

EndNet: Sparse AutoEncoder Network for Endmember Extraction and Hyperspectral Unmixing

Savas Ozkan, *Member, IEEE*, Berk Kaya, and Gozde Bozdagi Akar, *Senior Member, IEEE*

Abstract—Data acquired from multi-channel sensors is a highly valuable asset to interpret the environment for a variety of remote sensing applications. However, low spatial resolution is a critical limitation for previous sensors and the constituent materials of a scene can be mixed in different fractions due to their spatial interactions. Spectral unmixing is a technique that allows us to obtain the material spectral signatures and their fractions from hyperspectral data. In this paper, we propose a novel endmember extraction and hyperspectral unmixing scheme, so called *EndNet*, that is based on a two-staged autoencoder network. This well-known structure is completely enhanced and restructured by introducing additional layers and a projection metric (i.e., spectral angle distance (SAD) instead of inner product) to achieve an optimum solution. Moreover, we present a novel loss function that is composed of a Kullback-Leibler divergence term with SAD similarity and additional penalty terms to improve the sparsity of the estimates. These modifications enable us to set the common properties of endmembers such as non-linearity and sparsity for autoencoder networks. Lastly, due to the stochastic-gradient based approach, the method is scalable for large-scale data and it can be accelerated on Graphical Processing Units (GPUs). To demonstrate the superiority of our proposed method, we conduct extensive experiments on several well-known datasets. The results confirm that the proposed method considerably improves the performance compared to the state-of-the-art techniques in literature.

Index Terms—Hyperspectral Unmixing, Endmember Extraction, Sparse Autoencoder.

I. INTRODUCTION

IN remote sensing, hyperspectral data is an essential imaging sensory output by which we gain insight into the Earth system by utilizing information beyond the human visible spectrum. This sensor type has been widely used in a variety of remote sensing applications from environmental monitoring to military surveillance for several decades [1], [2]. However, spatial resolution of hyperspectral sensors is very limited compared to other optics operating in visible spectrum domain. Hence, the pixels can be composed of mixture of material spectra. For further process of hyperspectral applications, the constituent materials (endmembers) and their fractions (abundances) should be determined from data.

The mixture of constituent materials is generally formulated with a linear model [3] wherein each data point on an image $\mathbf{x} \in \mathbb{R}^D$ is a linear combination of endmembers $\mathbf{E} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_K] \in \mathbb{R}^{D \times K}$ with different fractions

$\mathbf{y} = [y_1, y_2, \dots, y_K] \in \mathbb{R}^K$. Here, K denotes the number of endmembers and D is the spectral bands of the hyperspectral data. In addition, this formulation has an additional term η (it is assumed to possess a zero mean Gaussian noise) to simulate possible noise sources in the process such as sensor readout noise or illumination variability due to surface topography:

$$\mathbf{x} = \sum_{k=1}^K \mathbf{e}_k y_k + \eta = \mathbf{E}\mathbf{y} + \eta, \quad s.t. \quad y_k \geq 0, \quad \sum_{k=1}^K y_k = 1 \quad (1)$$

Unmixing of hyperspectral data is separated into two major steps (note that we assume that the optimum number of endmembers is known) as endmember extraction and quantifying the abundances of these endmembers per pixel. These two unknown variables can be solved either simultaneously or individually depending on the approaches.

In literature, linear and nonlinear models have been vastly studied and several promising methods have been proposed as follows.

For linear models, geometrical volume-based algorithms are quite common to identify endmembers from hyperspectral data [4], [5], [6], [7]. They treat the distribution of data samples as a simplex set [5]. They ultimately embrace the fact that the vertices of this simplex set correspond to the endmembers since all of the data samples can be spanned with these points.

By exploiting this assumption (i.e., simplex set) and the notion of the presence of pure pixels, the most straightforward linear solution is to determine unmixed pixels from data for all materials [7], [8]. Vertex Component Analysis (VCA)-like methods [4], [6] can extend the assumption by projecting and imposing an orthogonality condition onto endmembers in their estimations. However, these methods have severe drawbacks [9]. Briefly, the mixture of materials in both macroscopic and microscopic levels [9], [10] can lead to erroneous spectral estimates. The assumption can collapse with the lack of the purities of one or more materials in data. To solve these drawbacks, several concepts have been introduced in the literature [11], [12], [13], [14], [15] that basically add a margin to the derivation with extra constraints or a kernel structure. This produces more reliable outputs even if the data does not directly satisfy the purest material condition. However, particularly for [14], [15], presence of outliers and noise can significantly reduce the performance of the suboptimal solution of the methods [16].

Lastly, methods utilizing a codebook of endmembers (collected with a spectrometry and available in a spectral library) [9], [17], [18] are another solution to identify fractional

S. Ozkan is with the Image Processing Department, TUBITAK Space Technologies Research Institute and the Department of Electrical and Electronics Engineering, Middle East Technical University, Ankara, Turkey e-mail: savas.ozkan@tubitak.gov.tr

B. Kaya and G.B. Akar are with the Department of Electrical and Electronics Engineering, Middle East Technical University, Ankara, Turkey.

abundances and corresponding endmembers from hyperspectral data. The observation throughout our paper, the importance of sparsity is particularly emphasized in these studies.

Although the linear mixture model is simple and provides practical advantages, it is not adequate to solve the mixing problem. Multiple scattering effects [19], microscopic-level material mixtures [9], [10] and water-absorbed environments [20] can be shown as severe natural cases that cannot be handled by the linear models. One way to deal with such scenarios is to use unsupervised nonlinear projections of data [21], [22]. However, this type of approaches requires a high computational workload, thus it eventually reduces the scalability of the methods for large-scale data. An alternative strategy is to replace the conventional operations with a nonlinear kernel function $f(\cdot)$ such that $\mathbf{x} = f(\mathbf{E}\mathbf{y}) + \eta$. This methodology simply introduces a distortion onto endmembers and abundance estimates to enhance the robustness against nonlinear interactions [23], [24], [25], [26], [27], [28], [29], [30], [31], [32].

Leveraging supervised data with neural-network architectures is another solution proposed for this problem [33], [34], [35], [36], [37]. Although critical discussions are presented for optimum training sample selection, the quality of training samples drastically influences the performance. Also, data labeling is costly and impractical compared to the unsupervised approaches. The endmembers and fractional abundances should be blind. A recent study [38] proposes a blind endmember extraction method based on neural networks by introducing a set of layer components to the encoder layer of an autoencoder. This approach particularly helps to increase the sparsity of fractional abundance estimates. However, no detailed experimental result is presented to validate the superiority of neural networks over the conventional methods and the model is currently suboptimal due to the tendency of overfitting.

The adaptation of neural network structure to hyperspectral unmixing has not been completely solved in previous works, especially for unsupervised setup. To the best of our knowledge, our proposed method will be the first successful attempt based on an autoencoder network that outperforms the conventional methods, all in an unsupervised manner.

Our Contribution: In this study, we propose a novel two-staged neural network autoencoder and an end-to-end learning scheme that are specialized to extract endmembers with their fractional abundances from hyperspectral data. Our method uses some of the major approaches previously introduced for hyperspectral unmixing, however it significantly differs from the traditional autoencoder pipeline with novel features.

Similar to [14], [15], we observed that rather than strictly parameterized all of the physical properties of the endmembers, solving the problem with some of the constraints can yield better results. Since the spectral angle distance (SAD) is used in our loss function which enforces the maximum angular similarity between samples rather than minimizing the Euclidian distance, the value range of the estimated endmembers is not constrained to [0,1]. However the negative values for endmembers are penalized by a SAD metric in the loss function, thus non-negativity constraint is automatically

applied in the learning process. Other constraints such as sum-to-one constraint (i.e., $\sum_{k=1}^K y_k = 1$) and greater than zero constraint (i.e., $y_k \geq 0$) are preserved because of the internal characteristics of the layer components in the architecture which will be explained in the following section.

In the proposed method, first, sparsity and non-linearity are effectively applied by a rectified activation function, ReLU (i.e., bounds the negative responses of hidden abstracts) [39] and a normalization layer [40] as in [38]. This adaptation enables us to set the bias terms to zero in the formulation and it eliminates the adverse effects of bias terms in the solution. However, the complete sparsity is still missing and it is prone to overfitting due to the fact that the network tends to generate hidden abstract (i.e., hidden representation of the neural network model) constantly for highly correlated materials. In other words, a small number of responses should contribute to the composition of a pixel for a practical solution. In our paper, this is achieved by hard-response selection (i.e., i.e., selection of only top responses of the hidden abstracts) and regularization techniques [41], [42] which allow only the highest responses of the fractional abundances contribute to the reconstruction stage. To this end, this leads to a better solution and improves the sparsity for fractional abundances.

Second, we replace the inner product operators at the encoder layer with spectral angle distance (SAD) to obtain more discriminative hidden abstracts. Third, we introduce an extra set of penalty terms to the loss function in which they can enforce closer angular similarity between the original and the reconstructed data along with the standard Euclidean (l_2) reconstruction penalty. Moreover, sparse hidden representations can be achieved with additional penalty terms. Fourth, in contrast to some conventional unmixing approaches [4], [7], [26], [29], [30], our proposed method solves the problem by optimizing endmembers and their corresponding fractional abundances concurrently with a stochastic gradient-based solver. This leads to the optimum solutions for both endmembers and fractions while it scales the method for large-scale data as stated in [81].

Lastly, we observed that usage of VCA-like [4], [7], [29] methods as an autoencoder filter initializer helps us to attain state-of-the-art performance. This makes the filters converge faster to the global optimum. Note that this step is practiced previously in different optimization-based algorithms [15] and it similarly avoids the parameters which leads to poor quality solutions.

The remainder of our paper is organized as follows. First, necessary information about the proposed architecture and the optimization step are presented in Section 2. Section 3 is reserved for the results and discussions on the experiments conducted on several publicly available datasets. Conclusion and overall discussions are presented in Section 4.

II. SPECTRAL UNMIXING

Motivation: A neural network autoencoder is simply composed of two consecutive information processing layers. Encoder layer transforms input samples to hidden abstracts (hidden abstracts correspond to the responses of hidden layers)

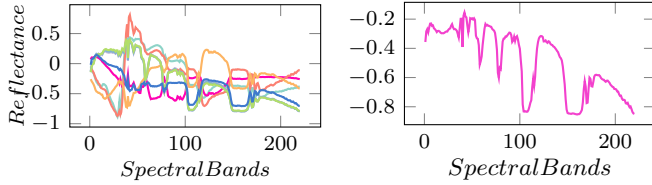


Fig. 1. Estimates of a conventional autoencoder for hyperspectral unmixing problem. First plot (left side) shows decoder layer parameters, $\mathbf{W}^{(e)}$, which intuitively correspond to the endmembers while second plot (right side) illustrates bias term response. As can be seen from the results, bias term acts like one of the endmembers exhibited from the scene. Moreover, the estimated signatures (left side) do not satisfy the physical conditions of endmembers, i.e., $y_k \geq 0$.

while the decoder layer tries to restore the original inputs from these hidden abstracts with high accuracy. Ultimately, latent correlations in data are unveiled in this chain of the transformations and conserved as trainable parameter sets in the model [43], [44], [45], [46].

For a given sample $\mathbf{x} \in \mathbb{R}^D$, an autoencoder initially maps it into a hidden representation $\mathbf{y}^1 \in \mathbb{R}^K$ with the inner product of a trainable parameter set \mathbf{W} , \mathbf{b} . Following this, the reconstruction of the input sample $\hat{\mathbf{x}} \in \mathbb{R}^D$ is recomputed:

$$\mathbf{y} = f(\mathbf{W}^{(e)}\mathbf{x} + \mathbf{b}^{(e)}), \quad (2a)$$

$$\hat{\mathbf{x}} = f(\mathbf{W}^{(d)}\mathbf{y} + \mathbf{b}^{(d)}). \quad (2b)$$

Here $f(\cdot)$ indicates an element-wise nonlinear activation function and a logistic activation function is frequently selected as either sigmoid or tanh.

Furthermore, some structural variations are observed for activation functions and/or parameter sets in literature [44], [45], [46], [47]. More precisely, the activation function $f(\cdot)$ at the decoder layer can be discarded or same/disjoint parameter sets can be used for better error propagation.

Finally, the parameters are optimized by minimizing the standard Euclidean (l_2) reconstruction error of the original and the reconstructed samples:

$$\mathcal{L} = \frac{1}{2} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2. \quad (3)$$

where \mathcal{L} is the loss function. Fig. 1 presents the possible parameter outcomes of a conventional autoencoder learned by Eqs.(2a), (2b) and (3) on the University of Pavia dataset [48].

Current Limitations For Unmixing Problem: The direct use of a traditional autoencoder might not yield a stable solution particularly for several reasons: 1) Logistic functions tend to generate redundant responses even if the input samples are not correlated with the parameter set [39], [49], [38]. This also violates and decreases the sparsity of true projections. 2) The inner product is not sufficiently discriminative and it can lead to erroneous fractional abundances for data. 3) Although the standard Euclidean (l_2) reconstruction term is the main penalty function used to minimize the estimation error, it should be carefully regularized with extra terms related to the domain in

order to reach to the optimum solution. Otherwise, the solution can be stuck to a local minimum.

For these reasons, these current limitations aggravate conventional autoencoders to attain state-of-the-art performance on hyperspectral unmixing problem.

A. Sparse Autoencoder for Hyperspectral Unmixing

The primary objective of this paper is to accomplish the adaptation of an autoencoder pipeline to the endmember extraction for hyperspectral data. In particular, we strive to preserve the estimation power of autoencoders on latent data correlations as well as the observations that are previously introduced for endmember extraction.

To enhance the performance and make the autoencoder effective for the problem, we made several modifications in the architecture and the optimization step² as follows. First, bias terms are removed from the formulations due to their adverse effects in the problem. Second, we replace and introduce additional layers to the architecture to improve sparsity as well as consistency. Later, inner product is replaced with spectral angle distance (SAD) to estimate more discriminative fractional abundances. Finally, we propose a novel loss function to optimize the model parameters effectively. In the following subsections, we describe these modifications in details.

Zero-biased Filters: In neural network architecture, bias term is one of the critical parameters, especially for linear regression models [50], since it determines the influence of parameters for next layers by thresholding the responses through activation functions. However, bias terms violate one of the crucial properties of endmembers by which none of the materials can be constantly expected in every pixel compositions, if a scene is not composed of one material.

More clearly, since bias terms constantly affect the composition of pixels, they behave like one of the endmembers observed from the data at the end of the optimization step. Fig. 1. plots a possible bias term response for a conventional autoencoder. For this reason, we removed the bias terms in both Eqs.(2a) and (2b) for our model. Also, this step partially retains the simplex set assumption (i.e., an affine projection) which constitutes a basis for the conventional methods in the literature [4], [5], [6].

Moreover, we discarded the activation function at the decoder layer to propagate the error more effectively and intuitively $\mathbf{W}^{(d)}$ parameter begins to represent endmembers \mathbf{E} estimated from the scene. Similarly, hidden abstract \mathbf{y} corresponds to the fractional abundances of a pixel.

Lastly, since abundance estimates possess non-linear relations and cannot be computed directly from endmembers $\mathbf{W}^{(d)}$, we used disjoint parameter sets for both encoder and decoder layers ($\mathbf{W}^{(e)} \in \mathbb{R}^{K \times D}$ and $\mathbf{W}^{(d)} \in \mathbb{R}^{D \times K}$).

Sparsity and Nonlinearity: Sparsity and nonlinearity are critical features for endmember extraction and hyperspectral unmixing [9], [17], [25], [26] in order to make a robust estimation. As previously mentioned limitations, logistic functions are insufficient to conserve the sparsity for the architecture

¹ \mathbf{y} intuitively corresponds to fractional abundances per pixel in hyperspectral unmixing.

²Source code and presented results will be available on <https://github.com/savasozkan/endnet>

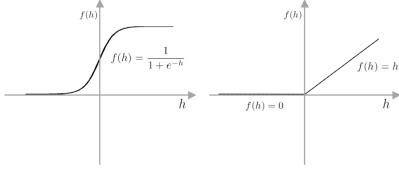


Fig. 2. Response characteristics of Sigmoid (left) and ReLU (right) activation functions.

even if they generate nonlinear responses. Therefore, we made three critical modifications in the architecture to achieve sparse, consistent and nonlinear hidden abstracts.

First, a rectified linear activation [39] is used at the encoder layer instead of logistic functions. Fig. 2 illustrates the response characteristics of the Sigmoid and Rectified Linear Unit (ReLU) activations for different input values. ReLU bounds the negative influence of filters and increases the selectivity of the representation. In addition, it boosts the parameters and allows them to saturate more quickly in the optimization step. In the scope of this paper, we tested our proposed method with the variants of rectified activation units [49], [51], [52] to be sure about the optimum solution. These activation units permit limited amount of negative responses to the subsequent layers. However, these negative values ultimately perturb the relation of endmembers (i.e. they become more correlated) and violates the greater than zero constraint (i.e., $y_k > 0$) for the abundance estimates.

However, trainable parameters can quickly become ill-posed (i.e. small changes in the parameters can lead to oscillation and instability in the model) with the removal of the logistic function [40], [53]. Therefore, we introduced a normalization layer [40] before the ReLU activation function at the encoder layer. For our model, this layer especially induces selectivity on the top activation responses (it has a similar objective as a bias term [50]) in addition to the mitigation of ill-posed effects. It normalizes the responses by scaling and reordering them. Fig. 3 illustrates the possible output distribution characteristics of this normalization layer. If we assume that shifting parameter ρ in Eq.(4) is discarded, by combining with ReLU layer, it would permit roughly 50% of the top positive responses to the next layer at each iteration:

$$\text{BN}(\mathbf{h}) = \frac{(\mathbf{h} - \boldsymbol{\mu})}{\sqrt{\boldsymbol{\sigma}^2 + \epsilon}} + \rho, \quad (4)$$

where $\text{BN}(\cdot)$ is the output of the normalization layer and it is fed into ReLU layer in Eq.(5). ϵ is a very small

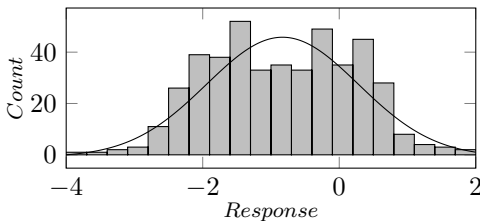


Fig. 3. Distribution characteristic of batch normalization layer outputs. Due to shifting parameter ρ , mean value of the distribution slightly moves to left in the plot.

constant (10^{-8}) to prevent zero-division in the formulation. $\mathbf{h} = \mathbf{W}^{(e)}\mathbf{x}$, $\mathbf{h} \in \mathbb{R}^K$ is the filter response for input \mathbf{x} (We should note that we will redefine this filter response in the following subsection. This formulation is given here to preserve the consistency). $\boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$, $\boldsymbol{\mu} \in \mathbb{R}^K$ and $\boldsymbol{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \boldsymbol{\mu})^2$, $\boldsymbol{\sigma}^2 \in \mathbb{R}^K$ are the mean and variance of mini-batch filter responses at t^{th} iteration. $\boldsymbol{\rho} \in \mathbb{R}^K$ are the trainable parameters to adjust the permutability of filter responses. We should remark that scaling factor in [40] is discarded due to better fractional abundances and endmember extraction:

$$f(\mathbf{h}) = \begin{cases} \text{BN}(\mathbf{h}), & \text{if } \text{BN}(\mathbf{h}) > 0 \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

For ReLU activation and normalization layers, similar observation (i.e., sparsity and ill-posed artifacts) can be found in the very recent work [38].

However, autoencoders might not still have the optimum solution in some cases due to parameter overfitting (i.e. lack of full parameter convergence and sparsity). In order to make improvements for reliable abundance estimates and endmembers, we utilized three additional steps at the end of the activation function. For this purpose, we first employed a regularization layer, i.e., Dropout [41], [42] with $\mathbf{r} \sim \text{Bernoulli}(p)$, after the activation function $f(\cdot)$ as $\mathbf{z} = \mathbf{r} * f(\mathbf{h})$. Ultimately, for each iteration, the network introduces a new solution for the problem by randomly ignoring some of the activations so that the generalization capacity of the network is enhanced. Here $*$ denotes an element-wise product and $0 < p \leq 1$ should be defined by the user based on the spectral correlations of materials in a scene (i.e if highly correlated materials exist in the scene, p should be close to 0.5). Note that the default value of p is set to 1.0.

Second, we hardly selected the top n activations (i.e., highest ones) from \mathbf{z} as in [54] (n is fixed to 2 due to the optimum spatial material mixture) to increase the selectivity of $\mathbf{W}^{(d)}$. Finally, we applied $l1$ normalization to satisfy the sum-to-one constraint (i.e., $\sum_{k=1}^K y_k = 1$) on abundance estimates:

$$\mathbf{y} = \frac{\mathbf{z}^*}{(\|\mathbf{z}^*\|_1 + \epsilon)}, \quad (6)$$

where \mathbf{z}^* denotes n -top activations. To learn the parameters in our model, we use backpropagation algorithm [55] and Eq.(6) is differentiable except at zero (For this case, partial derivative is directly set to zero). From chain rule, the partial derivative of the $l1$ normalization layer is computed as follows:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^*} = \frac{1}{\|\mathbf{z}^*\|_1} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{y}} - \mathbf{y} \sum_{k=1}^K \frac{\partial \mathcal{L}}{\partial y_k} \text{sign}(z_k^*) \right). \quad (7)$$

where \mathcal{L} is the loss function. The term $\frac{\partial \mathcal{L}}{\partial \mathbf{y}}$ defines the propagated error up to the hidden abstracts \mathbf{y} from the upper layer. Lastly, $\text{sign}(\cdot)$ denotes a function which returns the sign of its input.

Discriminative Hidden Abstracts: For neural network autoencoders, input data samples need to be mapped to hidden abstracts truly by storing their latent correlations, in order to recompute the original ones with high accuracy. This is why,

an accurate mapping is of a critical requirement to reveal these latent correlations from data.

Except in minor instances [56], the inner products of input samples with filter parameters have produced the best performance for almost all applications in the visible domain [44], [49], [53]. For the hyperspectral domain, there are various useful operators which here prove their success on several tasks such as classification and spectral signature comparison [57], [58], [59], [60], [61].

Fig. 4 plots the similarity score distributions of the normalized inner product (first row) and the spectral angle distance (SAD) (second row) of same (blue) and very coherent (red) material spectral signatures. From these plots, it is observed that inner product constantly produces high scores even if different materials are compared. However, the same material scores for spectral angle distance (SAD) can be separated from the coherent ones. In addition, [61] explains that SAD is better for separability, while spectral information divergence (SID) is more practical for preserving spectral patterns. Note that separability and discriminative power are especially important to estimate sparse responses at the encoder layer.

For these reasons, we replaced inner product with SAD at the encoder layer to obtain more discriminative and separable hidden abstracts as well as endmembers from hyperspectral data.

Spectral angle distance (SAD) measures the spectral angle between two input samples and the score closer to zero implies higher correlation. One of the concerns addressed for the SAD metric [58], [62] is that it is limitation to the nonlinear cases. As a remedy, several improvements are presented in the literature and most of which are based on kernel-based methods [58], [59], [60]. However, this limitation is not a drawback for our case, since nonlinearity is supplied with the nonlinear activation function [39] and the normalization layer [40] at the output of this metric.

SAD computes the similarity score of two samples, $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$, as follows:

$$S(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \cos^{-1}(\theta(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})), \quad (8)$$

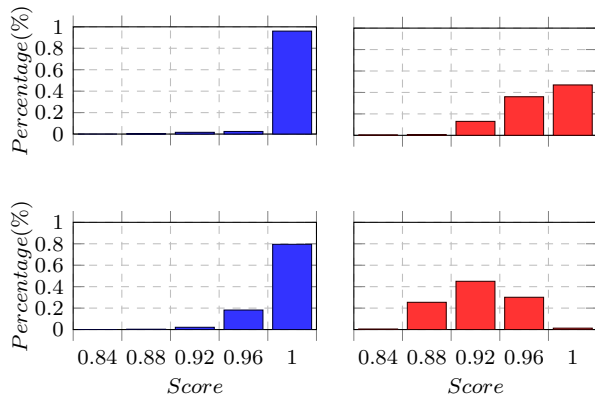


Fig. 4. Normalized similarity score distributions for same (blue) and very coherent (red) class comparisons with inner product (first row) and spectral angle distance (second row). As can be seen from the plots, spectral angle distance obtain more discriminative results.

where

$$\theta(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \frac{\mathbf{x}^{(i)} \mathbf{x}^{(j)}}{\|\mathbf{x}^{(i)}\|_2 \|\mathbf{x}^{(j)}\|_2}. \quad (9)$$

In Eq.(8), the SAD scores vary between $[0, \pi]$ and zero value means that two samples are identical. Therefore, we first need to normalize the score $S(\cdot, \cdot)$ to $[0, 1]$, since a higher score means a higher similarity at the decoder layer:

$$C(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = 1.0 - \frac{S(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})}{\pi}. \quad (10)$$

From this moment on, we reformulate the filter responses as $\mathbf{h} = C(\mathbf{x}, \mathbf{W}^{(e)})$. SAD is differentiable to optimize necessary parameters with the backpropagation algorithm by minimizing the loss function \mathcal{L} . By applying chain rule again, the gradient of k^{th} row \mathbf{w}_k of $\mathbf{W}^{(e)}$, $\mathbf{w}_k \in \mathbb{R}^{1 \times D}$ can be computed as:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_k} = \frac{\partial \mathcal{L}}{\partial C_k} \frac{\partial C_k}{\partial S_k} \frac{\partial S_k}{\partial \theta_k} \frac{\partial \theta_k}{\partial \mathbf{w}_k}, \quad k = 1, 2, \dots, K. \quad (11)$$

where C_k , S_k and θ_k are the simplifications of $C(\mathbf{x}, \mathbf{w}_k)$, $S(\mathbf{x}, \mathbf{w}_k)$ and $\theta(\mathbf{x}, \mathbf{w}_k)$ respectively to ease the formulation and understandability. Also, each derivation term in Eq.(11) can be written as:

$$\frac{\partial C_k}{\partial S_k} = \frac{-1}{\pi}, \quad (12a)$$

$$\frac{\partial S_k}{\partial \theta_k} = \frac{-1}{\sqrt{1 - \theta_k^2}}, \quad (12b)$$

$$\frac{\partial \theta_k}{\partial \mathbf{w}_k} = \frac{\mathbf{x}}{\|\mathbf{w}_k\|_2 \|\mathbf{x}\|_2} - \frac{\mathbf{w}_k * (\mathbf{w}_k \mathbf{x})}{\|\mathbf{w}_k\|_2^3 \|\mathbf{x}\|_2}. \quad (12c)$$

In the next subsection, we will analyze the optimization step of the proposed method by introducing modifications in the objective function.

Learning: Parameter optimization aims to minimize the error between input samples and their reconstructed versions by learning useful latent correlations from data. However, we found out that use of the Euclidean (l_2) norm as a primary objective term to unveil these correlations is not completely adequate for the problem. The main limitation of Euclidean norm is that it aggravates the method by estimating inappropriate/underestimated endmembers under severe illumination changes and nonlinearity (i.e., spectral variability) [63]. On the other hand, the spectral angle distance-like (SAD) operator overcomes these limitations and improves unmixing performance by exploiting geometric features of samples as explained in [30], [64]. In addition, smoothness and sparsity priors should be considered in the loss function for better parameter convergence. For this purpose, we made a critical set of modifications in the objective function related to this domain and a novel objective function \mathcal{L} is reformulated as:

$$\mathcal{L} = \frac{\lambda_0}{2} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 - \lambda_1 D_{KL}(1.0 \| C(\mathbf{x}, \hat{\mathbf{x}})) + \lambda_2 \|\mathbf{z}\|_1 + \lambda_3 \|\mathbf{W}^{(e)}\|_2 + \lambda_4 \|\mathbf{W}^{(d)}\|_2 + \lambda_5 \|\boldsymbol{\rho}\|_2 \quad (13)$$

where $C(\cdot, \cdot)$ is the normalized SAD score between the original and reconstructed version as in Eq.(10). λ_0 is the term that controls the influence of Euclidean norm. Additionally, we add a Kullback-Leibler divergence term (D_{KL}) [65], [66]



Fig. 5. Four real hyperspectral datasets used in the experiments: Urban, Samson, Jasper Ridge and Cuprite respectively. These datasets are broadly used to evaluate the performance of a method for endmember extraction and abundances unmixing of hyperspectral data.

to maximize SAD score distributions between the original and the reconstructed samples. λ_1 determines the effect of this cost term to the objective function. Note that Euclidean norm term is still critical for stable estimates and it cannot be directly set to zero, contrary to [38]. Thus, we empirically set λ_0 and λ_1 to 0.01 and 10 respectively.

As stated, the sparsity is important for effective hyperspectral unmixing. Therefore, we introduce a l_1 regularization term $\|\mathbf{z}\|_1$ as in [67] which penalizes the hidden layers that constantly generate responses for different input samples. Thus, the sparsity is enforced to the hidden abstracts and more distinct filter parameters are obtained for each material. The influence of l_1 sparsity term λ_2 is tuned to 0.1 which is a relatively large value. But, we should note that this value needs to be increased/decreased depending on the mixture level of the scenes.

We also add l_2 smoothing terms to the objective function for each trainable parameter and we define their values as 10^{-5} , 10^{-5} and 10^{-3} for λ_3 , λ_4 and λ_5 respectively.

The model parameters are optimized with a gradient-based stochastic Adam [68] optimizer by minimizing the loss function. We fixed the learning rate and number of training iterations to 0.001 and 400K respectively for all datasets (These terms can be still tuned). Additionally, mini-batch size N is set to 64 to balance the accuracy and the complexity at each iteration. Unlike the suggested values for momentum terms as in [68], we empirically found out that β_1 value should be defined as 0.7 to reduce the oscillation and instability of learning in the model.

Lastly, the parameters ($\mathbf{W}^{(e)}$ and $\mathbf{W}^{(d)}$) are initialized with the estimates of VCA-like methods [4], [7], [29] instead of a random initialization. This approach provides several advantages: First, it is practical to start the optimization from a more reliable initialization which can span all data points to decrease the possibility of parameter overshooting (if any geometrical volume-based algorithm is used). Second, it speeds up the convergence of the parameters.

Prevent Autoencoder From Overfitting: As mentioned, we utilized Dropout [41], [42] at the hidden layer to prevent hidden abstracts from overfitting to a poor quality solution. In addition, we employ denoising autoencoder scheme [43] in the parameter optimization step to be robust to noise exhibited on data (i.e. approximation error η in the process). Even though the primary principles for both methods are similar, we used these methods at different layers to generalize our model parameters.

This scheme is essentially based on the fact that data

samples are initially corrupted with an additive Gaussian noise on purpose while the reconstruction error is computed over the original ones. This helps to improve the generalization capacity of the encoder parameter set and it leads to better solution for robust representations by considering possible deformations (i.e. illumination changes, sensor noise etc.) in advance.

In practice, we opt to utilize a corruption process that does not hurt hidden abstracts significantly. This is particularly important since SAD is sensitive to large variations. For this purpose, isotropic Gaussian noise ($\tilde{\mathbf{x}}|\mathbf{x} \sim \mathcal{N}(\mathbf{x}|0, \sigma^2)$) and mask noise (i.e. it randomly chooses a portion of the elements) are jointly used to perturb the data. Additionally, we limit the mask noise to alter at most 40% of the original elements with the additive noise. However, the value of this parameter can be reduced for noisy data [48]. Fig. 6 visualizes the original spectral signatures and their corrupted versions under this assumption.

B. Fractional Abundance Estimation

After the endmember extraction, we need to find the fractional abundances of each material on data pixels. In our proposed method, we can obtain these values in two different ways.

- We can use the hidden abstracts \mathbf{y} of our model for each sample.
- Also, we can solve an inverse problem with the estimated endmembers $\mathbf{W}^{(d)}$ since the decoder layer is linear.

As explained in [69], due to the averaging operations for mean and variance values in batch normalization, hidden abstracts tend to generate different responses in test time and this can affect the accuracy of estimates (i.e., depending on the distribution of samples per material in the scene.). Therefore, we adapt the Simplex Projection Unmixing (SPU) [29], [70] algorithm with SAD kernel for abundance estimation. This algorithm estimates the abundances by using the mutual distances of material spectral signatures and data samples. It also assumes that nonnegativity and sum-to-one constraints are preserved in the nonlinear solution.

In particular, the estimation of the abundances by solving an inverse problem with $\mathbf{W}^{(d)}$ and SPU empirically introduces further improvements to the performance, even though estimated hidden abstracts \mathbf{y} still yield compatible results. Thus, the results for the fractional abundances are reported by the combination of these methods for our proposed method.

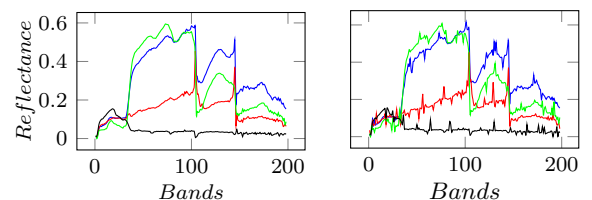


Fig. 6. Four samples corrupted with the denoising criterion. Different color is used for each hyperspectral signature pair.

TABLE I
SAD AND RMSE RESULTS ON URBAN DATASET. MEAN AND STANDARD DEVIATION ARE REPORTED. BEST RESULTS ARE ILLUSTRATED IN BOLD.

Endm.	Spectral Angle Distance (SAD) ($\times 10^{-2}$)									
	VCA	DMaxD	MVSA	SPICE	SCM	$l_{1 2}$ -NMF	DgS-NMF	EndNet-VCA	EndNet-DMaxD	
#1	21.56 \pm 3.1	13.16 \pm 0.0	19.04 \pm 0.1	45.04 \pm >9	9.38 \pm 0.0	6.06 \pm 0.2	5.86 \pm 0.1	6.72 \pm 0.2	6.88 \pm 0.2	
#2	39.04 \pm 5.1	>99 \pm 0.0	>99 \pm 1.3	>99 \pm >9	10.55 \pm 0.0	20.05 \pm 0.4	13.69 \pm 0.4	4.04 \pm 0.1	3.92 \pm 0.3	
#3	26.77 \pm 7.5	7.43 \pm 0.0	37.83 \pm 9.4	30.89 \pm >9	14.03 \pm 0.0	3.71 \pm 0.0	4.12 \pm 0.0	3.68 \pm 0.2	3.53 \pm 0.1	
#4	82.39 \pm 0.1	21.74 \pm 0.0	18.41 \pm 1.9	66.37 \pm >9	41.86 \pm 0.0	14.22 \pm 0.2	10.54 \pm 0.3	3.99 \pm 0.7	3.35 \pm 0.5	
Avg.	41.77 \pm 4.5	37.88 \pm 0.0	43.81 \pm 3.1	60.57 \pm >9	18.79 \pm 0.0	11.01 \pm 0.2	8.55 \pm 0.2	4.60 \pm 0.3	4.42 \pm 0.3	
	Root Mean Square Error (RMSE) ($\times 10^{-2}$)									
	VCA	DMaxD	MVSA	SPICE	SCM	$l_{1 2}$ -NMF	DgS-NMF	EndNet-VCA	EndNet-DMaxD	
#1	31.31 \pm 7.2	27.08 \pm 0.0	27.69 \pm 0.3	29.43 \pm 3.4	32.79 \pm 0.0	14.77 \pm 0.1	13.18 \pm 0.1	10.24 \pm 0.1	10.41 \pm 0.2	
#2	41.98 \pm 5.6	48.99 \pm 0.0	31.07 \pm 0.0	35.32 \pm 4.5	36.25 \pm 0.0	16.16 \pm 0.2	12.95 \pm 0.0	13.01 \pm 0.3	12.24 \pm 0.3	
#3	27.93 \pm 3.7	28.07 \pm 0.0	27.27 \pm 0.4	26.33 \pm 1.9	32.61 \pm 0.0	12.65 \pm 0.2	9.57 \pm 0.1	8.69 \pm 0.3	8.35 \pm 0.3	
#4	21.91 \pm 2.1	21.59 \pm 0.0	21.13 \pm 0.4	21.93 \pm 6.9	32.86 \pm 0.0	6.90 \pm 0.1	6.27 \pm 0.0	6.13 \pm 0.2	5.92 \pm 0.1	
Avg.	30.79 \pm 4.7	31.43 \pm 0.0	26.79 \pm 0.4	28.25 \pm 4.2	33.59 \pm 0.0	12.62 \pm 0.1	10.49 \pm 0.1	9.51 \pm 0.2	9.23 \pm 0.2	

Finally, the proposed method should be considered as a spectral unmixing method as well as endmember extractor, since it optimizes the loss function jointly. We conduct additional experiments to analyze the abundance performance and a discussion for these comparisons can be found in Section III. C.

III. EXPERIMENTS

Here, we demonstrate the performance of the proposed method on several datasets. For this purpose, we report quantitative and qualitative results of our proposed method on the hyperspectral unmixing problem. The estimated spectral signatures and the fractional abundances are compared with the corresponding ground truth.

In the beginning of this section, we will summarize the details of publicly available datasets, baseline methods and evaluation metrics that we used in the experiments. Lastly, we summarize the parameter settings determined in the experiments.

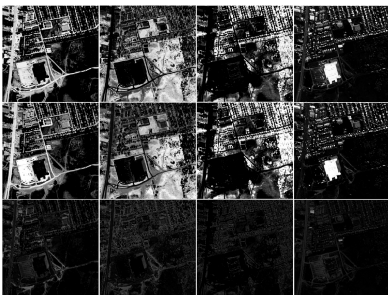


Fig. 7. Fractional abundances on Urban dataset with EndNet-DMaxD method. Each column corresponds to four materials: Asphalt, Grass, Tree and Roof respectively. Also, each row indicates estimated abundance results, ground truth and their absolute differences.

A. Hyperspectral Datasets

In order to make fair comparisons, we evaluate the proposed method both on synthetic data and on well-known

and extensively used real datasets for endmember extraction [71], [72], [73], [74], [48], [75], [76], [77], [8]. Even though synthetic data is able to measure the true fractional abundances/endmembers, it is quite hard to simulate some cases such as extreme non-linearity as in real data (even if Hapke Model or Bilinear Mixture Model (BMM) is used). In Table III, a base experiment is performed on a synthetic dataset (i.e., we utilized Spheric dataset from the IC Synthetic Hyperspectral Collections. To simulate the non-linear case, the noisy version with 40 db is used). It is clear that conventional models can even achieve ideal performance and the differences between the methods are quite small. Therefore, all the remaining experiments are conducted on real datasets throughout the paper to show the actual effectiveness of the methods particularly against the spectral variability. (Additional quantitative results are on the project website.)

Urban [71], [72]: The pixel resolution of data is 307×307 . Several channels (1-4, 76, 87, 101-111, 136-153 and 198-210) are removed due to water-vapor absorption and atmospheric effects. There are four constituent materials: Asphalt (#1), Grass (#2), Tree (#3), Roof (#4).

Samson [72]: Data is generated by the SAMSON sensor. It contains 156 channels and covers the spectral range from 0.4 to $0.9\mu\text{m}$. In order to make fair comparisons, a subimage of

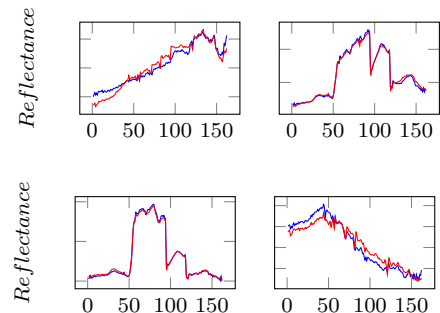


Fig. 8. Endmember signatures (Asphalt, Grass, Tree and Roof) on Urban dataset for ground truth (blue) and EndNet-DMaxD (red).

TABLE II

SAD AND RMSE RESULTS ON SAMSON DATASET. MEAN AND STANDARD DEVIATION ARE REPORTED. BEST RESULTS ARE ILLUSTRATED IN BOLD.

Endm.	Spectral Angle Distance (SAD) ($\times 10^{-2}$)								
	VCA	DMaxD	MVSA	SPICE	SCM	$l_{1 2}$ -NMF	DgS-NMF	EndNet-VCA	EndNet-DMaxD
#1	22.33 \pm 2.7	4.04 \pm 0.0	4.19 \pm 0.0	45.37 \pm >9	1.52 \pm 0.0	6.21 \pm 7.3	5.64 \pm 7.4	1.32 \pm 0.2	1.29 \pm0.1
#2	4.90 \pm 0.2	2.19 \pm0.0	5.61 \pm 0.0	12.19 \pm >9	3.79 \pm 0.0	5.23 \pm 0.3	4.80 \pm 0.3	4.72 \pm 0.2	4.69 \pm 0.1
#3	12.29 \pm 0.0	13.04 \pm 0.0	18.15 \pm 0.0	9.24 \pm >9	20.21 \pm 0.0	11.97 \pm 2.1	4.7 \pm 0.3	3.36 \pm 0.2	2.95 \pm0.3
Avg.	13.17 \pm 1.0	6.42 \pm 0.0	9.31 \pm 0.0	22.26 \pm >9	8.51 \pm 0.0	7.80 \pm 3.2	5.05 \pm 2.7	3.13 \pm 0.2	2.98 \pm0.2
	Root Mean Square Error (RMSE) ($\times 10^{-2}$)								
	VCA	DMaxD	MVSA	SPICE	SCM	$l_{1 2}$ -NMF	DgS-NMF	EndNet-VCA	EndNet-DMaxD
#1	16.92 \pm 6.3	15.08 \pm 0.0	19.66 \pm 0.0	27.76 \pm >9	18.65 \pm 0.0	8.58 \pm 3.3	7.77 \pm 3.8	5.86 \pm 0.0	5.72 \pm0.0
#2	16.64 \pm 3.3	22.08 \pm 0.0	23.81 \pm 0.0	28.16 \pm 2.8	21.65 \pm 0.0	7.44 \pm 3.7	7.74 \pm 3.6	4.15 \pm 0.1	3.84 \pm0.1
#3	25.42 \pm 0.6	26.41 \pm 0.0	31.78 \pm 0.0	35.57 \pm 3.4	34.61 \pm 0.0	5.55 \pm 0.9	2.70 \pm 0.9	2.07 \pm0.0	2.11 \pm 0.0
Avg.	19.66 \pm 3.2	21.19 \pm 0.0	25.08 \pm 0.0	30.49 \pm 5.4	24.97 \pm 0.0	7.19 \pm 2.4	6.07 \pm 2.8	4.01 \pm 0.0	3.88 \pm0.0

the original data (95×95) is considered as in [72]. Three material types are observed in the scene: Soil (#1), Tree (#2) and Water (#3).

Jasper Ridge [73]: Data is captured by the AVIRIS sensor. A subimage of 100×100 pixels of the original data is used. Some of the channels (1-3, 108-112, 154-166 and 220-224) are discarded due to the atmospheric effects and water-vapor absorption. Tree (#1), Water (#2), Soil (#3) and Road (#4) are observed in the scene.

Cuprite [71], [72], [74]: Data is captured with the AVIRIS over Cuprite, Nevada. It has 188 spectral reflectance bands covering the wavelength range from 0.4 to 2.5 μm . It has a nominal ground resolution of 20 m and a spectral resolution of 10 nm. A subimage of the original data [72] (250×190 pixels) is considered. Noisy (1-2 and 221-224) and water-vapor absorption (104-113 and 148-167) channels are also removed from data. This dataset hosts 12 unique mineral spectral signatures: Alunite (#1), Andradite (#2), Buddingtonite (#3), Dumortierite (#4), Kaolinite₁ (#5), Kaolinite₂ (#6), Muscovite (#7), Montmorillonite (#8), Nontronite (#9), Pyrope (#10), Sphene (#11) and Chalcedony (#12).

University of Pavia [48]: Data is recorded by the ROSIS sensor over Pavia, Italy. The spectral band number is 103 and the spectral range is varied from 0.43 to 0.86 μm . The spatial pixel resolution is 610×340 and the ground resolution is approximately 1.3 m. It comprises 9 labeled classes that covers different man-made structures and natural objects. Since bricks-gravel and asphalt-bitumen have similar spectral signatures, we considered these classes as joint classes [77]. In

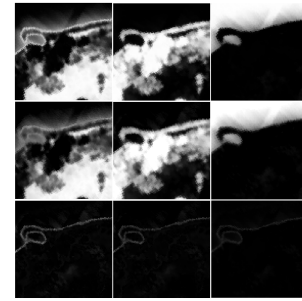


Fig. 9. Fractional abundances on Samson dataset with EndNet-DMaxD method. Results for three materials, Soil, Tree and Water, are reported in different columns. Estimated fractional abundances, ground truth and their absolute differences are illustrated in columns respectively.

the experiment, 7 labeled classes are used: Asphalt-Bitumen (#1), Meadows (#2), Trees (#3), Metal Sheets (#4), Bare Soil (#5), Gravel-Bricks (#6) and Shadow (#7).

Mississippi Gulfport [75]: Data is collected from the University of Southern Mississippi-Gulfport Campus. It has 64 spectral reflectance bands covering the wavelength range from 0.368 to 1.043 μm . The spatial resolution of the data is 1 m. It contains 11 man-made structures and natural materials. We eliminated sidewalk, yellow curb, cloth panels and water classes due to the lack of sufficient amounts of data samples. The remaining classes are used in the experiments: Trees (#1), Grass Pure (#2), Grass Ground (#3), Dirt & Sand (#4), Road (#5), Shadow (#6) and Building (#7).

B. Baselines, Metrics and Parameter Settings

Baselines: We compare our proposed method, EndNet, with several open-source hyperspectral unmixing algorithms

TABLE III
SAD RESULTS ON SYNTHETIC SPHERIC DATASET (40 DB NOISE).

	Spectral Angle Distance (SAD) ($\times 10^{-2}$)		
	VCA	DMaxD	EndNet-DMaxD
#1	0.59	3.91	1.29
#2	0.03	0.60	0.43
#3	0.11	2.12	1.38
#4	0.10	1.00	1.16
#5	0.14	0.84	0.61
Avg.	0.19	1.69	0.97

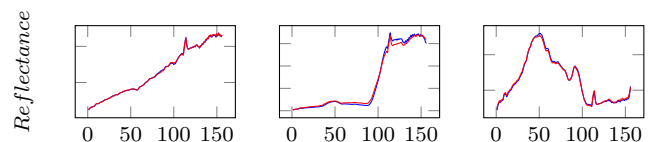


Fig. 10. Endmember signatures (Soil, Tree and Water) on Samson dataset for ground truth (blue) and EndNet-DMaxD (red).

TABLE IV

SAD AND RMSE RESULTS ON JASPER RIDGE DATASET. MEAN AND STANDARD DEVIATION ARE REPORTED. BEST RESULTS ARE ILLUSTRATED IN BOLD.

Endm.	Spectral Angle Distance (SAD) ($\times 10^{-2}$)								
	VCA	DMaxD	MVSA	SPICE	SCM	$l_{1 2}$ -NMF	DgS-NMF	EndNet-VCA	EndNet-DMaxD
#1	24.34 \pm 5.3	15.58 \pm 0.0	18.41 \pm 1.2	40.95 \pm >9	7.34 \pm 0.0	15.10 \pm 0.3	4.66 \pm0.2	8.23 \pm 3.1	4.99 \pm 0.4
#2	25.21 \pm 0.4	25.39 \pm 0.0	>99 \pm 1.0	59.83 \pm >9	7.80 \pm 0.0	4.60 \pm 0.0	4.60 \pm 0.0	3.62 \pm0.4	4.23 \pm 0.9
#3	34.81 \pm >9	13.35 \pm 0.0	26.27 \pm 0.7	28.03 \pm >9	5.61 \pm 0.0	6.16 \pm 0.5	5.66 \pm 0.2	6.23 \pm 0.7	4.47 \pm0.3
#4	50.56 \pm 9.3	10.69 \pm 0.0	7.01 \pm 0.8	86.19 \pm >9	3.36 \pm 0.0	9.81 \pm 0.1	6.73 \pm 0.1	67.06 \pm 0.5	1.96 \pm0.2
Avg.	33.73 \pm 6.2	16.25 \pm 0.0	37.92 \pm 0.9	53.75 \pm >9	6.02 \pm 0.0	7.19 \pm 2.4	5.41 \pm 0.1	21.29 \pm 1.2	3.91 \pm0.5
	Root Mean Square Error (RMSE) ($\times 10^{-2}$)								
	VCA	DMaxD	MVSA	SPICE	SCM	$l_{1 2}$ -NMF	DgS-NMF	EndNet-VCA	EndNet-DMaxD
#1	11.99 \pm 5.2	15.14 \pm 0.0	20.80 \pm 0.4	22.05 \pm 6.9	10.31 \pm 0.0	16.16 \pm 0.5	11.66 \pm 0.2	8.19 \pm 2.6	8.24 \pm0.4
#2	11.59 \pm 0.9	21.57 \pm 0.0	24.03 \pm 0.7	27.08 \pm >9	14.51 \pm 0.0	5.57 \pm 0.0	4.13 \pm0.0	26.18 \pm 0.5	6.17 \pm 0.3
#3	12.54 \pm 6.9	13.89 \pm 0.0	20.30 \pm 1.8	21.35 \pm 5.8	9.70 \pm 0.0	17.02 \pm 0.4	11.13 \pm 0.3	19.91 \pm 0.9	8.98 \pm0.2
#4	14.46 \pm 3.2	15.76 \pm 0.0	12.87 \pm 0.2	27.44 \pm >9	7.87 \pm 0.0	6.73 \pm 0.2	5.68 \pm0.1	30.65 \pm 0.3	8.55 \pm 0.1
Avg.	12.65 \pm 4.1	16.59 \pm 0.0	19.50 \pm 0.8	24.48 \pm 8.1	10.60 \pm 0.0	11.37 \pm 0.2	8.15 \pm 0.2	21.22 \pm 1.1	7.96 \pm0.3

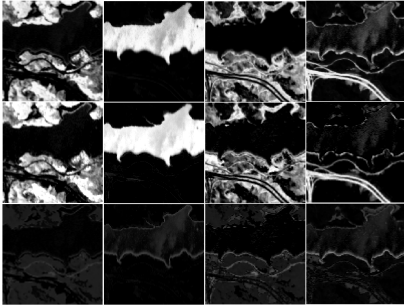


Fig. 11. Estimated fractional abundances and the absolute differences with the ground truth on Jasper Ridge dataset for EndNet-DMaxD method. Four materials, Tree, Water, Soil and Road, are illustrated in different columns.

throughout this study:

- Vertex Component Analysis (VCA) [4] (the code is available on <http://www.lx.it.pt/biucas/code.htm>)
- Minimum Volume Simplex Analysis (MVSA) [15] (the code is available on <http://www.lx.it.pt/biucas/code.htm>)
- Sparsity Promoting Iterated Constrained Endmembers (SPICE) [76] (the code is available on <https://github.com/GatorSense/SPICE>)
- Spatial Compositional Model (SCM) [77] (the code is available on <https://github.com/zhouyuanzxcv/Hyperspectral>)
- Distance-MaxD (DMaxD) [8] (the code is available on <https://sites.google.com/site/robhelenresearch/code>)

In addition to these open-source algorithms, we report the performance of recent methods for several datasets by referring directly to their reported scores, since their codes are not available: $l_{1|2}$ -NMF [71], [72], DgS-NMF [72]. Particularly, we should note that we compare our proposed method with the method like DgS-NMF [72] to show that our proposed method outperforms even the methods that exploit spatial priors in their solutions.

Metrics: To evaluate unmixing performance and compare with ground truth, we utilize two metrics: Spectral Angle Distance

(SAD) and Root Mean Square Error (RMSE):

$$\text{SAD}(\mathbf{e}, \hat{\mathbf{e}}) = \cos^{-1} \left(\frac{\mathbf{e} \hat{\mathbf{e}}}{\|\mathbf{e}\|_2 \|\hat{\mathbf{e}}\|_2} \right), \quad (14)$$

As mentioned earlier, SAD is used to evaluate the quality of estimated endmember with ground truth by measuring angle distance. Similarly, to assess the accuracy of the estimated abundances, we utilized the RMSE metric as:

$$\text{RMSE}(\mathbf{y}, \hat{\mathbf{y}}) = \sqrt{\frac{1}{N} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2}. \quad (15)$$

Note that smaller values indicate better performance for both metrics.

Parameter Settings: We fixed the default parameter values of the baseline methods in our experiments. We tuned the parameters of SCM algorithm for each dataset based on the observations in [77].

For fractional abundance estimation, we selected Multilinear Mixing Model (MLM) [78] for MVSA, VCA and DMaxD algorithms to preserve nonlinearity (Note that we also tested these algorithms with Fully-Constrained Least Square (FCLS), Generalized Bilinear Model (GBM) [79], Postnonlinear Mixing Model (PPNM) [32] and Simplex Projection Unmixing (SPU). MLM yielded the best performance). SPICE and SCM compute the optimum abundances for the scene in their algorithms.

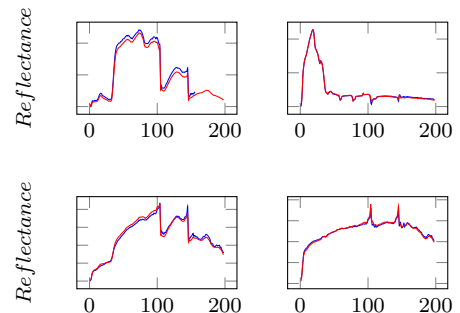


Fig. 12. Endmember signatures (Tree, Water, Soil and Road) on Jasper dataset for ground truth (blue) and EndNet-DMaxD (red).

TABLE V
SAD RESULTS ON CUPRITE DATASET. MEAN AND STANDARD DEVIATION ARE REPORTED. BEST RESULTS ARE ILLUSTRATED IN BOLD.

Endm.	Spectral Angle Distance (SAD) ($\times 10^{-2}$)								
	VCA	DMaxD	MVSA	SPICE	SCM	$l_{1 2}$ -NMF	DgS-NMF	EndNet-VCA	EndNet-DMaxD
#1	12.89 \pm 5.3	8.58 \pm 0.0	23.60 \pm 5.7	14.71 \pm 1.6	24.43 \pm 4.6	16.22 \pm 2.0	12.48 \pm 1.8	13.60 \pm 2.9	14.18 \pm 2.5
#2	8.21 \pm 1.9	7.15 \pm 0.0	>99 \pm >9	9.32 \pm 1.6	7.19 \pm 0.2	10.15 \pm 3.0	7.59 \pm 1.2	7.65 \pm 0.1	7.34 \pm 0.2
#3	9.04 \pm 2.0	11.58 \pm 0.0	>99 \pm >9	10.26 \pm 0.3	13.25 \pm 0.3	12.50 \pm 7.5	10.81 \pm 3.2	11.44 \pm 0.8	11.66 \pm 1.1
#4	9.51 \pm 3.0	8.01 \pm 0.0	38.36 \pm >9	11.91 \pm 1.3	12.22 \pm 0.9	13.07 \pm 5.3	11.01 \pm 2.1	8.38 \pm 0.6	7.81 \pm 0.7
#5	8.88 \pm 1.5	8.47 \pm 0.0	25.44 \pm 6.2	13.09 \pm 2.7	9.20 \pm 1.3	9.42 \pm 1.8	9.02 \pm 2.9	12.06 \pm 0.9	13.16 \pm 1.1
#6	7.08 \pm 5.2	8.52 \pm 0.0	>99 \pm >9	9.60 \pm 0.5	6.80 \pm 0.4	9.82 \pm 2.7	7.42 \pm 1.2	6.05 \pm 0.3	5.69 \pm 0.2
#7	17.13 \pm 4.0	10.26 \pm 0.0	41.73 \pm >9	9.97 \pm 0.5	17.15 \pm 2.5	29.86 \pm 7.4	20.51 \pm 6.0	15.93 \pm 1.9	15.45 \pm 1.1
#8	6.22 \pm 0.4	6.19 \pm 0.0	81.84 \pm >9	8.60 \pm 0.5	6.28 \pm 0.2	10.27 \pm 4.7	7.27 \pm 1.2	5.51 \pm 0.3	5.75 \pm 0.1
#9	7.95 \pm 1.0	7.36 \pm 0.0	65.95 \pm >9	11.45 \pm 1.6	7.61 \pm 0.1	12.80 \pm 4.1	8.88 \pm 1.7	8.36 \pm 0.7	7.78 \pm 0.1
#10	11.11 \pm 4.5	13.08 \pm 0.0	88.65 \pm >9	6.99 \pm 1.9	5.89 \pm 0.1	8.12 \pm 1.9	8.82 \pm 5.0	5.87 \pm 0.5	6.12 \pm 0.4
#11	8.40 \pm 5.9	28.56 \pm 0.0	52.02 \pm >9	8.99 \pm 2.4	7.11 \pm 1.7	11.03 \pm 3.3	8.15 \pm 2.1	9.03 \pm 1.4	8.38 \pm 1.2
#12	10.07 \pm 4.5	8.23 \pm 0.0	9.10 \pm 1.35	8.63 \pm 0.9	9.29 \pm 0.2	13.72 \pm 6.0	13.62 \pm 6.3	10.56 \pm 1.6	11.14 \pm 0.9
Avg.	9.71 \pm 2.9	10.54 \pm 0.0	71.92 \pm >9	10.30 \pm 1.3	10.54 \pm 1.0	13.07 \pm 4.1	10.46 \pm 2.9	9.53 \pm 1.0	9.54 \pm 0.8

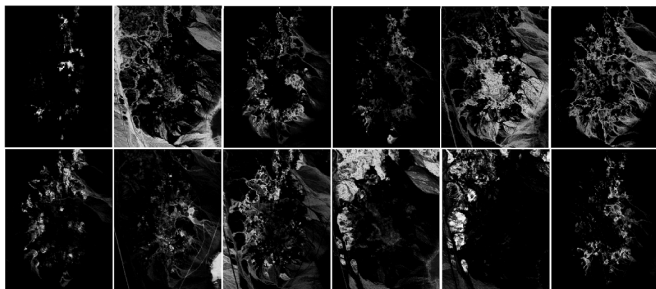


Fig. 13. Estimated fractional abundance results on Cuprite dataset with the EndNet-DMaxD method. First Row: Alunite, Andradite, Buddingtonite, Dumortierite, Kaolinite₁, Kaolinite₂. Second Row: Muscovite, Montmorillonite, Nontronite, Pyrope, Spinel, Chalcedony.

For our proposed method, three parameters should be tuned by the user for different scenes, namely p in the dropout layer, λ_2 in the optimization step and the percentage of mask noise. We will redefine the values of these parameters if the scene needs to be tuned by presenting reasonable explanations. Otherwise, the default values are used on each dataset ($p = 1.0$, $\lambda_2 = 0.1$ and 40% mask noise). Also, EndNet-VCA and EndNet-DMaxD in the experiments denote that which of the algorithms is used for the parameter initialization of our proposed method.

Finally, the corresponding ground truth material (i.e. most similar ground truth material) is determined by measuring its highest SAD similarity score with the estimated endmembers in the experiments.

Computational Complexity: Speed and memory requirement are two important parameters for the problem. Since a stochas-

TABLE VI
AVERAGE COMPUTATION TIME OF THE METHODS FOR THREE DATASETS.

	Computation Time (Sec)		
	Urban	Samson	Jasper
VCA-MLM	\approx 3600	\approx 349	\approx 352
SCM	\approx 1021	\approx 187	\approx 472
EndNet-DMaxD-SPU	\approx 914	\approx 789	\approx 855

tic gradient-based solver is utilized in the proposed method, it practically scales the problem for large-scale data [81]. Intuitively, the computation time and memory-requirement is independent from data since batch-based learning is used. Moreover, the computation time for the proposed method on each dataset is illustrated in Table VI (all codes are implemented on Matlab). Note that the computation time of our proposed method is independent from data size and even if data size increases, the computation time will be in the same order. Lastly, the proposed method can be easily parallelized on Graphical Processing Units (GPUs) due to the neural network architecture.

C. Experiments on Hyperspectral Unmixing Datasets

In this section, we compare our proposed method with the baselines on the Urban, Samson, Jasper Ridge and Cuprite datasets. These four datasets are illustrated in Fig. 5. For reliable assessments, tests are repeated 20 times for each method, thus mean and standard deviation of the results are reported. We also illustrate the qualitative results of the proposed method to provide visual comparisons on the estimated fractional abundances.

TABLE VII
COMPARISON OF ABUNDANCE ESTIMATION RESULTS FOR SPU (FIRST ROWS) AND ENDNET-DMAXD (SECOND ROWS).

	Root Mean Square Error (RMSE) ($\times 10^{-2}$)		
	Urban	Samson	Jasper
#1	SPU: 10.41 \pm 0.2 EndN: 13.04 \pm 0.3	SPU: 5.72 \pm 0.0 EndN: 9.41 \pm 0.1	SPU: 8.24 \pm 0.4 EndN: 10.12 \pm 0.6
#2	SPU: 12.24 \pm 0.3 EndN: 14.43 \pm 0.3	SPU: 3.84 \pm 0.1 EndN: 6.47 \pm 0.3	SPU: 6.17 \pm 0.3 EndN: 11.48 \pm 0.8
#3	SPU: 8.35 \pm 0.3 EndN: 8.71 \pm 0.5	SPU: 2.11 \pm 0.0 EndN: 3.93 \pm 0.1	SPU: 8.19 \pm 0.2 EndN: 9.53 \pm 0.3
#4	SPU: 5.92 \pm 0.2 EndN: 7.59 \pm 0.2	- -	SPU: 8.55 \pm 0.1 EndN: 12.29 \pm 0.4
Avg.	SPU: 9.23 \pm 0.2 EndN: 10.94 \pm 0.4	SPU: 3.88 \pm 0.0 EndN: 5.72 \pm 0.2	SPU: 7.96 \pm 0.3 EndN: 10.85 \pm 0.6

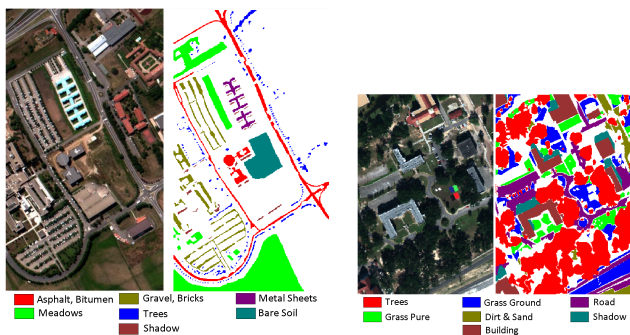


Fig. 14. Two real hyperspectral datasets with pixel-wise class labels: University of Pavia and Mississippi Gulfport respectively.

In Tables I, II and IV, each row corresponds to SAD or RMSE performance of different methods for a single material. The last row at each table denotes the average performance of all materials for the corresponding metric. In Table V, only the SAD performance is reported since Cuprite dataset does not have a quantitative fractional abundance ground truth. Fig. 7, 9, 11 and 13 visualize the qualitative abundance performance of our EndNet-DMaxD method on these datasets. In these figures, we provide the estimated fractional abundances, the ground truth (if it is available) and the absolute differences of the estimated abundances. In addition, the estimated endmember and ground truth signatures are illustrated for the Urban, Samson and Jasper datasets in Fig. 8, 10 and 12. (The project website provides the qualitative results on different datasets.)

Lastly, note that for SCM method, the parameters are tuned to $\beta_1 = 10$, $\beta_2 = 5$ and $\rho = 0.01$ for the Samson and the Jasper datasets while they are set to $\beta_1 = 10$, $\beta_2 = 10$ and $\rho = 0.01$ for Urban and Cuprite datasets.

Urban: Quantitative and qualitative results on Urban dataset are summarized in Table I and Fig. 7. From the results, our proposed method achieves the best overall scores for both SAD and RMSE metrics compared to other methods. EndNet-DMaxD is approximately 4.2% and 1.2% better than the second best result. The second best result is obtained by DgS-NMF method that exploits spatial priors in the estimation. This is important since our proposed method outperforms the methods that use spatial information about the datasets, without exploiting any extra prior.

Fig. 7 visualizes the estimated fractional abundance for the dataset with EndNet-DMaxD method. To compare the estimation accuracy, absolute differences of the estimated fractional abundance and ground truth are provided for each material (i.e rows). From these results, we can clearly observe that the estimation error is usually concentrated on the object boundaries. Despite the fact that our proposed method might yield false estimates where high spectral mixtures have occurred, we must emphasize that the ground truth can still exhibit some noise since the determination of perfect abundance labels of these locations can be challenging.

Samson: Results are shown in Table II and Fig. 9. As for the previous scenarios, our proposed methods obtain the best overall performance for both metrics, and the RMSE

results are prominently improved by the proposed method. Approximately 2% and 3% improvements are introduced for SAD and RMSE metrics respectively. From Fig. 9, it can be seen that the absolute difference is quite small. Only the water-ground intersection areas yield a small estimation error where all three materials can contribute.

Jasper Ridge: Remark that soil (#3) and road (#4) are highly mixed materials in this scene. Therefore, we redefine λ_2 as 0.4 to increase the sparsity of the material abundances in the proposed method. Table IV shows the SAD and RMSE performance of the methods. From the results, our proposed method introduces small improvements over DgS-NMF. However, it should also be noted that even if the use of spatial priors provides huge advantages, our method still yields the best results for highly mixed scenes. Fig. 11 illustrates the absolute error with respect to the ground truth. The error is intensified at the boundaries of water-ground and the regions designated road-soil.

Lastly, an important observation for one of our proposed methods is that the performance of EndNet-VCA drastically decreases on this dataset while EndNet-DMaxD maintains the exceptional performance. Our explanation of this result is that the DMaxD method is inclined to detect more uncorrelated spectral signatures from data than VCA (even if the estimates do not need to be the true results), since it maximizes the distances of the estimated spectral signatures with one another. Also, VCA can ignore the materials that are numerically small in a scene. Thus, considering the methods such as DMaxD for the proposed method initialization can be more advantageous.

Cuprite: Experimental results for this dataset are reported in Table V and Fig. 13. From these results, it is clear that the methods using the unmixed pixel assumption such as VCA, DMaxD etc. obtain high performance. Therefore, we decreased the value of the sparsity term λ_2 to 0.001. Another important feature of the dataset is that it consists of quite correlated materials. For this reason, we also set p to 0.8 to improve the convergence of the parameters.

The experimental results show that the proposed methods produce the lowest SAD results compared to the baseline methods. Fig. 13 visualizes the estimated fractional abundance for each material with the EndNet-DMaxD method. Since this dataset has no ground truth for the true abundances, we can only make visual comparisons to analyze the consistency with other methods [74]. When we compare our results with [74], we observed that the estimated fractional abundances are

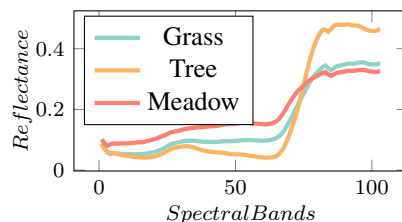


Fig. 15. Three highly correlated materials observed in University of Pavia dataset: Grass, Tree and Meadow. Even if Grass and Meadow are considered as one constituent material in the original dataset, they have two distinct spectral signatures as shown in the plot.

TABLE VIII
SAD RESULTS ON UNIVERSITY OF PAVIA DATASET. MEAN AND STANDARD DEVIATION ARE REPORTED. BEST RESULTS ARE ILLUSTRATED IN BOLD.

Endm.	Spectral Angle Distance (SAD) ($\times 10^{-2}$)					
	VCA	DMaxD	SPICE	SCM	EndNet-VCA	EndNet-DMaxD
#1	51.02 \pm 8.9	7.50 \pm 0.0	57.98 \pm >9	2.93 \pm 0.0	2.31 \pm1.1	2.38 \pm 0.6
#2	39.43 \pm >9	92.98 \pm 0.0	66.33 \pm >9	1.85 \pm0.0	4.19 \pm 1.1	6.07 \pm 1.8
#3	28.85 \pm >9	>99 \pm 0.0	12.50 \pm 5.3	4.70 \pm 0.0	1.62 \pm0.1	2.03 \pm 0.8
#4	42.47 \pm 4.6	30.60 \pm 0.0	44.45 \pm >9	11.67 \pm 0.0	10.11 \pm 1.2	8.57 \pm1.7
#5	85.19 \pm >9	70.95 \pm 0.0	>99 \pm >9	7.69 \pm 0.0	5.73 \pm0.6	7.04 \pm 0.8
#6	42.65 \pm 8.8	56.93 \pm 0.0	93.00 \pm >9	27.65 \pm 0.0	4.62 \pm 2.6	1.59 \pm0.5
#7	53.04 \pm 3.7	53.06 \pm 0.0	94.72 \pm >9	5.29 \pm0.0	55.80 \pm 0.3	13.40 \pm 1.9
Avg.	48.96 \pm >9	61.33 \pm 0.0	67.79 \pm >9	8.82 \pm 0.0	12.07 \pm 1.0	5.87 \pm1.2

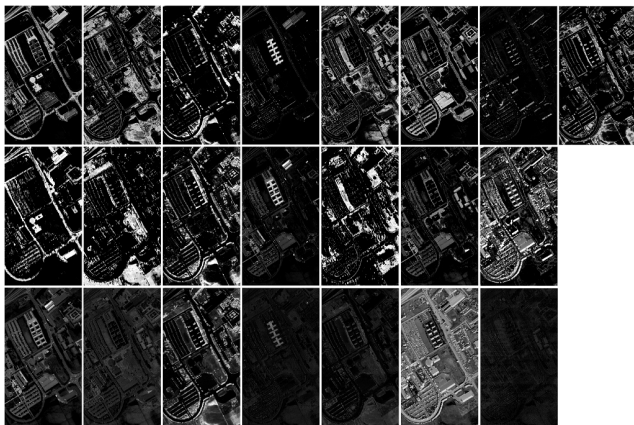


Fig. 16. Estimated fractional abundances on University of Pavia dataset. Rows indicate methods: EndNet-DMaxD, SCM and SPICE, while each column corresponds to the fractional abundance of a single material: Asphalt-Bitumen (#1), Meadows (#2), Trees (#3), Metal Sheets (#4), Bare Soil (#5), Gravel-Bricks (#6) and Shadow (#7) respectively. We also visualize the fractional abundance of Grass material obtained by EndNet-DMaxD method as the last image.

consisted with the reported visual results.

Comparison of Abundance Map Estimation: The detailed experimental results for the abundance estimation of SPU and EndNet-DMaxD are reported in Table VII. For three datasets, the SPU method introduces further improvements to the performance compared to EndNet-DMaxD abundance estimates. However, note that EndNet-DMaxD achieves similar or better performance compared to the baseline methods.

D. Experiments on Hyperspectral Classification Datasets

In this section, we analyze SAD and the qualitative performance of the baseline methods on the University of Pavia and Mississippi Gulfport datasets. In particular, these experimental results can provide a basis for future hyperspectral classification methods which could use the proposed method in these studies. Tests are repeated 20 times. Fig. 14 illustrates these datasets and their pixel-wise class labels. Note that no RMSE performance on abundance estimates are reported in the paper on the contrary to [77], since class labels are quite noisy and this information might not convey the true performance of the method.

Due to the fact that there is no endmember ground truth available for these datasets, we used the average spectra of

the pixels for each material as the ground truth in SAD comparisons as proposed in [77], [80].

Table VIII and IX show the SAD performance of these methods. Fig. 16 and 17 visualize the qualitative fractional abundance performance for three methods: SPICE, SCM and EndNet-DMaxD (VCA, DMaxD and MVSA methods all obtain insufficient abundance results). Remark that $l_{1|2}$ -NMF and DgS-NMF scores are not reported since neither their source codes nor previous performance are available on these datasets.

For University of Pavia dataset, the parameters are set to $\beta_1 = 10$, $\beta_2 = 10$ and $\rho = 0.01$ for the SCM method. These values are slightly different from the recommended settings in [77], since these configurations yield better visual and quantitative results in our experiments. Similarly, for Mississippi Gulfport, these values are tuned to $\beta_1 = 10$, $\beta_2 = 5$ and $\rho = 0.01$. Lastly, the prune threshold value is increased to 10^{-5} for SPICE algorithm on Mississippi Gulfport dataset.

University of Pavia: Table VIII and Fig. 16 show experimental results on this dataset. As mentioned previously, the scene contains 7 labeled classes including the associations of gravel-bricks and asphalt-bitumen materials. However, we empirically observed that the scene comprises another distinct material, 'Grass', which has a completely different spectral response from 'Meadow' and 'Tree'. These three materials are plotted in Fig. 15.

For this reason, we set the optimum endmember number in the scene to 8 just for the proposed method and preserved the default value (i.e. 7) for the baseline methods (these methods yield worse results when we increase the endmember number).

We redefine λ_2 and p as 0.01 and 0.8 respectively as for the Cuprite dataset, since the data contains a number of correlated materials. We decreased the mask noise to 10% due to the noise level of data and significantly spectral correlations. From the SAD results, it is clear that the proposed method achieves the best performance compared to the baselines and a 3% improvement is introduced over the second best result.

Of course, the fractional abundance estimation is especially important for this dataset to distinguish the classes. Fig. 16 visualizes the fractional abundance estimation for three methods. It can be seen that our proposed method, EndNet-DMaxD, outperforms the other baseline methods and obtains meaningful abundance results for the materials. This is a gratifying result since it confirms that 'Grass' material

TABLE IX
SAD RESULTS ON MISSISSIPPI GULFPORT DATASET. MEAN AND STANDARD DEVIATION ARE REPORTED. BEST RESULTS ARE ILLUSTRATED IN BOLD.

Endm.	Spectral Angle Distance (SAD) ($\times 10^{-2}$)					
	VCA	DMaxD	SPICE	SCM	EndNet-VCA	EndNet-DMaxD
#1	17.04 \pm >9	12.71 \pm 0.0	>99 \pm >9	15.31 \pm 0.0	3.47 \pm 1.0	2.75 \pm0.2
#2	45.73 \pm >9	12.21 \pm 0.0	>99 \pm >9	2.07 \pm 0.0	6.32 \pm 5.2	1.28 \pm0.1
#3	25.90 \pm 1.8	26.12 \pm 0.0	>99 \pm >9	1.93 \pm0.0	9.99 \pm 8.4	3.97 \pm 0.4
#4	42.09 \pm >9	34.12 \pm 0.0	>99 \pm >9	4.78 \pm0.0	13.55 \pm >9	8.31 \pm 0.7
#5	60.43 \pm 9.7	24.13 \pm 0.0	88.03 \pm >9	29.92 \pm 0.0	23.48 \pm >9	2.72 \pm0.7
#6	39.10 \pm 9.0	47.06 \pm 0.0	>99 \pm >9	21.03 \pm 0.0	13.69 \pm >9	11.78 \pm1.0
#7	37.42 \pm 4.7	9.28 \pm 0.0	46.60 \pm >9	6.40 \pm0.0	8.66 \pm 2.5	11.56 \pm 0.4
Avg.	38.25 \pm >9	23.67 \pm 0.0	>99 \pm >9	11.65 \pm 0.0	11.31 \pm 7.9	6.06 \pm0.5

in hyperspectral unmixing provides critical improvements to the scene where 'Trees' is repeatedly miscategorized with 'Grass'/'Meadow' by the baseline methods.

Mississippi Gulfport: Quantitative and qualitative results for this dataset are reported in Table IX and Fig. 17. As before, λ_2 and p values are tuned as 0.1 and 0.8, and the mask noise is set to 10%. Several materials (i.e. sidewalk, yellow curb, cloth panels and water) were not considered in the computation, since their samples are too scarce in the scene.

In Table IX, our proposed method yields the highest overall SAD score over ground truth materials and the quality of abundance estimates in Fig. 17 is significantly better compared to the baselines. In particular, for 'Tree', 'Building' and 'Shadow', misprediction rates are quite high for the baseline methods, whereas our proposed method obtains nearly optimum results for these materials in the scene.

IV. CONCLUSION

In this paper, we proposed a novel method, EndNet, for end-member extraction and hyperspectral unmixing. To this end, we improved and restructured the conventional autoencoder neural network by introducing additional layers and a novel loss function. These modifications enable us to address the common challenges of hyperspectral unmixing such as non-linearity, sparsity, some of the physical constraints etc. Also,

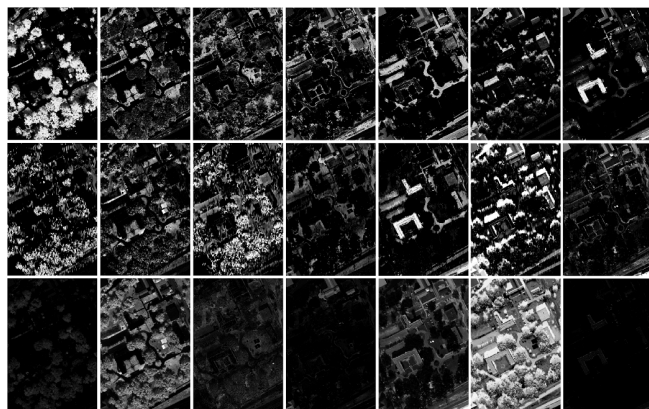


Fig. 17. Abundance maps for Mississippi Gulfport dataset. Similarly, the results of EndNet-DMaxD, SCM and SPICE are illustrated in different rows respectively. Trees (#1), Grass Pure (#2), Grass Ground (#3), Dirt & Sand (#4), Road (#5), Shadow(#6) and Building (#7) are identified in columns.

backpropagation with a stochastic-gradient based solver is used to optimize the problem and this scales the method for large-scale data, in contrast to more complex derivations/inferences in the literature.

We adapted the SPU algorithm to further improve the estimation of fractional abundances from a scene with the estimated endmembers, even when the estimates of hidden abstracts yield compatible results. This adaptation is achieved by replacing the standard l_2 norm kernel with a SAD kernel that we have shown is more efficacious.

From the experimental results, we have observed that our proposed method makes significant performance improvements to hyperspectral unmixing domain and achieves outstanding results on well-known hyperspectral datasets. In addition, we selected the DgS-NMF method as a baseline technique in the experiments and we conclusively demonstrated that our method outperforms these type of methods even though they use the spatial priors of a scene in their extraction.

Finally, we should emphasize that the proposed method is the first successful end-to-end learning algorithm based on a neural network that attains superior performance in an unsupervised manner for hyperspectral unmixing problem. That's why, we strongly believe that the findings of our study will provide a basis for further studies into neural network techniques in this domain.

V. ACKNOWLEDGMENT

We would like to thank Engin Tola and Ufuk Sakarya for their feedback on the text.

REFERENCES

- [1] J. M. Bioucas-Dias *et al.*, "Hyperspectral unmixing overview: Geometrical, statistical, and sparse regression-based approaches," *IEEE JSTARS*, pp. 354–379, 2012.
- [2] A. J. Brown *et al.*, "Hydrothermal formation of clay-carbonate alteration assemblages in the nili fossae region of mars," *Earth and Planetary Science Letters*, vol. 297, no. 1-2, pp. 174–182, 2010.
- [3] B. Hapke, "Bidirectional reflectance spectroscopy: 1. theory," *Journal of Geophysical Research: Solid Earth*, pp. 3039–3054, 1981.
- [4] J. M. Nascimento and J. M. Dias, "Vertex component analysis: A fast algorithm to unmix hyperspectral data," *IEEE TGRS*, pp. 898–910, 2005.
- [5] J. W. Boardman, "Geometric mixture analysis of imaging spectrometry data," in *IGARSS*, 1994, pp. 2369–2371.
- [6] J. C. Harsanyi and C.-I. Chang, "Hyperspectral image classification and dimensionality reduction: An orthogonal subspace projection approach," *IEEE TGRS*, pp. 779–785, 1994.

- [7] M. E. Winter, "N-findr: An algorithm for fast autonomous spectral end-member determination in hyperspectral data," in *SPIE*, 1999, pp. 266–275.
- [8] R. Heylen *et al.*, "Non-linear spectral unmixing by geodesic simplex volume maximization," *IEEE J-STSP*, pp. 534–542, 2011.
- [9] M.-D. Iordache *et al.*, "Sparse unmixing of hyperspectral data," *IEEE TGRS*, pp. 2014–2039, 2011.
- [10] R. B. Singer and T. B. McCord, "Mars-large scale mixing of bright and dark surface materials and implications for analysis of spectral reflectance," in *Lunar and Planetary Science Conference Proceedings*, 1979, pp. 1835–1848.
- [11] A. Ifarraguerri and C.-I. Chang, "Multispectral and hyperspectral image analysis with convex cones," *IEEE TGRS*, pp. 756–770, 1999.
- [12] F.-Y. Wang *et al.*, "Nonnegative least-correlated component analysis for separation of dependent sources by volume maximization," *IEEE PAMI*, pp. 875–888, 2010.
- [13] V. F. Haertel and Y. E. Shimabukuro, "Spectral linear mixing model in low spatial resolution image data," *IEEE TGRS*, pp. 2555–2562, 2005.
- [14] L. Miao and H. Qi, "Endmember extraction from highly mixed data using minimum volume constrained nonnegative matrix factorization," *IEEE TGRS*, pp. 765–777, 2007.
- [15] J. Li and J. M. Bioucas-Dias, "Minimum volume simplex analysis: A fast algorithm to unmix hyperspectral data," in *IGARSS*, 2008.
- [16] N. Dobigeon *et al.*, "Joint bayesian endmember extraction and linear unmixing for hyperspectral imagery," *IEEE TSP*, pp. 4355–4368, 2009.
- [17] M.-D. Iordache *et al.*, "Collaborative sparse regression for hyperspectral unmixing," *IEEE TGRS*, pp. 341–354, 2014.
- [18] A. J. Brown, "Spectral curve fitting for automatic hyperspectral data analysis," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 44, no. 6, pp. 1601–1608, 2006.
- [19] B. Somers *et al.*, "Nonlinear hyperspectral mixture analysis for tree cover estimates in orchards," *Remote Sensing of Environment*, vol. 113, no. 6, pp. 1183–1193, 2009.
- [20] N. Keshava and J. F. Mustard, "Spectral unmixing," *IEEE Signal Processing Magazine*, pp. 44–57, 2002.
- [21] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, pp. 2323–2326, 2000.
- [22] C. M. Bachmann *et al.*, "Improved manifold coordinate representations of large-scale hyperspectral scenes," *IEEE TGRS*, pp. 2786–2803, 2006.
- [23] G. Zhao *et al.*, "Multilayer unmixing for hyperspectral imagery with fast kernel archetypal analysis," *IEEE GRSL*, pp. 1532–1536, 2016.
- [24] J. Broadwater and A. Banerjee, "A comparison of kernel functions for intimate mixture models," in *IEEE WHISPERS*. IEEE, 2009, pp. 1–4.
- [25] J. Chen *et al.*, "Nonlinear unmixing of hyperspectral data based on a linear-mixture/nonlinear-fluctuation model," *IEEE Transactions on Signal Processing*, pp. 480–492, 2013.
- [26] C. Févotte and N. Dobigeon, "Nonlinear hyperspectral unmixing with robust nonnegative matrix factorization," *IEEE Transactions on Image Processing*, pp. 4810–4819, 2015.
- [27] R. Close *et al.*, "Using physics-based macroscopic and microscopic mixture models for hyperspectral pixel unmixing," in *SPIE*, 2012, p. 83901L.
- [28] Y. Altmann *et al.*, "Bilinear models for nonlinear unmixing of hyperspectral images," in *IEEE WHISPERS*. IEEE, 2011, pp. 1–4.
- [29] R. Heylen *et al.*, "Nonlinear unmixing by using different metrics in a linear unmixing chain," *IEEE JSTARS*, pp. 2655–2664, 2015.
- [30] F. Kizel *et al.*, "A stepwise analytical projected gradient descent search for hyperspectral unmixing and its code vectorization," *IEEE TGRS*, pp. 4925–4943, 2017.
- [31] N. Dobigeon *et al.*, "Nonlinear unmixing of hyperspectral images: Models and algorithms," *IEEE Signal Processing Magazine*, vol. 31, no. 1, pp. 82–94, 2014.
- [32] Y. Altmann *et al.*, "Supervised nonlinear spectral unmixing using a postnonlinear mixing model for hyperspectral imagery," *IEEE TIP*, vol. 21, no. 6, pp. 3017–3025, 2012.
- [33] J. Plaza *et al.*, "On the use of small training sets for neural network-based characterization of mixed pixels in remotely sensed hyperspectral images," *Pattern Recognition*, pp. 3032–3045, 2009.
- [34] K. J. Guilfoyle *et al.*, "A quantitative and comparative analysis of linear and nonlinear spectral mixture models using radial basis function neural networks," *IEEE TGRS*, pp. 2314–2318, 2001.
- [35] J. Plaza *et al.*, "Joint linear/nonlinear spectral unmixing of hyperspectral image data," in *IGARSS*. IEEE, 2007, pp. 4037–4040.
- [36] B. Ayerdi and M. Grana, "Hyperspectral image nonlinear unmixing and reconstruction by elm regression ensemble," *Neurocomputing*, vol. 174, pp. 299–309, 2016.
- [37] B. Pan *et al.*, "R-vcnet: A new deep-learning-based hyperspectral image classification method," *IEEE JSTAR*, pp. 1975–1986, 2017.
- [38] F. Palsson *et al.*, "Neural network hyperspectral unmixing with spectral information divergence objective," *IGARSS*, 2017.
- [39] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *ICML*, 2010, pp. 807–814.
- [40] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *ICML*, 2015, pp. 448–456.
- [41] N. Srivastava *et al.*, "Dropout: a simple way to prevent neural networks from overfitting," *JMLR*, no. 1, pp. 1929–1958, 2014.
- [42] S. Wager *et al.*, "Dropout training as adaptive regularization," in *NIPS*, 2013, pp. 351–359.
- [43] P. Vincent *et al.*, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion," *Journal of Machine Learning Research*, vol. 11, no. Dec, pp. 3371–3408, 2010.
- [44] Y. Bengio *et al.*, "Representation learning: A review and new perspectives," *IEEE PAMI*, pp. 1798–1828, 2013.
- [45] Y. LeCun *et al.*, "Deep learning," *Nature*, pp. 436–444, 2015.
- [46] S. Rifai *et al.*, "Contractive auto-encoders: Explicit invariance during feature extraction," in *ICML*, 2011, pp. 833–840.
- [47] A. Rasmus *et al.*, "Semi-supervised learning with ladder networks," in *NIPS*, 2015, pp. 3546–3554.
- [48] S. Holzwarth *et al.*, "Hysens-dais 7915/rosis imaging spectrometers at dlr," in *Proceedings of the 3rd EARSeL Workshop on Imaging Spectroscopy*, 2003, pp. 3–14.
- [49] K. He *et al.*, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *ICCV*, 2015, pp. 1026–1034.
- [50] R. Memisevic and D. Krueger, "Zero-bias autoencoders and the benefits of co-adapting features," *Stat*, p. 13, 2014.
- [51] D.-A. Clevert *et al.*, "Fast and accurate deep network learning by exponential linear units (elus)," *arXiv preprint arXiv:1511.07289*, 2015.
- [52] B. Xu *et al.*, "Empirical evaluation of rectified activations in convolutional network," *arXiv preprint arXiv:1505.00853*, 2015.
- [53] A. Krizhevsky *et al.*, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012, pp. 1097–1105.
- [54] A. Makhzani and B. Frey, "K-sparse autoencoders," *arXiv preprint arXiv:1312.5663*, 2013.
- [55] Y. LeCun *et al.*, "Backpropagation applied to handwritten zip code recognition," *Neural computation*, pp. 541–551, 1989.
- [56] L. Chunjie *et al.*, "Cosine normalization: Using cosine similarity instead of dot product in neural networks," *arXiv preprint arXiv:1702.05870*, 2017.
- [57] P. E. Dennison *et al.*, "A comparison of error metrics and constraints for multiple endmember spectral mixture analysis and spectral angle mapper," *Remote Sensing of Environment*, pp. 359–367, 2004.
- [58] G. Camps-Valls, "Kernel spectral angle mapper," *Electronics Letters*, pp. 1218–1220, 2016.
- [59] G. Camps-Valls *et al.*, "Composite kernels for hyperspectral image classification," *IEEE GRSL*, pp. 93–97, 2006.
- [60] M. Fauvel *et al.*, "Evaluation of kernels for multiclass classification of hyperspectral remote sensing data," in *ICASSP*. IEEE, 2006.
- [61] H. Li and L. Zhang, "A hybrid automatic endmember extraction algorithm based on a local window," *IEEE GRSL*, pp. 4223–4238, 2011.
- [62] X. Liu and C. Yang, "A kernel spectral angle mapper algorithm for remote sensing image classification," in *Image and Signal Processing*. IEEE, 2013, pp. 814–818.
- [63] F. Tsai and W. Philpot, "Derivative analysis of hyperspectral data," *Remote Sensing of Environment*, pp. 41–51, 1998.
- [64] Z. B. Rabah *et al.*, "A new method to change illumination effect reduction based on spectral angle constraint for hyperspectral image unmixing," *IEEE GRSL*, pp. 1110–1114, 2011.
- [65] J. R. Hershey and P. A. Olsen, "Approximating the kullback leibler divergence between gaussian mixture models," in *ICASSP*. IEEE, 2007, pp. IV–317.
- [66] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," *JMLR*, pp. 2579–2605, 2008.
- [67] K. Kavukcuoglu *et al.*, "Fast inference in sparse coding algorithms with applications to object recognition," *arXiv preprint arXiv:1010.3467*, 2010.
- [68] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [69] D. Ulyanov *et al.*, "Instance normalization: The missing ingredient for fast stylization," *arXiv:1607.08022*, 2016.
- [70] R. Heylen *et al.*, "Fully constrained least squares spectral unmixing by simplex projection," *IEEE TGRS*, pp. 4112–4122, 2011.

- [71] Y. Qian *et al.*, “Hyperspectral unmixing via $l_{1/2}$ sparsity-constrained nonnegative matrix factorization,” *IEEE TGRS*, pp. 4282–4297, 2011.
- [72] F. Zhu *et al.*, “Spectral unmixing via data-guided sparsity,” *IEEE TIP*, pp. 5412–5427, 2014.
- [73] F. Zhu *et al.*, “Structured sparse method for hyperspectral unmixing,” *ISPRS Journal of Photogrammetry and Remote Sensing*, pp. 101–118, 2014.
- [74] F. A. Kruse *et al.*, “Comparison of aviris and hyperion for hyperspectral mineral mapping,” in *11th JPL Airborne Geoscience Workshop*, 2002.
- [75] P. Gader *et al.*, “Muufi gulfport hyperspectral and lidar airborne data set,” *Univ. Florida, Gainesville, FL, USA, Tech. Rep. REP-2013-570*, 2013.
- [76] A. Zare and P. Gader, “Sparsity promoting iterated constrained endmember detection in hyperspectral imagery,” *IEEE GRSL*, pp. 446–450, 2007.
- [77] Y. Zhou *et al.*, “A spatial compositional model for linear unmixing and endmember uncertainty estimation,” *IEEE TIP*, pp. 5987–6002, 2016.
- [78] R. Heylen and P. Scheunders, “A multilinear mixing model for nonlinear spectral unmixing,” *IEEE TGRS*, pp. 240–251, 2016.
- [79] A. Halimi *et al.*, “Nonlinear unmixing of hyperspectral images using a generalized bilinear model,” *IEEE TGRS*, pp. 4153–4162, 2011.
- [80] A. Zare *et al.*, “Piecewise convex multiple-model endmember detection and spectral unmixing,” *IEEE TGRS*, pp. 2853–2862, 2013.
- [81] L. Bottou, “Large-scale machine learning with stochastic gradient descent,” in *Proceedings of COMPSTAT’2010*. Springer, 2010, pp. 177–186.

Savas Ozkan is a senior researcher at TUBITAK Space Technologies Research Institute and he is currently working toward the PhD degree from the Department of Electrical and Electronics Engineering, Middle East Technical University. His research interests include deep learning, image/video retrieval, hyperspectral image processing, generative adversarial networks and biomedical image processing.

Berk Kaya is a MSc student at ETH Zurich. His research interests are hyperspectral image processing and signal processing.

Gozde Bozdagi Akar is currently a Professor with the Department of Electrical and Electronics Engineering, Middle East Technical University. Her research interests are in face recognition, 2-D and 3-D video compression, multimedia streaming and hyperspectral image processing.