

Spectral Unmixing via Data-guided Sparsity

Feiyun Zhu, Ying Wang, Bin Fan, Shiming Xiang, Gaofeng Meng and Chunhong Pan

Abstract—Hyperspectral unmixing, the process of estimating a common set of spectral bases and their corresponding composite percentages at each pixel, is an important task for hyperspectral analysis, visualization and understanding. From an unsupervised learning perspective, this problem is very challenging—both the spectral bases and their composite percentages are unknown, making the solution space too large. To reduce the solution space, many approaches have been proposed by exploiting various priors. In practice, these priors would easily lead to some unsuitable solution. This is because they are achieved by applying an identical strength of constraints to all the factors, which does not hold in practice. To overcome this limitation, we propose a novel sparsity based method by learning a data-guided map to describe the individual mixed level of each pixel. Through this data-guided map, the ℓ_p ($0 < p < 1$) constraint is applied in an adaptive manner. Such implementation not only meets the practical situation, but also guides the spectral bases toward the pixels under highly sparse constraint. What’s more, an elegant optimization scheme as well as its convergence proof have been provided in this paper. Extensive experiments on several datasets also demonstrate that the data-guided map is feasible, and high quality unmixing results could be obtained by our method.

Index Terms—Data-guided Sparse (DgS), Data-guided Map (DgMap), Nonnegative Matrix Factorization (NMF), DgS-NMF, Mixed Pixel, Hyperspectral Unmixing (HU).

I. INTRODUCTION

HYPERSPECTRAL imaging, the process of capturing a 3D image cube at hundreds of contiguous and narrow spectral channels, has been used in a wide range of fields [1], [2]. Although this type of images contains substantial information, there are two underlying “problems”. One “problem” is that as a 3D image cube, it is very hard for computers to display [3], thus hampering human to understand this type of images. Another “problem” is called “mixed” pixels—due to the low spatial resolution of hyperspectral sensors, the spectra of different substances would unavoidably blend together [1], [4], [5], yielding a great number of mixed pixels as shown in Fig. 1. To address the above two problems, various Hyperspectral Unmixing (HU) methods have been proposed. What is more, HU is essential for various hyperspectral applications, such as sub-pixel mapping [6], hyperspectral enhancement [7], high-resolution hyperspectral imaging [8], detection and identification of ground targets [9].

Formally, the HU method takes in a hyperspectral image with L channels and assumes that each pixel spectrum \mathbf{y} is a composite of K spectral bases $\{\mathbf{m}_k\}_{k=1}^K \in \mathbb{R}_+^L$ [2], [10], [11]. Each spectral base is called an *endmember*, representing the pure spectrum, such as the spectra of “water”, “grass” etc.

Feiyun Zhu, Ying Wang, Bin Fan, Shiming Xiang, Gaofeng Meng and Chunhong Pan are with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences (e-mail: {fyzhu, ywang, bfan, smxiang, gfmeng and chpan}@nlpr.ia.ac.cn).

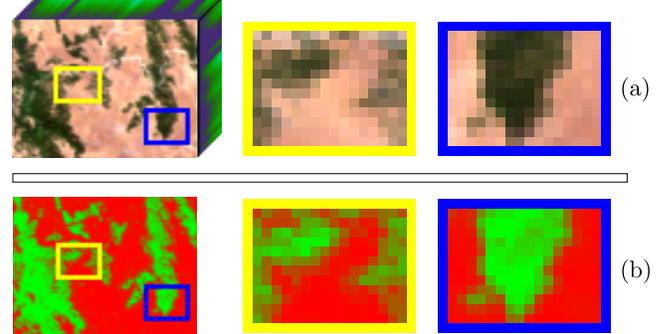


Figure 1. Two observations behind the figure: the mixed level of each pixel varies over image grids; the pixels in the transition area are more likely to be highly mixed. (a) Hyperspectral image and its close-ups. (b) Abundances of two substances in (a), indicated by the proportions of red and green inks.

Specifically, the pixel spectrum \mathbf{y} is generally approximated by a nonnegative linear combination as

$$\mathbf{y} = \sum_{k=1}^K \mathbf{m}_k a_k, \quad \text{s.t. } a_k \geq 0 \text{ and } \sum_{k=1}^K a_k = 1, \quad (1)$$

where a_k is the composite percentage (i.e. *abundance*) of the k^{th} *endmember*. In the unsupervised setting, both *endmembers* $\{\mathbf{m}_k\}_{k=1}^K$ and *abundances* $\{a_k\}_{k=1}^K$ are unknown. Such case makes the solution space really large [12]. Prior knowledge is required to restrict the solution space, or even to bias the solution toward good results.

To shrink the solution space, many methods have been proposed by exploiting various constraints on *abundances* [9], [11] and *endmembers* [13], [14]. Specifically, the sparse constraints [9], [12] and the spatial constraints [11], [15] are the most popular ones. Unfortunately, all these methods exploit an identical strength of constraints on all the factors, which may not meet the practical situation. An example is illustrated in Fig. 1, where the mixed level¹ of each pixel varies over image grids. Such an example indicates that it is better to impose the sparse constraint of adaptive strengths for the pixels.

In this paper, we propose a Data-guided Sparsity regularized Nonnegative Matrix Factorization (DgS-NMF) method for the HU task. The basic motivation is that the mixed level of each pixel might be different from each other, as shown in Fig. 1. To give a more accurate model, a data-guided map (DgMap) is incorporated into the NMF framework so as to adaptively impose the sparse constraint for each pixel. First, via a two step strategy, the DgMap is learned from the hyperspectral image, describing the mixed level of each pixel. Given this DgMap, the ℓ_p ($0 < p < 1$)-norm based sparsity constraint is individually imposed. For each pixel, the choice of p is totally

¹Note that a pixel with higher mixed levels should own the *abundance* vectors of lower sparse levels, and vice versa.

dependent on the corresponding DgMap value. Such case is better suited to the practical situation, thus expected to achieve better HU results. Besides, this adaptive sparsity constraint would influence the estimation of *endmembers*, potentially, guiding the *endmembers* toward the pixels under highly sparse constraints. Extensive empirical results verify that our method is highly promising for the HU task.

The rest of this paper is organized as follows: in Section II, we briefly review several recent HU methods. Section III presents how to learn DgMaps from the original hyperspectral image cube. The DgS-NMF method as well as its properties are given in Section IV. Then, extensive experiments and detailed comparisons are provided in Section V. Finally, the conclusion of this work is drawn in Section VI.

II. PREVIOUS WORK

The HU methods could be typically categorized into two types: geometric methods [16], [17], [18] and statistical ones [10], [19], [20], [21]. Usually, the geometric methods utilize a simplex to describe the distribution of hyperspectral pixels. The vertices of this simplex are viewed as the *endmembers*. Perhaps, N-FINDR [22] and Vertex Component Analysis (VCA) [16] are the most popular geometric methods. In N-FINDR, the *endmembers* are identified by inflating a simplex inside the hyperspectral pixel distribution and treating the vertices of a simplex with the largest volume as *endmembers* [22]. While VCA [16] projects all pixels onto a direction orthogonal to the simplex spanned by the chosen *endmembers*; the new *endmember* is identified as the extreme of the projection. Although these methods are simple and fast, they suffer from the requirement of pure pixels for each *endmember*, which is usually unavailable in practice.

Accordingly, a number of statistical methods have been proposed for or applied to the HU task, among which the Non-negative Matrix Factorization (NMF) [23] and its variants are the most popular ones. As an unsupervised method, the goal of NMF is to find two nonnegative matrices to approximate the original matrix with their product [24]. Specifically, the nonnegative constraint on the two factor matrices only allows additive combinations, not subtractions, resulting in a parts-based representation. This parts-based property could ensure the representation results to be more intuitive and interpretable, since psychological and physiological evidences have shown that human brain works in a parts-based way [25], [26].

Although the NMF method is well suited to many applications, such as face analysis [27], [28] and documents clustering [29], [30], the objective function of NMF is non-convex, inherently resulting in large solution space [31]. Many extensions have been proposed by exploiting various priors to restrict the solution space. For the HU problem, these priors are either imposed to the *abundance* matrix or to the *endmember* matrix. For example, the Local Neighborhood Weights regularized NMF method (W-NMF) [15] assumes that the hyperspectral pixels are on a manifold structure, which could be transferred to the *abundance* space through a Laplace graph constraint. Actually, this constraint has a smooth influence, and eventually weaken the parts-based property of NMF.

Inspired by the MVC-NMF [13] method, Wan et al. [14] proposed the EDC-NMF method. The basic assumption is that due to the high spectral resolution of sensors, the *endmember* spectra should be smooth itself and different as much as possible from each other. However, in their algorithm, they take a derivative of *endmembers*, introducing negative values to the updating rules. To make up this drawback, the elements in the *endmember* matrix are required to project to a given nonnegative value after each iteration. Consequently, the regularization parameters could not be chosen freely, limiting the efficacy of this method.

Other algorithms assume that in hyperspectral images most pixels are mixed by only a few *endmembers*, hence exploiting various kinds of sparse constraints on the *abundance* [4]. Specifically, the $\ell_{1/2}$ -NMF [9] is a very popular sparsity regularized NMF method. It is an improvement from Hoyer's lasso regularized NMF method [32]. There are two advantages of the $\ell_{1/2}$ -NMF over the lasso regularized NMF. One advantage is that the lasso constraint [33], [34] could not enforce further sparse when the full additivity constraint is used, limiting the effectiveness of this method [9]. Another advantage is that Fan et al. [35] has proven that the ℓ_p ($0 < p < 1$) constraint could obtain sparser solutions than the ℓ_1 norm does.

Our method is also derived from the sparse assumption on the *abundance*. Different from the existing methods, the strength of sparse constraints is learned from the data itself and applied in an adaptive way. Such improvement not only meets the practical situation better, but also help the optimization process to reach a more suitable local minimum.

III. DATA-GUIDED MAP (DGMAP)

Generally, the Data-guided Map (DgMap) is a map learnt from the hyperspectral image that describes the strength of priors (constraints) for each factor. In this work, the DgMap depicts the mixed level of each pixel. It comes from two observations that: 1) in the local image window, the mixed level of each pixel might be more or less different from each other as shown in Fig. 1; 2) in the whole image, the pixels in the transition area are very likely to be highly mixed (c.f. Footnote 1). For the second idea, Fig. 1 illustrates an example, where there are two targets (i.e. "tree" and "soil") in the scene. The pixels in the transition area are very likely to be mixed by spatially neighboring pixels from these two targets, thus, yielding a great number of mixed pixels. Therefore, these pixels in the transition area should receive weaker sparse constraint than pixels in the other areas. In the following, we would elaborate how to learn such a DgMap from the hyperspectral image via a two step strategy.

A. Initial Data-guided Map

Suppose we are given a hyperspectral image $\{\mathbf{y}_n\}_{n=1}^N \in \mathbb{R}_+^{L \times N}$ with N pixels and L channels. It is reasonable to assume that the pixels in the transition area are more or less different from their spatial neighbors as shown in Fig. 1. For this reason, the initial DgMap $\mathbf{h}^{(0)} \in \mathbb{R}_+^N$ could be learnt by measuring the uniformity of neighboring pixels over the entire image, i.e. $\mathbf{h}^{(0)} = f(\mathbf{y}_1, \dots, \mathbf{y}_N)$. In this way, the inhomogeneous areas

are treated as the transition ones. For the i^{th} pixel, its value in the DgMap could be estimated by measuring the similarity between spatially neighboring pixels as follows:

$$h_i^{(0)} = \sum_{j \in \mathcal{N}_i} s_{ij}, \quad (2)$$

where \mathcal{N}_i is the neighborhood of the i^{th} pixel that includes four neighbors; s_{ij} is the similarity between the i^{th} pixel and its neighboring pixel \mathbf{y}_j by the dot-product metric

$$s_{ij} = \frac{\mathbf{y}_i^T \mathbf{y}_j}{\|\mathbf{y}_i\| \cdot \|\mathbf{y}_j\|}, \quad (3)$$

which is a classic measure in the HU study [9], [12], or by the heat kernel similarity metric

$$s_{ij} = \exp\left(-\frac{\|\mathbf{y}_j - \mathbf{y}_i\|_2^2}{\sigma}\right). \quad (4)$$

The value of σ controls the contrast of DgMaps. Generally, a smaller σ results in a DgMap with higher contrast. In all the experiments, σ is set as $\sigma \in [0.005, 0.08]$.

To evaluate the effectiveness of the definition of DgMaps in (2), we collect a set of 36 hyperspectral images² and calculate their DgMaps according to the similarity measures in (3) and (4). The results are plotted in Fig. 2, showing the histogram of DgMap values over all the 36 images. As Fig. 2a shows, it is of little effect by using the dot-product measure—up to 99.61% guided values are located in a narrow range of [0.95, 1]. Such case suggests that almost all the pixels have similar DgMap values, lacking of guided information. Contrarily, the DgMaps from the heat kernel measure contain more information, as shown in Figs. 2b. Therefore, we choose the heat kernel measure to learn the initial DgMap.

B. Fine Tuned Data-guided Map

Through the local uniformity assumption and the heat kernel measure, the learned DgMap does not have the global consistency over the entire image. Therefore, we further propose a fine tuning step to refine the initial DgMap. For this purpose, the closed-form method [36], [37] is adopted. The advantages are in two folds: 1) this fine tuning process not only propagates the guidance information over the entire image, but also maintains the structures latent in the original hyperspectral data cube [38], [39]; 2) according to our experiments, the fine tuned DgMap further improves the HU performance although the initial DgMap already outperforms the state-of-the-art.

Specifically, the closed-form method is based on the assumption that in each small window, the data-guided values come from the same projection as [36]:

$$h_j = \mathbf{w}_i^T \mathbf{y}_j + b_i, \quad \text{for } j \in \mathcal{N}_i, i \in \{1, \dots, N\},$$

where \mathbf{y}_i is the i^{th} pixel; \mathcal{N}_i is the neighborhood of \mathbf{y}_i ³; \mathbf{w}_i is the projection vector and b_i is a bias term. The local adjustment

²This image set includes hyperspectral scenes of urban areas, suburbs areas, farmland areas, mine areas, airports and so on. In average, there are 369×369 pixels and 193 channels in a hyperspectral image.

³Note that the neighborhood \mathcal{N}_i defined here is different from the neighborhood \mathcal{N}_i used in Section III-A.

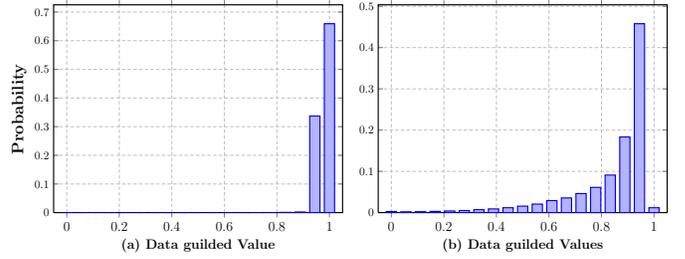


Figure 2. Histogram of DgMap values learnt from 36 hyperspectral images respectively by: (a) dot-product measure (3) and (b) heat kernel measure (4).

from the initial DgMap $\{h_j^{(0)}\}_{j=1}^N$ is formulated as

$$h_j^* \leftarrow \arg \min_{h_j} \left(\alpha \left(h_j - h_j^{(0)} \right)^2 + \left(h_j - \mathbf{w}_i^T \mathbf{y}_j - b_i \right)^2 \right), \\ \forall j \in \mathcal{N}_i, i \in \{1, 2, \dots, N\}.$$

The local window is placed in an overlapping manner. This case ensures the property of propagating guidance information between neighboring pixels [36]. The bigger the local window is, the wider the propagation could spread. Besides, in each local window, the gradient field of the DgMap is linearly related to the corresponding image gradient field as $\nabla h_j = \mathbf{w}_i^T \nabla \mathbf{y}_j, \forall j \in \mathcal{N}_i, i \in \{1, \dots, N\}$, transferring the gradient distributions as well as the transition information latent in the original hyperspectral image into the newly learnt DgMap [40], [39]. As a result, we could refine the DgMap according to the structures latent in the original image cube.

Considering all the local minimizing problems together as well as the numerical stability, we can fine tune the DgMap by minimizing the following quadratic function [40], [36]:

$$E(\mathbf{h}, \mathbf{w}, b) = \alpha \left\| \mathbf{h} - \mathbf{h}^{(0)} \right\|_2^2 + \sum_{i=1}^N \left(\epsilon \|\mathbf{w}_i\|_2^2 + \sum_{j \in \mathcal{N}_i} \left(h_j - \mathbf{w}_i^T \mathbf{y}_j - b_i \right)^2 \right), \quad (5)$$

where $\mathbf{h}^{(0)} = [h_1^{(0)}, \dots, h_N^{(0)}]^T \in \mathbb{R}_+^N$ is the initial DgMap; $\epsilon \in [10^{-7}, 10^{-4}]$ controls the smooth level of the refined DgMap \mathbf{h} ; $\alpha \in [10^{-6}, 10^{-4}]$ controls the strength of the fine tuning process. A smaller α corresponds to a stronger refinement.

The objective function above could be further simplified by setting $\frac{\partial E}{\partial \mathbf{w}} = 0, \frac{\partial E}{\partial b} = 0$ and substituting their solutions into (5), yielding a compact objective function as:

$$E(\mathbf{h}) = \alpha \left\| \mathbf{h} - \mathbf{h}^{(0)} \right\|_2^2 + \mathbf{h}^T \mathbf{L} \mathbf{h}, \quad (6)$$

where \mathbf{L} is a highly sparse matrix that has been proven to be a graph Laplacian by [41]. It is defined as

$$\mathbf{L} = \sum_{n=1}^N \mathbf{S}_n^T \mathbf{L}_i \mathbf{S}_n, \quad (7)$$

where \mathbf{S}_i^T is the i^{th} column selection matrix that selects the pixels in the i^{th} local window from the whole hyperspectral image as $\mathbf{Y}_i = \mathbf{Y} \mathbf{S}_i^T, \mathbf{L}_i = \mathbf{G}_i \mathbf{G}_i$, in which $\mathbf{G}_i = \left(\mathbf{P} - \bar{\mathbf{Y}}_i^T (\bar{\mathbf{Y}}_i \bar{\mathbf{Y}}_i^T + \epsilon \mathbf{I})^{-1} \bar{\mathbf{Y}}_i \right), \mathbf{P} = \mathbf{I} - \frac{1}{|\mathcal{N}_i|} \mathbf{1} \mathbf{1}^T$ is

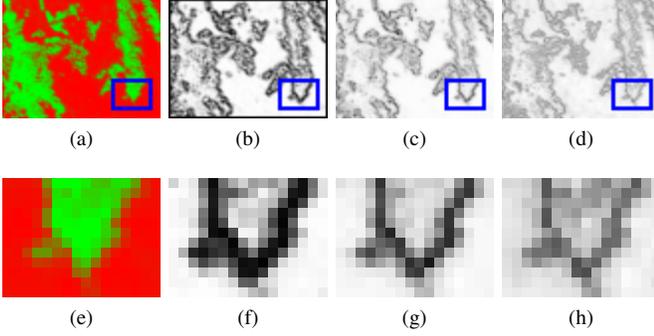


Figure 3. (a) The *abundance* map of the hyperspectral image in Fig. 1, where the proportions of red and green inks represent the *abundances* of two targets. (b) The initial DgMap from the heat kernel measure. (c) Fine tuned DgMap learnt by 3×3 window. (d) Fine tuned DgMap learnt by 7×7 window. (e)-(h) are the close ups of (a)-(d) respectively. (Best viewed in color)

the centering matrix with $|\mathcal{N}_i| \times |\mathcal{N}_i|$ elements and $\bar{\mathbf{Y}}_i = \mathbf{Y}_i \mathbf{P}$ contains the zero mean pixels in \mathcal{N}_i [42], [43].

Since the objective function (6) is quadratic in \mathbf{h} , it can be solved by setting the derivative to zero as $\nabla E(\mathbf{h}) = \mathbf{0}$, yielding a highly sparse linear equation

$$(\mathbf{L} + \alpha \mathbf{I}) \mathbf{h} = \alpha \mathbf{h}^{(0)}, \quad (8)$$

which could be efficiently solved [36]. In order to simplify the incorporation of the learnt DgMap into the ℓ_p ($0 < p < 1$) norm, the data-guided values are resized into the range of (0, 1) as

$$h_n \leftarrow \frac{h_n - \min(\mathbf{h})}{\max(\mathbf{h}) - \min(\mathbf{h}) + \beta}, \quad n = 1, 2, \dots, N,$$

where $\beta = 10^{-8}$ is a small value used to prevent $\{h_n\}_{n=1}^N$ being equal to 1, ensuring the numerical stability for the sparse constraint in the next section.

To study the influence of the local window size, we conduct an experiment as illustrated in Fig. 3, where Fig. 3a shows the reference *abundance* map, Fig. 3b shows the initial DgMap, Fig. 3c shows the refined DgMap with 3×3 local window, followed by the fine tuned result with 7×7 local window in Fig. 3d. As Fig. 3 shows, the 3×3 local window is sufficient to get suitable result at low computational costs. Therefore, the 3×3 local window size is chosen in this work.

IV. DATA-GUIDED SPARSE NMF (DGS-NMF)

A. Data-guided Regularization and DgS-NMF Model

Based on the linear combination model in (1), a hyperspectral image $\mathbf{Y} \triangleq [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N] \in \mathbb{R}_+^{L \times N}$, with L channels and N pixels, could be approximated by two factor matrices:

$$\mathbf{Y} = \mathbf{M}\mathbf{A} + \mathbf{E}, \quad (9)$$

where $\mathbf{M} \triangleq [\mathbf{m}_1, \dots, \mathbf{m}_K] \in \mathbb{R}_+^{L \times K}$ is the *endmember* matrix including K spectral vectors, $K \ll \min\{L, N\}$; $\mathbf{A} \triangleq [\mathbf{a}_1, \dots, \mathbf{a}_N] \in \mathbb{R}_+^{K \times N}$ is the corresponding *abundance* matrix, whose n^{th} column vector \mathbf{a}_n contains all the K *abundances* at pixel \mathbf{y}_n ; \mathbf{E} is a residual term. Specifically, (9) could be naturally translated into the Nonnegative Matrix Factorization [23] (NMF) problem by strictly constraining the

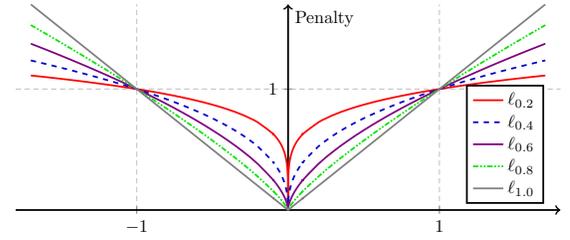


Figure 4. The shape of the ℓ_p -norm with different $p \in (0, 1]$, indicating that a smaller p tends to find a sparser solution [35].

nonnegative property of both factors, i.e. $\mathbf{M} \geq \mathbf{0}$, $\mathbf{A} \geq \mathbf{0}$, which agrees with the nonnegative requirement on both *endmembers* and *abundances*. Such case suggests that NMF is physically suitable for the HU task.

Suppose we are given the fine tuned DgMap $\mathbf{h} \in \mathbb{R}_+^N$. Different from the traditional sparse regularization [32], [9] that constrains all factors $\{\mathbf{a}_n\}_{n=1}^N$ at the same sparse level as

$$\mathcal{J}(\mathbf{A}) = \sum_{n=1}^N \|\mathbf{a}_n\|_1, \quad \mathcal{J}(\mathbf{A}) = \sum_{n=1}^N \|\mathbf{a}_n\|_{1/2}^{1/2}, \quad (10)$$

where $\|\mathbf{a}\|_p^p = \sum_k |a_k|^p$ ($\forall 0 < p < 1$), this paper proposes a novel data-guided constraint as

$$\mathcal{J}(\mathbf{A}) = \sum_{n=1}^N \|\mathbf{a}_n\|_{1-h_n}^{1-h_n} = \sum_{n=1}^N \left(\sum_{k=1}^K |A_{kn}|^{1-H_{kn}} \right), \quad (11)$$

where H_{kn} is the (k, n) -th element in the matrix $\mathbf{H} = \mathbf{1}_K \mathbf{h}^T$. The elements in the same column of \mathbf{H} are identical to each other, i.e. $H_{1,n} = H_{2,n} \dots = H_{K,n} = h_n, \forall n \in \{1, \dots, N\}$.

In this way, all the *abundance* factors $\{\mathbf{a}_n\}_{n=1}^N$ are constrained in the ℓ_p ($0 < p < 1$)-norm. For each factor, the level (strength) of sparse constraint is closely related to the choice of p —a smaller p corresponds to a sparser constraint [35] (cf. Fig. 4). This amounts to the dependence on the DgMap value h_n , as shown in (11). So, for instance, a pixel for which $h_n = 0.2$ will be constrained by a weak sparsity regularization in the $\ell_{0.8}$ -norm, whereas one for which $h_n = 0.8$ will be constrained by the $\ell_{0.2}$ regularization and so will enjoy a heavy sparsity constraint (cf. Fig. 4). Additionally, the DgMap values in the transition areas are generally small (cf. Fig. 3c). As a result, they will be constrained at relatively low levels of sparsity constraints, conforming to their mixed properties.

Compared with the traditional regularization (10), the advantages of our constraint (11) lie in three aspects: 1) as the fine tuned DgMap describes the mixed level over the entire image, our constraint is more agreeable with the practical mixed property of each pixel; 2) with the careful constraint in (11), the non-convex objective function (12) is more likely to converge to some suitable local minima; 3) although the adaptive sparsity regularization is constrained on the *abundance* factors, it would explicitly influence the estimation of *endmembers*, guiding the *endmembers* toward the pixels with highly sparse constraint. This doesn't mean that the pixels with highly sparse constraints are *endmembers*. Many pixels with highly sparse constraints compete for the *endmember*, and some trade-off spectra could also be *endmembers*.

Apart from the advantages above, it is easy to find that the traditional sparse constraints (10) are special cases of our adaptive sparse constraint (11). Given a constant DgMap with each pixel $\{h_n\}_{n=1}^N$ equal to zero, the adaptive sparse constraint degrades into the ℓ_1 regularization, i.e. $\mathcal{J}(\mathbf{A}) = \sum_n \|\mathbf{a}_n\|_{1-h_n}^{1-h_n} = \sum_n \|\mathbf{a}_n\|_1$; whereas if each element in the DgMap is equal to 1/2, the adaptive sparse constraint turns into the $\ell_{1/2}$ regularization. Moreover, for the HU task, all the elements in \mathbf{A} are within the range of (0, 1) [12], [9]. Thus, once allowed the limit $h_n \rightarrow -\infty$, the adaptive sparse constraint would degrade into the non-regularization case, i.e. $\mathcal{J}(\mathbf{A}) \rightarrow 0$, since for any $a \in (0, 1)$, we have $a^{1+\infty} \rightarrow 0$.

To obtain the optimal factor matrices, we model the matrix representation problem (9) as the Data-guided Sparsity regularized Nonnegative Matrix Factorization (DgS-NMF) objective

$$\begin{aligned} \mathcal{O}(\mathbf{M}, \mathbf{A}) &= \frac{1}{2} \|\mathbf{Y} - \mathbf{MA}\|_F^2 + \lambda \sum_{n=1}^N \sum_{k=1}^K |A_{kn}|^{1-H_{kn}} \quad (12) \\ \text{s.t. } &\mathbf{M} \geq \mathbf{0}, \mathbf{A} \geq \mathbf{0}, \end{aligned}$$

where $\lambda \geq 0$ is a balancing parameter that controls the average sparsity of the factor representation. In the next subsection, the optimization for the DgS-NMF problem will be analyzed.

B. Updating Rules for DgS-NMF

Akin to NMF [31] and EM [44], the objective function in (12) is non-convex for \mathbf{M} and \mathbf{A} together. No global minima could be reached. Alternatively, we propose an iterative algorithm that alternately updates \mathbf{M} and \mathbf{A} at each iteration. It has the ability to arrive at some local minima after finite iterations, which will be proved in Section IV-C.

Specifically, the Lipschitz constant [45] of the data-guided constraint (11) will be infinity for $A_{kn} = 0, \forall k, n$. To ensure the Lipschitz condition, we reformulate our model (12) as

$$\begin{aligned} \mathcal{O}(\mathbf{M}, \mathbf{A}) &= \frac{1}{2} \|\mathbf{Y} - \mathbf{MA}\|_F^2 + \lambda \sum_{n=1}^N \sum_{k=1}^K (A_{kn} + \xi)^{1-H_{kn}} \\ \text{s.t. } &\mathbf{M} \geq \mathbf{0}, \mathbf{A} \geq \mathbf{0}, \quad (13) \end{aligned}$$

where ξ is a small positive value to ensure the numerical condition. It is obvious that the objective (13) is reduced to (12) when $\xi \rightarrow 0$. For simplicity, we use $\mathbf{A} + \xi = [A_{kn} + \xi]$ to express the idea of adding ξ to every entry $A_{kn}, \forall k, n$.

Considering the constraints of $\mathbf{M} \geq \mathbf{0}, \mathbf{A} \geq \mathbf{0}$, the objective function (13) could be rewritten as the Lagrange Multiplier:

$$\begin{aligned} \mathcal{L} &= \frac{1}{2} \|\mathbf{Y} - \mathbf{MA}\|_F^2 + \lambda \sum_{n=1}^N \sum_{k=1}^K (A_{kn} + \xi)^{1-H_{kn}} \\ &\quad + \text{Tr}(\Psi \mathbf{M}^T) + \text{Tr}(\Gamma \mathbf{A}^T), \quad (14) \end{aligned}$$

where ψ_{lk}, γ_{kn} are the lagrange multipliers for the inequality constraints $M_{lk} \geq 0$ and $A_{kn} \geq 0$ respectively, and $\Psi = [\psi_{lk}] \in \mathbb{R}_+^{L \times K}$, $\Gamma = [\gamma_{kn}] \in \mathbb{R}_+^{K \times N}$ are the lagrange multipliers in matrix format. To find the local minima, one intuitive approach is to differentiate (14) and set the partial derivatives to zero. This amounts to solving the following linear equations

$$\nabla_{\mathbf{M}} \mathcal{L} = \mathbf{MAA}^T - \mathbf{YA}^T + \Psi = \mathbf{0} \quad (15)$$

Algorithm 1 for DgS-NMF

Input: the hyperspectral image $\mathbf{Y} \in \mathbb{R}_+^{L \times N}$, the number of *endmembers* (i.e. K) and the penalty parameters λ .

Output: two factor matrices $\mathbf{M} \in \mathbb{R}_+^{L \times K}$ and $\mathbf{A} \in \mathbb{R}_+^{K \times N}$.

- 1: Calculate initial DgMap $\mathbf{h}^{(0)} \in \mathbb{R}_+^N$ according to Eq. (2).
 - 2: Get the fine tuned DgMap \mathbf{h} by solving the highly sparse linear equation (8). Calculate $\mathbf{H} = \mathbf{1}_K \mathbf{h}^T \in \mathbb{R}^{K \times N}$.
 - 3: Initialize the factor matrices \mathbf{M} and \mathbf{A} .
 - 4: **repeat**
 - 5: update \mathbf{A} by the updating rule (20).
 - 6: update \mathbf{M} by the updating rule (19).
 - 7: scale \mathbf{M} and \mathbf{A} by Eq. (21) after each iteration.
 - 8: **until** convergence
 - 9: Output \mathbf{M} and \mathbf{A} as the final unmixing result.
-

$$\begin{aligned} \nabla_{\mathbf{A}} \mathcal{L} &= \mathbf{M}^T \mathbf{MA} - \mathbf{M}^T \mathbf{Y} + \Gamma + \\ &\quad \lambda (\mathbf{1} - \mathbf{H}) \circ (\mathbf{A} + \xi)^{-\mathbf{H}} = \mathbf{0}, \quad (16) \end{aligned}$$

where \circ is the Hadamard product between matrices; $\mathbf{A}^{\mathbf{H}} = \left[(A_{kn})^{H_{kn}} \right] \in \mathbb{R}_+^{K \times N}$ is an elementwise exponential operation. Based on the Karush-Kuhn-Tucker conditions $\psi_{lk} M_{lk} = 0$ and $\gamma_{kn} A_{kn} = 0$, we could simplify (15) and (16) by multiplying both sides with M_{lk} and A_{kn} respectively, yielding

$$(\mathbf{MAA}^T)_{lk} M_{lk} - (\mathbf{YA}^T)_{lk} M_{lk} = 0 \quad (17)$$

$$\begin{aligned} (\mathbf{M}^T \mathbf{MA})_{kn} A_{kn} - (\mathbf{M}^T \mathbf{Y})_{kn} A_{kn} + \\ \lambda \left((\mathbf{1} - \mathbf{H}) \circ (\mathbf{A} + \xi)^{-\mathbf{H}} \right)_{kn} A_{kn} = 0. \quad (18) \end{aligned}$$

Solving Eqs. (17) and (18), we get the updating rules as

$$M_{lk} \leftarrow M_{lk} \frac{(\mathbf{YA}^T)_{lk}}{(\mathbf{MAA}^T)_{lk}} \quad (19)$$

$$A_{kn} \leftarrow A_{kn} \frac{(\mathbf{M}^T \mathbf{Y})_{kn}}{\left(\mathbf{M}^T \mathbf{MA} + \lambda (\mathbf{1} - \mathbf{H}) \circ (\mathbf{A} + \xi)^{-\mathbf{H}} \right)_{kn}}. \quad (20)$$

However, if \mathbf{M} and \mathbf{A} form the solution of NMF, \mathbf{DU} and $\mathbf{U}^{-1} \mathbf{A}$ are the solution for any positive diagonal matrix \mathbf{U} [24], [4]. To get rid of this kind of uncertainty, one intuitive method is to scale each row of \mathbf{A} or each column of \mathbf{M} to be unit ℓ_1 -norm or ℓ_2 -norm [29] as follows

$$M_{lk} \leftarrow M_{lk} \left(\sum_{n=1}^N |A_{kn}| \right), \quad A_{kn} \leftarrow \frac{A_{kn}}{\sum_{n=1}^N |A_{kn}|}. \quad (21)$$

Similarly, we scale \mathbf{M} and \mathbf{A} by (21) after each iteration.

The algorithm for DgS-NMF is summarized in Algorithm 1. For the updating rules in (19) and (20), we have the following theorem, which will be proven in the next section, as

Theorem 1. *The objective function (12) is non-increasing under the updating rules (19) and (20).*

C. Convergence Proof for DgS-NMF

To ensure the reliability of (19) and (20), the convergence proofs of both updating rules are discussed. Fortunately, the convergence proof of (19) could be eliminated since it has been

analyzed in [31]. A common skill used in EM [44], [46] and NMF [31] is employed by introducing an auxiliary function:

Definition 2. $G(\mathbf{A}, \mathbf{A}')$ is an auxiliary function of $\mathcal{O}(\mathbf{A})$ if the following properties are satisfied,

$$G(\mathbf{A}, \mathbf{A}') \geq \mathcal{O}(\mathbf{A}), \quad G(\mathbf{A}, \mathbf{A}) = \mathcal{O}(\mathbf{A}). \quad (22)$$

Lemma 3. By minimizing the energy of $G(\mathbf{A}, \mathbf{A}')$ given by

$$\mathbf{A}^{(t+1)} = \arg \min_{\mathbf{A}} G(\mathbf{A}, \mathbf{A}^{(t)}),$$

we can obtain a solution $\mathbf{A}^{(t+1)}$ that makes $\mathcal{O}(\mathbf{A})$ non-increasing at each iteration, i.e. $\mathcal{O}(\mathbf{A}^{(t+1)}) \leq \mathcal{O}(\mathbf{A}^{(t)})$. Finally, $\mathcal{O}(\mathbf{A})$ will converge after finite iterations.

Proof: This is because of the following inequalities:

$$\begin{aligned} \mathcal{O}(\mathbf{A}^{(\min)}) &\leq \dots \leq \mathcal{O}(\mathbf{A}^{(t+1)}) \leq G(\mathbf{A}^{(t+1)}, \mathbf{A}^{(t)}) \\ &\leq \mathcal{O}(\mathbf{A}^{(t)}) \leq \dots \leq \mathcal{O}(\mathbf{A}^{(0)}) \end{aligned}$$

Now we consider the objective function (13) with \mathbf{A} as the only variable:

$$\mathcal{O}(\mathbf{A}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{M}\mathbf{A}\|_F^2 + \lambda \sum_{n=1}^N \sum_{k=1}^K (A_{kn}^{(t)} + \xi)^{1-H_{kn}}. \quad (23)$$

Specifically, it is approximately a quadratic function as follows

$$\begin{aligned} \mathcal{O}(\mathbf{A}) &\approx \mathcal{O}(\mathbf{A}^{(t)}) + \text{Tr}(\mathbf{C}^T \nabla \mathcal{O}(\mathbf{A}^{(t)})) \\ &\quad + \frac{1}{2} [\text{Tr}(\mathbf{C}^T (\mathbf{M}^T \mathbf{M}) \mathbf{C}) - \lambda F(\mathbf{A})], \end{aligned} \quad (24)$$

where $\mathbf{C} = (\mathbf{A} - \mathbf{A}^{(t)})$ and

$$F(\mathbf{A}) = \sum_{n,k} H_{kn} (1 - H_{kn}) (A_{kn}^{(t)} + \xi)^{-(H_{kn}+1)} C_{kn}^2.$$

To prove the convergence property of (20), we have to find an auxiliary function of (24), by which the updating rule (20) could be obtained by differentiating this auxiliary function and setting the derivatives to zero. Conversely, a function constituted based on the updating rule (20) is given by

$$\begin{aligned} G(\mathbf{A}, \mathbf{A}^{(t)}) &= \mathcal{O}(\mathbf{A}^{(t)}) + \text{Tr}(\mathbf{C}^T \nabla \mathcal{O}(\mathbf{A}^{(t)})) \\ &\quad + \frac{1}{2} \sum_{k=1}^K \sum_{n=1}^N Q_{kn} C_{kn}^2, \end{aligned} \quad (25)$$

where

$$Q_{kn} = \frac{(\mathbf{M}^T \mathbf{M} \mathbf{A}^{(t)} + \lambda (\mathbf{1} - \mathbf{H}) \circ (\mathbf{A}^{(t)} + \xi)^{-\mathbf{H}})_{kn}}{A_{kn}^{(t)}}.$$

It could be separated into two parts $Q_{kn} = Q_{kn}^{(1)} + \lambda Q_{kn}^{(2)}$ as

$$\begin{aligned} Q_{kn}^{(1)} &= \left(\sum_{l=1}^K \frac{(\mathbf{M}^T \mathbf{M})_{kl} A_{ln}^{(t)}}{A_{kn}^{(t)}} \right) \\ Q_{kn}^{(2)} &= \frac{(1 - H_{kn}) (A_{kn}^{(t)} + \xi)^{-H_{kn}}}{A_{kn}^{(t)}}. \end{aligned}$$

Since $\frac{(A_{kn}^{(t)} + \xi)^{-H_{kn}}}{A_{kn}^{(t)}} > (A_{kn}^{(t)} + \xi)^{-(H_{kn}+1)}$, we have

$$Q_{kn}^{(2)} > (1 - H_{kn}) (A_{kn}^{(t)} + \xi)^{-(H_{kn}+1)}. \quad (26)$$

Specifically, we have to prove the following lemma:

Lemma 4. The function $G(\mathbf{A}, \mathbf{A}^{(t)})$ defined in (25) is an auxiliary function for $\mathcal{O}(\mathbf{A})$ defined in (24).

Proof: On the one hand, the equation $\mathcal{O}(\mathbf{A}) = \mathcal{O}(\mathbf{A}^{(t)}) = G(\mathbf{A}, \mathbf{A}^{(t)})$ holds for any $\mathbf{A} = \mathbf{A}^{(t)}$, i.e. $\mathbf{C} = \mathbf{0}$. On the other hand, when $\mathbf{A} \neq \mathbf{A}^{(t)}$, i.e. $\mathbf{C} \neq \mathbf{0}$ we have to prove $\mathcal{O}(\mathbf{A}^{(t)}) \leq G(\mathbf{A}, \mathbf{A}^{(t)})$.

Since the constant term and linear term in (24) are identical to their counterparts in (25), Lemma 4 could be proven by only comparing the quadratic terms as

$$\sum_{k=1}^K \sum_{n=1}^N Q_{kn} C_{kn}^2 \geq \text{Tr}(\mathbf{C}^T (\mathbf{M}^T \mathbf{M}) \mathbf{C}) - F(\mathbf{A}). \quad (27)$$

The inequality above could be expressed as two terms

$$\underbrace{\sum_{k,n} Q_{kn}^{(1)} C_{kn}^2 - \text{Tr}(\mathbf{C}^T \mathbf{M}^T \mathbf{M} \mathbf{C})}_{\text{first term}} + \underbrace{\lambda \left(\sum_{k,n} Q_{kn}^{(2)} C_{kn}^2 + F(\mathbf{A}) \right)}_{\text{second term}} \geq 0. \quad (28)$$

We could prove the inequality (28) by verifying that both terms are greater than or equal to zero. Therefore, the inequality (28) could be proven by comparing the first term [4]

$$\begin{aligned} f_1 &= \sum_{k,n,l} \left(\frac{(\mathbf{M}^T \mathbf{M})_{kl} A_{ln}^{(t)}}{A_{kn}^{(t)}} C_{kn}^2 - C_{kn} C_{ln} (\mathbf{M}^T \mathbf{M})_{lk} \right) \\ &= \sum_{k,n,l} \frac{(\mathbf{M}^T \mathbf{M})_{kl}}{2 A_{kn}^{(t)} A_{ln}^{(t)}} \left(A_{ln}^{(t)} C_{kn} - A_{kn}^{(t)} C_{ln} \right)^2 \geq 0. \end{aligned} \quad (29)$$

Then considering the inequality (26), the second term becomes

$$\begin{aligned} f_2 &> \sum_{k,n} (1 - H_{kn}) (A_{kn}^{(t)} + \xi)^{-(H_{kn}+1)} C_{kn}^2 + \\ &\quad \sum_{k,n} H_{kn} (1 - H_{kn}) (A_{kn}^{(t)} + \xi)^{-(H_{kn}+1)} C_{kn}^2 \\ &= \sum_{k,n} (1 - H_{kn}^2) (A_{kn}^{(t)} + \xi)^{-(H_{kn}+1)} C_{kn}^2, \end{aligned} \quad (30)$$

where C_{kn}^2 is undoubtedly nonnegative. Since any element H_{kn} lies in the range of $(0, 1)$, this ensures the nonnegative property of $(1 - H_{kn}^2)$. The expression $(A_{kn}^{(t)} + \xi)^{-(H_{kn}+1)}$ is positive as $A_{kn}^{(t)} + \xi$ is always positive. Therefore, $f_2 \geq 0$ holds for any condition. We have proven the inequality (27) or (28) by proving $f_1 \geq 0$ and $f_2 \geq 0$. In consequence, $G(\mathbf{A}, \mathbf{A}^{(t)})$ is an auxiliary function of $\mathcal{O}(\mathbf{A})$. ■

Through the theoretical analyses above, we have proven Theorem 1. In addition, the empirical convergence property of DgS-NMF will be analyzed in Section V-G.

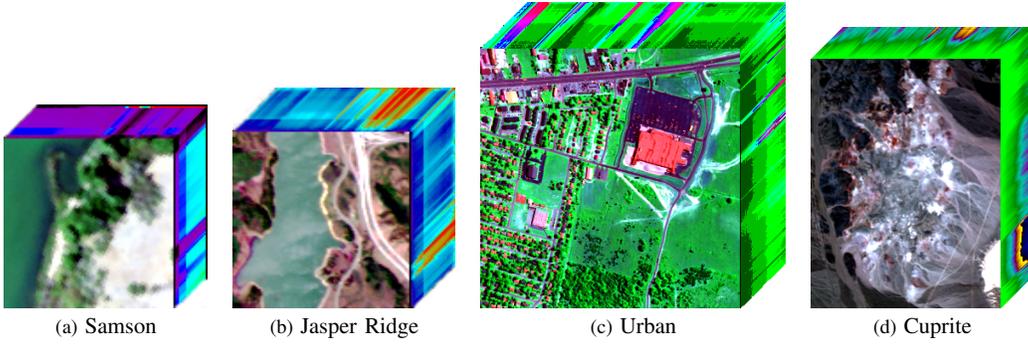


Figure 5. The four real hyperspectral images, i.e. Samson, Jasper Ridge, Urban and Cuprite respectively, used in the experiments.

Table I
COMPUTATIONAL OPERATION COUNTS FOR NMF AND DGS-NMF AT EACH ITERATION.

Methods	Arithmetic Operations in float-point format				Overall
	Addition	Multiplication	Division	Exponent	
NMF	$2LNK + 2K^2(L + N) - 2K(N + L) - 2K^2$	$2LNK + 2K^2(L + N) + K(L + N)$	$K(L + N)$	-	$O(KLN)$
DgS-NMF	$2LNK + 2K^2(L + N) - K(N + 2L) - 2K^2$	$2LNK + 2K^2(L + N) + K(L + 3N)$	$K(L + N)$	KN	$O(KLN)$

Table II
PARAMETERS USED IN COMPUTATIONAL COMPLEXITY ANALYSIS.

Parameters	Description
K	number of <i>endmembers</i>
L	number of channels
N	number of pixels in hyperspectral image
t	number of iterations
$q (= 9)$	number of pixels in the local window

D. Computational Complexity Analysis for DgS-NMF

Speed is important for algorithms [47], [48]. For this reason, the computational complexity of DgS-NMF is thoroughly analyzed by comparing with that of NMF. Since both algorithms are iteratively updated, the complexity is analyzed by summarizing the arithmetic operations at each iteration, then considering the iteration steps. For convenience, the parameters used here are listed in Table II.

In the updating rules (19) and (20), there are four kinds of arithmetic operations, i.e. addition, multiplication, division and exponent respectively. Table I summaries the counts of each arithmetic operation as well as the overall cost. In terms of the four operations, the differences between DgS-NMF and NMF are limited: DgS-NMF requires KN more additions, $2KN$ more multiplications and KN more exponents. Nevertheless, both methods have a $O(KLN)$ overall cost at one iteration step, as shown in the last column of Table I.

Apart from the updating costs, the DgS-NMF method requires $O(qLN)$ to obtain the initial DgMap and $O(q^2LN + (2\sqrt{q} - 1)^2N)$ [49], [50], [51] to get the fine tuned one. Thus, if both methods needs t iterations, the total computational complexities are $O(tKLN)$ for NMF and

$$O(tKLN + qLN + q^2LN + (2\sqrt{q} - 1)^2N)$$

for DgS-NMF. For the HU task, we have $N \gg \max(K, L, t, q)$, thus, indicating that the computational complexity of DgS-NMF is a only bit more than that of NMF, but still in the same order of magnitude.

V. EVALUATION

In this section, we evaluate the performance of the proposed method for the HU task. Several experiments are carried out to show that DgS-NMF is successfully adapted to the HU task.

A. Real Hyperspectral Images

This section introduces the information of four hyperspectral data used in the experiment. Specifically, the ground truth is achieved via the method introduced in [52], [53], [54].

Samson, as shown in Fig. 5a, is an simple data available on <http://opticks.org/confluence/display/opticks/Sample+Data>. There are 952×952 pixels in it. Each pixel is observed at 156 channels covering the wavelength from 0.401 to $0.889\mu m$. As a result, the spectral resolution is highly up to $3.13nm$. The original image is very large, which could be computationally expensive for the HU study. A region of 95×95 pixels is considered, whose first pixel corresponds to the (252, 332)-th pixel in the original image. There are three *endmembers* in this image, i.e. ‘#1 Soil’, ‘#2 Tree’ and ‘#3 Water’.

Jasper Ridge, as shown in Fig. 5b, is a popular hyperspectral data used in [55], [4]. There are 512×614 pixels in it. Each pixel is recorded at 224 channels ranging from 0.38 to $2.5\mu m$. The spectral resolution is up to $9.46nm$. Since this hyperspectral image is too complex to get the ground truth, we consider a subimage of 100×100 pixels. The first pixel starts from the (105, 269)-th pixel in the original image. After removing the channels 1–3, 108–112, 154–166 and 220–224 (due to dense water vapor and atmospheric effects), we remain

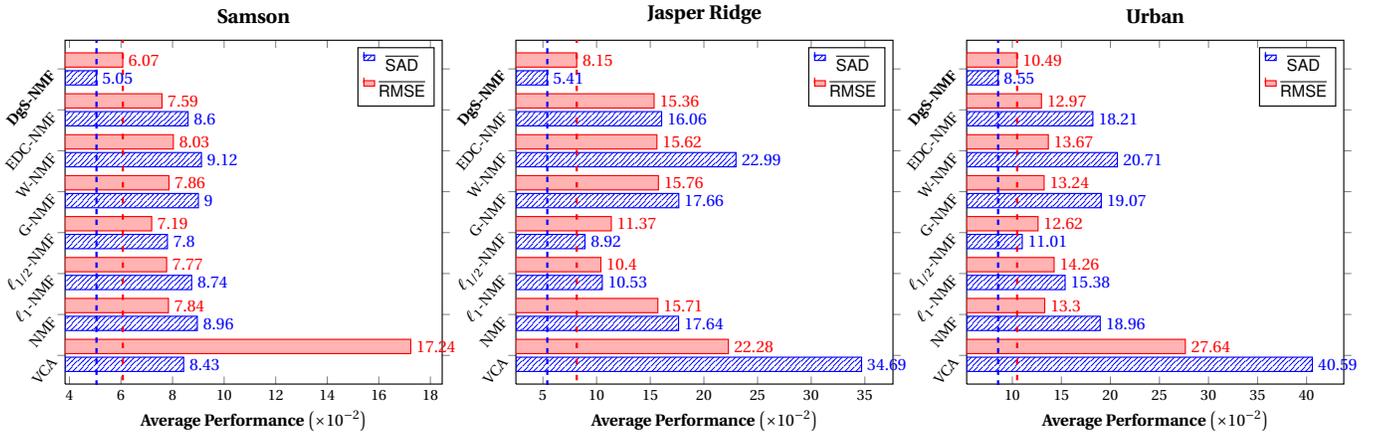


Figure 6. The average performances (i.e. $\overline{\text{SAD}}$ and $\overline{\text{RMSE}}$) of eight methods on the three datasets: Samson, Jasper Ridge and Urban, respectively.

198 channels (this is a common preprocess for HU analyses). There are four *endmembers* latent in this data: ‘#1 Tree’, ‘#2 Soil’, ‘#3 Water’ and ‘#4 Road’, as shown in Fig. 5b.

Urban is one of the most widely used hyperspectral data used in the HU area [12], [9], [54]. There are 307×307 pixels in it, each of which corresponds to a $2 \times 2 m^2$ area. In this image, there are 210 wavelengths ranging from 0.4 to $2.5 \mu m$, resulting in a spectral resolution of 10 nm. After the channels 1–4, 76, 87, 101–111, 136–153 and 198–210 are removed (due to dense water vapor and atmospheric effects), we remain 162 channels. There are four *endmembers*: ‘#1 Asphalt’, ‘#2 Grass’, ‘#3 Tree’ and ‘#4 Roof’ as shown in Fig. 5c.

Cuprite is the most benchmark dataset for the HU research [11], [12], [9], [14], [16] that covers the Cuprite in Las Vegas, NV, U.S. There are 224 channels, ranging from 0.37 to $2.48 \mu m$. After removing the noisy channels (1–2 and 221–224) and water absorption channels (104–113 and 148–167) [11], [9], we remain 188 channels. In this paper, a region (cf. Fig. 5d) of 250×190 pixels is considered, where there are 14 types of minerals [16]. Since there are minor differences between variants of the same mineral, we reduce the number of *endmembers* to 12. Note that there are small differences in the setting of *endmembers* among the papers [11], [12], [9], [14], [16]. Thus, the results of the same method in their papers might be slightly different from each other, as well slightly different from ours.

B. Compared Algorithms

To verify the performance, the proposed method is compared with seven related methods. The information of all these methods are summarized as follows:

- 1) **Our algorithm:** Data-guided Sparse regularized NMF (DgS-NMF) is a new method proposed in this paper.
- 2) Vertex Component Analysis [16] (VCA) is a classic geometric method. The code is available on <http://www.lx.it.pt/biucas/code.htm>.
- 3) Nonnegative Matrix Factorization [23] (NMF) is a benchmark statistical method. The code is obtained from <http://www.cs.helsinki.fi/u/phoyer/software.html>.

- 4) Nonnegative sparse coding [32] (ℓ_1 -NMF) is a classic sparse regularized NMF method. The code is available from <http://www.cs.helsinki.fi/u/phoyer/software.html>.
- 5) $\ell_{1/2}$ sparsity-constrained NMF [9] ($\ell_{1/2}$ -NMF) is a state-of-the-art method that could get sparser results than ℓ_1 -NMF. Since the code is unavailable, we implement it.
- 6) Graph regularized NMF [24] (G-NMF) is a good algorithm that transfer graph information latent in data to the new representation. The code is obtained from <http://www.cad.zju.edu.cn/home/dengcai/Data/GNMF.html>.
- 7) Local Neighborhood Weights regularized NMF [15] (W-NMF) is a graph based NMF method. It integrates the spectral information and spatial information when constructing the weighted graph. Since the code is unavailable from the author, we implement it.
- 8) *Endmember* Dissimilarity Constrained NMF [14] (EDC-NMF) urges the *endmember* to be smooth and different from each other. The code is implemented by ourself.

There is no parameter in VCA and NMF. For the other six methods, there is mainly one parameter. In the next subsection, we will introduce how to set the parameter for each algorithm.

C. Parameter Settings

Similar to ℓ_1 -NMF and $\ell_{1/2}$ -NMF, there is one essential parameter λ in DgS-NMF controlling the average sparsity of the new representation. To estimate an optimal parameter, two steps are required. First, an parameter range of $[\lambda_{\min}, \lambda_{\max}]$ is carefully determined by trying the values at very large steps. Second, given this parameter range, we search the best parameter by densely searching the range of $[\lambda_{\min}, \lambda_{\max}]$ at a number of equally spaced values. The parameter value that helps to achieve the best result is treated as the optimal parameter setting. For the other methods, the parameters are determined similarly. Specifically, for our method, the optimal λ is located in the range of $[0.005, 0.9]$ on all the datasets.

D. Evaluation Metrics

To assess the quantitative HU performance, two benchmark metrics are introduced, i.e. the Spectral Angle Distance

Table III

THE SADs AND RMSEs, AS WELL AS THEIR STANDARD DERIVATIONS, ON THE SAMSON DATA. FOR EACH TARGET, THE RESULTS ARE ARRANGED IN ROWS, WHERE THE RED VALUE CORRESPONDS TO THE BEST RESULT, WHILE THE BLUE VALUE IS THE SECOND BEST ONE. (BEST VIEWED IN COLOR)

End.	Spectral Angle Distance SAD ($\times 10^{-2}$)							
	VCA	NMF	ℓ_1 -NMF	$\ell_{1/2}$ -NMF	G-NMF	W-NMF	EDC-NMF	DgS-NMF
#1	7.47 \pm 12.50	6.47 \pm 7.34	6.40 \pm 7.33	6.21 \pm 7.31	6.55 \pm 7.58	6.71 \pm 7.75	6.29 \pm 7.27	5.64 \pm 7.36
#2	4.96 \pm 0.04	5.57 \pm 0.33	5.44 \pm 0.27	5.23 \pm 0.28	5.66 \pm 0.35	5.92 \pm 0.40	5.56 \pm 0.34	4.80 \pm 0.27
#3	12.87 \pm 0.55	14.85 \pm 1.84	14.39 \pm 1.72	11.97 \pm 2.11	14.79 \pm 1.85	14.75 \pm 1.86	13.94 \pm 1.70	4.70 \pm 0.34
Avg.	8.43 \pm 4.20	8.96 \pm 1.76	8.74 \pm 1.82	7.80 \pm 2.13	9.00 \pm 1.83	9.12 \pm 1.86	8.60 \pm 1.79	5.05 \pm 2.42
End.	Root Mean Square Error RMSE ($\times 10^{-2}$)							
	VCA	NMF	ℓ_1 -NMF	$\ell_{1/2}$ -NMF	G-NMF	W-NMF	EDC-NMF	DgS-NMF
#1	17.06 \pm 2.53	9.05 \pm 3.16	9.06 \pm 3.11	8.58 \pm 3.33	9.08 \pm 3.24	9.30 \pm 3.23	8.82 \pm 3.16	7.77 \pm 3.78
#2	13.35 \pm 4.51	7.64 \pm 3.62	7.63 \pm 3.57	7.44 \pm 3.66	7.67 \pm 3.70	7.97 \pm 3.71	7.59 \pm 3.64	7.74 \pm 3.63
#3	21.31 \pm 3.78	6.82 \pm 0.45	6.63 \pm 0.43	5.55 \pm 0.88	6.82 \pm 0.46	6.84 \pm 0.44	6.36 \pm 0.73	2.70 \pm 0.91
Avg.	17.24 \pm 3.54	7.84 \pm 2.18	7.77 \pm 2.23	7.19 \pm 2.40	7.86 \pm 2.24	8.03 \pm 2.23	7.59 \pm 2.06	6.07 \pm 2.76

Table IV

THE SADs AND RMSEs, AS WELL AS THEIR STANDARD DERIVATIONS, ON THE JASPER RIDGE DATA. FOR EACH TARGET, THE RESULTS ARE ARRANGED IN ROWS, WHERE THE RED VALUE CORRESPONDS TO THE BEST RESULT, WHILE THE BLUE VALUE IS THE SECOND BEST ONE. (BEST VIEWED IN COLOR)

End.	Spectral Angle Distance SAD ($\times 10^{-2}$)							
	VCA	NMF	ℓ_1 -NMF	$\ell_{1/2}$ -NMF	G-NMF	W-NMF	EDC-NMF	DgS-NMF
#1	68.87 \pm 5.13	21.40 \pm 0.28	21.29 \pm 0.28	15.10 \pm 0.33	21.51 \pm 0.23	25.53 \pm 0.45	19.71 \pm 0.25	4.66 \pm 0.21
#2	22.70 \pm 0.01	12.09 \pm 0.41	7.41 \pm 0.47	6.16 \pm 0.50	12.03 \pm 0.43	25.15 \pm 0.78	10.84 \pm 0.48	5.66 \pm 0.24
#3	24.08 \pm 1.87	18.58 \pm 0.23	6.37 \pm 0.04	4.60 \pm 0.05	18.61 \pm 0.21	20.28 \pm 0.28	15.54 \pm 0.16	4.60 \pm 0.01
#4	23.09 \pm 2.02	18.48 \pm 0.13	7.05 \pm 0.10	9.81 \pm 0.08	18.50 \pm 0.13	21.00 \pm 0.35	18.16 \pm 0.08	6.73 \pm 0.08
Avg.	34.69 \pm 1.38	17.64 \pm 0.16	10.53 \pm 0.08	8.92 \pm 0.14	17.66 \pm 0.16	22.99 \pm 0.43	16.06 \pm 0.10	5.41 \pm 0.10
End.	Root Mean Square Error RMSE ($\times 10^{-2}$)							
	VCA	NMF	ℓ_1 -NMF	$\ell_{1/2}$ -NMF	G-NMF	W-NMF	EDC-NMF	DgS-NMF
#1	34.39 \pm 3.45	18.80 \pm 0.44	14.93 \pm 0.64	16.16 \pm 0.46	18.89 \pm 0.41	17.69 \pm 0.32	19.16 \pm 0.44	11.66 \pm 0.36
#2	19.40 \pm 1.27	20.14 \pm 0.39	16.41 \pm 0.73	17.02 \pm 0.44	20.23 \pm 0.36	18.89 \pm 0.31	20.04 \pm 0.44	11.13 \pm 0.32
#3	13.75 \pm 0.98	11.69 \pm 0.06	5.41 \pm 0.01	5.57 \pm 0.03	11.71 \pm 0.06	12.45 \pm 0.10	10.46 \pm 0.05	4.13 \pm 0.01
#4	21.60 \pm 3.08	12.22 \pm 0.20	4.84 \pm 0.11	6.73 \pm 0.21	12.24 \pm 0.18	13.45 \pm 0.29	11.78 \pm 0.16	5.68 \pm 0.10
Avg.	22.28 \pm 1.87	15.71 \pm 0.18	10.40 \pm 0.33	11.37 \pm 0.23	15.76 \pm 0.17	15.62 \pm 0.24	15.36 \pm 0.20	8.15 \pm 0.18

(SAD) [12], [16], [11] and the Root Mean Square Error (RMSE) [12], [9], [56]. SAD is used to evaluate the estimated *endmembers*. It is defined as

$$\text{SAD}(\mathbf{m}, \hat{\mathbf{m}}) = \arccos\left(\frac{\mathbf{m}^T \hat{\mathbf{m}}}{\|\mathbf{m}\| \cdot \|\hat{\mathbf{m}}\|}\right), \quad (31)$$

where $\hat{\mathbf{m}}$ is the estimated *endmember* and \mathbf{m} is the corresponding ground truth. As the metric above describes the angel distance between two vectors, a smaller SAD corresponds to a better performance. To assess the estimated *abundance*, we employ the RMSE metric, which is given by

$$\text{RMSE}(\mathbf{z}, \hat{\mathbf{z}}) = \left(\frac{1}{N} \|\mathbf{z} - \hat{\mathbf{z}}\|_2^2\right)^{1/2}, \quad (32)$$

where N is the number of pixels in the image, $\hat{\mathbf{z}}$ (a row vector in the *abundance* matrix $\hat{\mathbf{A}}$) is the estimated *abundance* map, and \mathbf{z} is the corresponding ground truth. In general, a smaller RMSE corresponds to a better result.

E. Performance Evaluation

To verify the performance of our method, eight experiments are carried out. Each experiment is repeated 20 times. The mean results as well as their standard deviations are reported. The evaluation includes two parts: quantitative comparisons and visual comparisons.

1) *Quantitative Comparisons*: The quantitative results are summarized in Tables III, IV, V, VI and plotted in Fig. 6. In Table III, there are two sub-tables that show SADs and RMSEs respectively on Samson. In the sub-table, each row shows the performances of one *endmember*, i.e. ‘#1 Soil’, ‘#2 Tree’ and ‘#3 Water’ in sequence. The last row shows the average performance. In each category, the value in the red ink is the best, while the blue value is the second best. As Table III shows, our method generally achieves the best results, and in a few cases it achieves comparable results with the best results of other methods. Such case is better illustrated in the 1st subfigure of Fig. 6, where DgS-NMF is the best method that reduces 35.3% for SAD and 15.6% for RMSE according

Table V

THE SADs AND RMSEs, AS WELL AS THEIR STANDARD DERIVATIONS, ON THE URBAN DATA. FOR EACH TARGET, THE RESULTS ARE ARRANGED IN ROWS, WHERE THE RED VALUE CORRESPONDS TO THE BEST RESULT, WHILE THE BLUE VALUE IS THE SECOND BEST ONE. (BEST VIEWED IN COLOR)

Endm.	Spectral Angle Distance SAD ($\times 10^{-2}$)							
	VCA	NMF	ℓ_1 -NMF	$\ell_{1/2}$ -NMF	G-NMF	W-NMF	EDC-NMF	DgS-NMF
#1	21.10 \pm 2.34	17.87 \pm 0.27	16.98 \pm 0.36	6.06 \pm 0.22	17.98 \pm 0.21	20.68 \pm 0.16	14.88 \pm 0.34	5.86 \pm 0.09
#2	35.35 \pm 4.42	21.89 \pm 0.40	23.45 \pm 0.35	20.05 \pm 0.37	22.05 \pm 0.37	23.40 \pm 0.46	22.35 \pm 0.37	13.69 \pm 0.23
#3	28.95 \pm 6.54	10.43 \pm 0.04	4.05 \pm 0.02	3.71 \pm 0.03	10.53 \pm 0.03	11.45 \pm 0.04	10.16 \pm 0.03	4.12 \pm 0.01
#4	76.98 \pm 0.09	25.67 \pm 0.21	17.06 \pm 0.35	14.22 \pm 0.19	25.72 \pm 0.22	27.30 \pm 0.26	25.45 \pm 0.20	10.54 \pm 0.28
Avg.	40.59 \pm 1.82	18.96 \pm 0.10	15.38 \pm 0.07	11.01 \pm 0.09	19.07 \pm 0.11	20.71 \pm 0.14	18.21 \pm 0.08	8.55 \pm 0.08
Root Mean Square Error RMSE ($\times 10^{-2}$)								
#1	28.22 \pm 5.50	14.16 \pm 0.09	18.97 \pm 0.09	14.77 \pm 0.14	14.09 \pm 0.08	14.60 \pm 0.09	13.38 \pm 0.07	13.18 \pm 0.06
#2	35.24 \pm 8.38	16.22 \pm 0.12	18.17 \pm 0.13	16.16 \pm 0.20	16.11 \pm 0.10	16.73 \pm 0.09	15.97 \pm 0.09	12.95 \pm 0.05
#3	28.35 \pm 3.33	13.34 \pm 0.14	14.02 \pm 0.18	12.65 \pm 0.16	13.28 \pm 0.12	13.52 \pm 0.11	13.14 \pm 0.09	9.57 \pm 0.12
#4	18.73 \pm 3.32	9.49 \pm 0.03	5.89 \pm 0.06	6.90 \pm 0.12	9.49 \pm 0.02	9.85 \pm 0.01	9.40 \pm 0.02	6.27 \pm 0.06
Avg.	27.64 \pm 3.62	13.30 \pm 0.09	14.26 \pm 0.07	12.62 \pm 0.11	13.24 \pm 0.07	13.67 \pm 0.07	12.97 \pm 0.06	10.49 \pm 0.06

Table VI

THE SADs AND THEIR STANDARD DERIVATIONS OF 6 METHODS ON THE CUPRITE DATASET. THERE ARE 12 KINDS OF *endmembers*. FOR EACH *endmember*, THE RESULTS ARE ARRANGED IN ROWS, WHERE THE RED VALUE IS THE BEST ONE. (BEST VIEWED IN COLOR)

Endmembers	Spectral Angle Distance SAD ($\times 10^{-2}$)					
	NMF	ℓ_1 -NMF	$\ell_{1/2}$ -NMF	EDC-NMF	W-NMF	DgS-NMF
#1 Alunite	16.00 \pm 2.19	16.22 \pm 2.05	14.64 \pm 1.80	14.32 \pm 4.58	15.88 \pm 3.08	12.48\pm1.85
#2 Andradite	10.52 \pm 3.13	10.15 \pm 3.02	7.86 \pm 1.37	9.32 \pm 0.84	9.75 \pm 3.35	7.59\pm1.19
#3 Buddingtonite	12.50 \pm 7.73	12.50 \pm 7.51	11.84 \pm 5.62	11.20 \pm 4.84	11.68 \pm 6.81	10.81\pm3.18
#4 Dumortierite	13.83 \pm 5.56	13.07 \pm 5.31	11.53 \pm 3.22	12.82 \pm 2.33	12.59 \pm 5.49	11.01\pm2.15
#5 Kaolinite ₁	9.52 \pm 1.98	9.42 \pm 1.82	9.68 \pm 2.30	8.74\pm2.02	9.75 \pm 2.20	9.02 \pm 2.86
#6 Kaolinite ₂	10.17 \pm 2.76	9.84 \pm 2.72	7.24\pm0.96	7.80 \pm 0.77	8.78 \pm 2.17	7.42 \pm 1.22
#7 Muscovite	28.89 \pm 8.56	29.86 \pm 7.37	25.12 \pm 5.51	24.31 \pm 7.18	20.89 \pm 10.43	20.51\pm6.02
#8 Montmorillonite	10.48 \pm 4.41	10.27 \pm 4.70	7.46 \pm 1.50	8.42 \pm 1.36	8.42 \pm 2.70	7.27\pm1.25
#9 Nontronite	12.69 \pm 4.16	12.80 \pm 4.06	9.54 \pm 2.35	9.73 \pm 1.69	12.69 \pm 4.56	8.88\pm1.71
#10 Pyrope	8.84 \pm 3.44	8.12\pm1.87	8.94 \pm 5.22	8.64 \pm 6.18	10.75 \pm 4.47	8.82 \pm 5.01
#11 Sphene	10.97 \pm 3.35	11.03 \pm 3.34	7.91\pm2.54	8.88 \pm 6.70	10.45 \pm 4.58	8.15 \pm 2.14
#12 Chalcedony	13.48 \pm 6.08	13.72 \pm 6.00	13.76 \pm 6.37	12.84\pm5.86	13.85 \pm 7.94	13.62 \pm 6.33
Avg.	13.16 \pm 1.25	13.08 \pm 1.16	11.29 \pm 0.99	11.42 \pm 0.88	12.12 \pm 1.39	10.46\pm0.58

to the results of the second best method, i.e. $\ell_{1/2}$ -NMF.

Table IV summarizes the performances of eight methods on Jasper Ridge. The rows show the results of four targets, i.e. ‘#1 Road’, ‘#2 Soil’, ‘#3 Water’ and ‘#4 Tree’ respectively. Generally, the sparsity constrained methods, i.e. ℓ_1 -NMF, $\ell_{1/2}$ -NMF and DgS-NMF, achieve better results than other methods. This is since sparse constraints tend to find expressive *endmembers* [27], which might be more reliable for the HU task. The average performances (i.e. SAD and RMSE) are illustrated in the 2nd subfigure of Fig. 6. As we shall see, our method obtains extraordinary advantages—compared with the second best methods, i.e. $\ell_{1/2}$ -NMF and ℓ_1 -NMF, DgS-NMF reduces 39.3% and 21.6% respectively for SAD and RMSE.

The results on Urban are illustrated in Table V, where the rows contain results of ‘#1 Asphalt’, ‘#2 Grass’, ‘#3 Tree’ and ‘#4 Roof’ respectively. It can be seen that apart from our method, $\ell_{1/2}$ -NMF, ℓ_1 -NMF and EDC-NMF generally achieve better results than the others. However, in general, DgS-NMF

obtains the best performance. In Fig. 6, the 3rd subfigure shows the average performances. Compared with the second best methods, i.e. $\ell_{1/2}$ -NMF, our method reduces 22.3% and 16.9% for SAD and RMSE respectively.

For the former three datasets, the number of *endmembers* is small, i.e. K is small. To verify the performance of our method on a dataset with large K , we carry out an experiment on the Cuprite dataset. It is worth mentioning that the Cuprite is the most important benchmark dataset for the HU research [9], [14]. As the ground truth for the *abundance* is unavailable, only the results of *endmembers* are reported in Table VI. As we shall see, our method generally obtains the best performance. Besides, the sparsity constrained methods, i.e. $\ell_{1/2}$ -NMF and ℓ_1 -NMF, usually achieve relatively good results.

2) *Visual Comparisons* : In order to give an intuitive HU comparison, we illustrate the *abundance* maps in two ways: in pseudo color and in gray scale. Fig. 8a illustrates an example of the pseudo color manner, where there are mainly four color

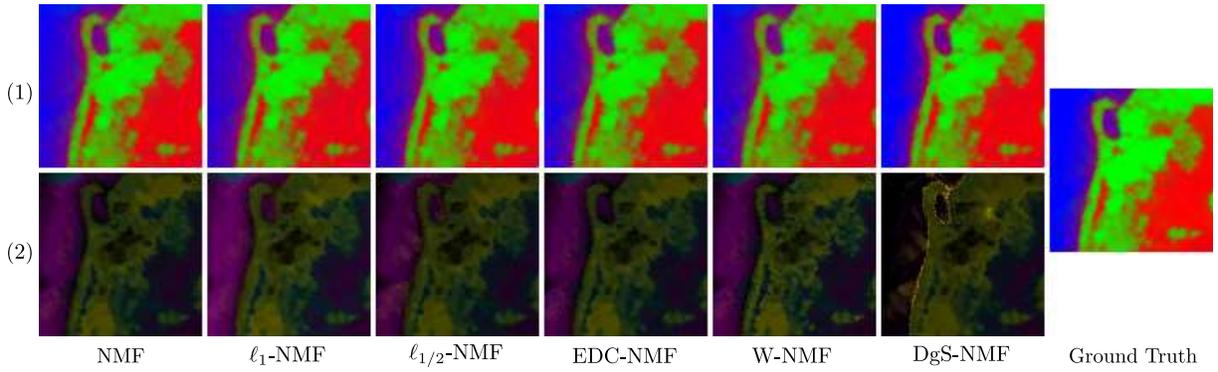


Figure 7. The *abundance* maps in pseudo color on the Samson data. There are seven columns and two rows in this figure. From the 1st to the 6th column, each column shows the result of one algorithm. The last column shows the ground truth. The second row shows the absolute difference between the estimated result $\hat{\mathbf{A}}$ and the ground truth \mathbf{A} , i.e. $|\mathbf{A} - \hat{\mathbf{A}}| \in \mathbb{R}_+^{K \times N}$. For each subfigure, the proportions of Red, Green and Blue inks associated with each pixel represent the *abundances* of ‘Soil’, ‘Tree’ and ‘Water’ in the corresponding pixel. (Best viewed in color)

inks. Through these colors, we could represent the fractional *abundances* A_{kn} associated with pixel \mathbf{y}_n by plotting the corresponding pixel using the proportions of red, blue, green and black inks given by A_{kn} for $k = 1, 2, 3, 4$, respectively. So, for instance, a pixel for which $A_{2n} = 1$ will be colored blue, whereas one for which $A_{1n} = A_{2n} = 0.5$ will be colored with equal proportions of red and blue inks and so will appear purple. Figs. 7, 8a and 9a are obtained in this way.

For the Samson dataset, because of the high quality of all the estimated *abundances*, the *abundance* maps in gray scale might be very similar. For this reason, we only illustrate the pseudo color version. The results are illustrated in Fig. 7. The top row shows the *abundance* maps in pseudo color, and the bottom row shows the absolute difference between the estimated results $\hat{\mathbf{A}}$ and the ground truth \mathbf{A} , i.e. $|\mathbf{A} - \hat{\mathbf{A}}| \in \mathbb{R}_+^{K \times N}$. As Fig. 7 shows, in general, the DgS-NMF method achieves the minimal difference according to the ground truth.

For the Jasper Ridge data, the *abundance* maps in pseudo color and in gray scale are both provided in Fig. 8. There are four targets, i.e. ‘#1 Tree’, ‘#2 Soil’, ‘#3 Water’ and ‘#4 Road’ respectively, the fractional *abundances* of which are illustrated by the proportions of red, blue, green and black inks associated with each pixel, as shown in Fig. 8a. As can be seen, the sparse constraint methods, i.e. ℓ_1 -NMF, $\ell_{1/2}$ -NMF and DgS-NMF, get better results than the other methods. Specifically, DgS-NMF achieves extraordinary results—the absolute difference map in the (2, 6)-th subfigure is the minimal one.

In Fig. 9, the *abundance* maps in pseudo color and in gray scale are shown for the Urban data. The four targets are as follows: ‘#1 Asphalt’, ‘#2 Grass’, ‘#3 Tree’ and ‘#4 Roof’. The *abundances* of these targets are equal to the proportions of red, green, blue and black inks at each pixel. Similar to the results in Figs. 7 and 8, our method achieves the best result in terms of the absolute difference map as shown in the 6th subfigure in the second row in Fig. 9a.

F. Influences of Varying Parameters

To test the stability of our method, the influences of parameters are evaluated. Nine experiments have been conducted with respect to nine varying parameters: $\lambda = 0.2\lambda_0, \dots, 1.8\lambda_0$. Here,

λ_0 is the optimal parameter for each algorithm; it might be different either for different algorithms or on different datasets. To reduce the randomness, each experiment is repeated ten times and the mean results are reported.

The quantitative performances are summarized in Fig. 10, where there are two rows and three columns. The top row shows the average SADs, while the bottom row displays the average RMSEs. Each column shows the results on one dataset. As can be seen, the curves of NMF and EDC-NMF are plain. For the former method, there is no parameter in it. For the latter one, we fix λ at the optimal parameter λ_0 . This is because the parameter in EDC-NMF can not be set freely; too big parameter value would lead to failure updating. In general, the sparse constraint methods achieve better results for all tested parameter values. Additionally, for most cases, DgS-NMF achieves great advantages.

G. Convergence Study

In Section IV-C, it has been proven that the objective (12) could converge to a minimum by using the updating rules (19) and (20). To verify this conclusion, we study the empirical convergence property of DgS-NMF by comparing its convergence curves with that of NMF (a benchmark method). As shown in Fig. 11, there are three subfigures, each of which shows the results on one dataset. In each subfigure, the X-axis shows the number of iteration t , and the Y-axis illustrates the relative decrement of the objective energy, i.e. $\frac{(\mathcal{O}_t - \mathcal{O}_{t+1})}{\mathcal{O}_t}$, of NMF and DgS-NMF. All values in Fig. 11 are nonnegative, indicating that the objective energy of both methods decrease at each iteration. Besides, DgS-NMF converges to a local minimum with comparable iteration steps as NMF. In this way, we’ve proven Theorem 1 via empirical results.

H. Influences of DgMaps

This section gives two kinds of evaluations: 1) to evaluate the estimated DgMaps, 2) to evaluate the contribution of the fine tuning step proposed in Section III-B. Both visual and quantitative comparisons have been introduced.

Obviously, the mixed level of each pixel is closely related to the sparse level of the corresponding *abundance* vector. It is

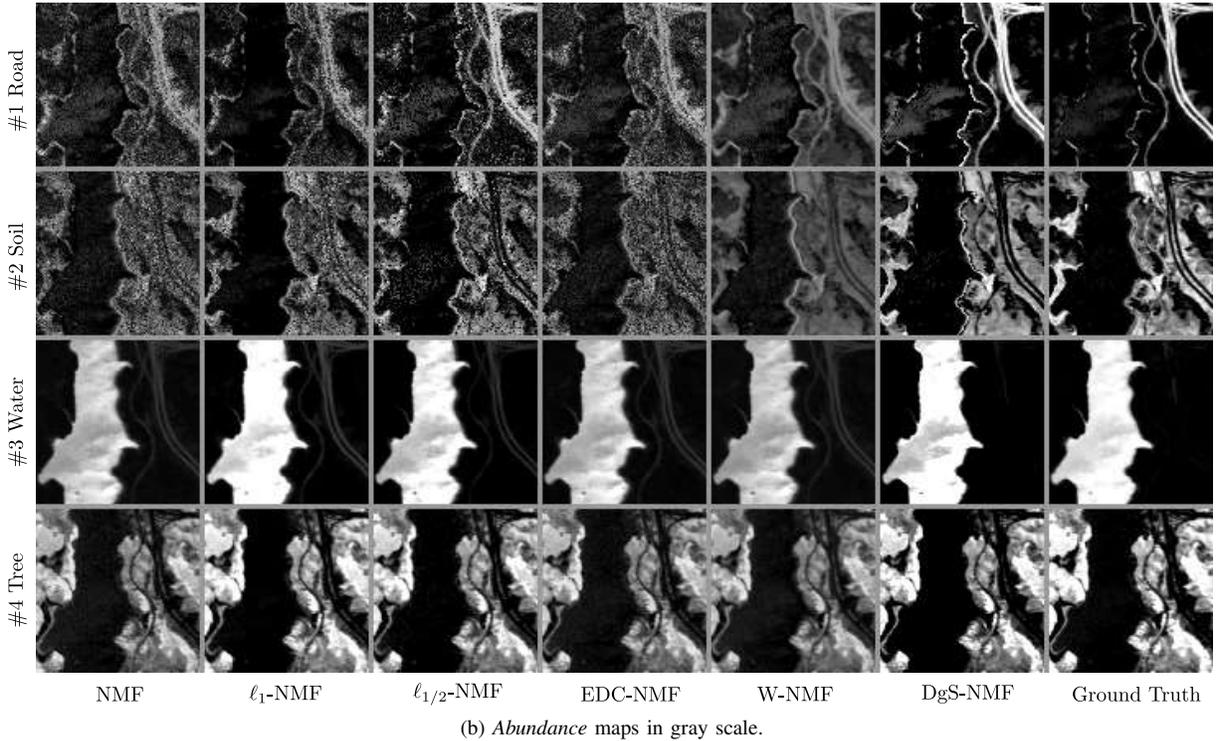
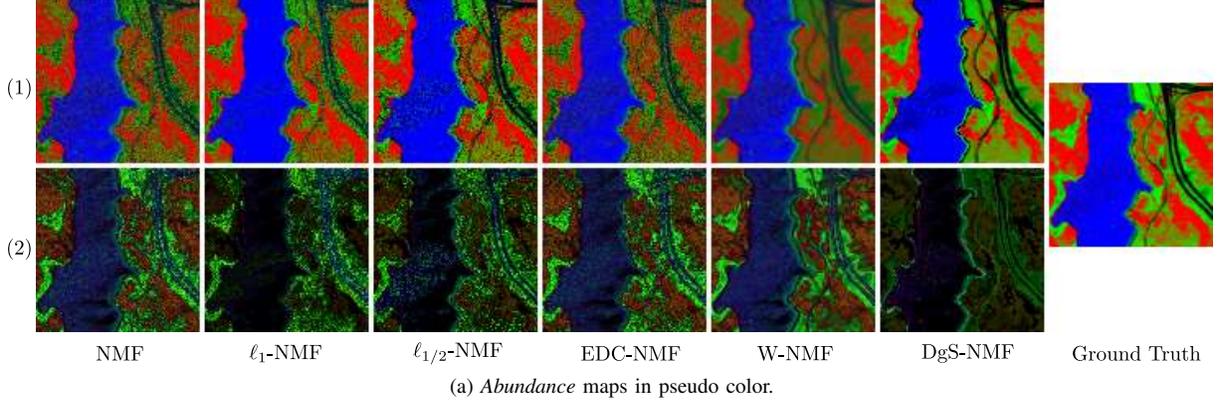


Figure 8. The *abundance maps* on the Jasper Ridge data: (a) in pseudo color and (b) in gray scale. There are two rows in (a). The second row shows the absolute difference between the estimated result $\hat{\mathbf{A}}$ and the ground truth \mathbf{A} , i.e. $|\mathbf{A} - \hat{\mathbf{A}}| \in \mathbb{R}_+^{K \times N}$. For each subfigure in (a), the proportions of Red, Blue, Green and Black inks associated with each pixel represent the fractional *abundances* of ‘Tree’, ‘Water’, ‘Soil’ and ‘Road’ in the corresponding pixel. There are four rows and seven columns in (b). Each row shows the *abundance maps* of one target. From the 1st to the 6th column, each column illustrates the results of one algorithm. The last column shows the ground truth. (Best viewed in color)

reasonable to assess an estimated DgMap by comparing with the corresponding sparse map of *abundances* from the ground truth. Specifically, given *abundance vectors* $\{\mathbf{a}_n\}_{n=1}^N \in \mathbb{R}_+^K$, the n^{th} value in the sparse map is obtained by measuring the sparsity [57], [9] of \mathbf{a}_k :

$$S_n = \frac{\sqrt{K} - \|\mathbf{a}_n\|_1 / \|\mathbf{a}_n\|_2}{\sqrt{K} - 1}, \quad \forall n \in \{1, 2, \dots, N\}, \quad (33)$$

where K is the number of elements in the *abundance vector*.

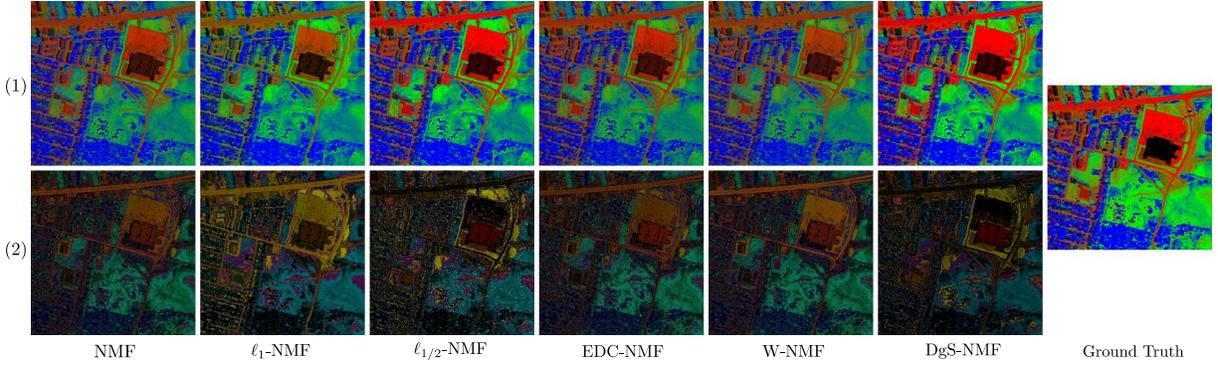
The visual comparisons of the fine tuned DgMap and the sparse map from ground truths are illustrated in Fig. 12. As we shall see, the estimated DgMap is generally good. It achieves very good results in the sudden change areas, while in the smooth areas our method fails to capture the mixed

information.

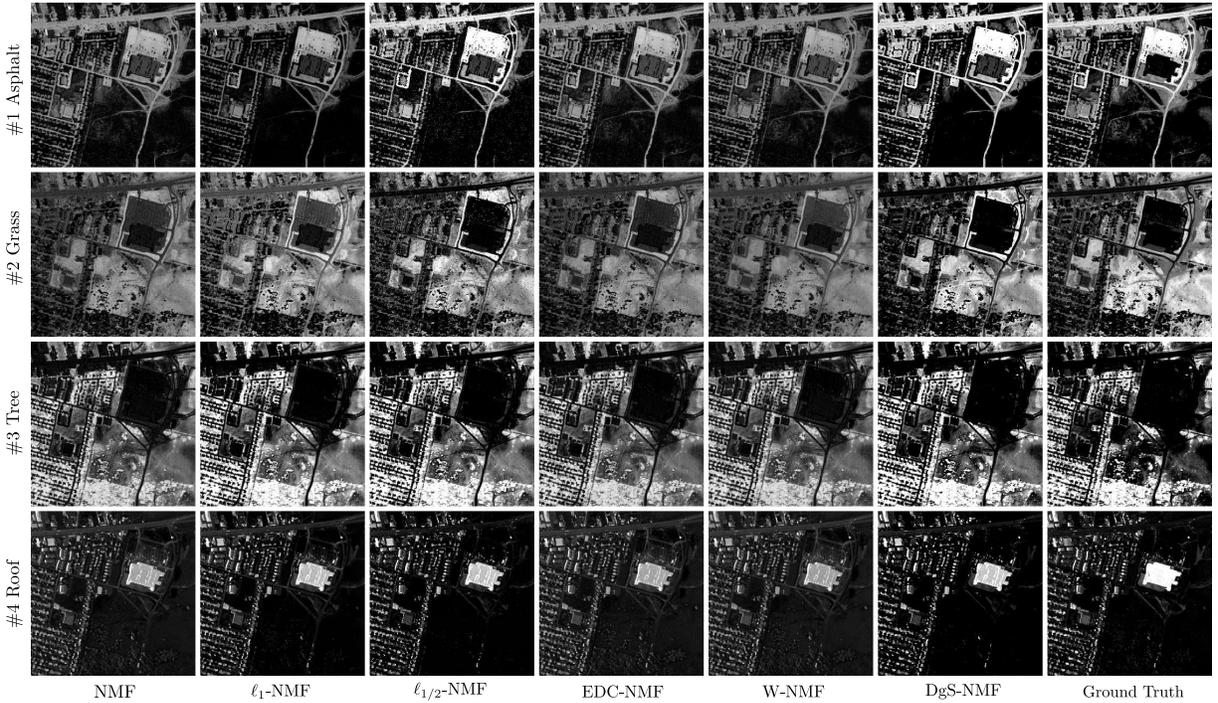
To study the quantitative evaluations, the HU performances⁴ are summarized in Table VII and visualized in Fig. 13. There are three kinds of results of DgS-NMF with respect to three maps: 1) “map₁” is the initial DgMap; 2) “map₂” denotes the fine tuned DgMap; and 3) “map₃” means the sparse map from ground truths. As we shall see, in most cases, the results of “map₃” are the best, and the results of “map₂” are the second best. Such observations are better illustrated in Fig. 13. These observations above imply that:

- the results of the proposed data guided sparse model (DgS-NMF) is quite promising. One can expect an even

⁴Since the standard variation of each method is similar, only the average HU performances are provided.



(a) *Abundance maps in pseudo color.*



(b) *Abundance maps in gray scale.*

Figure 9. The *abundance maps* on the Urban data: (a) in pseudo color and (b) in gray scale. There are two rows in (a). The second row shows the absolute difference between the estimated result $\hat{\mathbf{A}}$ and the ground truth \mathbf{A} , i.e. $|\mathbf{A} - \hat{\mathbf{A}}| \in \mathbb{R}_+^{K \times N}$. For each subfigure in (a), the proportions of Red, Blue, Green and Black inks associated with each pixel represent the fractional *abundances* of ‘Asphalt’, ‘Tree’, ‘Grass’ and ‘Roof’ in the corresponding pixel. There are four rows and seven columns in (b). Each row shows the *abundance maps* of one target. From the 1st to the 6th column, each column illustrates the results of one algorithm. The last column shows the ground truths. (Best viewed in color)

better result with a better estimation of DgMap.

- although the initial DgMap helps DgS-NMF to achieve good HU results, the fine tuning process could further improve the HU performances very much.

There are mainly two contributions of this paper. First, we propose a data-guided sparsity model for the HU task. We have verified its effectiveness by a heuristic DgMap estimation method. If we can obtain a more accurate DgMap, the result can be further improved. Second, our work introduces a new and open problem for the hyperspectral image: how to effectively estimate a DgMap from a hyperspectral image cube? This problem has never been considered in this area. Owing to the encouraging result obtained by introducing the data guided sparsity, we would like to do some further research to make

it sound. The learning based methods might be exploited to estimate better DgMaps. Besides, the accelerating techniques used in [10] will be considered as well.

VI. CONCLUSIONS

In this paper, we have provided a novel Data-guided Sparse NMF (DgS-NMF) method by deriving a data-guided map from the original hyperspectral image. Through this data-guided map, the sparse constraint could be applied in an adaptive manner. Such case not only agrees with the practical situation but also leads the *endmember* toward some spectra resembling the highly sparse regularized pixel. What is more, experiments on the four datasets demonstrate the advantages of DgS-NMF: 1) under the optimal parameter setting, DgS-

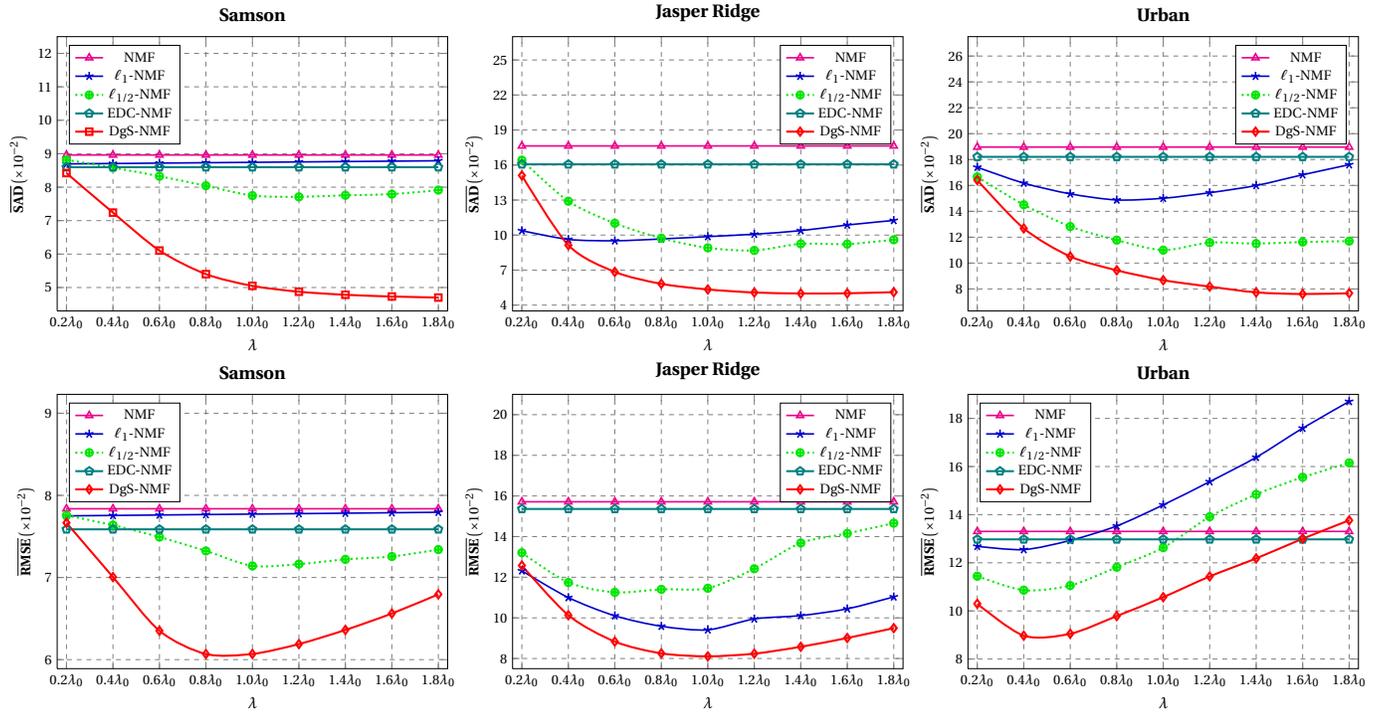


Figure 10. Performance vs. parameter λ . There are two rows and three columns. The top row shows $\overline{\text{SAD}}$ s and the bottom row shows $\overline{\text{RMSE}}$ s. Each column shows the performance on one dataset. In each subfigure, λ_0 on the X-axis donates the best parameter setting for each algorithm. (Best viewed in color)

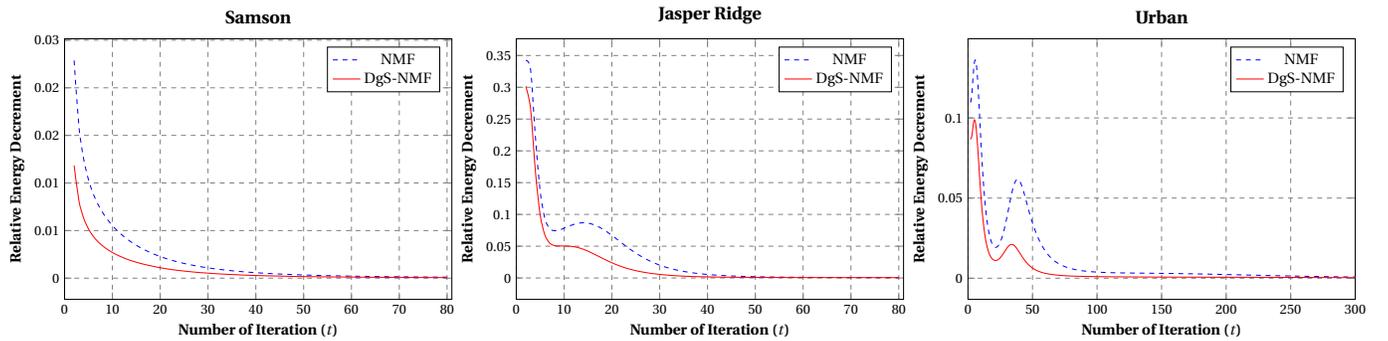


Figure 11. Relative decrement of the objective energy, i.e. $\frac{(\mathcal{O}_t - \mathcal{O}_{t+1})}{\mathcal{O}_t}$, of NMF and DgS-NMF on the three datasets: Samson, Jasper Ridge and Urban.

NMF achieves better results than all the other methods in terms of both quantitative and visual performances; 2) when the parameter varies, in most cases, our method achieves remarkable advantages over its competitors. Besides, both theoretic proof and empirical results verify the convergence ability of our method.

ACKNOWLEDGEMENTS

The authors would like to thank the editor and reviewers for their valuable comments and suggestions. This work is supported by the projects (Grant No. 61331018, 91338202, 61305049 and 61375024) of the National Natural Science Foundation of China.

REFERENCES

- [1] N. Keshava, "A survey of spectral unmixing algorithms," *Lincoln Lab. J.*, vol. 14, no. 1, pp. 55–78, Jan 2003.
- [2] Bioucas-Dias *et al.*, "Hyperspectral unmixing overview: Geometrical, statistical, and sparse regression-based approaches," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens. (JSTARS)*, vol. 5, no. 2, pp. 354–379, april 2012.
- [3] S. Cai, Q. Du, and R. Moorhead, "Hyperspectral imagery visualization using double layers," *IEEE Trans. Geosci. Remote Sens.*, vol. 45, no. 10, pp. 3028–3036, Oct 2007.
- [4] F. Zhu, Y. Wang, S. Xiang, B. Fan, and C. Pan, "Structured sparse method for hyperspectral unmixing," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 88, pp. 101–118, 2014.
- [5] N. Keshava and J. F. Mustard, "Spectral unmixing," *IEEE Signal Processing Magazine*, vol. 19, no. 1, pp. 44–57, Jan 2002.
- [6] K. C. Mertens, L. P. C. Verbeke, E. I. Ducheyne, and R. R. D. Wulf, "Using genetic algorithms in sub-pixel mapping," *International Journal of Remote Sens.*, vol. 24, no. 21, pp. 4241–4247, 2003.
- [7] Z. Guo, T. Wittman, and S. Osher, "L1 unmixing and its application to hyperspectral image enhancement," vol. 7334, no. 1, 2009, p. 73341M.
- [8] R. Kawakami, J. Wright, Y.-W. Tai, Y. Matsushita, M. Ben-Ezra, and K. Ikeuchi, "High-resolution hyperspectral imaging via matrix factorization," *IEEE CVPR*, vol. 0, pp. 2329–2336, 2011.
- [9] Y. Qian, S. Jia, J. Zhou, and A. Robles-Kelly, "Hyperspectral unmixing via $\ell_{1/2}$ sparsity-constrained nonnegative matrix factorization," *IEEE Trans. Geosci. Remote Sens.*, vol. 49, no. 11, pp. 4282–4297, nov 2011.

Table VII

THE COMPARISON OF HU PERFORMANCES OF DGS-NMF VS. THREE MAPS, ON THE THREE DATASETS. “MAP₁” INDICATES THE INITIAL DGMAP; “MAP₂” MEANS THE FINE-TUNED DGMAP; “MAP₃” DENOTES THE SPARSE MAP FROM GROUND TRUTHS, DEFINED BY THE SPARSE METRIC (33). THE RED VALUE CORRESPONDS TO THE BEST RESULT, WHILE THE BLUE VALUE IS THE SECOND BEST RESULT. (BEST VIEWED IN COLOR)

End.	Spectral Angle Distance SAD ($\times 10^{-2}$)									Root Mean Square Error RMSE ($\times 10^{-2}$)								
	Samson			Jasper Ridge			Urban			Samson			Jasper Ridge			Urban		
	map ₁	map ₂	map ₃	map ₁	map ₂	map ₃	map ₁	map ₂	map ₃	map ₁	map ₂	map ₃	map ₁	map ₂	map ₃	map ₁	map ₂	map ₃
#1	5.99	5.64	4.36	5.94	4.66	3.11	5.78	5.86	6.34	8.07	7.77	5.47	11.69	11.66	5.91	12.11	13.18	12.26
#2	5.08	4.80	3.94	8.76	5.66	3.64	16.04	13.69	6.15	7.85	7.74	5.92	12.06	11.13	6.34	12.57	12.95	11.80
#3	6.67	4.70	4.72	3.64	4.60	3.81	4.04	4.12	3.71	3.30	2.70	2.44	4.20	4.13	4.20	10.99	9.57	5.89
#4	–	–	–	8.48	6.73	6.02	11.71	10.54	9.96	–	–	–	5.95	5.68	5.64	6.28	6.27	6.28
Avg.	5.92	5.05	4.34	6.71	5.41	4.15	9.39	8.55	6.54	6.41	6.07	4.61	8.47	8.15	5.52	10.49	10.49	9.06

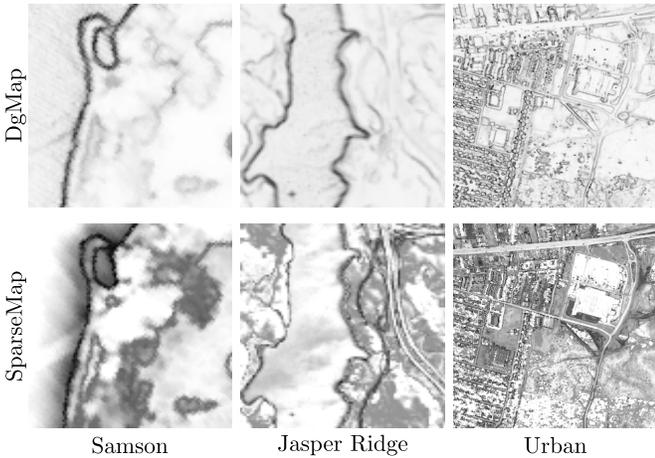


Figure 12. The comparison of DgMaps vs. Sparse Maps on the three datasets. DgMaps are obtained by Section III-B; the Sparse maps are achieved by measuring the sparsity of the abundances from the ground truth by (33).

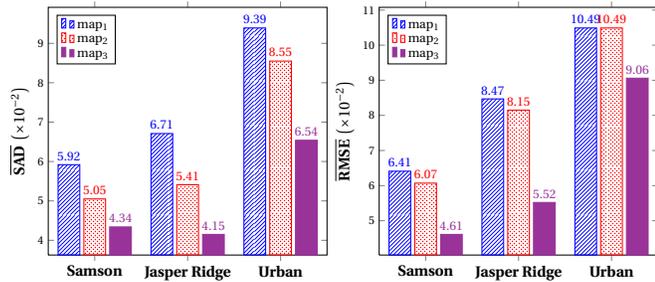


Figure 13. The comparison of average performances (i.e. \overline{SAD} and \overline{RMSE}) of DgS-NMF vs. three maps, on the three datasets. “map₁” indicates the initial DgMap; “map₂” means the fine-tuned DgMap; “map₃” denotes the sparse map from ground truths, defined by the sparse metric (33).

[10] C. Li, T. Sun, K. Kelly, and Y. Zhang, “A compressive sensing and unmixing scheme for hyperspectral data processing,” *IEEE TIP*, vol. 21, no. 3, pp. 1200–1210, March 2012.

[11] X. Lu, H. Wu, Y. Yuan, P. Yan, and X. Li, “Manifold regularized sparse nmf for hyperspectral unmixing,” *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 5, pp. 2815–2826, 2013.

[12] X. Liu, W. Xia, B. Wang, and L. Zhang, “An approach based on constrained nonnegative matrix factorization to unmix hyperspectral data,” *IEEE TGRS*, vol. 49, no. 2, pp. 757–772, 2011.

[13] L. Miao, H. Qi, and H. Qi, “Endmember extraction from highly mixed data using minimum volume constrained nonnegative matrix factorization,” *IEEE TGRS*, vol. 45, no. 3, pp. 765–777, 2007.

[14] N. Wang, B. Du, and L. Zhang, “An endmember dissimilarity constrained non-negative matrix factorization method for hyperspectral

unmixing,” *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens. (JSTARS)*, vol. 6, no. 2, pp. 554–569, 2013.

[15] J. Liu, J. Zhang, Y. Gao, C. Zhang, and Z. Li, “Enhancing spectral unmixing by local neighborhood weights,” *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 5, no. 5, pp. 1545–1552, 2012.

[16] J. M. P. Nascimento and J. M. B. Dias, “Vertex component analysis: a fast algorithm to unmix hyperspectral data,” *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 4, pp. 898–910, 2005.

[17] C.-I. Chang, C.-C. Wu, W. Liu, and Y. C. Ouyang, “A new growing method for simplex-based endmember extraction algorithm,” *IEEE Trans. Geosci. Remote Sens.*, vol. 44, no. 10, pp. 2804–2819, 2006.

[18] G. Martin and A. Plaza, “Spatial-spectral preprocessing prior to end-member identification and unmixing of remotely sensed hyperspectral data,” *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens. (JSTARS)*, vol. 5, no. 2, pp. 380–395, 2012.

[19] N. Dobigeon, S. Moussaoui, M. Coulon, J.-Y. Tourneret, and A. O. Hero, “Joint bayesian endmember extraction and linear unmixing for hyperspectral imagery,” *IEEE TSP*, vol. 57, no. 11, pp. 4355–4368, 2009.

[20] J. M. Bioucas-Dias, “A variable splitting augmented lagrangian approach to linear spectral unmixing,” in *WHISPERS*, 2009, pp. 1–4.

[21] J. M. P. Nascimento and J. M. Bioucas-Dias, “Hyperspectral unmixing based on mixtures of dirichlet components,” *IEEE Trans. Geosci. Remote Sens.*, vol. 50, no. 3, pp. 863–878, 2012.

[22] M. E. Winter, “N-findr: an algorithm for fast autonomous spectral end-member determination in hyperspectral data,” in *Proc. SPIE*, vol. 3753, 1999, pp. 266–275.

[23] D. D. Lee and H. S. Seung, “Learning the parts of objects with nonnegative matrix factorization,” *Nature*, vol. 401, no. 6755, pp. 788–791, Oct 1999.

[24] D. Cai, X. He, J. Han, and T. S. Huang, “Graph regularized nonnegative matrix factorization for data representation,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 8, pp. 1548–1560, aug 2011.

[25] S. E. Palmer, “Hierarchical structure in perceptual representation,” *Cognitive Psychology*, vol. 9, no. 4, pp. 441–474, 1977.

[26] N. K. Logothetis and D. L. Sheinberg, “Visual object recognition,” *Annu. Rev. Neurosci.*, vol. 19, no. 1, pp. 577–621, 1996.

[27] S. Z. Li, X. Hou, H. Zhang, and Q. Cheng, “Learning spatially localized, parts-based representation,” in *IEEE CVPR*, 2001, pp. 207–212.

[28] R. Sandler *et al.*, “Nonnegative matrix factorization with earth mover’s distance metric for image analysis,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 8, pp. 1590–1602, 2011.

[29] W. Xu, X. Liu, and Y. Gong, “Document clustering based on non-negative matrix factorization,” in *Int. Conf. on Res. and Development in Inform. Retrieval (SIGIR)*, 2003, pp. 267–273.

[30] F. Shahnaz, M. W. Berry, V. P. Pauca, and R. J. Plemmons, “Document clustering using nonnegative matrix factorization,” *Inform. Proc. & Management*, vol. 42, no. 2, pp. 373–386, 2006.

[31] D. D. Lee and H. S. Seung, “Algorithms for non-negative matrix factorization,” in *NIPS*. MIT Press, 2000, pp. 556–562.

[32] P. O. Hoyer, “Non-negative sparse coding,” in *Proc. 12th IEEE Workshop Neural Netw. Signal Process.*, 2002, pp. 557–565.

[33] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society*, vol. 58, no. 1, pp. 267–288, 1996.

[34] D. L. Donoho, “Compressed sensing,” *IEEE Trans. Inform. Theory*, vol. 52, no. 4, pp. 1289–1306, 1996.

[35] J. Fan and R. Li, “Variable selection via nonconcave penalized likelihood and its oracle properties,” *Journal of the American Statistical Association*, vol. 96, no. 456, pp. 1348–1360, 2001.

- [36] A. Levin, D. Lischinski, Y. Weiss, and Y. Weiss, "A closed-form solution to natural image matting," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 2, pp. 228–242, 2008.
- [37] K. He, J. Sun, and X. Tang, "Single image haze removal using dark channel prior," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 12, pp. 2341–2353, 2011.
- [38] A. Levin, D. Lischinski, and Y. Weiss, "A closed form solution to natural image matting," in *IEEE CVPR*, 2006, pp. 61–68.
- [39] Q. Gu, Z. Li, and J. Han, "Learning a kernel for multi-task clustering," in *Proc. Assoc. Adv. Artif. Intell. (AAAI)*, 2011.
- [40] Q. Shan, J. Jia, and M. Brown, "Globally optimized linear windowed tone mapping," *IEEE TVCG*, vol. 16, no. 4, pp. 663–675, 2010.
- [41] D. M. Cvetkovic, M. Doob, and H. Sachs, *Spectral of Graphs-Theory and Applications*. New York: Academic Press, 1980.
- [42] S. Xiang, C. Pan, F. Nie, and C. Zhang, "Turbopixel segmentation using eigen-images," *IEEE TIP*, vol. 19, no. 11, pp. 3024–3034, 2010.
- [43] S. Xiang, F. Nie, C. Pan, and C. Zhang, "Regression reformulations of l1e and l1sa with locally linear transformation," *IEEE Trans. on Systems, Man, and Cybernetics*, vol. 41, no. 5, pp. 1250–1262, 2011.
- [44] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *Royal Statistical Society*, vol. 39, no. 1, pp. 1–38, 1977.
- [45] H. Li, S. Tak, and J. C. Ye, "Lipschitz-killing curvature based expected euler characteristics for p-value correction in fnirs," *Journal of neuroscience methods*, vol. 204, no. 1, p. 61–67, February 2012.
- [46] L. Saul and F. Pereira, "Aggregate and mixed-order markov models for statistical language proc." in *Conf. on Empirical Methods in Natural Language Process*, 1997, pp. 81–89.
- [47] J. Nocedal and S. J. Wright, *Numerical Optimization*, 2nd ed. New York: Springer, 2006.
- [48] F. Zhu, B. Fan, X. Zhu, Y. Wang, S. Xiang, and C. Pan, "10,000+ times accelerated robust subset selection (ARSS)," in *Proc. Assoc. Adv. Artif. Intell. (AAAI)*, 2015. [Online]. Available: <http://arxiv.org/abs/1409.3660>
- [49] T. A. Davis, *Direct Methods for Sparse Linear Systems*. Philadelphia, PA: SIAM, 2006.
- [50] C. M. Bishop, *Pattern Recognition and Machine Learning*, ser. Inform. Science and Statistics. Springer, 2006.
- [51] MathWorks, "Reference available on <http://www.mathworks.cn/help/matlab/math/sparse-matrix-operations.html#brc1y0o>."
- [52] F. Zhu, "Hyperspectral unmixing datasets & ground truths on http://www.escience.cn/people/feiyunZHU/Dataset_GT.html," 2014.
- [53] F. Zhu, Y. Wang, B. Fan, G. Meng, and C. Pan, "Effective spectral unmixing via robust representation and learning-based sparsity," *arXiv*, vol. :1409.0685, 2014.
- [54] S. Jia and Y. Qian, "Constrained nonnegative matrix factorization for hyperspectral unmixing," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 1, pp. 161–173, 2009.
- [55] Envi-Tutorials, "An overview of hyperspectral remote sensing <http://www.cossa.csiro.au/hswwww/Overview.htm>," 2013.
- [56] K. Canham, A. Schlamm, A. Ziemann, B. Basener, and D. W. Messinger, "Spatially adaptive hyperspectral unmixing," *IEEE Trans. Geosci. Remote Sens.*, vol. 49, no. 11, pp. 4248–4262, 2011.
- [57] P. O. Hoyer and P. Dayan, "Non-negative matrix factorization with sparseness constraints," *Journal of Machine Learning Res. (JML)*, vol. 5, no. 12, pp. 1457–1469, 2004.