[a]*College of Computer Science and Technology, Qingdao University*
[b]*Department of Computer Science, Harbin Institute of Technology*
[c]*School of Computer Science and Engineering, University of Electronic Science and Technology of China*
[d]*Department of Computer Science, University of Kentucky*
[e]*Institute of Biomedical Informatics, University of Kentucky*

**Abstract**

Nonnegative matrix factorization (NMF) has been widely studied in recent years due to its effectiveness in representing nonnegative data with parts-based representations. For NMF, a sparser solution implies better parts-based representation. However, current NMF methods do not always generate sparse solutions. In this paper, we propose a new NMF method with log-norm imposed on the factor matrices to enhance the sparseness. Moreover, we propose a novel column-wisely sparse norm, named $\ell_{2,\log}$-(pseudo) norm to enhance the robustness of the proposed method. The $\ell_{2,\log}$-(pseudo) norm is invariant, continuous, and differentiable. For the $\ell_{2,\log}$ regularized shrinkage problem, we derive a closed-form solution, which can be used for other general problems. Efficient multiplicative updating rules are developed for the optimization, which theoretically guarantees the convergence of the objective value sequence. Extensive experimental results confirm the effectiveness of the proposed method, as well as the enhanced sparseness and robustness.

*Keywords:* Nonnegative matrix factorization, sparse, robust, convergence

*Corresponding author: Chenglizhao Chen
Email address:* `cclz123@163.com` (Chenglizhao Chen)

## 1. Introduction

It has been increasingly ubiquitous to use high-dimensional data in various areas such as machine learning, data mining, and multimedia data processing, which makes the task of learning from examples challenging [1, 2]. A widely used technique to handle such data is dimension reduction, among which matrix factorization methods have drawn significant attention. Matrix factorization seeks two or more low-dimensional matrices to approximate the original data such that the high-dimensional data can be represented with reduced dimensions [3, 4]. Typical matrix factorization techniques include principal component analysis (PCA), nonnegative matrix factorization (NMF), singular value decomposition (SVD), eigenvalue decomposition (EVD), etc.

For some types of data, the entries are naturally nonnegative. For example, the pixel values of images or the frequencies of words in a document are naturally nonnegative. For such data, parts-based representation is believed to commonly exist in human brain with both psychological and physiological evidence [5, 6, 7]. The parts-based representation inspires us to seek two nonnegative factor matrices to approximate the original nonnegative data, which leads to the NMF. Among the two factor matrices, one is considered as the basis while another is treated as the representation or soft indicator matrix. NMF only allows additive combinations of the basis vectors, which enables NMF to learn a parts-based representation [8].

NMF has been extensively studied in recent years [9, 10, 11, 12, 13, 14, 15], which has found applications in various areas, such as pattern recognition [16], multimedia analysis [17], and text mining [18]. Basically, the goal of the original NMF is to approximate a nonnegative matrix with two nonnegative factor matrices with physical meanings [8]. For the constrained optimization problem, multiplicative updating-based strategy has been developed and commonly adopted for its optimization [19]. The original NMF seeks the factorization in Euclidean space and thus it only accounts for the linear relationship of data while omitting the nonlinear ones. To tackle this issue, variants such as the

2

graph regularized NMF (GNMF) [20], the robust manifold NMF (RMNMF) [21] are developed based on manifold learning. Other than manifold technique, kernel method is also used in NMF methods, which relies on the convex NMF framework [22, 23]. It assumes that the basis can be represented as a combination of the data and thus it is doable to calculate the similarity of examples and basis in kernel space.

For NMF methods, it is pointed out that sparser solutions reveal better parts-based representation [20]. However, recent studies show that NMF does not always generate sparse factorization, which implies failure in learning parts-based representation [24, 25]. To combat this issue, various approach have been attempted for sparse solutions, such as Bayesian sparse learning [26, 27]. For example, a maximum a posteriori (MAP) estimation framework is developed to address the sparse nonnegative matrix factorization problems [27], which is built upon the sparse Bayesian learning. The sparse Bayesian learning framework places a sparsity-promoting prior on the data [28], which has been shown to give rise to many models in literature [29]. Bayesian group sparse learning introduces Laplacian scale mixture distribution for sparse coding given a sparseness control parameter [26]. It is natural to impose the $\ell_0$-norm for sparseness, which counts the number of nonzero elements in a matrix. However, $\ell_0$-norm is generally hard to solve and the $\ell_1$-norm has been widely used as a relaxation for sparseness property in the machine learning community. Unfortunately, the $\ell_1$ may be inaccurate to approximate $\ell_0$ if there are large entries in the input matrix or vector. Recently, nonconvex approximations to the rank function have drawn significant attention in various applications such as subspace learning [30] and robust PCA [31]. It is shown that nonconvex approximations such as log-determinant rank approximation have been successful in low-rank matrix recovery problems. The reason is that the nonconvex approximations can better reveal the behavior of the true rank function than the convex approach. In fact, the low-rankness of a matrix is closely related to the sparsity of its singular values, where the rank function is equivalent to the $\ell_0$-norm of the vector of singular values. Thus, the success of nonconvex approximations to the

3

rank function inspires us to design nonconvex approximations to the $\ell_0$-norm for enhanced sparse property.

Moreover, nonnegative data such as images often have noise, which has adverse effects and degrades the learning performance. Thus, there is a crucial need to tackle noise effects from data [21, 32]. The $\ell_{2,1}$-norm, which is defined as the summation of $\ell_2$-norms of all column vectors in a matrix, has been widely used to enforce example-wise sparsity to deal with noise effects [33, 21, 34, 35]. Different from $\ell_1$-norm that treats all entries independently, $\ell_{2,1}$-norm measures input matrix in an example-wise way, which allows it to preserve spatial information of examples [21]. For nonnegative data such as images, such property of $\ell_{2,1}$-norm is indeed desirable. It is noted that the calculation of $\ell_{2,1}$-norm is closely related to the $\ell_1$-norm in that it adds the $\ell_2$-norms of all columns with equal weights. Thus, the $\ell_{2,1}$ may have similar issues to $\ell_1$-norm in resulting column-wise sparse property. In this paper, to better tackle noise issues and enhance column-wise sparsity, we propose a novel $\ell_{2,\log}$-(pseudo) norm with column-wisely sparser property. For the $\ell_{2,\log}$-based shrinkage problem, we provide a closed-form solution, which can be generally used in various other problems.

We summarize the key contributions of this paper as follows: 1) We propose a novel NMF model with log-based sparsity constraints. The new model generates sparser solutions to the factorization, which reveals better parts-based representation; 2) Multiplicative updating rules are developed for efficient optimization; 3) Regarding the updating rules, we provide theoretical analysis that guarantees the convergence of our algorithm; 4) We propose a novel $\ell_{2,\log}$-(pseudo) norm to restrict column-wise sparsity. The $\ell_{2,\log}$-(pseudo) norm is invariant. Similar to the soft-thresholding problem, we formally provide the $\ell_{2,\log}$-shrinkage operator, which is the solution to the $\ell_{2,\log}$-(pseudo) norm associated thresholding problem. The $\ell_{2,\log}$-shrinkage operator can be generally used in other problems; 5) The $\ell_{2,\log}$-shrinkage operator guarantees that the data with noise subtraction are nonnegative at each iteration, which ensures the nonnegativity and convergence of the factorization; 6) Extensive experiments confirm the effectiveness of

4

our method in clustering and data representation.

The rest of this paper is organized as follows: In Section 2, we briefly review some methods that are closely related to our research. Then we introduce our method, including its formulation, optimization, and convergence analysis in Section 3. To enhance the robustness of the new method to noise effects, we present the robust model in Section 4, including its formulation, optimization, and convergence analysis. We conduct extensive experiments and present the detailed results in Section 5. Finally, we conclude the paper in Section 6.

**Remark**: In this paper, the proposed log-based sparse approximations do not satisfy the definition of norms. They are named as (pseudo) norms for simplicity of representation.


## 2. Related Work

In this section, we briefly review some closely related works, including the original NMF and graph Laplacian.


### 2.1. Original NMF

Given nonnegative data $X = [x_1, \cdots, x_n] \in \mathcal{R}^{p \times n}$ with $p$ being the dimension and $n$ sample size, NMF is to factor $X$ into $U \in \mathcal{R}^{p \times k}$ (basis) and $V \in \mathcal{R}^{n \times k}$ (coefficients) with the following optimization problem [8]:

$$\min_{u_{ij} \geq 0, v_{ij} \geq 0} \|X - UV^T\|_F^2, \tag{1}$$

where $u_{ij}$ and $v_{ij}$ are the $ij$-th elements of $U$ and $V$, respectively, and $k \ll n$ enforces a low-rank approximation of the original data.


### 2.2. Graph Laplacian

To account for nonlinear relationships of data in a mapped low-dimensional space, graph Laplacian is a powerful technique that has been widely used [36]. It is defined as

$$\mathbf{Tr}(V^T D V) - \mathbf{Tr}(V^T W V) = \mathbf{Tr}(V^T L V), \tag{2}$$

where $W = [w_{ij}]$ is a similarity matrix that measures pair-wise nonlinear similarities of the examples, $D = [d_{ij}]$ is a diagonal matrix with $d_{ii} = \sum_j w_{ij}$, and $L = D - W$ is the graph Laplacian matrix.

## 3. Log-norm Regularized Sparse NMF

In this section, we will present the log-norm regularized sparse NMF model, including its formulation, optimization, and convergence analysis.

### 3.1. Formulation

For nonnegative data such as images, the parts-based representation has been shown effective in representing their latent structures. This inspires us to develop parts-based representation to represent such data. The NMF methods have been extensively studied for parts-based representation, which are assumed to generate sparse representation due to the fact that they only involve positive combination of the basis. However, recent studies show that NMF does not always result in sparse factorization [24, 25], which implies that NMF does not always succeed in finding good parts-based representation. It is crucial for NMF methods to have sparse solutions, since sparser solutions imply better parts-based representation [20], which reveals the underlying true structures of the data. To enhance the sparseness and reveal the nature of parts-based representation, various approaches have been developed in the literature, such as Bayesian NMF [14] and regularization technique [24]. In this paper, we adopt the most widely used approach and impose the sparsity constraints with the $\ell_1$-norm regularization on the basis and representation matrices $U$ and $V$, respectively, which leads to

$$
\min_{U,V} \|X - UV^T\|_F^2 + \alpha \|U\|_1 + \beta \|V\|_1
$$
$$
s.t. \quad u_{ij} \geq 0, v_{ij} \geq 0,
$$
(3)

where $\alpha \geq 0$ and $\beta \geq 0$ are two balancing parameters, and $\|M\|_1 = \sum_{ij} |m_{ij}|$ is the $\ell_1$-norm for a matrix $M = [m_{ij}]$. In practice, the $\ell_1$-norm is often used

as a tight convex relaxation of the $\ell_0$-norm, whereas the latter is rarely used in practice due to the hard optimization. However, the $\ell_1$ sparsity measure becomes an increasingly bad proxy to the $\ell_0$ norm if any of the elements are large [37]. When the input matrix has large values, the $\ell_1$-norm may approximate the $\ell_0$-norm with significant error, which may lead to inaccurate approximation and sub-optimal solution. Although greater values of $\alpha$ and $\beta$ may lead to sparser solutions, such solutions do not necessarily have good interpretations for the data. To obtain sparsity with a learning model, it is natural to require that the model first approximates the sparsity accurately. Unfortunately, the $\ell_1$-norm may be inaccurate in approximating sparsity and more accurate approximation is admirable for sparsity learning. To withdraw this drawback and better restrict the sparsity, we propose to impose the following regularization instead of the $\ell_1$-norm:

$$\|M\|_{\log} = \sum_{ij} \log(1 + |m_{ij}|).$$

In this paper, we call the above term log-norm ($\ell_{\log}$-(pseudo) norm), which can be considered as a special case in [37]. Particularly, the $\ell_{\log}$-(pseudo) norm admits the followings properties:

- $\|M\|_{\log} \geq 0$ always holds.

- For large $|m_{ij}|$, we have $\log(1+|m_{ij}|) \ll |m_{ij}|$. This reveals that $\|M\|_{\log} \ll \|M\|_1$, implying closer approximation to the true sparsity if a matrix contains large values.

We replace the $\ell_1$-norm with the $\ell_{\log}$-(pseudo) norm in (3) and obtain the following model:

$$\min_{U,V} \|X - UV^T\|_F^2 + \alpha\|U\|_{\log} + \beta\|V\|_{\log} \tag{4}$$
$$s.t. \quad u_{ij} \geq 0, v_{ij} \geq 0.$$

It is seen that the $\ell_{\log}$-(pseudo) norm imposed on the basis and representation matrices renders the model to restrict more accurate sparse constraints, which may lead to sparser solutions. For NMF methods, multiplicative optimization

7

strategies are often adopted. With the $\ell_{\log}$-(pseudo) norm regularization, it is seen that it is more challenging to design efficient multiplicative optimization algorithm for (4). It is noted that model (4) only considers the linear relationships of data in Euclidean space. However, this might be less sufficient since there often exist nonlinear relationships of data, which are omitted in the above model. To account for nonlinear structures of the data, we extend the above model and seek the representation on manifold, which leads to the log-norm regularized sparse NMF (LS-NMF):

$$\min_{U,V} \|X - UV^T\|_F^2 + \lambda \mathbf{Tr}(V^T L V) + \alpha \|U\|_{\log} + \beta \|V\|_{\log}$$

$$s.t. \quad u_{ij} \geq 0, v_{ij} \geq 0,$$

(5)

where $\lambda \geq 0$ is a balancing parameter and $L$ is the graph Laplacian matrix. It is seen that with the graph Laplacian, the new model (5) enforces the smoothness of the data representation from nonlinear kernel space to the low-dimensional space, such that the learned representation $V$ represents nonlinear structures of the data. In rest of this section, we will develop an efficient multiplicative updating rule for the optimization of (5) and provide the corresponding convergence analysis. We will further extend (5) to its robust version in the next section.

*3.2. Optimization*

In this subsection, we will design an efficient optimization algorithm for (5) and present the detailed derivations. The objective in (5) is equivalent to

$$\mathcal{O} = \mathbf{Tr}(XX^T) - 2\mathbf{Tr}(XVU^T) + \mathbf{Tr}(UV^T VU^T)$$

$$+ \lambda \mathbf{Tr}(V^T LV) + \alpha \|U\|_{\log} + \beta \|V\|_{\log}.$$

(6)

The Lagrangian function of $\mathcal{O}$ is

$$\mathcal{L} = \mathbf{Tr}(XX^T) - 2\mathbf{Tr}(XVU^T) + \mathbf{Tr}(UV^T VU^T)$$

$$+ \lambda \mathbf{Tr}(V^T LV) + \alpha \|U\|_{\log} + \beta \|V\|_{\log} + \mathbf{Tr}(\Psi U^T) + \mathbf{Tr}(\Phi V^T),$$

(7)

where $\Psi = [\psi_{ij}]$ and $\Phi = [\phi_{ij}]$ are two matrices with $\psi_{ij}$ and $\phi_{ij}$ being the Lagrangian multipliers of constraints $u_{ij} \geq 0$ and $v_{ij} \geq 0$, respectively. We take

8

the partial derivatives w.r.t $U$ and $V$, respectively, and obtain

$$
\begin{aligned}
\frac{\partial \mathcal{L}}{\partial U} &= -2XV + 2UV^T V + \alpha \mathbf{1}_{d \times k} \oslash (\mathbf{1}_{d \times k} + U) + \Psi \\
\frac{\partial \mathcal{L}}{\partial V} &= -2X^T U + 2VU^T U + 2\lambda LV + \beta \mathbf{1}_{n \times k} \oslash (\mathbf{1}_{n \times k} + V) + \Phi
\end{aligned}
\tag{8}
$$

where $\oslash$ is the element-wise division operation of two matrices, $\mathbf{1}$ is matrix of 1's with size being clarified in the corresponding subscript. Using the Karush-Kuhn-Tucker (KKT) conditions of $\psi_{ij} u_{ij} = 0$ and $\phi_{ij} v_{ij} = 0$, we have

$$
\begin{aligned}
&- 2(XV)_{ij} u_{ij} + 2(UV^T V)_{ij} u_{ij} + \alpha(\mathbf{1}_{d \times k} \oslash (\mathbf{1}_{d \times k} + U))_{ij} u_{ij} + \psi_{ij} u_{ij} \\
={}& - 2(XV)_{ij} u_{ij} + 2(UV^T V)_{ij} u_{ij} + \alpha \frac{u_{ij}}{1 + u_{ij}} = 0,
\end{aligned}
\tag{9}
$$

$$
\begin{aligned}
&- 2(X^T U)_{ij} v_{ij} + 2(VU^T U)_{ij} v_{ij} + 2\lambda(LV)_{ij} v_{ij} + \beta(\mathbf{1}_{n \times k} \oslash (\mathbf{1}_{n \times k} + V))_{ij} v_{ij} + \phi_{ij} v_{ij} \\
={}& - 2(X^T U)_{ij} v_{ij} + 2(VU^T U)_{ij} v_{ij} + 2\lambda(DV)_{ij} v_{ij} - 2\lambda(WV)_{ij} v_{ij} + \beta \frac{v_{ij}}{1 + v_{ij}} = 0.
\end{aligned}
\tag{10}
$$

The above equations (9) and (10) lead to the following updating rules of the proposed model:

$$
u_{ij} \leftarrow u_{ij} \frac{(XV)_{ij}}{(UV^T V)_{ij} + \alpha \frac{1}{2(1 + u_{ij})}}
\tag{11}
$$

$$
v_{ij} \leftarrow v_{ij} \frac{(X^T U)_{ij} + \lambda(WV)_{ij}}{(VU^T U)_{ij} + \lambda(DV)_{ij} + \beta \frac{1}{2(1 + v_{ij})}},
\tag{12}
$$

or equivalently

$$
u_{ij} \leftarrow u_{ij} \frac{(2XV)_{ij}}{(2UV^T V + \alpha \mathbf{1}_{d \times k} \oslash (\mathbf{1}_{d \times k} + U))_{ij}}
\tag{13}
$$

$$
v_{ij} \leftarrow v_{ij} \frac{2(X^T U + \lambda WV)_{ij}}{(2VU^T U + 2\lambda DV + \beta \mathbf{1}_{n \times k} \oslash (\mathbf{1}_{n \times k} + V))_{ij}}.
\tag{14}
$$

For the convergence of the above updating rules, we will provide the detailed theoretical analysis in the following subsection. It is noted that the detailed derivations of optimization and analysis of convergence for (5) are crucial because they are necessary for the robust model in the next section.

### 3.3. Convergence Analysis

In this section, we theoretically analyze the convergence of the updating rules provided in (13) and (14). Regarding the updating rules, we have the

following theorem.

**Theorem 1.** *The objective $\mathcal{O}$ in* (5) *is non-increasing under the updates in* (13) *and* (14)*. The objective function is invariant under these updates if and only if $U$ and $V$ are at a stationary point.*

In the following, we will provide the proof regarding the updates of $U$ and $V$, respectively. To begin the proof, we need to introduce the definition of auxiliary function, which is described below.

**Definition 1.** *For the functions $G(v, v')$ and $F(v)$, if the following conditions*

$$G(v, v') \geq F(v), \quad G(v, v) = F(v)$$

*are satisfied, then $G(v, v')$ is an auxiliary function of $F(v)$.*

Regarding auxiliary function, it has a useful property, which is provided in the following lemma.

**Lemma 1** ([20])**.** *If $G(v, v')$ is an auxiliary function of $F(v)$, then $F(v)$ is non-increasing under the updating rule of*

$$v^{(t+1)} = \underset{v}{\operatorname{argmin}} \, G(v, v^{(t)}). \tag{15}$$

*Proof.* The statement is easily seen according to the following chain of inequality:

$$F(v^{(t+1)}) \leq G(v^{(t+1)}, v^{(t)}) \leq G(v^{(t)}, v^{(t)}) = F(v^{(t)}). \tag{16}$$

$\square$

It is seen that Lemma 1 guarantees that the objective function value of $F(v)$ non-increasing if the updating rule of (15) is adopted. Next, we will show that the updates of (13) and (14) are exactly the updating rule of (15) with proper auxiliary functions of $U$ and $V$, respectively.

Since the updating rules are essentially performed in element-wise manner, it is sufficient to show that the objective function is non-increasing with respect to each element of the matrix variable. In the following, we first consider the

updating of $V$. For $V$, we denote each of its elements by $v_{ab}$. Correspondingly, we use $F_{ab}$ to denote the $v_{ab}$-associated part in $\mathcal{O}$. Then it is straightforward to obtain the first and second derivatives of $F_{ab}$ with respect to $v_{ab}$ as follows:

$$F'_{ab} = (-2X^T U + 2VU^T U + 2\lambda L V)_{ab} + \frac{\beta}{1 + v_{ab}}, \tag{17}$$

$$F''_{ab} = 2(U^T U)_{bb} + 2\lambda L_{aa} - \frac{\beta}{(1 + v_{ab})^2}. \tag{18}$$

Then, we formally have the following lemma, which defines an auxiliary function for $F_{ab}$.

**Lemma 2.** *The function*

$$
\begin{aligned}
G(v, v_{ab}^{(t)}) =& F_{ab}(v_{ab}^{(t)}) + F'_{ab}(v_{ab}^{(t)})(v - v_{ab}^{(t)}) \\
&+ \frac{(VU^T U)_{ab} + \lambda(DV)_{ab} + \frac{\beta}{2(1 + v_{ab}^{(t)})}}{v_{ab}^{(t)}} (v - v_{ab}^{(t)})^2
\end{aligned} \tag{19}
$$

*is an auxiliary function for $F_{ab}(v)$.*

*Proof.* It is easy to check the second condition for auxiliary function, which is seen as below:

$$
\begin{aligned}
G(v, v) =& F_{ab}(v) + F'_{ab}(v)(v - v) \\
&+ \frac{(VU^T U + \lambda DV)_{ab} + \frac{\beta}{2(1+v)}}{v} (v - v)^2 \tag{20} \\
=& F_{ab}(v).
\end{aligned}
$$

To show that the first condition, i.e., $G(v, v_{ab}^{(t)}) \geq F_{ab}(v)$ holds, we compare $G(v, v_{ab}^{(t)})$ with the Tylor expansion series of $F_{ab}(v)$, which is expanded as follows:

$$
\begin{aligned}
F_{ab}(v) =& F_{ab}(v_{ab}^{(t)}) + F'_{ab}(v_{ab}^{(t)})(v - v_{ab}^{(t)}) \\
&+ \left[ (U^T U)_{bb} + \lambda L_{aa} - \frac{\beta}{2(1 + v_{ab}^{(t)})^2} \right] (v - v_{ab}^{(t)})^2.
\end{aligned} \tag{21}
$$

Thus, to show $G(v, v_{ab}^{(t)}) \geq F_{ab}(v)$, we only need to show (19) $\geq$ (21), which is equivalent to show

$$\frac{(VU^T U)_{ab} + \lambda(DV)_{ab} + \frac{\beta}{2(1+v_{ab}^{(t)})}}{v_{ab}^{(t)}} \geq (U^T U)_{bb} + \lambda L_{aa} - \frac{\beta}{2(1 + v_{ab}^{(t)})^2}. \tag{22}$$

11

With straightforward algebra, it is easy to see that

$$(VU^TU)_{ab} = \sum_{l=1}^{k} v_{al}^{(t)} (U^TU)_{lb} \geq v_{ab}^{(t)} (U^TU)_{bb}, \tag{23}$$

and

$$\lambda(DV)_{ab} = \lambda \sum_{l=1}^{n} D_{al} v_{lb}^{(t)} \geq \lambda D_{aa} v_{ab}^{(t)}$$

$$\geq \lambda(D-W)_{aa} v_{ab}^{(t)} = \lambda L_{aa} v_{ab}^{(t)}. \tag{24}$$

Thus,

$$\frac{(VU^TU)_{ab} + \lambda(DV)_{ab} + \frac{\beta}{2(1+v_{ab}^{(t)})}}{v_{ab}^{(t)}}$$

$$\geq \frac{(VU^TU)_{ab} + \lambda(DV)_{ab}}{v_{ab}^{(t)}} \tag{25}$$

$$\geq (U^TU)_{bb} + \lambda L_{aa}$$

$$\geq (U^TU)_{bb} + \lambda L_{aa} - \frac{\beta}{2(1+v_{ab}^{(t)})^2}.$$

Thus, $G(v, v_{ab}^{(t)})$ is an auxiliary function of $F_{ab}(v)$. $\qquad\square$

Next, with the above definition 1 and Lemmas 1 and 2, we will prove Theorem 1 in the following.

*Proof of Theorem 1.* To obtain $v_{ab}^{(t+1)}$, we need to solve the following problem

$$v_{ab}^{(t+1)} = \operatorname*{argmin}_{v} G(v, v_{ab}^{(t)}). \tag{26}$$

It is seen that $G(v, v_{ab}^{(t)})$ defined in (19) is quadratic and convex. Thus, (26) admits solution with first-order optimal condition:

$$F_{ab}'(v_{ab}^{(t)}) + 2\frac{(VU^TU)_{ab} + \lambda(DV)_{ab} + \frac{\beta}{2(1+v_{ab}^{(t)})}}{v_{ab}^{(t)}}(v - v_{ab}^{(t)}) = 0. \tag{27}$$

It is seen that (27) leads to

$$2\frac{(VU^TU)_{ab} + \lambda(DV)_{ab} + \frac{\beta}{2(1+v_{ab}^{(t)})}}{v_{ab}^{(t)}} \cdot v$$

$$= 2(VU^TU)_{ab} + 2\lambda(DV)_{ab} + \frac{\beta}{(1+v_{ab}^{(t)})} - F_{ab}'(v_{ab}^{(t)}) \tag{28}$$

Hence,

$$
\begin{aligned}
v_{ab}^{(t+1)} =& v_{ab}^{(t)} - F_{ab}'(v_{ab}^{(t)}) \frac{v_{ab}^{(t)}}{2(VU^TU)_{ab}+2\lambda(DV)_{ab}+\frac{\beta}{(1+v_{ab}^{(t)})}} \\
=& v_{ab}^{(t)} \frac{2(VU^TU)_{ab}+2\lambda(DV)_{ab}+\frac{\beta}{(1+v_{ab}^{(t)})}-F_{ab}'(v_{ab}^{(t)})}{2(VU^TU)_{ab}+2\lambda(DV)_{ab}+\frac{\beta}{(1+v_{ab}^{(t)})}} \\
=& v_{ab}^{(t)} \frac{2(X^TU)_{ab}+2\lambda(WV)_{ab}}{2(VU^TU)_{ab}+2\lambda(DV)_{ab}+\frac{\beta}{(1+v_{ab}^{(t)})}},
\end{aligned}
\tag{29}
$$

which essentially results in the updating rule of (14). Since $G(v, v_{ab}^{(t)})$ is an auxiliary function for $F_{ab}(v)$, (29) guarantees the non-increasing property of $F_{ab}(v)$. Hence, the objective $\mathcal{O}$ is non-increasing under the update rule of (14).

Mathematically, the matrices $U$ and $V$ are playing similar rules in the model and thus the proof regarding (13) follows (14). We only need to replace $X$ with $X^T$ and set $\lambda = 0$, then the above analysis applies to (13), which concludes the proof.

$\square$

*Remark.* In the above analysis, it is been shown that the value of objective function is decreasing with the alternative updating rules of $U$ and $V$. We define $\Upsilon = [U^T, V^T]^T \in \mathcal{R}^{(d+n)\times k}$ and treat the updates of (13) and (14) as a mapping $\Upsilon^{(t+1)} = \mathcal{M}(\Upsilon^{(t)})$. Then, clearly we have $\Upsilon^* = \mathcal{M}(\Upsilon^*)$ at convergence. Following [22, 38], with non-negativity constraint enforced, we expand $\Upsilon \approxeq \mathcal{M}(\Upsilon^*) + (\partial\mathcal{M}/\partial\Upsilon)(\Upsilon - \Upsilon^*)$, which indicates that $\|\Upsilon^{(t+1)} - \Upsilon^*\| \leq \|\partial\mathcal{M}/\partial\Upsilon\| \cdot \|\Upsilon^{(t)} - \Upsilon^*\|$ under an appropriate matrix norm. In fact, $\|\partial\mathcal{M}/\partial\Upsilon\| \neq 0$ generally holds. Thus, (13) and (14) roughly have a first-order convergence rate.

## 4. Robust Log-norm Regularized Sparse NMF

In this section, we further develop a robust model based on the LS-NMF, which is named robust log-norm regularized sparse NMF model (RLS-NMF). In particular, with the fundamentals of the LS-NMF in Section 3, we will present

the detailed formulation, optimization, and convergence analysis for RLS-NMF in this section.

### 4.1. Formulation of RLS-NMF

It is noted that the LS-NMF model raised in Section 3 seeks the nonnegative representation with original data. Unfortunately, data are often observed and collected with noise, which severely degrades the learning performance. Thus, there is a demanding need to develop more robust model to handle noise effects and promote the learning performance. To enhance the robustness, we adopt the more robust measure $\ell_{2,1}$-norm to minimize the residual instead of the Frobenius-norm, which leads to

$$\min_{U,V} \|X - UV^T\|_{2,1} + \lambda \mathbf{Tr}(V^T LV) + \alpha \|U\|_{\log} + \beta \|V\|_{\log} \tag{30}$$
$$s.t. \quad u_{ij} \geq 0, v_{ij} \geq 0,$$

where, for a matrix $M$, $\|M\|_{2,1} = \sum_j \|m_j\|_2$ is the $\ell_{2,1}$-norm with column-wise sparsity. Here, the $\ell_{2,1}$-norm is invariant and helps keep spatial information of the examples. However, the optimization of $\ell_{2,1}$ with nonnegative constraints is difficult. To facilitate the optimization, we further decompose the data as $X = UV^T + S$, where the matrix $S$ with column-wise sparsity is introduced to account for the noise. With the above assumption, we relax model (30) to the following

$$\min_{U,V,S} \|X - S - UV^T\|_F^2 + \gamma \|S\|_{2,1}$$
$$+ \lambda \mathbf{Tr}(V^T LV) + \alpha \|U\|_{\log} + \beta \|V\|_{\log} \tag{31}$$
$$s.t. \quad u_{ij} \geq 0, v_{ij} \geq 0,$$

where $\gamma \geq 0$ is a balancing parameter. It is seen that the relaxed model (31) is easier to solve and the balancing parameter $\gamma$ allows the model to have more freedom. It is noted that the $\ell_{2,1}$-norm is defined as the summation of $\ell_2$-norms of all column vectors in a matrix, where the summation actually performs in a way similar to $\ell_1$-norm. As pointed out in earlier section that the $\ell_1$-norm might be less efficient in approximating the true sparsity, we design the following novel

$\ell_{2,\log}$-(pseudo) norm to better restrict column-wise sparsity:

$$\|M\|_{2,\log} = \sum_j \log(1 + \|m_j\|_2). \tag{32}$$

For any type of noise, the expectation of the sparsity measurement by (32) is less than the $\ell_{2,1}$-based measurement. Our explanation of the above statement is as follows. Let $c$ be a column of $S$ and we denote the elements of $c$ by $c_1, c_2, \cdots, c_d$. For any types of distribution of $c_i$ for $i = 1, \cdots, d$, the expectation of the log-based approximation is generally less than the $\ell_2$-based approximation. Let $\mathbf{E}(\cdot)$ be the expectation and $f_{\sum_{i=1}^d c_i^2}(y)$ be the probability density function for $y = \sum_{i=1}^d c_i^2$, then the above conclusion can be formally analyzed in the following way:

$$\begin{aligned}
\mathbf{E}\left( \log\left(1 + \sqrt{\sum_{i=1}^d c_i^2}\right)\right) &= \int_0^{+\infty} \log(1 + \sqrt{y}) f_{\sum_{i=1}^d c_i^2}(y)\mathrm{d}y \\
&< \int_0^{+\infty} \sqrt{y} f_{\sum_{i=1}^d c_i^2}(y)\mathrm{d}y = \mathbf{E}\left(\sqrt{\sum_{i=1}^d c_i^2}\ \right).
\end{aligned} \tag{33}$$

The above inequality generally holds for all columns of $S$, i.e., $s_1, \cdots, s_n$, thus it is straightforward that

$$\begin{aligned}
\mathbf{E}\Big(\|S\|_{2,\log}\Big) &= \mathbf{E}\Big(\sum_{i=1}^n \log(1 + \|s_i\|_2)\Big) = \sum_{i=1}^n \mathbf{E}(\log(1 + \|s_i\|_2)) \\
&< \sum_{i=1}^n \mathbf{E}(\|s_i\|_2) = \mathbf{E}\Big(\sum_{i=1}^n \|s_i\|_2\Big) = \mathbf{E}(\|S\|_{2,1}).
\end{aligned} \tag{34}$$

Moreover, if $s_i$ only contains essentially small values, then it is natural that $s_i$ contains noise and is indeed sparse. Thus, for such a column of $S$, it is essentially important that the approximation is close to 0 rather than 1 to distinguish noise effects and useful information. It is noted that $\log(1 + \sqrt{x}) < \sqrt{x}$ holds for small $x > 0$, which indicates that the log-based approximation is closer to 0 than the $\ell_2$-based approach and thus is more accurate in approximating the real sparsity. Thus, it is expected that the log-based approximation is more accurate in approximating the real sparse indicator of the columns than the $\ell_2$-based approach.

We incorporate the $\ell_{2,\log}$-(pseudo) norm into (31), which leads to the robust

log-norm regularized sparse NMF model (RLS-NMF):

$$\min_{U,V,S}\|X - S - UV^T\|_F^2 + \gamma\|S\|_{2,\log}$$

$$+ \lambda\mathbf{Tr}(V^T LV) + \alpha\|U\|_{\log} + \beta\|V\|_{\log} \tag{35}$$

$$s.t. \quad u_{ij} \geq 0, v_{ij} \geq 0.$$

For the optimization and convergence analysis of the RLS-NMF model in (35), we will present them in details in rest of this section.

### 4.2. Optimization of RLS-NMF

For $S$-minimization, the sub-problem is

$$\min_S\|X - S - UV^T\|_F^2 + \gamma\|S\|_{2,\log}. \tag{36}$$

We formally provide the following theorem to solve this type of optimization problem.

**Theorem 2** ($\ell_{2,\log}$-shrinkage operator)**.** *Given matrix $Y$ and a nonnegative parameter $\tau$, the following problem*

$$\min_W \frac{1}{2}\|Y - W\|_F^2 + \tau\|W\|_{2,\log} \tag{37}$$

*admits closed-form solution in a column-wise manner:*

$$w_i = \begin{cases} \frac{\xi}{\|y_i\|_2}y_i, & if\ f_i(\xi) \leq \frac{1}{2}\|y_i\|_2^2,\ (1 + \|y_i\|_2)^2 > 4\tau,\ \xi > 0 \\ 0, & otherwise, \end{cases} \tag{38}$$

*where $f_i(x) = \frac{1}{2}(x - \|y_i\|_2)^2 + \tau\log(1 + x)$, and $\xi = \frac{\|y_i\|_2 - 1}{2} + \sqrt{\frac{(1 + \|y_i\|_2)^2}{4} - \tau}$.*

*Proof.* The objective of (37) can be rewritten in a column-wise manner as

$$\min_{w_i} \sum_{i=1}^{n} \left\{\frac{1}{2}\|y_i - w_i\|_2^2 + \tau\log(1 + \|w_i\|_2)\right\}, \tag{39}$$

such that each $w_i$ can be obtained by column independently. For $w_i$, the sub-problem is

$$\min_{w_i} \frac{1}{2}\|y_i - w_i\|_2^2 + \tau\log(1 + \|w_i\|_2). \tag{40}$$

16

We may treat $w_i$ as a special matrix and perform thin SVD to it. Then it is seen that $w_i$ has exactly one singular value, which is $\sigma(w_i) = \sqrt{w_i^T w_i} = \|w_i\|_2$, where $\sigma(\cdot)$ is the singular value of the input vector. Thus, (40) is equivalent to

$$\min_{w_i} \frac{1}{2}\|y_i - w_i\|_2^2 + \tau \log(1 + \sigma(w_i)). \tag{41}$$

Hence, according to [30], the solution to (41) is obtained with

$$w_i = u_i \sigma^*(w_i) v_i^T, \tag{42}$$

where $u_i$ and $v_i$ are left and right singular vectors of $y_i$ and $\sigma^*(w_i) = \operatorname{argmin}_{x \geq 0} \frac{1}{2}(\sigma(y_i) - x)^2 + \tau \log(1 + x)$. Thus, by solving the equation, we have

$$\sigma^*(w_i) = \begin{cases} \xi, & \text{if } f_i(\xi) \leq f_i(0), \ (1 + \sigma(y_i))^2 > 4\tau, \ \xi > 0, \\ 0, & \text{otherwise,} \end{cases} \tag{43}$$

with $f_i(x) = \frac{1}{2}(x - \sigma(y_i))^2 + \tau \log(1 + x)$, and $\xi = \frac{\sigma(y_i) - 1}{2} + \sqrt{\frac{(1 + \sigma(y_i))^2}{4} - \tau}$. It is straightforward that $y_i = \frac{y_i}{\|y_i\|_2}\|y_i\|_2[1]$ is a thin SVD of $y_i$, where $[1]$ is a special matrix with 1 being the only one element. We substitute $u_i = \frac{y_i}{\|y_i\|_2}$, $\sigma(y_i) = \|y_i\|_2$, and $v_i = [1]$ into above equations, which leads to (38) and concludes the proof. $\square$

*Remark.* Regarding the problem (37), it is easy to verify that for a given $Y$, a larger $\tau$ generally leads to a potentially sparser solution for $W$. To see this, we consider the three conditions given in (38). We only consider the case $0 \leq \tau < \frac{(1 + \|y_i\|_2)^2}{4}$, since $\tau \geq \frac{(1 + \|y_i\|_2)^2}{4}$ directly returns a zero matrix as the solution for $W$. Then, it is straightforward to see that for $\tau' > \tau$, the corresponding $\xi' = \frac{\|y_i\|_2 - 1}{2} + \sqrt{\frac{(1 + \|y_i\|_2)^2}{4} - \tau'} < \xi$. Thus, given $Y$, $Prob(\xi' > 0|Y) < Prob(\xi > 0|Y)$, implying that the third condition is more difficult to satisfy for $\tau'$. For the first condition, our analysis is as follows. Let $e = \frac{1 + \|y_i\|_2}{2}$, then

$$\begin{aligned} f_i(\xi) &= \frac{1}{2}(\xi - \|y_i\|_2)^2 + \tau \log(1 + \xi) \\ &= \frac{1}{2}(-e + \sqrt{e^2 - \tau})^2 + \tau \log(e + \sqrt{e^2 - \tau}) \\ &= (e^2 - e\sqrt{e^2 - \tau} - \frac{\tau}{2}) + \tau \log(e + \sqrt{e^2 - \tau}). \end{aligned} \tag{44}$$

17

We treat (44) as a function of $\tau$ and let $g(\tau) = f_i(\xi)$, then it is seen that

$$
\begin{aligned}
g'(\tau) &= \frac{e}{2\sqrt{e^2 - \tau}} - \frac{1}{2} + \log(e + \sqrt{e^2 - \tau}) + \frac{\tau}{e + \sqrt{e^2 - \tau}} \cdot \frac{-1}{2\sqrt{e^2 - \tau}} \\
&= \frac{e(e + \sqrt{e^2 - \tau}) - \sqrt{e^2 - \tau}(e + \sqrt{e^2 - \tau}) - \tau}{2\sqrt{e^2 - \tau}(e + \sqrt{e^2 - \tau})} + \log(e + \sqrt{e^2 - \tau}) \\
&= \frac{e^2 + e\sqrt{e^2 - \tau} - e\sqrt{e^2 - \tau} - e^2 + \tau - \tau}{2\sqrt{e^2 - \tau}(e + \sqrt{e^2 - \tau})} + \log(e + \sqrt{e^2 - \tau}) \quad (45) \\
&= \log(e + \sqrt{e^2 - \tau}) \\
&= \log\left((1 + \|y_i\|_2)/2 + \sqrt{(1 + \|y_i\|_2)^2/4 - \tau}\right) = \log(1 + \xi) > 0.
\end{aligned}
$$

Thus, it is straightforward that $g(\tau') > g(\tau)$ for $\tau' > \tau$, which implies that $Prob(g(\tau') \leq \frac{\|y_i\|_2^2}{2}|Y) < Prob(g(\tau) \leq \frac{\|y_i\|_2^2}{2}|Y)$. Thus, the first condition is also more difficult to satisfy for $\tau'$. In summary, it is seen that the conditions in (38) are more difficult to satisfy for a larger value of $\tau$, which suggests that a larger $\tau$ potentially leads to a larger number of zero columns for $w$ and thus leads to a potentially sparser solution.

For ease of representation, we denote the $\ell_{2,\log}$-shrinkage operator in (38) as $\mathcal{S}_\tau(Y)$. Then $S$ admits a closed-form solution with the $\ell_{2,\log}$-shrinkage operator:

$$
S = \mathcal{S}_{\frac{\gamma}{2}}(X - UV^T). \tag{46}
$$

To solve $U$ and $V$, the associated sub-problem is

$$
\begin{aligned}
\min_{U,V} \|X - S - UV^T\|_F^2 + \lambda \mathbf{Tr}(V^T LV) + \alpha \|U\|_{\log} \\
+ \beta \|V\|_{\log} \quad s.t. \quad u_{ij} \geq 0, v_{ij} \geq 0.
\end{aligned} \tag{47}
$$

It is seen that the above problem is similar to (5) except that the factorization is performed on $X - S$ in (47) instead of $X$. To derive updating rules for $U$ and $V$ in a similar way to (13) and (14), we first provide the following theorem to guarantee the non-negativity of $X - S$, which is essential for the optimization and nonnegativity of $U$ and $V$.

**Theorem 3.** *Given nonnegative data $X$ and values of $U$ and $V$, the matrix $X - S$ is nonnegative under the updating rule of (46).*

18

*Proof.* We denote $M = X - UV^T$, then it is easy to see that the optimal $s_j$ is either a zero vector or scaled $m_j$ with a positive scaling factor $\xi/\|m_i\|_2$. We consider the following two cases.

1) For the columns that $s_j = 0$, it is easy to verify that $x_j - s_j = x_j$, which is nonnegative.

2) For the columns that $s_j \neq 0$, it is easy to see that

$$
\begin{aligned}
\xi &= \frac{\|m_j\|_2 - 1}{2} + \sqrt{\frac{(1 + \|m_j\|_2)^2}{4} - \tau} \\
&\leq \frac{\|m_j\|_2 - 1}{2} + \sqrt{\frac{(1 + \|m_j\|_2)^2}{4}} = \frac{\|m_j\|_2 - 1}{2} + \frac{1 + \|m_j\|_2}{2} = \|m_j\|_2.
\end{aligned}
\tag{48}
$$

Thus, it is seen that the corresponding $s_j$ is obtained by scaling $m_j$ with a factor $\frac{\xi}{\|m_j\|_2} \leq 1$, which indicates that $x_j - s_j = m_j - s_j + (UV^T)_j = \frac{\|m_j\|_2 - \xi}{\|m_j\|_2} + (UV^T)_j$ is nonnegative.

It is seen that all columns of $X - S$ are nonnegative and thus the matrix $X - S$ is nonnegative. $\square$

Theorem 3 is important in that it guarantees the nonnegativity of $U$ and $V$ with the following updating rules, which is essential to the nature of parts-based representation:

$$
u_{ij} \leftarrow u_{ij} \frac{(2(X - S)V)_{ij}}{(2UV^TV + \alpha \mathbf{1}_{d \times k} \oslash (\mathbf{1}_{d \times k} + U))_{ij}}
\tag{49}
$$

$$
v_{ij} \leftarrow v_{ij} \frac{2((X - S)^TU + \lambda WV)_{ij}}{(2VU^TU + 2\lambda DV + \beta \mathbf{1}_{n \times k} \oslash (\mathbf{1}_{n \times k} + V))_{ij}}.
\tag{50}
$$

*4.3. Convergence Analysis for RLS-NMF*

For the updating rules of (46), (49) and (50), it is guaranteed that the objective function value sequence converges. We formally provide the theoretical result with the following theorem.

**Theorem 4.** *Given nonnegative initial values of $U$ and $V$, the objective function of (35) is monotonally decreasing under the updating rules of (46), (49) and (50).*

*Proof.* Given nonnegative initial values of $U$ and $V$, $X - S$ is nonnegative. Then the proof of $U$ and $V$ follows Theorem 1 by replacing $X$ with $X - S$. Thus, the objective function is non-increasing under the updating rules of (49) and (50). For the updating rule of (46), since it is the optimal solution to (36), the objective function is guaranteed to be non-increasing. Due to the nonnegativity of the objective function (35), the value sequence must converge under the updating rules (46), (49) and (50). $\qquad\square$

## 5. Experiments

In this section, we conduct extensive experiments to testify the effectiveness of the proposed method. In particular, we compare our method with several state-of-the-art NMF methods, including NMF [8], weighted NMF (WNMF) [39], orthogonal NMF (ONMF) [40], convex NMF (CNMF) [22], graph regularized NMF (GNMF) [20], robust manifold NMF (RMNMF) [21], and semi NMF (Semi-NMF) [22]. We testify all methods on 8 widely used benchmark data sets, including Yale [41], Jaffe [42], ORL [43], AR [44], Extended Yale B (EYaleB) [45], COIL20 [46], Pendigits [47], and Semeion [48]. All examples of these data sets are scaled to have a unit $\ell_2$ norm. Three evaluation metrics are used in the experiment, including clustering accuracy, normalized mutual information (NMI), and purity. All these metrics have values ranging within $[0, 1]$, where the higher values represent better clustering results. For these metrics, we briefly introduce them in the following. Clustering accuracy measures the extent to which each cluster contains data points from the same class. It is defined as

$$\text{Accuracy} = \frac{\sum_{i=1}^{n} \delta(map(r_i), l_i)}{n}, \tag{51}$$

where $n$ is the total number of data points, $r_i$ and $l_i$ are the predicted and true labels of the data point $x_i$, respectively, $\delta(a, b)$ is a delta function that returns 1 when $a = b$ otherwise 0, and $map(r_i)$ is a mapping function that maps $r_i$ to an equivalent label by permutation such that (51) is maximized. Normalized
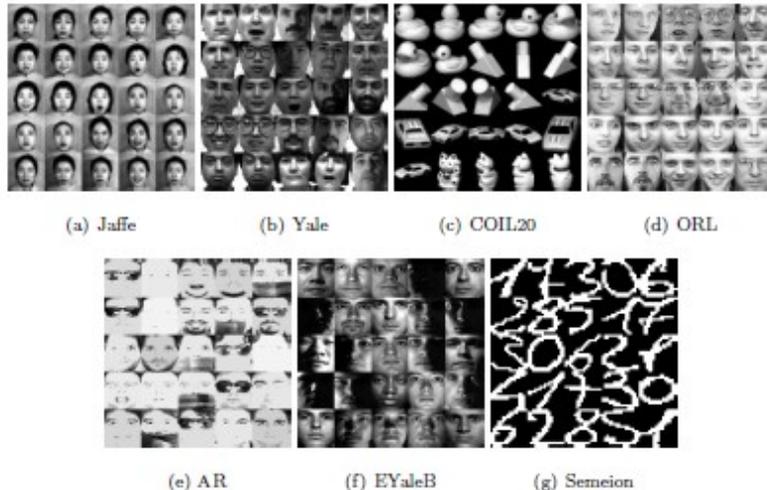
20

Figure 1: Examples of the data sets used in our experiments. Because Pendigits data set has low resolution and it is hard to observe visual feature details, we do not show examples from this data set.

mutual information measures the quality of the clusters, which is defined as

$$\text{Normalized mutual information} = \frac{\sum_{i=1}^{N} \sum_{i=1}^{N} n_{i,j} \log \frac{n_{i,j}}{n_i \hat{n}_j}}{\sqrt{(\sum_{i=1}^{N} n_i \log \frac{n_i}{n})(\sum_{j=1}^{N} \hat{n}_j \log \frac{\hat{n}_j}{n})}}, \quad (52)$$

where $N$ is the number of clusters, $n_i$ and $\hat{n}_j$ denote the sizes of the $i$-th cluster and $j$-th class, respectively, and $n_{i,j}$ denotes the number of data points in the intersection between them. Purity is a simple and transparent evaluation measure, which measures the extent to which each cluster contains data points from primarily one class. It is defined as

$$\text{Purity} = \frac{1}{n} \sum_{i=1}^{N} \max_{j}(n_{i,j}). \quad (53)$$

For all methods in comparison, we follow the following settings for parameters. We tune all the balancing parameters within the set $\{0.001, 0.01, 0.1, 1, 10, 100, 1000\}$ for all methods. For the graph Laplacian that is used in methods such as GNMF, RMNMF, and RLS-NMF, without loss of generality, we use the binary weighting strategy with 5 neighbors kept in the similarity graph ma-

21

trix. For all methods in comparison, the exact number of clusters of the data is provided to determine $k$, which follows a common setting in literature [23, 20]. After we obtain the factorization from each method, K-means clustering is performed on the representation matrix $V$ to obtain the final clustering result. For all methods, we tune the parameters with all possible combinations and report the best performance.

Table 1: Clustering Performance on 8 Benchmark Data Sets

| Data | Accuracy (%) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| Semeion | 52.86 | 63.03 | 43.57 | 55.56 | 45.20 | 50.09 | 57.94 | 68.49 | 68.86 |
| EYaleB | 12.34 | 16.16 | 16.86 | 14.79 | 9.61 | 11.43 | 24.61 | 18.64 | 37.20 |
| ORL | 53.50 | 55.75 | 56.25 | 52.25 | 23.00 | 49.00 | 57.25 | 62.25 | 68.50 |
| AR | 10.46 | 22.69 | 27.54 | 26.00 | 11.92 | 22.85 | 32.62 | 26.85 | 29.23 |
| Jaffe | 85.45 | 97.65 | 95.77 | 90.61 | 69.95 | 82.63 | 97.65 | 98.59 | 98.59 |
| Yale | 21.82 | 43.03 | 44.85 | 44.85 | 40.00 | 41.21 | 40.00 | 48.48 | 48.48 |
| COIL20 | 67.43 | 82.71 | 56.39 | 63.40 | 56.87 | 57.36 | 61.74 | 85.97 | 85.97 |
| Pendigits | 77.97 | 79.20 | 49.55 | 73.62 | 60.12 | 58.58 | 72.78 | 88.16 | 88.26 |
| Data | Normalized Mutual Information (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| Semeion | 43.50 | 58.88 | 35.44 | 44.82 | 37.96 | 40.12 | 47.91 | 61.69 | 62.98 |
| EYaleB | 21.59 | 25.86 | 28.46 | 26.73 | 14.39 | 16.32 | 43.32 | 29.12 | 44.11 |
| ORL | 74.51 | 74.72 | 73.08 | 72.78 | 43.51 | 69.76 | 75.39 | 76.41 | 81.44 |
| AR | 25.71 | 43.49 | 44.04 | 43.43 | 27.12 | 41.23 | 49.46 | 44.53 | 44.72 |
| Jaffe | 85.40 | 96.50 | 93.54 | 89.44 | 70.65 | 84.46 | 96.48 | 98.13 | 98.13 |
| Yale | 29.25 | 48.34 | 48.20 | 51.09 | 41.37 | 43.99 | 47.52 | 51.50 | 51.50 |
| COIL20 | 76.00 | 90.59 | 66.65 | 72.73 | 70.53 | 70.94 | 73.92 | 90.32 | 90.32 |
| Pendigits | 71.09 | 73.02 | 40.48 | 66.07 | 60.29 | 58.84 | 67.27 | 83.88 | 84.06 |
| Data | Purity (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| Semeion | 54.05 | 65.29 | 45.95 | 56.56 | 45.20 | 54.74 | 58.82 | 68.49 | 68.86 |
| EYaleB | 13.01 | 18.39 | 17.94 | 15.29 | 09.32 | | 25.43 | 19.26 | 38.19 |
| ORL | 60.50 | 62.25 | 61.00 | 58.50 | 25.00 | 47.25 | 63.25 | 65.75 | 72.25 |
| AR | 11.08 | 25.23 | 29.31 | 28.92 | 13.00 | 20.38 | 35.38 | 29.00 | 31.62 |
| Jaffe | 85.45 | 97.65 | 95.77 | 90.61 | 74.18 | 82.36 | 97.65 | 98.59 | 98.59 |
| Yale | 26.06 | 44.24 | 44.85 | 47.27 | 40.61 | 40.61 | 41.21 | 48.48 | 48.48 |
| COIL20 | 69.24 | 84.44 | 58.13 | 64.65 | 60.07 | 60.14 | 63.61 | 86.25 | 86.25 |
| Pendigits | 77.97 | 79.20 | 49.57 | 73.62 | 65.78 | 65.01 | 72.78 | 88.16 | 88.26 |

The top three performances are highlighted in red, blue, and green, respectively.

## 5.1. Comparison of Clustering Performance

In this test, we conduct extensive experiments to testify the effectiveness of the proposed method. For all the methods, we follow the settings as described above and report the detailed clustering performance in Table 1. It is observed that the proposed method generally achieves the best performance with significant improvements. In particular, the RLS-NMF achieves the best performances with significant improvements on 7 out of 8 data sets in clustering accuracy and purity, and 6 out of 8 data sets in NMI, respectively. On EYaleB and ORL data sets, the RLS-NMF improves the performance by about 10% in accuracy and purity. In the other cases, the RLS-NMF achieves the top second

Table 2: Clustering Performance on Corrupted Jaffe Data Set

| Corruption Level | Accuracy (%) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 20% | 77.00 | 89.20 | 76.53 | 84.51 | 75.59 | 68.08 | 90.61 | 89.20 | 92.96 |
| 40% | 61.50 | 77.93 | 65.73 | 67.14 | 56.34 | 58.69 | 63.38 | 76.53 | 80.75 |
| 60% | 52.11 | 57.75 | 55.40 | 50.23 | 48.83 | 43.66 | 50.23 | 68.54 | 69.95 |
| Corruption Level | Normalized Mutual Information (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 20% | 79.30 | 87.52 | 75.51 | 81.80 | 76.80 | 70.93 | 86.28 | 86.82 | 89.31 |
| 40% | 58.66 | 73.20 | 62.11 | 69.84 | 57.71 | 54.51 | 60.30 | 73.94 | 78.29 |
| 60% | 49.25 | 59.07 | 50.11 | 47.97 | 44.75 | 41.26 | 47.96 | 62.63 | 66.41 |
| Corruption Level | Purity (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 20% | 80.75 | 89.20 | 79.34 | 84.51 | 76.53 | 72.77 | 90.61 | 89.20 | 92.96 |
| 40% | 62.44 | 77.93 | 67.61 | 71.83 | 61.50 | 59.62 | 63.38 | 76.53 | 80.75 |
| 60% | 54.93 | 60.56 | 56.34 | 54.46 | 50.23 | 46.01 | 52.11 | 69.48 | 69.95 |

Table 3: Clustering Performance on Corrupted Yale Data Set

| Corruption Level | Accuracy (%) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 20% | 36.36 | 36.36 | 38.18 | 38.18 | 34.55 | 30.30 | 39.39 | 44.85 | 46.67 |
| 40% | 31.52 | 39.39 | 38.18 | 38.79 | 27.27 | 34.55 | 27.88 | 40.00 | 43.64 |
| 60% | 30.91 | 32.73 | 30.91 | 31.52 | 26.06 | 32.73 | 31.52 | 39.39 | 40.00 |
| Corruption Level | Normalized Mutual Information (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 20% | 39.55 | 42.21 | 41.38 | 43.43 | 39.71 | 36.98 | 46.60 | 46.62 | 49.15 |
| 40% | 35.96 | 43.03 | 40.84 | 44.11 | 33.55 | 38.08 | 34.39 | 44.33 | 44.62 |
| 60% | 35.51 | 35.90 | 34.48 | 35.39 | 32.03 | 34.99 | 35.22 | 42.59 | 44.09 |
| Corruption Level | Purity (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 20% | 37.58 | 38.79 | 40.00 | 40.00 | 37.58 | 32.73 | 43.64 | 44.85 | 48.48 |
| 40% | 34.55 | 40.61 | 38.18 | 40.00 | 30.91 | 35.15 | 29.09 | 40.00 | 43.64 |
| 60% | 33.33 | 32.73 | 34.55 | 32.73 | 29.70 | 34.55 | 33.94 | 40.00 | 41.21 |

performance with comparable performance to the best method. For example, on COIL20 data set, the RLS-NMF has the top performance in accuracy and purity, respectively. In NMI, the RSL-NMF achieves the top second performance, which is slightly inferior to the GNMF by 0.27%. It is noted that among all baseline methods, GNMF and Semi-NMF are among the most competing ones. In particular, Semi-NMF has the best performance on AR data set in all metrics while GNMF has the best performance in NMI on COIL20 data set. In all other cases, these two methods have inferior performances to the RLS-NMF. Generally, we observe that the RLS-NMF has better performance than the LS-NMF. However, on some data sets, such as Jaffe and COIL20, the RLS-NMF has the same performance as the LS-NMF. We explain this as follows: In such data sets as Jaffe and COIL20, the noise effects are not strong, which can be observed

Table 4: Clustering Performance on Corrupted COIL20 Data Set

| Corruption | Accuracy (%) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Level | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 20% | 60.76 | 79.31 | 56.32 | 64.79 | 53.75 | 60.49 | 59.51 | 80.97 | 83.82 |
| 40% | 56.53 | 67.85 | 58.89 | 56.53 | 51.32 | 58.47 | 56.94 | 76.25 | 80.49 |
| 60% | 61.46 | 67.99 | 56.74 | 61.60 | 54.48 | 56.39 | 61.25 | 71.32 | 71.32 |
| Corruption | Normalized Mutual Information (%) | | | | | | | | |
| Level | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 20% | 70.79 | 85.60 | 68.89 | 72.60 | 66.88 | 71.60 | 70.91 | 87.13 | 88.79 |
| 40% | 69.36 | 74.60 | 69.48 | 69.36 | 63.75 | 68.83 | 69.60 | 79.99 | 83.14 |
| 60% | 69.32 | 75.54 | 67.75 | 71.19 | 62.87 | 66.77 | 68.12 | 77.30 | 77.30 |
| Corruption | Purity (%) | | | | | | | | |
| Level | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 20% | 61.46 | 80.76 | 57.99 | 65.90 | 58.82 | 66.11 | 61.53 | 82.78 | 84.10 |
| 40% | 58.75 | 71.46 | 61.18 | 58.75 | 56.04 | 62.78 | 57.92 | 76.32 | 80.49 |
| 60% | 62.78 | 69.03 | 60.35 | 62.92 | 56.87 | 58.54 | 61.81 | 73.26 | 73.26 |

The top three performances are highlighted in red, blue, and green, respectively.

from Figure 1. Thus, we believe that the noise term and the $\ell_{2,\log}$-(pseudo) norm is not essential in this case. However, on other data sets with heavy noise effects, such as EYaleB data set, we can see that the RLS-NMF has significantly improved performance than the LS-NMF, which verifies the effectiveness and necessity of the robust model. In the next test, we will further testify the RLS-NMF on data sets with artificial noise to confirm its effectiveness.

## 5.2. Comparison on Noisy Data with Random Corruptions

To further testify the effectiveness and robustness to noise effects of the RLS-NMF model, we conduct experiments on noisy data. In particular, we consider two types of noise in our experiments, including random corruption and Gaussian noise. In this subsection, we first conduct experiments on randomly corrupted data. Throughout this subsection, we keep the same settings for all methods as in Section 5.1. Without loss of generality, we conduct experiments on Yale, Jaffe, and COIL20 data sets, where we randomly remove 20%, 40%, and 60% elements from these data sets, respectively. We show some examples of the corrupted data examples in Figure 2. It is seen that the images are severely damaged with the 60% level of corruption, which makes the clustering of such data more challenging. We report the comparison results in Tables 2 to 4. It is observed that the RLS-NMF obtains the best performances in all cases with significant improvements over the baseline methods. It should be noted that on

24

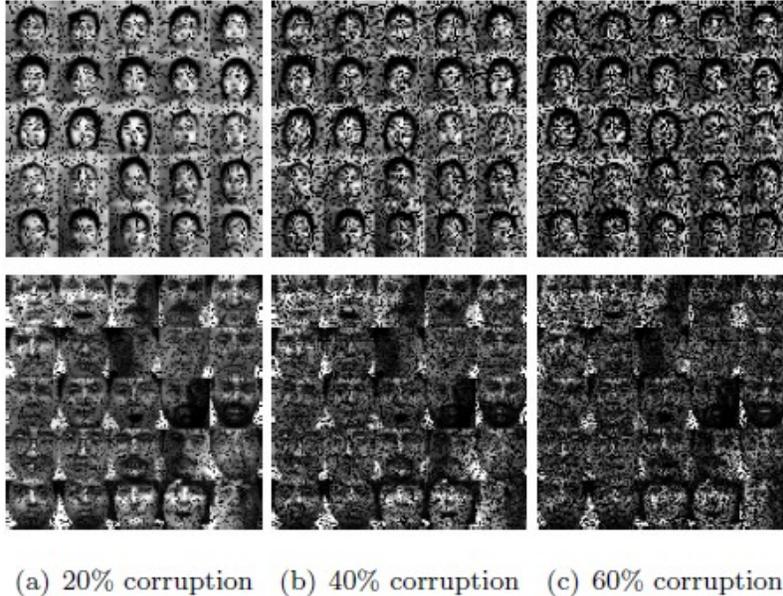(a) 20% corruption    (b) 40% corruption    (c) 60% corruption

Figure 2: Examples of the corrupted images from Jaffe (on top) and Yale (on bottom) data sets. From left to right are images with 20%, 40%, and 60% level corruptions, respectively.

original Yale and Jaffe data sets, the RLS-NMF obtains the same performance with the LS-NMF. However, with heavy noise effects, the RLS-NMF shows superior performance to the LS-NMF. These observations show the improved robustness of the robust model on randomly corrupted data sets.

*5.3. Comparison on Noisy Data with Gaussian Noise*

In this subsection, we further evaluate the proposed method on data sets with Gaussian noise. In particular, we test all methods under different noise level conditions. Without loss of generality, we conduct experiments using Yale, ORL, COIL20, and Semeion data sets. In our test, we add zero mean Gaussian noise to the data sets with variance varies within $\{0.005^2, 0.01^2, 0.015^2\}$, which are referred as light, moderate, and heavy noise conditions in this test. For each data set, we map the data to ensure nonnegativity by subtracting the smallest value. We show some examples of ORL and COIL20 data sets for illustration of

Table 5: Clustering Performance on ORL Data Set with Gaussian Noise

| Variance | Accuracy (%) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 0.005 | 56.50 | 56.00 | 57.75 | 55.50 | 18.00 | 39.50 | 47.75 | 61.00 | 63.50 |
| 0.010 | 49.75 | 52.75 | 52.16 | 49.75 | 15.75 | 28.50 | 32.75 | 54.25 | 54.75 |
| 0.015 | 33.50 | 39.25 | 37.96 | 33.50 | 17.25 | 23.00 | 25.50 | 40.75 | 44.00 |
| Variance | Normalized Mutual Information (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 0.005 | 73.61 | 73.86 | 62.18 | 73.99 | 37.39 | 59.44 | 65.01 | 75.86 | 76.52 |
| 0.010 | 67.32 | 69.67 | 69.10 | 67.32 | 34.09 | 49.39 | 51.49 | 69.97 | 73.17 |
| 0.015 | 56.62 | 59.26 | 53.42 | 56.62 | 34.68 | 43.61 | 44.53 | 58.25 | 61.45 |
| Variance | Purity (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 0.005 | 60.75 | 60.00 | 57.76 | 61.00 | 19.00 | 43.50 | 52.00 | 64.50 | 66.00 |
| 0.010 | 54.75 | 58.05 | 55.83 | 54.75 | 17.00 | 32.25 | 36.25 | 57.75 | 60.25 |
| 0.015 | 37.75 | 43.25 | 42.85 | 37.75 | 17.50 | 25.00 | 28.00 | 44.75 | 48.25 |

Table 6: Clustering Performance on Yale Data Set with Gaussian Noise

| Variance | Accuracy (%) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 0.005 | 40.61 | 43.03 | 43.89 | 41.82 | 24.85 | 38.79 | 41.82 | 48.48 | 49.70 |
| 0.010 | 40.00 | 44.24 | 39.77 | 40.00 | 32.12 | 35.15 | 36.36 | 43.03 | 48.48 |
| 0.015 | 36.36 | 41.21 | 37.56 | 40.00 | 23.64 | 32.73 | 36.36 | 44.85 | 46.06 |
| Variance | Normalized Mutual Information (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 0.005 | 45.68 | 45.03 | 44.66 | 45.87 | 32.26 | 43.55 | 47.74 | 50.59 | 49.04 |
| 0.010 | 41.76 | 49.16 | 41.06 | 41.76 | 37.92 | 38.28 | 45.41 | 44.32 | 50.92 |
| 0.015 | 40.02 | 43.34 | 39.98 | 42.47 | 26.36 | 34.47 | 42.42 | 49.71 | 49.03 |
| Variance | Purity (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 0.005 | 41.21 | 45.45 | 43.89 | 43.03 | 29.70 | 40.00 | 46.06 | 49.09 | 49.70 |
| 0.010 | 41.82 | 45.45 | 39.97 | 41.82 | 33.33 | 35.76 | 38.79 | 44.24 | 49.09 |
| 0.015 | 38.79 | 43.03 | 37.56 | 40.61 | 24.24 | 33.33 | 36.97 | 46.67 | 49.09 |

the noise effects in Figure 3. The other experimental settings remain the same as those in Sections 5.1 and 5.2. We report the clustering results in Tables 5 to 7. As the noise becomes heavier, it is observed that all methods have significantly reduced performances, which confirms the adverse effects of noise. Generally, we can see that the LS-NMF and RLS-NMF obtain the best performances among all methods. Moreover, the RLS-NMF has relatively better performance than the LS-NMF. For example, the RLS-NMF has slightly improved performance over LS-NMF on COIL20 data set under the light noise condition. However, the RLS-NMF has significantly improved performance over LS-NMF on COIL20 data set under moderate and heavy noise conditions, which, again, confirms the enhanced robustness of the RLS-NMF model and the effectiveness of using the $\ell_{2,\log}$-norm.

Table 7: Clustering Performance on COIL20 Data Set with Gaussian Noise

| Variance | Accuracy (%) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 0.005 | 58.75 | 79.17 | 77.35 | 63.40 | 47.29 | 59.65 | 60.28 | 83.54 | 83.96 |
| 0.010 | 64.24 | 77.64 | 73.33 | 64.24 | 26.94 | 61.32 | 65.59 | 80.76 | 82.01 |
| 0.015 | 67.64 | 72.22 | 70.09 | 59.93 | 19.79 | 57.50 | 61.32 | 74.24 | 77.50 |
| Variance | Normalized Mutual Information (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 0.005 | 73.76 | 85.85 | 85.07 | 74.30 | 56.95 | 70.90 | 70.91 | 90.35 | 90.63 |
| 0.010 | 76.03 | 84.98 | 85.59 | 76.03 | 38.89 | 70.72 | 72.21 | 87.39 | 88.43 |
| 0.015 | 75.95 | 79.03 | 77.74 | 72.72 | 28.28 | 69.76 | 68.40 | 79.47 | 80.98 |
| Variance | Purity (%) | | | | | | | | |
| | NMF | GNMF | RMNMF | WNMF | CNMF | ONMF | Semi-NMF | LS-NMF | RLS-NMF |
| 0.005 | 62.78 | 80.03 | 79.17 | 65.35 | 50.07 | 62.64 | 62.99 | 86.39 | 86.45 |
| 0.010 | 65.63 | 80.00 | 75.56 | 65.63 | 30.90 | 62.78 | 66.87 | 80.97 | 82.15 |
| 0.015 | 69.10 | 75.00 | 73.21 | 62.64 | 22.36 | 59.24 | 61.60 | 74.31 | 78.54 |



(a) $0.005^2$ variance    (b) $0.01^2$ variance    (c) $0.015^2$ variance

Figure 3: Examples of the images with Gaussian noise from ORL (on top) and COIL20 (on bottom) data sets. From left to right are images with Gaussian noise with variance being $0.005^2$, $0.01^2$, and $0.015^2$, respectively.

### 5.4. Parameter Sensitivity

For unsupervised learning methods, it is difficult to determine optimal parameters in real-world applications. Thus, it is important that unsupervised method performs well insensitively to parameters. In this test, we show the

Figure 4: Performance changes of RLS-NMF with respect to $\gamma$, where $\alpha$ and $\beta$ are fixed to be the optimal ones used in Section 5.1, respectively.



Figure 5: Performance changes of RLS-NMF with respect to $\alpha$ and $\beta$ on different data sets, where $\gamma$ is fixed to be the optimal one used in Section 5.1.

sensitivity of RLS-NMF to the balancing parameters. Without loss of generality, we show the results on four data sets, including COIL20, Jaffe, ORL, and Semeion. On other data sets, we can observe similar patterns. Specifically, we conduct experiments from two perspectives. That is, we first show the effects of parameter $\gamma$ and then show how the combination of parameters $\{\alpha, \beta\}$ affects the final clustering performance of RLS-NMF, respectively. For

Figure 6: Examples of convergence curves of the LS-NMF (on top) and RLS-NMF (on bottom) on different data sets.

$\gamma$, we fix $\alpha$ and $\beta$ to be the ones used in Section 5.1 and vary $\gamma$ within the set $\{0.001, 0.01, 0.1, 1, 10, 100, 1000\}$. We show the results in Figure 4. It is observed that with a broad range of values for $\gamma$, the RLS-NMF can achieve good performance. We can also observe that the RLS-NMF tends to perform better with larger $\gamma$ values, which might be explained that larger $\gamma$ values render the RLS-NMF better account for noise effects.

For $\{\alpha, \beta\}$, we fix $\gamma$ to be the ones used in Section 5.1 and vary $\alpha$ and $\beta$ within the set $\{0.001, 0.01, 0.1, 1, 10, 100, 1000\}$. We show the results in Figure 5. It is seen that the RLS-NMF obtains relatively high performance within a broad range of parameter selections. In particular, the RLS-NMF is insensitive to $\alpha$ and small $\beta$ values tend to be more effective. We observe similar patterns on other data sets, which suggests us to set small values for $\beta$ in real-world applications.

### 5.5. Convergence Analysis and Time Comparison

In Sections 3.3 and 4.3, we have provided theoretical analysis of the convergence of the objective value of the LS-NMF and RLS-NMF methods. In this section, we further provide experimental results to verify the convergent property of the optimization algorithms. Without loss of generality, we conduct experiments on four data sets, including Yale, Semeion, Jaffe, and EYaleB, where we plot how the value of objective function changes with respect to the it-
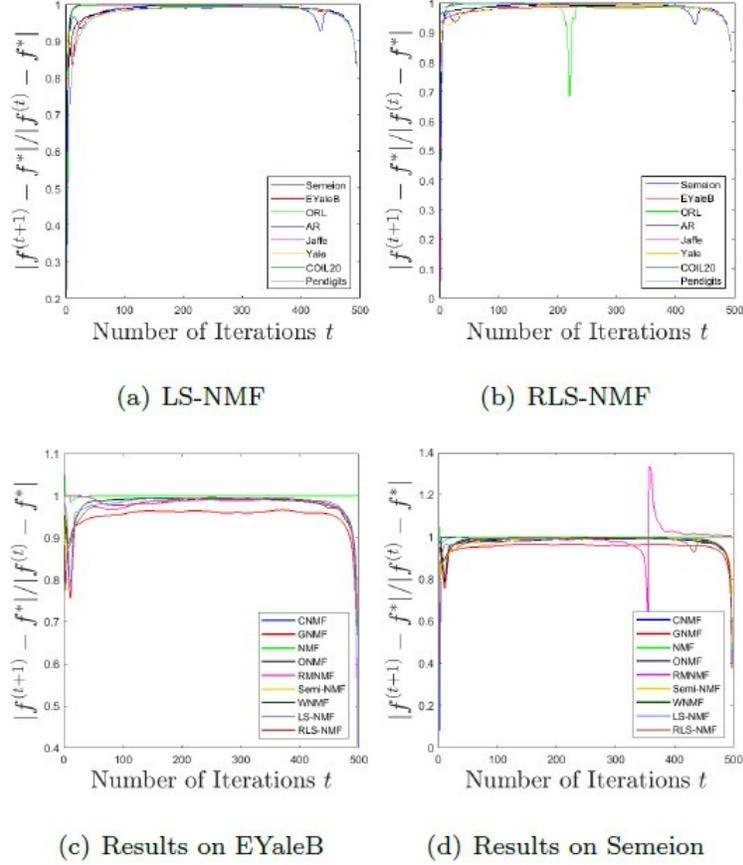
29

Figure 7: Empirical results on convergence rates of different methods: (a)-(b) curves of LS-NMF and RLS-NMF on all data sets, respectively; (c)-(d) comparison of all methods on EYaleB and Semeion data sets, respectively.

eration numbers for both LS-NMF and RLS-NMF. For each data set, we fix the parameters to be the ones used in Section 5.1, which lead to the highest clustering performance. We show the curves of the first 200 iterations in Figure 6. It is observed that the objective value sequences on these data sets indeed converge, which empirically verifies the convergent property of the proposed algorithms. Moreover, we observe that the proposed algorithms generally converge within about 100 iterations, which implies their fast convergence and efficiency.

Figure 8: Comparison of time cost of all methods on different data sets.

Then, we further empirically investigate the convergence rate of the proposed algorithms. In particular, we plot the values of $\left\{\frac{|f^{(t+1)}-f^*|}{|f^{(t)}-f^*|}\right\}$ in our experiment, where $\{f^{(t)}\}$ denotes the objective value sequence and $f^*$ denotes the convergent value of $\{f^{(t)}\}$. As seen in above test, the proposed LS-NMF and RLS-NMF algorithms converge within about 200 iterations. Thus, in this test we treat $f^{(500)}$ as the empirical value of $f^*$, where 500 is sufficiently large to guarantee the convergence of $\{f^{(t)}\}$, and show the results in ????. It is seen that, as far as can be observed, the values of $\frac{|f^{(t+1)}-f^*|}{|f^{(t)}-f^*|}$ are always less than 1, which indicates that the sequence $\{f^{(t)}\}$ is convergent with a linear convergence rate. Without loss of generality, we test all algorithms on EYaleB and Semeion data sets to compare their convergence rates and show the results in ????. It is seen that the curves generally imply that these methods have a linear convergence rate, which basically suggests that all the methods have comparable convergence rates.

Moreover, we test the time cost for all methods and show the results in Figure 8, where all algorithms are terminated after 500 iterations. This test is conducted with MATLAB 2018a on a workstation with Intel(R) Xeon(R) W-2133 CPU. It is seen that the LS-NMF and RLS-NMF have at least comparable efficiency to other methods, such as ONMF. Considering that the LS-NMF and RLS-NMF have superior performance to the other baseline methods in clustering, such comparable speed, though not the fastest, is indeed acceptable.
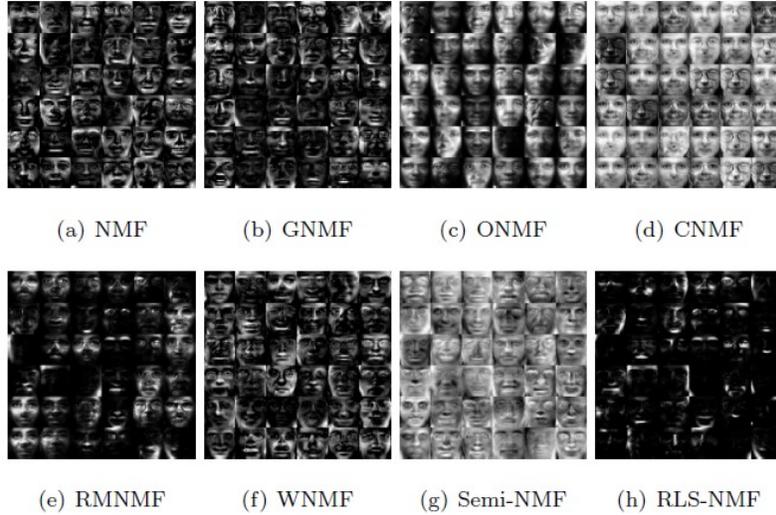
31

Figure 9: Visual examples of the basis vectors learned by various methods from ORL data set by different methods.

*5.6. Sparsity Learning*

In this test, we investigate the ability of the proposed method in sparsity learning with the $\ell_{\log}$-(pseudo) norm regularization. Throughout this test, we fix the parameters to those used in Section 5.1 for all methods, where these parameters lead to the best clustering performance. As pointed out in [20], sparser basis learned by NMF models implies better parts-based representation. In this test, we first show some examples of the basis learned by different methods. Without loss of generality, we conduct this test on ORL data set. The learned basis vectors have 1024 dimensions and we reshape them to size of 32×32 for visual illustration. Then we visually show the reshaped basis vectors as gray scale images in Figure 9. It is observed that the basis images learned by Semi-NMF are not sparse, which is explained by the mixed signs of the basis. Among all methods, RLS-NMF generates the sparsest basis vectors as can be observed from the results. The sparser basis vectors suggest that RLS-NMF can learn better parts-based representation than the baseline methods, which implies its effectiveness in finding a low-dimensional representation of the data.
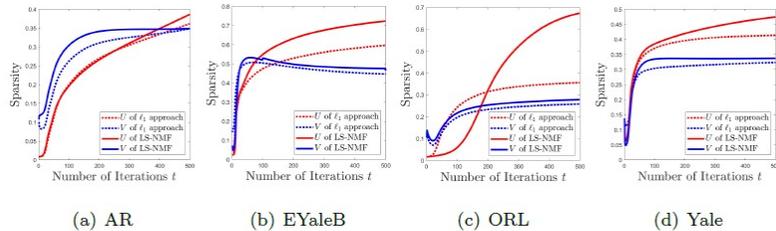
Figure 10: Comparison of LS-NMF and the $\ell_1$ approach in sparse learning.



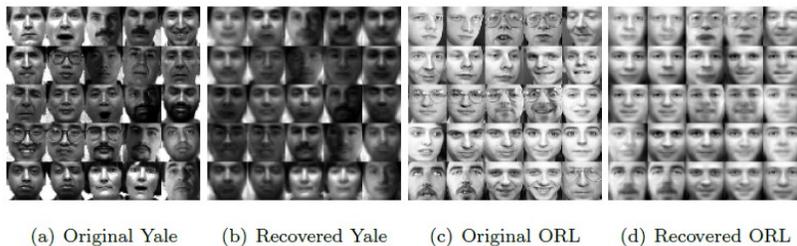(a) Original Yale    (b) Recovered Yale    (c) Original ORL    (d) Recovered ORL

Figure 11: Visual examples of the original and reconstructed images. (a)-(b) are from Yale while (c)-(d) are from ORL, respectively.

Moreover, to verify the effectiveness of the $\ell_{\log}$-(pseudo) norm in restricting sparseness of the factor matrices, we compare the LS-NMF with the $\ell_1$-norm approach. The $\ell_1$-norm approach is obtained by replacing the $\ell_{\log}$-(pseudo) norm in (4) with the $\ell_1$-norm, where multiplicative updating rules are used in a way similar to LS-NMF. For both approaches, we show the results on several data sets, including AR, EYaleB, Yale, and ORL, where the parameters are tuned such that best clustering performance is observed. We show how the sparsity of $U$ and $V$ matrices changes with respect to the iterations in Figure 10, where the sparsity follows the definition given in [24]. It is seen that the LS-NMF has superior performance in restricting the sparsity of the basis and representation vectors, which verifies the effectiveness of the proposed approach in learning parts-based representation.

### 5.7. Data Reconstruction

To better understand the factorization results of the proposed method, we further show some examples of reconstructed data by RLS-NMF. Without loss of generality, we show some reconstructed example images from Yale and ORL data sets in Figure 11. The parameters are set to be the ones used in Section 5.1, which leads to the best clustering performance. It is observed that the reconstructed images well capture the key features of original images. Moreover, some outliers like glasses are significantly removed in the reconstructed images. These observations verify the effectiveness of the RLS-NMF in finding parts-based representations.

## 6. Conclusion

In this paper, we propose a new type of sparse NMF methods, including the LS-NMF and RLS-NMF, with the latter being the robust version of the former. The RLS-NMF learns the basis and representation matrices with $\ell_{\log}$-(pseudo) norm, which enhances the sparseness of the learned parts-based representation. Moreover, to enhance the robustness of the RLS-NMF model to noise effects, a noise term is introduced, which is restricted to be example-wisely sparse with the novel $\ell_{2,\log}$-(pseudo) norm. Efficient multiplicative updating rules are developed for optimization, which have theoretical convergence guarantee. Extensive experiments verify the effectiveness of the RLS-NMF in clustering and data representation.

## References

[1] R. O. Duda, P. E. Hart, D. G. Stork, Pattern classification, John Wiley & Sons, 2012.

[2] C. Peng, Q. Cheng, Discriminative ridge machine: A classifier for high-dimensional data or imbalanced data, IEEE Transactions on Neural Networks and Learning Systems (2020) 1–15.

[3] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, Y. Ma, Robust recovery of subspace structures by low-rank representation, Pattern Analysis and Machine Intelligence, IEEE Transactions on 35 (1) (2013) 171–184.

[4] C. Peng, Q. Zhang, Z. Kang, C. Chen, Q. Cheng, Kernel two-dimensional ridge regression for subspace clustering, Pattern Recognition 113 (2021) 107749.

[5] S. E. Palmer, Hierarchical structure in perceptual representation, Cognitive psychology 9 (4) (1977) 441–474.

[6] E. Wachsmuth, M. Oram, D. Perrett, Recognition of objects and their component parts: responses of single units in the temporal cortex of the macaque, Cerebral Cortex 4 (5) (1994) 509–522.

[7] N. K. Logothetis, D. L. Sheinberg, Visual object recognition, Annual review of neuroscience 19 (1) (1996) 577–621.

[8] D. D. Lee, H. S. Seung, Learning the parts of objects by non-negative matrix factorization, Nature 401 (6755) (1999) 788–791.

[9] Yifeng, Pan, Youlian, Liu, Ziying, Multiclass nonnegative matrix factorization for comprehensive feature pattern discovery, Neural Networks & Learning Systems IEEE Transactions on 30 (2) (2019) 615–629.

[10] X. Wang, T. Zhang, X. Gao, Multiview clustering based on non-negative matrix factorization and pairwise measurements, IEEE Transactions on Cybernetics 49 (9) (2019) 3333–3346.

[11] Z. Li, J. Tang, X. He, Robust structured nonnegative matrix factorization for image representation, IEEE Transactions on Neural Networks & Learning Systems PP (5) (2018) 1–14.

[12] J. Yin, L. Gao, Z. Zhang, Scalable distributed nonnegative matrix factorization with block-wise updates, IEEE Transactions on Knowledge & Data Engineering PP (99) (2018) 1–1.

[13] C. Peng, Z. Zhang, C. Chen, Z. Kang, Q. Cheng, Two-dimensional semi-nonnegative matrix factorization for clustering, Information Sciences 590 (2022) 106–141. doi:https://doi.org/10.1016/j.ins.2021.12.098.

[14] M. N. Schmidt, O. Winther, L. K. Hansen, Bayesian non-negative matrix factorization, in: International Conference on Independent Component Analysis and Signal Separation, Springer, 2009, pp. 540–547.

[15] O. Dalhoumi, N. Bouguila, M. Amayri, W. Fan, Bayesian matrix factorization for semibounded data, IEEE Transactions on Neural Networks and Learning Systems (2021) 1–13doi:10.1109/TNNLS.2021.3111824.

[16] S. Z. Li, X. W. Hou, H. J. Zhang, Q. S. Cheng, Learning spatially localized, parts-based representation, in: CVPR 2001. IEEE Conference on, Vol. 1, IEEE, 2001, pp. I–207.

[17] M. Cooper, J. Foote, Summarizing video using non-negative similarity matrix factorization, in: Multimedia Signal Processing, 2002 IEEE Workshop on, IEEE, 2002, pp. 25–28.

[18] W. Xu, X. Liu, 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 informaion retrieval, ACM, 2003, pp. 267–273.

[19] D. D. Lee, H. S. Seung, Algorithms for non-negative matrix factorization, in: Advances in neural information processing systems, 2001, pp. 556–562.

[20] D. Cai, X. He, J. Han, T. S. Huang, Graph regularized nonnegative matrix factorization for data representation, IEEE Transactions on Pattern Analysis and Machine Intelligence 33 (8) (2011) 1548–1560.

[21] J. Huang, F. Nie, H. Huang, C. Ding, Robust manifold nonnegative matrix factorization, ACM Transactions on Knowledge Discovery from Data (TKDD) 8 (3) (2014) 11.

[22] C. H. Ding, T. Li, M. I. Jordan, Convex and semi-nonnegative matrix factorizations, IEEE transactions on pattern analysis and machine intelligence 32 (1) (2010) 45–55.

[23] C. Peng, Z. Zhang, Z. Kang, C. Chen, Q. Cheng, Nonnegative matrix factorization with local similarity learning, Information Sciences 562 (2021) 325–346.

[24] P. O. Hoyer, Non-negative matrix factorization with sparseness constraints, The Journal of Machine Learning Research 5 (2004) 1457–1469.

[25] P. O. Hoyer, Non-negative sparse coding (2002) 557–565.

[26] J.-T. Chien, H.-L. Hsieh, Bayesian group sparse learning for nonnegative matrix factorization, in: Thirteenth Annual Conference of the International Speech Communication Association, 2012.

[27] I. Fedorov, A. Nalci, R. Giri, B. D. Rao, T. Q. Nguyen, H. Garudadri, A unified framework for sparse non-negative least squares using multiplicative updates and the non-negative matrix factorization problem, Signal processing 146 (2018) 79–91.

[28] D. P. Wipf, B. D. Rao, Sparse bayesian learning for basis selection, IEEE Transactions on Signal processing 52 (8) (2004) 2153–2164.

[29] D. Wipf, S. Nagarajan, Iterative reweighted $\ell_1$ and $\ell_2$ methods for finding sparse solutions, IEEE Journal of Selected Topics in Signal Processing 4 (2) (2010) 317–329.

[30] C. Peng, Z. Kang, H. Li, Q. Cheng, Subspace clustering using log-determinant rank approximation, in: Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM, 2015, pp. 925–934.

[31] C. Peng, Y. Chen, Z. Kang, C. Chen, Q. Cheng, Robust principal component analysis: A factorization-based approach with linear complexity, Information Sciences 513 (2019).

[32] C. Peng, Z. Kang, Y. Hu, J. Cheng, Q. Cheng, Robust graph regularized nonnegative matrix factorization for clustering, Acm Transactions on Knowledge Discovery from Data 11 (3) (2017) 33.

[33] F. Nie, H. Huang, X. Cai, C. Ding, Efficient and robust feature selection via joint $\ell_{2,1}$-norms minimization, in: Proceedings of the 23rd International Conference on Neural Information Processing Systems - Volume 2, NIPS'10, Curran Associates Inc., Red Hook, NY, USA, 2010, p. 1813–1821.

[34] H. Xu, C. Caramanis, S. Sanghavi, Robust pca via outlier pursuit, in: Advances in Neural Information Processing Systems, 2010, pp. 2496–2504.

[35] M. McCoy, J. A. Tropp, et al., Two proposals for robust pca using semidefinite programming, Electronic Journal of Statistics 5 (2011) 1123–1160.

[36] F. R. Chung, Spectral graph theory, Vol. 92, American Mathematical Soc., 1997.

[37] E. J. Candes, M. B. Wakin, S. P. Boyd, Enhancing sparsity by reweighted L1 minimization, Journal of Fourier analysis and applications 14 (5) (2008) 877–905.

[38] L. Xu, M. I. Jordan, On convergence properties of the em algorithm for gaussian mixtures, Neural computation 8 (1) (1996) 129–151.

[39] Y.-D. Kim, S. Choi, Weighted nonnegative matrix factorization, in: Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on, IEEE, 2009, pp. 1541–1544.

[40] C. Ding, T. Li, W. Peng, H. Park, Orthogonal nonnegative matrix t-factorizations for clustering, in: Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, 2006, pp. 126–135.

[41] P. N. Belhumeur, J. P. Hespanha, D. J. Kriegman, Eigenfaces vs. fisherfaces: Recognition using class specific linear projection, IEEE Transactions on pattern analysis and machine intelligence 19 (7) (1997) 711–720.

[42] M. J. Lyons, S. Akamatsu, M. Kamachi, J. Gyoba, J. Budynek, The japanese female facial expression (jaffe) database (1998).

[43] F. S. Samaria, A. C. Harter, Parameterisation of a stochastic model for human face identification, in: Applications of Computer Vision, 1994., Proceedings of the Second IEEE Workshop on, IEEE, 1994, pp. 138–142.

[44] A. M. Martinez, The ar face database, CVC Technical Report 24 (1998).

[45] A. S. Georghiades, P. N. Belhumeur, D. J. Kriegman, From few to many: Illumination cone models for face recognition under variable lighting and pose, IEEE transactions on pattern analysis and machine intelligence 23 (6) (2001) 643–660.

[46] S. A. Nene, Columbia object image library(coil-20), Technical Report 5 (1996).

[47] K. Bache, M. Lichman, UCI machine learning repository (2013).
URL http://archive.ics.uci.edu/ml

[48]  Tactile, Semeion data set, Semeion Research Center of Sciences of Communication, via Sersale 117, 00128 Rome, Italy. Tattile Via Gaetano Donizetti, 1-3-5,25030 Mairano (Brescia), Italy.