

# Tensor-Based Classification Models for Hyperspectral Data Analysis

Konstantinos Makantasis, Anastasios Doulamis, Nikolaos Doulamis and Antonis Nikitakis

**Abstract**—In this paper, we present tensor-based linear and nonlinear models for hyperspectral data classification and analysis. By exploiting principles of tensor algebra, we introduce new classification architectures, the weight parameters of which satisfy the *rank-1* canonical decomposition property. Then, we propose learning algorithms to train both linear and non-linear classifiers. The advantages of the proposed classification approach are that i) it significantly reduces the number of weight parameters required to train the model (and thus the respective number of training samples), ii) it provides a physical interpretation of model coefficients on the classification output and iii) it retains the spatial and spectral coherency of the input samples. The linear tensor-based model exploits principles of logistic regression assuming the *rank-1* canonical decomposition property among its weights. For the non-linear classifier, we propose a modification of a feedforward neural network (FNN), called *rank-1* FNN, since its weights satisfy again the *rank-1* canonical decomposition property. An appropriate learning algorithm is also proposed to train the network. Experimental results and comparisons with state of the art classification methods, either linear (e.g., Linear SVM) or non-linear (e.g., deep learning) indicates the outperformance of the proposed scheme, especially in cases where a small number of training samples is available.

**Index Terms**—tensor-based classification, hyperspectral data analysis, dimensionality reduction, non-linear modeling, *rank-1* feedforward neural networks

## I. INTRODUCTION

THE recent advances in optics and photonics have stimulated the deployment of hyperspectral imaging sensors of high spatial and spectral resolution. These sensors are now placed on satellite, unmanned aerial vehicle, and ground acquisition platforms used for material, object and terrain land detection and classification [1]. Although high spatial and spectral resolution improves classification accuracy, it also imposes several research challenges derived as a consequence of the so-called "curse of dimensionality"; the difficulties arising when we need to analyze and organize data in high dimensional spaces. Hyperspectral data have their own unique characteristics, though being applied for a wide variety of applications, such as agriculture, surveillance, astronomy and biomedical imaging [2]; i) high dimensional data, ii) limited number of labeled samples and iii) large spatial variability of spectral signatures [3].

Konstantinos Makantasis is with the KIOS Research and Innovation Center of Excellence, Cyprus.

Anastasios Doulamis and Nikolaos Doulamis are with the National Technical University of Athens, Greece.

Antonis Nikitakis is with Althexis Solutions Ltd, Cyprus.

This work was supported by the European Union funded H2020 STOP-IT project, Strategic, Tactical, Operational Protection of water Infrastructure against cyber-physical Threats, grand agreement, 740610.

Most of existing works, concerning hyperspectral image classification, follow the conventional workflow of pattern recognition process, consisting of two separate steps. First, features are extracted from the raw data, creating labelled training datasets. Second, classifiers, linear or non-linear, such as Support Vector Machines (SVM) and Neural Networks (NNs) [4], are used to map the extracted features to the target (desired) outputs. The key problem, however, in applying such conventional processes in classifying high dimensional hyperspectral data is that a large number of labelled training samples is required to model the statistical input diversities and consequently to well train the classifier. In remote sensing applications, collection of a large number of labelled data is an expensive and time consuming process. Another drawback is that classifiers are often used as "black boxes" [5]. This means that there is no a direct interpretation of how spatial and spectral bands contribute to the final classification outcome.

One way to address issues deriving from the high dimensionality and heterogeneity of the data is to employ statistical learning methods [6]. However, even in this case, the problem of extracting a set of appropriate features remains. Features significantly affect the classification outcome and for data laying in high dimensional spaces is really an arduous task to estimate a suitable set of discriminant features so as to increase the accuracy of the classifier.

For this reason, recently deep learning paradigms have been investigated for classifying hyperspectral data [7], [8], [9], [10]. Deep learning machines receive as inputs, instead of features, the raw sensory data. Then, they non-linearly transform the raw inputs to hierarchies of representations which are used, in the following, as "the most suitable features" in a supervised mapping phase. Thus, deep learning tackles feature-related issues. This is also proven by the current research outcomes [11], [12], [13], [14] indicating the outperformance of deep learning machines in accurately detecting various objects in hyperspectral imaging data. Examples include the detection of man-made constructions rather than natural ones [15], [16], vehicles' detection [17], object tracking [18], land cover mapping [19] and critical infrastructure assessment [20].

However, a typical deep learning architecture contains a huge number of tunable parameters implying that a large number of samples is also needed to accurately train the network. In addition, deep learning processes present high computational complexity.

Tensor-based machine learning models are promising alternatives for hyperspectral data classification [21], [22]. In particular, Zhou et al. [21] and Tan et. al [22] have introduced a linear model using tensor-based regression with applications in neuro-imaging data analysis [21] and for classification [22].

These approaches are considered as the first works discussed the statistical inference procedure for the general raw tensor regression. In conventional learning models, usually the inputs are vector data. Therefore, in case of multidimensional input arrays, first tensor vectorization is carried out. However, vectorization destroys the inherent spatial and spectral structure of the input which can offer a physical interpretation of how spatial information and spectral bands contribute to the classification outcome. Furthermore, tensor vectorization fails to address issues stem from the high dimensionality of the data since again a large number of tunable parameters is required. To handle these limitations, we need to consider the input data as tensors, keeping the spatial and spectral structure of the data, and then, using principles of tensor algebra, to find out ways to reduce the number of parameters needed to be estimated during training.

### A. Our Contribution

In this work we propose a new machine learning model which receives as inputs multidimensional tensor signals as they are the hyperspectral images, i.e., 3D tensor cubes. The model weight parameters satisfy the *rank-1* canonical decomposition property meaning that the weight parameters are decomposed as a linear combination of a minimal number of possibly non-orthogonal rank-1 terms [23]. Using such decomposition of the model weights we are able to significantly reduce the number of parameters required to train the classifier. Thus, a smaller labelled dataset is needed than in conventional learning approaches where the tensor inputs are first vectorized.

Another advantage of the *rank-1* canonical decomposition property for the model weights is that it retains the structure of the spatial and spectral band information, which is a very important aspect for hyperspectral data classification. This is due to the fact that it actually permits the extraction of valuable information regarding the contribution of each of the hyperspectral bands to the classification. Thus, the proposed canonical decomposition provides a physical interpretation of the classification outcome, i.e., how the location of the pixels (spatial information) and the spectral bands (spectral information) influence the final classification performance.

Our work is motivated from [21] which proposes a single linear output tensor regression model for binary classification. In contrast to [21], in this paper, a multi-class classification problem is investigated using tensor-based logistic regression models of multiple outputs. In addition, the *rank-1* canonical decomposition property is also applied, apart from high-order linear, to non-linear classifiers, which is not a straightforward process. The proposed high order non-linear model is relied on a modification of a feedforward neural network (FNN), while it retains the universal approximation principles; capability of the network to approximate any unknown function, under some assumptions of continuity, within any degree of accuracy. The main difference is that the model weights satisfy the *rank-1* canonical decomposition property. Therefore, the number of parameters (and consequently the number of training samples) are significantly reduced, especially for cases where tensor

inputs are considered. A new learning algorithm is also introduced to train the network without spoiling the canonical decomposition assumption. We call the proposed high-order nonlinear model as *rank-1-FNN*.

To sum up, the main contribution of this paper is threefold. First, we introduce new learning models, linear and non-linear, that consider a *rank-1* canonical decomposition of their weight parameters. This way, a quite smaller number of weights is required compared to conventional learning paradigms. This, in the sequel, requires a smaller number of samples used to train the model, which is in the line of remote sensing applications, where a limited number of samples is available. The new models are suitable for tensor input data of high dimensions. Second, the *rank-1* canonical decomposition property allows for a physical interpretation of how each spatial and spectral band of the tensor inputs affects the classification outcome. This is an important attribute in analyzing hyperspectral data. Third, the introduction of a *rank-1* FNN non-linear classifier allows for modelling of complex relationships, due to the universal approximation property of neural nets [24], while simultaneously keeping the aforementioned advantages.

The rest of this paper is organized as follows; Section II presents the problem formulation, as well as, the notation and tensor algebra operations that will be used throughout the paper. Sections III and IV present the development of the high-order linear and nonlinear models. Experimental evaluation of the developed models is presented in Section V and Section VI concludes this work.

## II. PROBLEM FORMULATION AND TENSOR ALGEBRA NOTATION

### A. Problem Formulation

Let us denote as  $\mathbf{X}_i$  the  $i$ -th patch of a hyper-spectral image that we would like to classify into one of  $C$  available classes. Each class expresses, for example, the type of vegetation or soil properties for the patch  $\mathbf{X}_i$ . In this paper, we consider  $\mathbf{X}_i$  as a 3D tensor, i.e.,  $\mathbf{X}_i \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ , where variables  $p_1$  and  $p_2$  refer to the spatial dimensions of the hyperspectral patch and  $p_3$  to the number of spectral bands. In a more general case, variable  $\mathbf{X}_i$  can be seen as a  $D$ -dimensional tensor, that is,  $\mathbf{X}_i \in \mathbb{R}^{p_1 \times \dots \times p_D}$ .

Let us also denote as  $p_w^k(\mathbf{X}_i)$ , with  $k = 1, \dots, C$  a relationship (linear or non-linear) that expresses the probability of the observation  $\mathbf{X}_i$  to belong to the  $k$ -th class. Subscript  $w$  indicates dependence of the relationship  $p_w^k(\cdot)$  on weight parameters. Aggregating the values  $p_w^k(\cdot)$  over all  $C$  classes, we form a classification vector, say  $\mathbf{y}_i$ , the elements of which  $y_{i,k} \equiv p_w^k(\cdot)$ . Then, the maximum probability value over all  $k$  classes indicates the class to which the hyperspectral image patch belongs to.

$$\hat{k} = \arg \max_{\forall k} \{y_{i,k} \equiv p_w^k(\mathbf{X}_i)\}. \quad (1)$$

The values of  $y_{i,k}$  are estimated using machine learning algorithms. For this reason, a training dataset consisting  $N$  samples is considered

$$\mathcal{S} = \{(\mathbf{X}_i, \mathbf{t}_i)\}_{i=1}^N, \quad (2)$$

where  $\mathbf{t}_i$  is a  $C$ -dimensional vector, the elements of which  $t_{i,j}$  are all zero except for one which equals unity indicating the class to which  $\mathbf{X}_i$  belongs to. That is,  $\mathbf{t}_i \in \{0,1\}^C$  and  $\sum_{j=1}^C t_{i,j} = 1$ . In the following we omit subscript  $i$  for simplicity purposes if we refer to an input sample.

Multidimensional arrays are also known as tensors. Since tensors is a key concept of the proposed high-order learning model (linear and non-linear), in the following some basic notations and definitions regarding tensor algebra are presented that will be used through out this work.

### B. Tensor Algebra Notations and Definitions

In this paper tensors are denoted with bold uppercase letters, vectors with bold lowercase letters and scalars with lowercase letters.

**Tensor vectorization.** The  $\text{vec}(\mathbf{B})$  operator stacks the entries of a  $D$ -dimensional tensor  $\mathbf{B} \in \mathbb{R}^{p_1 \times \dots \times p_D}$  on a column vector. Specifically, an entry  $\mathbf{B} = [\dots b_{i_1, \dots, i_D} \dots]$  maps to the  $j^{\text{th}}$  entry of  $\text{vec}(\mathbf{B})$ , in which  $j = 1 + \sum_{d=1}^D (i_d - 1) \prod_{d'=1}^{d-1} p_{d'}$ .

**Tensor products.** Given two tensors  $\mathbf{A} = [\mathbf{a}_1 \dots \mathbf{a}_n] \in \mathbb{R}^{m \times n}$  and  $\mathbf{B} = [\mathbf{b}_1 \dots \mathbf{b}_q] \in \mathbb{R}^{p \times q}$  the following products can be defined;

- **Kronecker product.** The Kronecker product is the  $mp \times nq$  matrix

$$\mathbf{A} \otimes \mathbf{B} = [\mathbf{a}_1 \otimes \mathbf{B} \dots \mathbf{a}_n \otimes \mathbf{B}] = \begin{bmatrix} a_{11}\mathbf{B} & \dots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \dots & a_{mn}\mathbf{B} \end{bmatrix}$$

- **Khatri-Rao product.** The Khatri-Rao product  $\mathbf{A} \odot \mathbf{B}$  is defined as the  $mp \times n$  columnwise Kronecker product, i.e.,  $\mathbf{A} \odot \mathbf{B} = [\mathbf{a}_1 \otimes \mathbf{b}_1 \quad \mathbf{a}_2 \otimes \mathbf{b}_2 \dots \mathbf{a}_n \otimes \mathbf{b}_n]$ , if  $\mathbf{A}$  and  $\mathbf{B}$  have the same number of columns,  $n = q$ .
- **Outer product.** The outer product  $\mathbf{b}^{(1)} \circ \mathbf{b}^{(2)} \circ \dots \circ \mathbf{b}^{(N)}$  of the vectors  $\{\mathbf{b}^{(i)}\}_{i=1}^N$  forms a tensor whose element at  $(i_1, i_2, \dots, i_N)$  position equals  $\prod_{j=1}^N b_{i_j}^{(j)}$ .

**Tensor matricization.** The mode- $d$  matricization,  $\mathbf{B}_{(d)}$ , maps a tensor  $\mathbf{B}$  to a  $p_d \times \prod_{d' \neq d} p_{d'}$  matrix by arranging the mode- $d$  fibers to be the columns of the resulting matrix. In particular, the  $(i_1, \dots, i_D)$  element of  $\mathbf{B}$  maps to the  $(i_d, j)$  element of  $\mathbf{B}_{(d)}$ , where  $j = 1 + \sum_{d' \neq d} (i_{d'} - 1) \prod_{d'' < d', d'' \neq d} p_{d''}$ .

**Rank- $R$  decomposition.** A tensor  $\mathbf{B} \in \mathbb{R}^{p_1 \times \dots \times p_D}$  admits a rank- $R$  decomposition if  $\mathbf{B} = \sum_{r=1}^R \mathbf{b}_1^{(r)} \circ \dots \circ \mathbf{b}_D^{(r)}$ , where the symbol " $\circ$ " represents the vector outer product and  $\mathbf{b}_d^{(r)} \in \mathbb{R}^{p_d}$ ,  $d = 1, \dots, D$ ,  $r = 1, \dots, R$ . The decomposition can be represented by  $\mathbf{B} = [\mathbf{B}_1, \dots, \mathbf{B}_D]$ , where  $\mathbf{B}_d = [\mathbf{b}_d^{(1)}, \dots, \mathbf{b}_d^{(R)}] \in \mathbb{R}^{p_d \times R}$ ,  $d = 1, \dots, D$ . When a tensor  $\mathbf{B}$  admits a rank- $R$  decomposition the following results hold;

$$\mathbf{B}_{(d)} = \mathbf{B}_d (\mathbf{B}_D \odot \dots \odot \mathbf{B}_{d+1} \odot \mathbf{B}_{d-1} \odot \dots \odot \mathbf{B}_1)^T \quad (3)$$

and

$$\text{vec}(\mathbf{B}) = (\mathbf{B}_D \odot \dots \odot \mathbf{B}_1) \mathbf{1}_R, \quad (4)$$

where  $\mathbf{1}_R$  is a vector of  $R$  ones. For more information regarding tensor algebra please refer to [25].

## III. HIGH ORDER LINEAR MODELLING

### A. Vector Logistic Regression

Let us first assume for simplicity that the input samples are of vector forms. We denote these input vectors as  $\mathbf{x} \in \mathbb{R}^{p_1}$ , where variable  $p_1$  expresses vector dimension. Using the the logistic regression framework [26], one can model the probability function  $p_w^k(\cdot)$  as

$$p_w^k(\mathbf{x}) = \frac{\exp \mathbf{w}^{(k)T} \mathbf{x}}{\sum_{i=1}^C \exp \mathbf{w}^{(i)T} \mathbf{x}}, \quad (5)$$

where  $\mathbf{w}^{(k)} \in \mathbb{R}^{p_1}$  with  $k = 1, 2, \dots, C$  stands for the weights with respect to the  $k$ -th class. Matrix  $\mathbf{W} = [\mathbf{w}^{(1)} \mathbf{w}^{(2)} \dots \mathbf{w}^{(C)}]$  includes all the weights involved in the model. The rational meaning of the weight parameters  $\mathbf{w}^{(k)}$  is that they express the degree of confidence of the vector input  $\mathbf{x}$  to belong to the  $k$ -th out of  $C$  available classes. In addition, the elements  $w_j^{(k)}$  of  $\mathbf{w}^{(k)} = [\dots w_j^{(k)} \dots]^T$  express the degree of significance of each element of the input vector  $\mathbf{x}$  with respect to the  $k$ -th class.

One simple way to extend Eq. (5) in case that the inputs are tensors. i.e.,  $\mathbf{X}_i \in \mathbb{R}^{p_1 \times \dots \times p_D}$ , is to take the vectorized forms of them (see the respective text on algebra notation of Section II-*tensor vectorization*). The main limitation of such an approach is that (i) a large number of parameters is needed to be estimated, particularly  $(C \prod_{l=1}^D p_l)$  and (ii) vectorization spoils the spatial structure of tensor inputs, that is, pixels belonging to a neighboring region frequently present similar properties (spatial coherency). A large number of parameters also implies a large number of available labelled observations in order to successfully complete the training procedure. However, usually the available labeled samples are limited due to manual effort required to collect and annotate them.

### B. Matrix Logistic Regression

In case of matrix input observations,  $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2}$ , one can reduce the number of model parameters by taking into consideration concepts of [27], applied for electroencephalogram data classification. Then, the logistic regression model is given by

$$p_w^k(\mathbf{X}) = \frac{\exp \mathbf{w}_1^{(k)T} \mathbf{X} \mathbf{w}_2^{(k)}}{\sum_{i=1}^C \exp \mathbf{w}_1^{(i)T} \mathbf{X} \mathbf{w}_2^{(i)}}, \quad (6)$$

where  $\mathbf{w}_1^{(i)} \in \mathbb{R}^{p_1}$  and  $\mathbf{w}_2^{(i)} \in \mathbb{R}^{p_2}$ . We also define as  $\mathbf{W}_1 = [\mathbf{w}_1^{(1)} \mathbf{w}_1^{(2)} \dots \mathbf{w}_1^{(C)}]$  and  $\mathbf{W}_2 = [\mathbf{w}_2^{(1)} \mathbf{w}_2^{(2)} \dots \mathbf{w}_2^{(C)}]$  the aggregate model parameters.

The key advantage of Eq. (6) is that the number of required parameters is  $C(p_1 + p_2)$  instead of  $C(p_1 \cdot p_2)$  that would be needed if one vectorize matrix observation input data  $\mathbf{X}$ . This means that the weight parameters  $\mathbf{w}^{(k)}$ , in relation (6), for the  $k$ -th class are rank-1 canonically decomposed into the weights of  $\mathbf{w}_1^{(k)} \in \mathbb{R}^{p_1}$  and  $\mathbf{w}_2^{(k)} \in \mathbb{R}^{p_2}$ , that is,

$$\mathbf{w}^{(k)} = \mathbf{w}_2^{(k)} \otimes \mathbf{w}_1^{(k)}, \quad (7)$$

where operator  $\otimes$  is the Kronecker product as defined in Section II-*Kronecker product*.

Using Eq. (7), one can write Eq. (6) as

$$p_w^k(\mathbf{X}) = \frac{\exp[(\mathbf{w}_2^{(k)} \otimes \mathbf{w}_1^{(k)})^T \text{vec}(\mathbf{X})]}{\sum_{i=1}^C \exp[(\mathbf{w}_2^{(i)} \otimes \mathbf{w}_1^{(i)})^T \text{vec}(\mathbf{X})]}. \quad (8)$$

### C. Tensor-based Logistic Regression

The aforementioned concept can be extended in case of tensor inputs,  $\mathbf{X} \in \mathbb{R}^{p_1 \times \dots \times p_D}$ , [21]. In this way, we have that

$$p_w^k(\mathbf{X}) = \frac{\exp[(\mathbf{w}_D^{(k)} \otimes \dots \otimes \mathbf{w}_1^{(k)})^T \text{vec}(\mathbf{X})]}{\sum_{i=1}^C \exp[(\mathbf{w}_D^{(i)} \otimes \dots \otimes \mathbf{w}_1^{(i)})^T \text{vec}(\mathbf{X})]}, \quad (9)$$

where  $\mathbf{w}_l^{(k)} \in \mathbb{R}^{p_l}$  with  $l = 1, 2, \dots, D$  is the  $l$ -th *rank*-1 canonical decomposition of the weight parameters  $\mathbf{w}^{(k)}$  for the  $k$ -th class. This property is referred as CANDECOMP, stem from CANonical DECOMPosition, or PARAFAC, stem from PARAllel FACTors, in the literature [28], [29]. That is,

$$\mathbf{w}^{(k)} = \mathbf{w}_D^{(k)} \otimes \dots \otimes \mathbf{w}_1^{(k)}. \quad (10)$$

Eq. (9) can be seen as an expression of the Khatri-Rao product, denoted as operator  $\odot$  (see Section II), which is the column-wise Kronecker product of the *rank*-1 canonical decomposition weight parameters  $\mathbf{w}_l^{(k)}$ ,

$$p_w^k(\mathbf{X}) = \frac{\exp[(\mathbf{w}_D^{(k)} \odot \dots \odot \mathbf{w}_1^{(k)})^T \text{vec}(\mathbf{X})]}{\sum_{i=1}^C \exp[(\mathbf{w}_D^{(i)} \odot \dots \odot \mathbf{w}_1^{(i)})^T \text{vec}(\mathbf{X})]}, \quad (11)$$

In Eqs. (9) and (11), we can aggregate the total weight parameters as

$$\mathbf{W}_l = [\mathbf{w}_l^{(1)} \mathbf{w}_l^{(2)} \dots \mathbf{w}_l^{(C)}], \quad (12)$$

with  $l = 1, 2, \dots, D$ . In other words, matrix  $\mathbf{W}_l$  contains all the weight parameters with respect to the  $l$ -th dimension of the tensor  $\mathbf{X}$  overall classes, i.e, related with the  $p_l$  tensor dimension.

The representation of Eqs. (9) and (11) significantly reduces the number of model parameters needed for classifying the tensor inputs  $\mathbf{X}$ . In particular, the vector-based logistic regression model derived through vectorization of tensor  $\mathbf{X}$  requires the estimation of  $C \prod_{l=1}^D p_l$  parameters, while the tensor-based representation of (9) and (11) reduces this number to  $C \sum_{l=1}^D p_l$ .

The advantages of the aforementioned proposed representation for the classification of hyperspectral images are twofold. First, as we reduce the number of weight parameters, a smaller number of labelled data samples is required to train the logistic regression model. This is an important factor for developing classifiers that can better generalize to unseen hyperspectral data. Usually, the manual effort for the annotation is laborious and therefore a small number of labelled training data is available. Second, the *rank*-1 canonical decomposition of the model weights provides a physical interpretation of how the spatial and spectral information of the hyperspectral tensor input  $\mathbf{X}$  contributes to the classification. In particular, according to the statements made right after Eq. (5), the weights  $\mathbf{w}^{(k)}$  express

the degree of importance of the tensor input  $\mathbf{X}$  to belong to the  $k$ -th class. Since these weights are canonically decomposed into  $D$  separate weights  $\mathbf{w}_l^{(k)}$ , each indicating the contribution of  $p_l$  tensor dimension to the belong to the  $k$ -th class, the proposed model provides a quantitative representation of how the elements of each tensor dimension tunes the classification performance.

More specifically, in case of hyperspectral imaging, the dimension of tensor inputs equals 3, i.e.,  $D = 3$ . The first two dimensions refer to the spatial properties of the pixel data, while the third one to the spectral bands. Therefore, the first two decomposed weights, i.e.,  $\mathbf{w}^{(1)}$  and  $\mathbf{w}^{(2)}$ , express how the pixel spatial coherency affects the classification outcome. On the other hand, the third decomposed weight vector  $\mathbf{w}^{(3)}$  indicates how the spectral bands influence the classification and which of the spectral bands are the most salient. This property of the proposed model is very important towards the analysis of hyperspectral data since it facilitates the interpretation of the results and quantifies the importance of the spectral bands on the classification compressing the influence of the bands that are of less importance.

### D. Estimation of the Model Weights

The model weight parameters are estimated through a training set  $\mathcal{S} = \{(\mathbf{X}_i, t_i)\}_{i=1}^N$  of  $N$  samples as Eq. (2) indicates. We recall that if  $\mathbf{X}_i$  belongs to the  $k^{\text{th}}$  class, then  $t_i$  is a vector with all elements zero except the element  $k$ , which equals one, i.e.,  $t_{i,j} = 0, \forall j \neq k$  and  $t_{i,k} = 1$ . By minimizing the negative log-likelihood function,

$$L(\mathbf{W}_1, \dots, \mathbf{W}_D; \mathcal{S}) = - \sum_{i=1}^N \sum_{k=1}^C t_{i,k} \log p_w^k(\mathbf{X}_i), \quad (13)$$

the weight parameters can be estimated as

$$\{\hat{\mathbf{W}}_1, \hat{\mathbf{W}}_2, \dots, \hat{\mathbf{W}}_D\} = \arg \min_{\mathbf{W}_l} L(\mathbf{W}_1, \dots, \mathbf{W}_D; \mathcal{S}). \quad (14)$$

In Eq. (14), matrices  $\hat{\mathbf{W}}_l, l = 1, \dots, D$  refer to the optimal estimates of the  $\mathbf{W}_l$  weight parameters, expressing the contribution of the  $l$ -dimension of the tensor input to the classification for all  $C$  available classes.

Based on the statements of Section II, presenting some basic notations on tensor algebra (*Rank-R decomposition and tensor matricization*), it is held that

$$\begin{aligned} \langle \mathbf{w}_D^{(k)} \odot \dots \odot \mathbf{w}_1^{(k)}, \mathbf{X} \rangle = \\ \langle \mathbf{w}_l^{(k)}, \mathbf{X}_{(l)} (\mathbf{w}_D^{(k)} \odot \dots \odot \mathbf{w}_{l+1}^{(k)} \odot \mathbf{w}_{l-1}^{(k)} \odot \dots \odot \mathbf{w}_1^{(k)}) \rangle. \end{aligned} \quad (15)$$

The proof of Eq. (15) is given in Appendix A. In Eq. (15),  $\mathbf{X}_{(l)}$  denotes the mode- $l$  matricization of tensor  $\mathbf{X}$  obtained by keeping the  $l$ -dimension intact and concatenating the slices of the rest dimensions into a long matrix [25]. We also recall that the operator  $\odot$  refers to the Khatri-Rao product, while the  $\langle \cdot, \cdot \rangle$  to the inner product between two tensors.

The left hand of Eq. (15) expresses the arguments of  $\exp(\cdot)$  involved in the linear logistic regression model of Eq. (11). Therefore, (11) can be rewritten by taking into account the

right-hand of (15). In the sequel, generation of (15) over all available classes is obtained as

$$\langle \mathbf{W}_D \odot \cdots \odot \mathbf{W}_1, \mathbf{X} \rangle = \text{diag}(\langle \mathbf{W}_l, \mathbf{X}_{(l)}(\mathbf{W}_D \odot \cdots \odot \mathbf{W}_{l+1} \odot \mathbf{W}_{l-1} \odot \cdots \odot \mathbf{W}_1) \rangle). \quad (16)$$

As we can see Eq. (15) is actually the inner product of two vectors. The first is the weight parameters  $w_l^{(k)}$  expressing the contribution of the  $l$  dimension of the tensor input as far as the  $k$ -th class is concerned. On the other hand, the second vector is independent from the values of  $w_l^{(k)}$ . Therefore, one approach for optimally estimating the weights  $w_l^{(k)}$  and consequently  $w^{(k)}$  [see Eq. (10)] is one to consider all the weight parameters  $w_q^{(k)}$  with  $q \neq l$  apart from the  $l^{th}$  fixed and then solving with respect to  $w_l^{(k)}$ . This procedure is iteratively applied for all weight parameters  $w_l^{(k)}$  until some termination criteria are met. Similar statements can be concluded for the matrix representation of Eq. (16).

Actually, the learning algorithm used to optimally estimate the model weights in case that Eq. (10) is held, i.e., the *rank*-1 canonical decomposition, simulates a regression learning where we use as inputs the data of the right-hand of Eq. (16), that is,

$$\mathbf{X}_{(l)}(\mathbf{W}_D \odot \cdots \odot \mathbf{W}_{l+1} \odot \mathbf{W}_{l-1} \odot \cdots \odot \mathbf{W}_1). \quad (17)$$

Therefore, the parameter estimation problem can be solved by using conventional software such as scikit-learn in python<sup>1</sup>. The proposed learning procedure, considering the *rank*-1 canonical decomposition of the weight parameters is described in Algorithm 1.

---

**Algorithm 1:** Optimal Estimation of Model Weights using *rank*-1 Canonical Decomposition

---

**Initialization:**

1. Set Iteration Index  $n \rightarrow 0$
2. Randomize all the weight parameters  $\mathbf{W}_l(n) \in \mathbb{R}^{p_l \times C}$  for all  $l = 1, \dots, D$

**3. repeat**

**for**  $l = 1, \dots, D$  **do**

3.1 Calculate the matrix of Eq.(17), that is,  
 $\mathbf{X}_{(l)}(\mathbf{W}_D(n) \odot \cdots \odot \mathbf{W}_{l+1}(n) \odot \mathbf{W}_{l-1}(n) \odot \cdots \odot \mathbf{W}_1(n))$

3.2 Update matrix  $\mathbf{W}_l(n)$  using a regression learning algorithm through minimization of Eq. (14) such as

$$\mathbf{W}_l(n+1) \rightarrow \arg \min_{\mathbf{W}_l} L(\mathbf{W}_1(n+1), \dots, \mathbf{W}_{l-1}(n+1), \mathbf{W}_l(n), \dots, \mathbf{W}_D(n))$$

**end**

Set  $n \rightarrow n + 1$

**until** termination criteria are met;

---

Although, the high-order linear model can provide physically interpretable results, due to its structure, it is restricted to produce decision boundaries that are linear in the input space.

<sup>1</sup><http://scikit-learn.org/>

This implies that this model is not able to cope with more complicated problems, where non-linear decision boundaries are necessary to provide classification results of high accuracy. This motivates us to extend the previous linear regression model to a non-linear one. The proposed non-linear model should assume again a *rank*-1 canonical decomposition of the weight parameters in order to retain the aforementioned advantages in hyperspectral classification.

#### IV. HIGH-ORDER NON-LINEAR MODELLING

The proposed high-order non-linear model is based on the concepts of the previously discussed linear logistic regression filters with the difference that a nonlinear transformation is applied on the input data. This means, in other words, that the probability  $p_w^k(\cdot)$  of an input observation  $\mathbf{X}$  to belong to one of the  $C$  available classes is non-linearly interwoven with respect to the input tensor data and the weight parameters that influence the importance of these data on classification performance through a function  $f_w(\cdot)$ , i.e.,

$$p_w^k(\cdot) = f_w(\mathbf{X}). \quad (18)$$

The main difficulty in implementing Eq. (18) is that function  $f_w(\cdot)$  is actually unknown. One way to parameterize the unknown function  $f_w(\cdot)$  is to exploit the principles of the universal approximation theorem, stating that a function, under some assumptions of continuity, can be approximated by a feedforward neural network with a finite number of neurons within any degree of accuracy [30]. Feedforward neural networks have been proven as a reliable framework for image classification [31].

However, the main difficulty in applying a feedforward neural network (FNN) for hyperspectral classification problems is twofold. First, a large number of weight parameters is required to be learned, proportionally to  $Q \prod_{l=1}^D p_l + QC$ , where variable  $Q$  refers to the number of hidden neurons of the network. This outcome derives as a consequence of the structure of the network as is briefly described below (see Section IV-A). This, in the sequel, implies that a large number of labelled samples needed to successfully train the network. Second, the weights of the network are not directly related to the physical properties of hyperspectral data and how these properties affect the classification performance, since the inputs are vectorized and thus they do not preserve their structure. Networks are often treated as "black boxes" when they are applied for classification of hyperspectral data. To overcome these problems, we propose a modification of conventional feedforward neural network structures so that network weights from the input to the hidden layer satisfy the *rank*-1 canonical decomposition according to statements of the previous section. In addition, we introduce a learning algorithm to train the so-called *rank*-1 canonical decomposition feedforward neural network - *rank*-1 FNN. Before proposing *rank*-1 FNN, we briefly describe how  $p_w^k(\cdot)$  is modelled through a FNN.

##### A. FeedForward Neural Network Modelling

Fig. 1 illustrates a feedforward neural network nonlinearly approximating the probability  $p_w^k(\cdot)$ . Initially, the tensor input

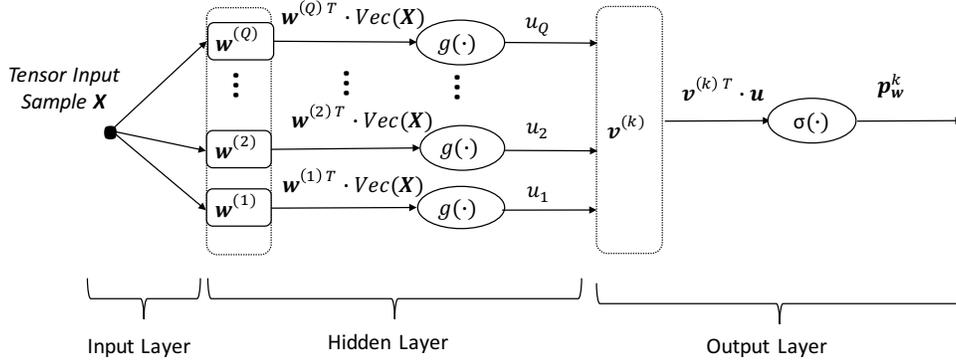


Fig. 1. The structure of a feedforward neural network.

$\mathbf{X}$  is vectorized creating a vector,  $\text{vec}(\mathbf{X})$ , of size  $\prod_{l=1}^D p_l$ . The network is assumed to have  $Q$  hidden neurons. Each neuron is associated with an activation function  $g(\cdot)$ . In this paper, the sigmoid function  $g(x) = 1/(1 + \exp(-x))$  is selected, where factor  $a$  regulates the steepness of the function. The activation function of the  $i$ -th out of  $Q$  hidden neurons receives as input the inner product of  $\text{vec}(\mathbf{X})$  and a weight vector  $\mathbf{w}^{(i)}$  associated with the  $i$ -th neuron and produces as output a scalar  $u_i$  given by [32]

$$u_i = g(\mathbf{w}^{(i)T} \text{vec}(\mathbf{X})) \equiv g(\langle \mathbf{w}^{(i)}, \text{vec}(\mathbf{X}) \rangle), \quad (19)$$

where we recall that the operator  $\langle \cdot, \cdot \rangle$  expresses the inner product. It should be mentioned that in the current notation the superscript  $i$  of the weights  $\mathbf{w}^{(i)}$  refers to the  $i$ -th hidden neuron. Gathering the responses of all hidden neurons in one vector  $\mathbf{u} = [u_1, u_2, \dots, u_Q]^T$ , we have that

$$\mathbf{u} = g(\langle \mathbf{W}, \mathbf{X} \rangle), \quad (20)$$

where  $\mathbf{W} = [\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(Q)}]^T$  is a matrix containing the weights  $\mathbf{w}^{(i)}$  for all hidden neurons,  $i = 1, \dots, Q$ . Thus, the output of the network is given as

$$\mathbf{p}_w^k = \sigma(\langle \mathbf{v}^{(k)}, \mathbf{u} \rangle) \equiv \sigma(\mathbf{v}^{(k)T} \mathbf{u}), \quad (21)$$

where  $\sigma(\cdot)$  stands for the softmax function,  $\mathbf{v}^{(k)}$  the weights between the hidden and the output layer and the superscript of the weights for the  $k$ -th class. The softmax function corresponds to the following conditional probability and is defined as follows for a class  $i$

$$P(y = i|x) = \frac{\exp(\mathbf{x}^T \cdot \mathbf{w}_i)}{\sum_{j=1}^{j=C} \exp(\mathbf{x}^T \cdot \mathbf{w}_j)}. \quad (22)$$

### B. Rank-1 Canonical Decomposition Feedforward Neural Networks - Rank-1 FNN

To reduce the number of parameters of the network and to relate the classification results to the spatial and spectral properties of hyperspectral input data, we  $\text{rank-1}$  canonically decomposed the weight parameters  $\mathbf{w}^{(i)}$  as in Eq. (10). We should stress that vector  $\mathbf{w}_l^{(i)}$  relates the input tensor data with the  $i$ -th hidden neuron and therefore  $i = 1, 2, \dots, Q$ .

Then, taking into account the properties of Eq. (15), the output of the  $i$ -th hidden neuron  $u_i$  can be written as

$$\begin{aligned} u_i &= g(\langle \mathbf{w}^{(i)}, \mathbf{X} \rangle) \\ &= g(\langle \mathbf{w}_D^{(i)} \otimes \dots \otimes \mathbf{w}_1^{(i)}, \mathbf{X} \rangle) \\ &= g(\langle \mathbf{w}_D^{(i)} \odot \dots \odot \mathbf{w}_1^{(i)}, \mathbf{X} \rangle) \\ &= g(\langle \mathbf{w}_l^{(i)}, \boldsymbol{\tau}_{\neq l}^{(i)} \rangle). \end{aligned} \quad (23)$$

In Eq. (23), vector  $\boldsymbol{\tau}_{\neq l}^{(i)}$  is a transformed version of tensor input  $\mathbf{X}$ . Vector  $\boldsymbol{\tau}_{\neq l}^{(i)}$  is independent from  $\mathbf{w}_l^{(i)}$ . That is,

$$\boldsymbol{\tau}_{\neq l}^{(i)} = \mathbf{X}_{(l)}(\mathbf{w}_D^{(i)} \odot \dots \odot \mathbf{w}_{l+1}^{(i)} \odot \mathbf{w}_{l-1}^{(i)} \odot \dots \odot \mathbf{w}_1^{(i)}), \quad (24)$$

where we recall that  $\mathbf{X}_{(l)}$  is the mode- $l$  matricization of  $\mathbf{X}$ .

Eq. (23) actually resembles the operation of a single perceptron with inputs the weights  $\mathbf{w}_l^{(i)}$  and the transformed version  $\boldsymbol{\tau}_{\neq l}^{(i)}$  of the input data. That is, if the  $\text{rank-1}$  canonically decomposed weights  $\mathbf{w}_r^{(i)}$  with  $r \neq l$  are known then  $\boldsymbol{\tau}_{\neq l}^{(i)}$  will be also known. Fig. 2 shows the structure of the proposed  $\text{rank-1}$  FNN. The model consists of one input layer, one hidden and one output layer. The main modification of this structure compared to a conventional FNN is in the hidden layer. More specifically, the weights of a hidden neuron are first decomposed into  $D$  canonical factors, each expresses the spatial and spectral band contribution to the classification performance. In this figure, we have shaded the weight vector  $\mathbf{w}_l^{(i)}$  since all the rest weight vectors are used to estimate the transformed input  $\boldsymbol{\tau}_{\neq l}^{(i)}$ .

### C. The Learning Algorithm

To train the proposed  $\text{rank-1}$  FNN model the set  $\mathcal{S} = \{(\mathbf{X}_i, \mathbf{t}_i)\}_{i=1}^N$  containing  $N$  samples is used. The learning algorithm minimizes the negative log-likelihood (see relation (13)) with respect to network responses  $\mathbf{y}_i = [\dots y_{i,k} \dots]^T$ , with  $y_{i,k} \equiv p_w^k(\mathbf{X}_i)$ , and targets  $\mathbf{t}_i$  over all training samples.

In case of using a conventional neural network training algorithm, the estimated weights do not satisfy the  $\text{rank-1}$  canonical decomposition assumption expressed by Eq. (10). For this reason, in the proposed learning algorithm we initially fix all the weights  $\mathbf{w}_r^{(i)}$  with  $r \neq l$ . This way, we are able to estimate the transformed input vector  $\boldsymbol{\tau}_{\neq l}^{(i)}$ . Therefore, the only unknown parameters of the hidden layer is the vector  $\mathbf{w}_l^{(i)}$ .

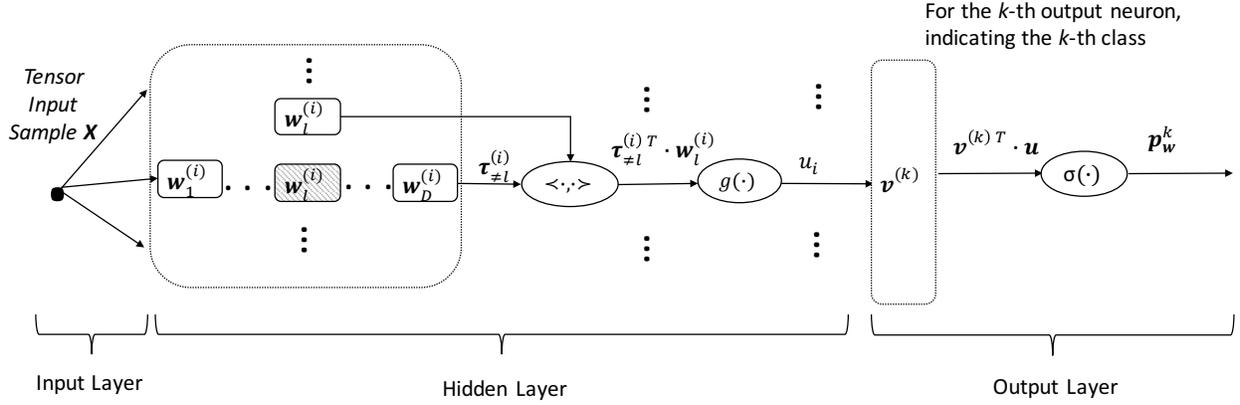


Fig. 2. The structure of the proposed *rank-1* feedforward neural network.

This can be derived through a gradient based optimization algorithm, assuming that the derivative  $\partial L/\partial w_l^{(i)}$  is known. Then, the weights are updated towards the negative direction of the partial derivative.

One way to estimate the partial derivative  $\partial L/\partial w_l^{(i)}$  is to exploit principles of the back-propagation algorithm which actually implements the chain rule property for estimating the derivative of the error with respect to the network weights. In particular, by using the back-propagation algorithm, we compute the partial derivatives

$$\begin{aligned} \partial L/\partial w_l^{(i)} \text{ for } l = 1, \dots, D \text{ and } i = 1, \dots, Q \\ \partial L/\partial v^{(k)} \text{ for } k = 1, \dots, C. \end{aligned} \quad (25)$$

The partial derivative  $\partial L/\partial w_l^{(i)}$  can be estimated if we assume all the weights  $w_r^{(i)}$  with  $r \neq l$  fixed, since in this case we can estimate the vector  $\tau_{neq l}^{(i)}$ . Therefore, estimation of the parameters of the *Rank-1* FNN is obtained by iteratively solving with respect to one of the  $D$  canonical decomposed weight vectors, assuming the remaining fixed. Algorithm 2 presents the main steps of the proposed algorithm, applying for the calculation of the hidden layer weights of the network.

## V. EXPERIMENTAL RESULTS ON HYPERSPECTRAL DATA

A natural question that arises is whether such a reduction in the number of parameters would limit the descriptive power of the models in (9) and in (21). To answer this question we present quantitative results regarding classification accuracy. Furthermore, we compare the performance of the linear high-order model (tensor-based logistic regression) against two other well-known linear models; vector-based logistic regression and linear SVMs, and the performance of the high-order nonlinear model against Fully Connected Feed Forward Neural Nets (FCFNN), nonlinear SVMs and two deep learning approaches for classifying hyperspectral data; an approach based on Stacked-Autoencoders [11] and an approach based on the exploitation of Convolutional Neural Networks [12]. These methods are the current state of the art in the literature.

The model in (9) is linear in the feature space and the estimated parameters can be used to estimate the importance of the spectral bands. In particular, the most important feature

### Algorithm 2: Optimal Estimation of Hidden Layer Model Weights of the *Rank-1* FNN

**Initialization:**

1. Set Iteration Index  $n \rightarrow 0$
2. Randomize all the weight vectors  $w_l^{(i)}(n) \in \mathbb{R}^{p_l \times 1}$  for all  $l = 1, \dots, D$  and for all  $i = 1, 2, \dots, Q$
3. Randomize all the weight vectors  $v^{(k)}(n) \in \mathbb{R}^{Q \times 1}$  for all  $k = 1, \dots, C$

**repeat**

**for**  $l = 1, \dots, D$  **do**

**for**  $i = 1, \dots, Q$  **do**

- 4.1 Estimate the transformed input vector  $\tau_{neq l}^{(i)}$   
 $\tau_{neq l}^{(i)} = \mathbf{X}_{(l)}(w_D^{(i)}(n) \odot \dots \odot w_{l+1}^{(i)}(n) \odot w_{l-1}^{(i)}(n+1) \odot \dots \odot w_1^{(i)}(n+1)),$
- 4.2 Estimate the derivative error of  $\partial L/\partial w_l^{(i)}$
- 4.3 Update the weights  $w_l^{(i)}(n)$  towards the negative direction of derivative

**end**

**end**

**for**  $k = 1, \dots, C$  **do**

- 4.4 Estimate the derivative error of  $\partial E/\partial v^{(k)}$
- 4.5 Update the weights  $v^{(k)}(n)$  towards the negative direction of derivative

**end**

Set  $n \rightarrow n + 1$

**until** termination criteria are met;

elements should have the highest, in absolute value, weight coefficients, while feature elements uncorrelated with the output variables should have coefficient values close to zero. This way, we can reduce the dimensionality of raw data keeping feature elements with the highest weight coefficients.

#### A. Datasets

In our study we used AVIRIS and ROSIS sensors hyperspectral datasets <sup>2</sup>. In particular, we used i) the Indian Pines dataset, which depicts a test site in North-western Indiana and

<sup>2</sup>The code for the Rank-1 FNN is available at <https://bit.ly/2GRMd85>

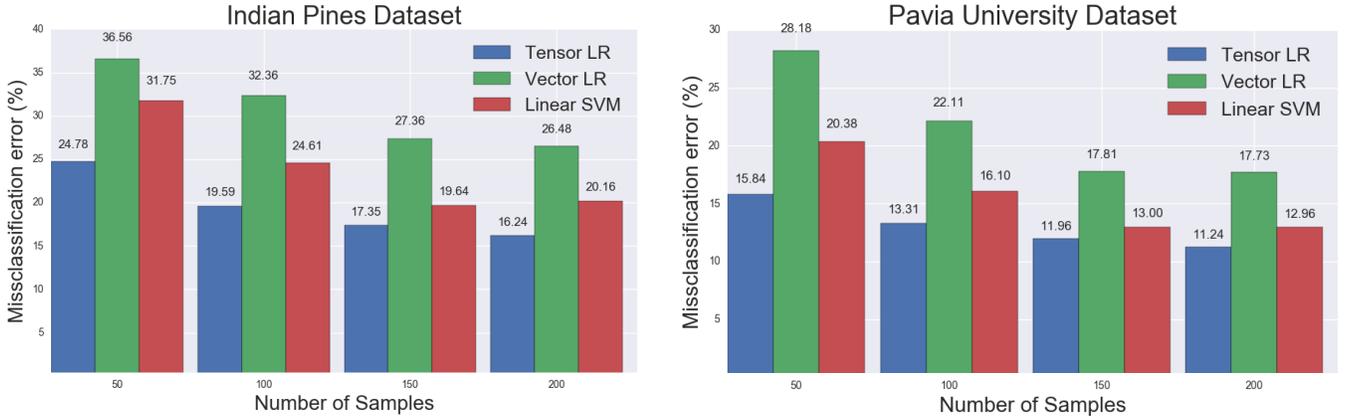


Fig. 3. Misclassification error on test set for tensor-based logistic regression, vector-based logistic regression and linear SVMs.

consists of 224 spectral reflectance bands and ii) the Pavia university datasets, whose number of spectral bands are 103.

A hyperspectral image is represented as a 3D tensor of dimensions  $p_1 \times p_1 \times p_3$ , where  $p_1$  and  $p_2$  correspond to the height and width of the image and  $p_3$  to the spectral bands. In order to classify a pixel  $I_{x,y}$  at location  $(x, y)$  on image plane and successfully fuse spectral and spatial information, we use a square patch of size  $s \times s$  centered at pixel  $I_{x,y}$ . Let us denote as  $t_{x,y}$  the class label of the pixel at location  $(x, y)$  and as  $\mathbf{X}_{x,y}$  the tensor patch centered at pixel  $I_{x,y}$ . Then, we can form a dataset  $S = \{(\mathbf{X}_{x,y}, t_{x,y})\}$  for  $x = 1, 2, \dots, p_2$  and  $y = 1, 2, \dots, p_1$ . Each one of the patches  $\mathbf{X}_{x,y}$  is also a 3D tensor with dimension  $s \times s \times p_3$ , which contains spectral and spatial information for the pixel located at  $(x, y)$ . The dataset  $S$  is used to train the classifiers. The Pavia University and the Indian Pines datasets contain 42,776 and 10,086 labeled pixels respectively.

### B. Classification Accuracy of the Tensor-based Logistic Regression

The performance of the tensor-based logistic regression method is evaluated into four different experiments, each of different number of training samples to assess classification accuracy even in cases where a small number of samples is used. The evaluation is compared against vector-based logistic regression and linear SVMs. For each of the four experiments, we use as training dataset a subset of  $S$  that contains 50, 100, 150 and 200 samples from each class respectively. The remaining samples are used as test data.

The same training data are used to train all models. The vector-based logistic regression as well as the linear SVMs are trained by using the vectorized versions of patches  $\mathbf{X}_{x,y}$  that belong to the training datasets. The tensor-based logistic regression model is trained by using the same patches without applying any transformation on them, and thus, keeps intact their spatial structure.

Figure 3 presents the classification accuracy, on the test set, for all tested methods. In both datasets and in all cases the tensor-based model outperforms linear SVMs and vector-based logistic regression, despite the fact that it employs the smallest number of parameters. These results quantitatively answer the

question regarding the capacity of the model in (9) under a small sample setting framework.

### C. Tensor-based Dimensionality Reduction

In the following we evaluate the quality of the dimensionality reduction that can be conducted by the tensor-based logistic regression model. Towards this direction, we utilize the model presented in [12], where a deep Convolutional Neural Network (CNN) is used for classifying hyperspectral data. The authors reduce the dimensionality of the data along the spectral dimension by applying principal components analysis and by utilizing the  $n$  principal components that preserve at least 99.9% of the total dataset variance.

In this work we utilize the same CNN structure, but we reduce the dimensionality of the raw data by selecting the  $n$  most significant spectral bands, i.e., spectral bands to which the tensor-based logistic regression model assigns the larger, in terms of absolute value, coefficients. Due to the fact that this method does not take the variation of estimation into account, we firstly normalized the data, so as to suppress the effect of variance. The results in terms of classification accuracy on the test set, using the same CNN for both dimensionality reduction methods are presented in Table I. Each experiment has been executed at 10 different runs and the standard deviation across all different runs is also depicted in Table I. In this table, we denote as TB-CNN the tensor-based dimensionality reduction followed by a CNN classifier and as PCA-CNN the PCA dimensionality reduction followed again by a CNN model. We conducted three experiments where we use 10%, 20% and 40% of the datasets as the training samples.

PCA-CNN performs slightly better than TB-CNN. However, our proposed method presents an advantage over PCA. Although, PCA maps the raw data to a lower dimensional feature space, the resulted features are not interpretable. Using the tensor-based dimensionality reduction, features at the lower dimensional space correspond to the most significant spectral bands providing a physical meaning on how spectral information affects the classification performance.

Furthermore, we compare the features extracted by the proposed tensor-based logistic regression model (high order linear modelling - Section III) with the features extracted using

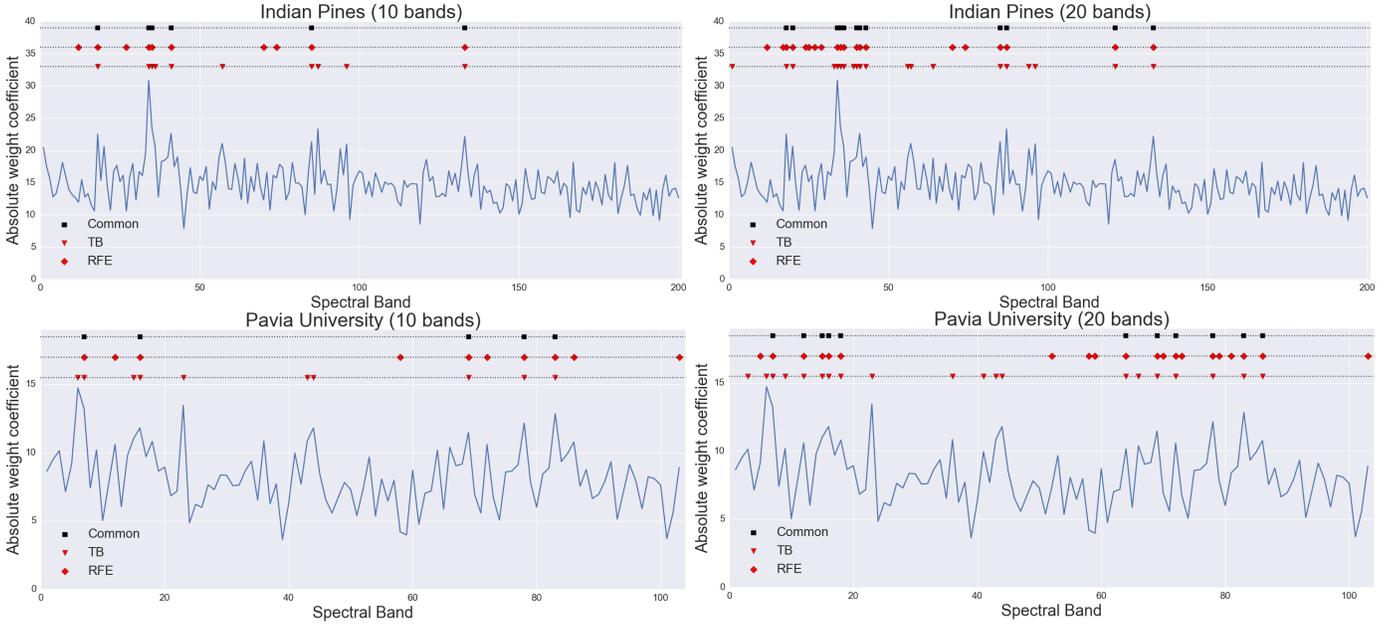


Fig. 4. Selected spectral bands using the tensor-based (TB) logistic regression model and the RFE technique. Figures on the left depict the selection of the ten (10) most important bands, while figures on the right depict the selection of the twelve (20) most important bands. Common spectral bands of both techniques are also depicted, denoting as Common in the figure.

TABLE I  
OVERALL CLASSIFICATION ACCURACY RESULTS (%) OF THE TENSOR-BASED LOGISTIC REGRESSION MODEL FOR BOTH DATASETS.

Pavia University			
Splitting ratio	10%	20%	40%
TB-CNN	97.33 ± 0.5	98.41 ± 0.3	99.31 ± 0.2
PCA-CNN	97.58 ± 0.3	98.53 ± 0.3	99.42 ± 0.1
Indian Pines			
Splitting ratio	10%	20%	40%
TB-CNN	82.29 ± 1.1	90.09 ± 0.7	95.98 ± 0.6
PCA-CNN	84.34 ± 0.9	91.72 ± 0.7	96.37 ± 0.5

the Recursive Feature Elimination (RFE) technique [33]. RFE is an iterative procedure. At each iteration, it trains a model using all available spectral bands and then it discards the less informative ones until a predefined number of features is reached. Fig. 4 presents the 10 and 20 most important spectral bands selected from both techniques (i.e., the tensor-based, called TB in the figure and the RFE). In this figure, the common bands detected by both techniques are also identified for clarification purposes. The results have been conducted for the India Pines and Pavia University dataset. To verify the quality of the selected spectral bands as important, we measure the classification accuracy obtained if we keep only these bands. Specifically, the extracted bands of both techniques are used to train the CNN model presented in [12]. Please note that the training set used for both feature selection techniques (i.e., TB and RFE) is the same and it was built by selecting 100 random samples per class. The results are presented in Table III for Pavia dataset. In addition, Fig. 5 shows the confusion matrix for the Pavia dataset when 10 or 20 salient bands are selected for the TB and RFE approach.

The overall classification accuracy for both dimensionality reduction techniques is almost the same, indicating that the proposed TB method is also appropriate for dimensionality reduction apart from classification capabilities - see Section V-B.

The features extracted by the tensor-based logistic regression model are also compared against the features selected by Tensor Discriminative Locality Alignment (TDLA) [34], a state-of-the-art method for tensor data dimensionality reduction (see Table II). At this point, we should mention that TDLA has been designed for dimensionality reduction, while our method is a tensor-based classification modeling framework. Feature selection is an interesting side effect of our proposed methodology.

Again the quality of the data dimensionality reduction is measured through classification accuracy. In this experiment we have used different classification models, instead of the CNN adopted previously, to indicate the robustness of the proposed method in selecting salient spectral bands under different classification frameworks. Particularly, the experiments have been conducted using a FCFFNN and a SVM classifier which receives as inputs the selected salient bands of the TB and TDLA approach. We have used different classification models (e.g., CNN, FCFFNN, SVM) to indicate the robustness of the proposed method in selecting salient spectral bands under different classification frameworks. Three different experiments, for 10, 15, 20 most important spectral bands selection, are conducted. The experiments have been executed using 10 different runs and the standard deviation across all runs is depicted to indicate the robustness of the classification accuracy. As is observed, the best overall classification accuracy is achieved by TDLA and an SVM classifier. However, our approach gives better results when non-linear models (such as FCFFNN) are adopted.

TABLE II  
OVERALL CLASSIFICATION ACCURACY RESULTS (%) WHEN TDLA AND TENSOR-BASED LOGISTIC REGRESSION ARE USED TO REDUCE THE DIMENSION OF THE DATA.

DR + Classifier	TDLA + FCFNN	TDLA + SVM	TB + FCFNN (10 bands)	TB + SVM (10 bands)	TB + FCFNN (15 bands)	TB + SVM (15 bands)	TB + FCFNN (20 bands)	TB + SVM (20 bands)
OA (%)	71.22	89.54	64.53 ± 1.4	68.14 ± 0.6	73.28 ± 0.8	73.77 ± 0.7	74.62 ± 0.7	73.95 ± 0.6

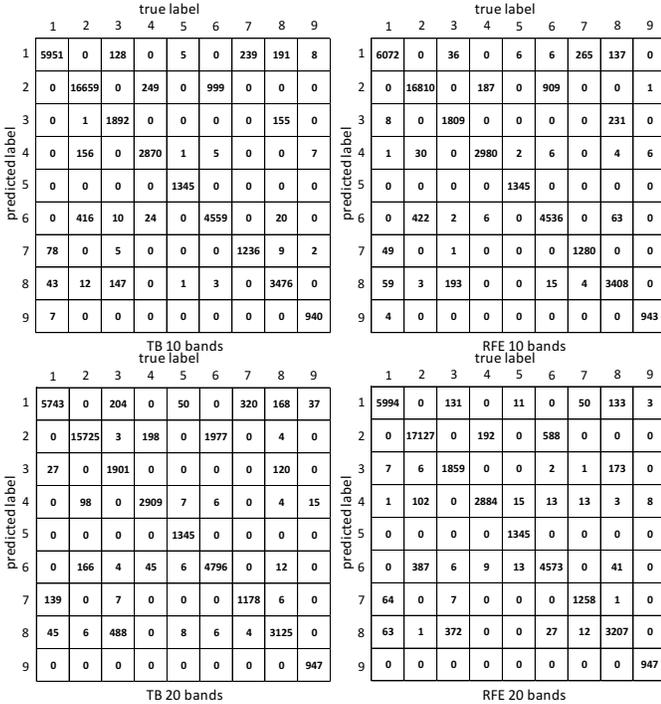


Fig. 5. Classification accuracy per class on Pavia University dataset when different feature selection techniques are used.

TABLE III  
CLASSIFICATION ACCURACY RESULTS (%) WHEN RFE AND TENSOR-BASED LOGISTIC REGRESSION ARE USED TO REDUCE THE DIMENSION OF THE DATA.

Feature selection	Number of bands	Accuracy (%) Pavia University	Accuracy (%) Indian Pines
Tensor-based	10	92.82	77.48
RFE	10	92.61	83.42
Tensor-based	20	93.02	77.46
RFE	20	93.25	83.56

The conclusions of the aforementioned experiments are that, although the proposed method has been designed for classification purposes (see the respective experiments of Section V-D), it also performs quite well for data dimensionality reduction, compared with state of the art methods that have been appropriately designed for this purpose.

D. High-Order Nonlinear Classification Model

For evaluating the performance of the high-order nonlinear classification model, that is, the rank-1 FNN, we also use the Pavia University and the Indian Pines datasets. A similar procedure as in Section V-B is followed to form the training

sets. We conduct different experiments using training sets of 50, 100, 150 and 200 random samples from each class respectively. A similar training map selection was made in [35] for the Indian Pines and Pavia University dataset. Please note that if for a specific class the number of samples are less than the number required for the experiment, then we select randomly 50% of the total available class samples. Fig. 6 depicts the effect of different training samples on classification accuracy for different window size  $s$  and assuming a constant number of neurons in the hidden layer. Particularly, we set  $Q = 100$  for the Pavia dataset and  $Q = 75$  for the Indian Pines. These values are derived from descriptions below (see also Fig. 7). Additionally, Fig. 7 presents the effect of the number of training samples on misclassification error using different number of neurons  $Q$  for both datasets.

Then, we evaluate the performance of proposed tensor-based high-order nonlinear classifier in relation to the size  $s$  of the spatial window around a pixel. We conduct experiments with a window size to be equal to 3, 5, 7, 9, and 11 as is depicted in Figure 6. These results suggest that the best overall classification accuracy is achieved when  $s = 5$ , i.e., the classification accuracy is mostly affected by the 24 closest neighbors of a pixel. When  $s < 5$  the spatial information is not adequate for achieving highly accurate classification results, while for  $s > 5$  the neighbourhood region is very probable to contain pixels that belong to different classes thus deteriorating the classification accuracy.

Moreover, we evaluate the performance of the high-order nonlinear classifier in relation to its complexity, i.e. the value of parameter  $Q$ , indicating the number of hidden neurons. Particularly, we set  $Q$  to be equal to 50, 75, 100 and 125 respectively. By increasing the value of  $Q$ , the complexity and the capacity of the model are also increased. The results of this evaluation are presented in Fig.7.

Regarding the Indian Pines dataset, we observe that the model with  $Q = 75$  outperforms all other models. When the training set size is very small, i.e. 50 samples per class, the model with  $Q = 50$  does not have the capacity to capture the statistical relation between the input and the output and thus underfits the data. On the other hand, the models with  $Q = 100$  and  $Q = 125$ , due to their high complexity, they slightly overfit the data. As training set increases, the misclassification error for all model decreases. This happens because increasing the training set size it also increases the amount of information that can be used during training to estimate the coefficients of the models.

As far as, the Pavia University dataset is concerned, we observe that the model with  $Q = 100$  outperforms all other models, when the size of dataset is larger than 50 samples per class. When the training dataset size is 50 samples per class

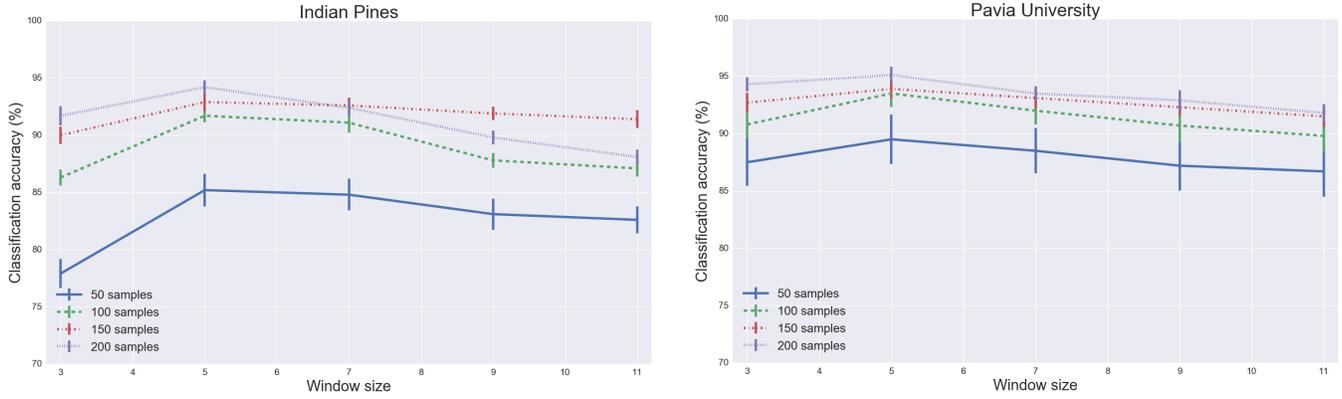


Fig. 6. Overall classification accuracy on test set versus the size of spatial window for the high-order nonlinear classification model.

model with  $Q = 75$  outperforms all other models. The model with  $Q = 125$  overfits the data, while the model with  $Q = 50$  underfits them. When the size of the training dataset increases, the classification accuracy of all models also increases.

In the following, we compare the performance of the proposed *rank-1* FNN against FCFNN, Radial Basis Function SVM (RBF-SVM), and two deep learning approaches that have been proposed for classifying hyperspectral data; the first one is based on Stacked-Autoencoders (SAE) [11], while the second one on the exploitation of Convolutional Neural Networks (CNN) [12]. The FCFNN receives as input the vectorized version of the Rank-1 FNN inputs. The same applies for the input of RBF-SVM. Moreover, FCFNN has one hidden layer that has the same number of hidden neuron as the Rank-1 FNN. SAE and CNN were developed as proposed in the original papers. The number of hidden neurons of FCFNN is 75 for the Indian Pines dataset and 100 the Pavia University dataset (as derived from Fig.7). The architecture of the network that exploits Stacked-Autoencoders consists of three hidden layers, while each hidden layer contains 10% less neurons than its input. We choose to gradually reduce the number of hidden neurons from one hidden layer to the next, in order not to permit for the network to learn the identity function during pre-training. Regarding CNN, we utilize exactly the same architecture as the one presented in [12]. The performance of all these models is evaluated on varying size training sets; training sets that contain 50, 100, 150 and 200 samples from each one of the available classes.

Table IV presents the outcome of this comparison. In this table, we also depict the standard deviation across 10 different experiment runs to better reveal the classification accuracy. When the training set size is small, our approach outperforms all other models. This stems from the fact that the proposed high-order nonlinear model exploits tensor algebra operations to reduce the number of coefficients that need to be estimated during training, while at the same time it is able to retain the spatial information of the input. Although, the FCFNN utilizes the same number of hidden neurons as our proposed model, it seems to overfit training sets when a small size of data is used, due to the fact that it employs a larger number of coefficients. RBF-SVM performs better than the FCFNN on the Pavia University dataset, but slightly worse on the Indian

TABLE IV  
OVERALL CLASSIFICATION ACCURACY RESULTS (%) OF HIGH-ORDER NONLINEAR CLASSIFICATION MODEL FOR BOTH DATASETS.

Pavia University				
Samples per class	50	100	150	200
<i>rank-1</i> FNN ( $Q=100$ )	<b>89.95</b> ± 0.7	<b>93.50</b> ± 0.4	93.89 ± 0.3	95.11 ± 0.3
FCFNN	67.79 ± 3.6	76.53 ± 2.7	78.48 ± 2.2	82.59 ± 2.1
RBF-SVM	86.98 ± 1.6	88.99 ± 1.3	89.86 ± 1.2	91.82 ± 1.2
SAE	86.54 ± 2.1	91.90 ± 1.8	92.38 ± 1.4	93.29 ± 1.1
CNN	88.89 ± 0.6	92.74 ± 0.3	<b>94.68</b> ± 0.2	<b>95.89</b> ± 0.2

Indian Pines				
Samples per class	50	100	150	200
<i>rank-1</i> FNN ( $Q=75$ )	<b>85.20</b> ± 1.2	<b>91.63</b> ± 0.8	<b>92.82</b> ± 0.5	94.15 ± 0.4
FCFNN	73.88 ± 4.2	81.10 ± 3.9	84.14 ± 2.7	85.86 ± 2.5
RBF-SVM	73.18 ± 2.3	77.86 ± 2.0	82.11 ± 1.5	84.99 ± 1.5
SAE	65.51 ± 4.4	70.66 ± 3.6	74.03 ± 3.6	76.49 ± 3.2
CNN	82.43 ± 0.9	85.48 ± 0.7	92.28 ± 0.2	<b>94.81</b> ± 0.2

TABLE V  
OVERALL CLASSIFICATION ACCURACY RESULTS (%) OF STM AND RANK-1 FNN FOR BOTH DATASETS.

Indian Pines			
Method	STM	STM-MPCA	Rank-1 FNN
OA	62.6 ± 2.6	80.6 ± 1.9	73.9 ± 1.4

Pavia University			
Method	STM	STM-MPCA	Rank-1 FNN
OA	79.5 ± 1.4	89.4 ± 0.5	88.2 ± 0.8

Pines dataset. The deep learning architecture based on the exploitation of SAE, is actually a FCFNN with three hidden layers, which employs an unsupervised pre-training phase to initialize its coefficients. The fully connectivity property of this architecture implies very high complexity, which is responsible

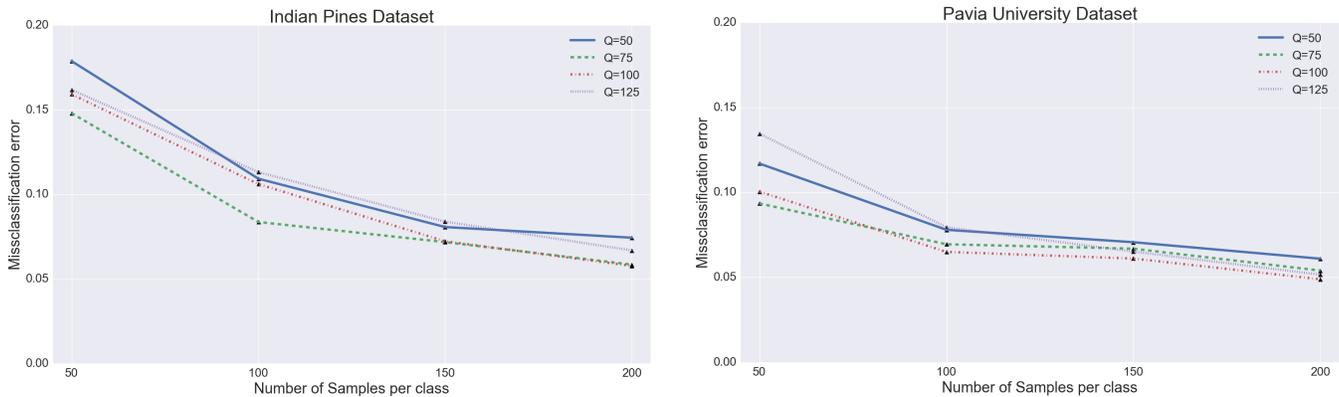


Fig. 7. Misclassification error on test set versus the complexity, determined by  $Q$ , of the high-order nonlinear classification model.

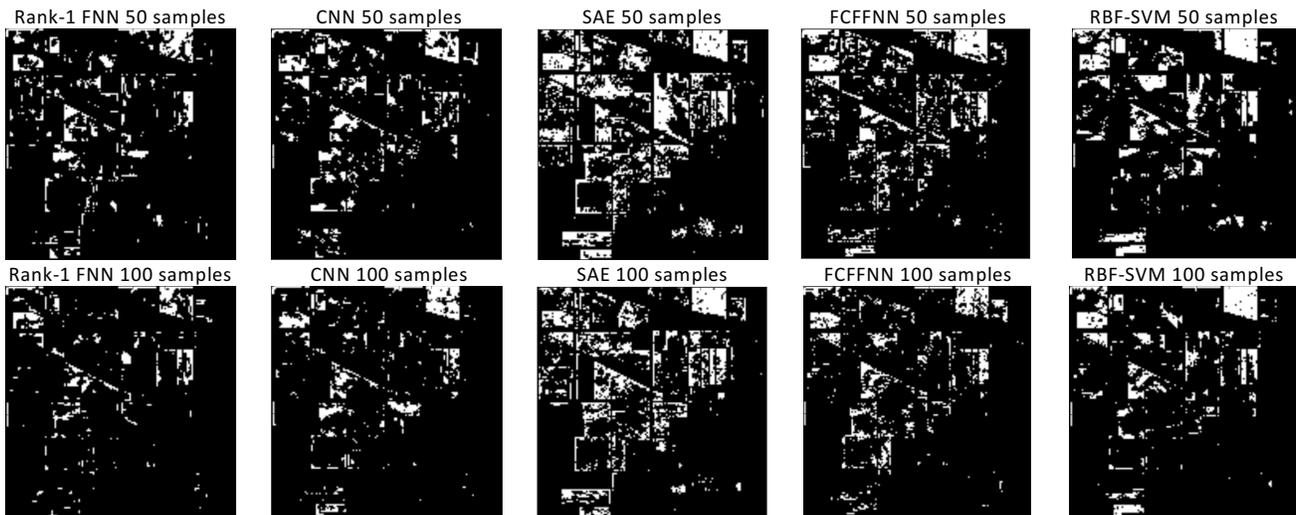


Fig. 8. Visualization of classification accuracy for all tested approaches on Indian Pines dataset. White pixels were misclassified.

for the poor performance, due to overfitting issues, on the Indian Pines dataset. Finally, the deep learning architecture based on CNN performs better than FCFNN, RBF-SVM and SAE mainly due to its sparse connectivity (low complexity) and the fact that it can inherently exploit the spatial information of the input. When the training set consists of 150 and 200 samples per class, for the Pavia University dataset, and 200 samples for the Indian Pines dataset, the deep learning architecture based on CNN seems to outperform even the proposed high-order nonlinear classification model. This happens because the CNN-based deep learning model has higher capacity than the proposed model, which implies that it is capable of capturing better the statistical relation between the input and the output, when the training set contains sufficient information. However, when the size of the training set is small, our proposed model, due to its lower complexity, outperforms the CNN-based deep learning model. Fig. 8 and Fig. 10 depict the classification accuracy for all tested models when the training set contains 50 and 100 samples per class (white labeled pixels correspond to misclassified samples), while Fig. 9 presents the convergence curves for the Rank-1 FNN. In these figures, the false images (misclassified pixels) are depicted.

Finally, we compare our proposed method against the Support Tensor Machine (STM) proposed in [36] using In-

Indian Pines and Pavia University datasets. The work of [36] evaluates the performance of STM using both raw data and data of a reduced dimension derived through a Multilinear Principle Component Analysis (MPCA) [37]. The results of this comparison are presented in Table V. In this table, we also depict the standard deviation of the results. As is observed, the proposed Rank-1 FNN outperforms the STM model when raw hyperspectral data are used. Instead, MPCA-STM works slightly better than our approach. This is mainly due to the fact that better features are used as inputs.

## VI. CONCLUSIONS

In this work, we present a linear and a non-linear tensor-based scheme for hyperspectral data classification and analysis. The advantages of the presented model is that i) it reduces the number of parameters required for the classification and thus it reduces the respective number of training samples, ii) it provides a physical interpretation regarding the model coefficients on the classification output and iii) it retains the spatial and spectral coherency of the input samples. By utilizing tensor algebra operations, we introduce learning algorithms to train the tensor-based classifiers (linear and non-linear). The linear classifier is based on a modification of a logistic regression model, while the non-linear one on

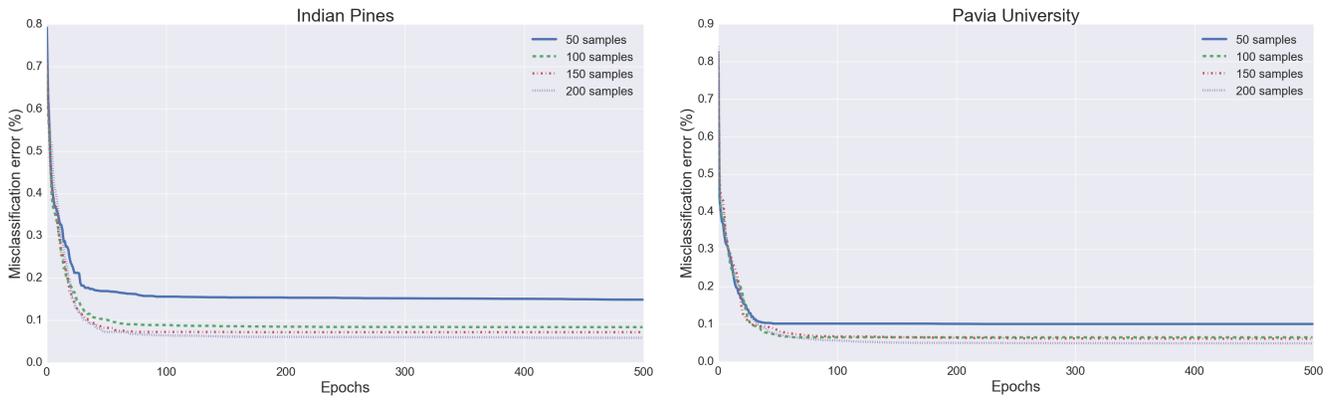


Fig. 9. Convergence curves for the Rank-1 FNN for both datasets and for different sizes of training set.

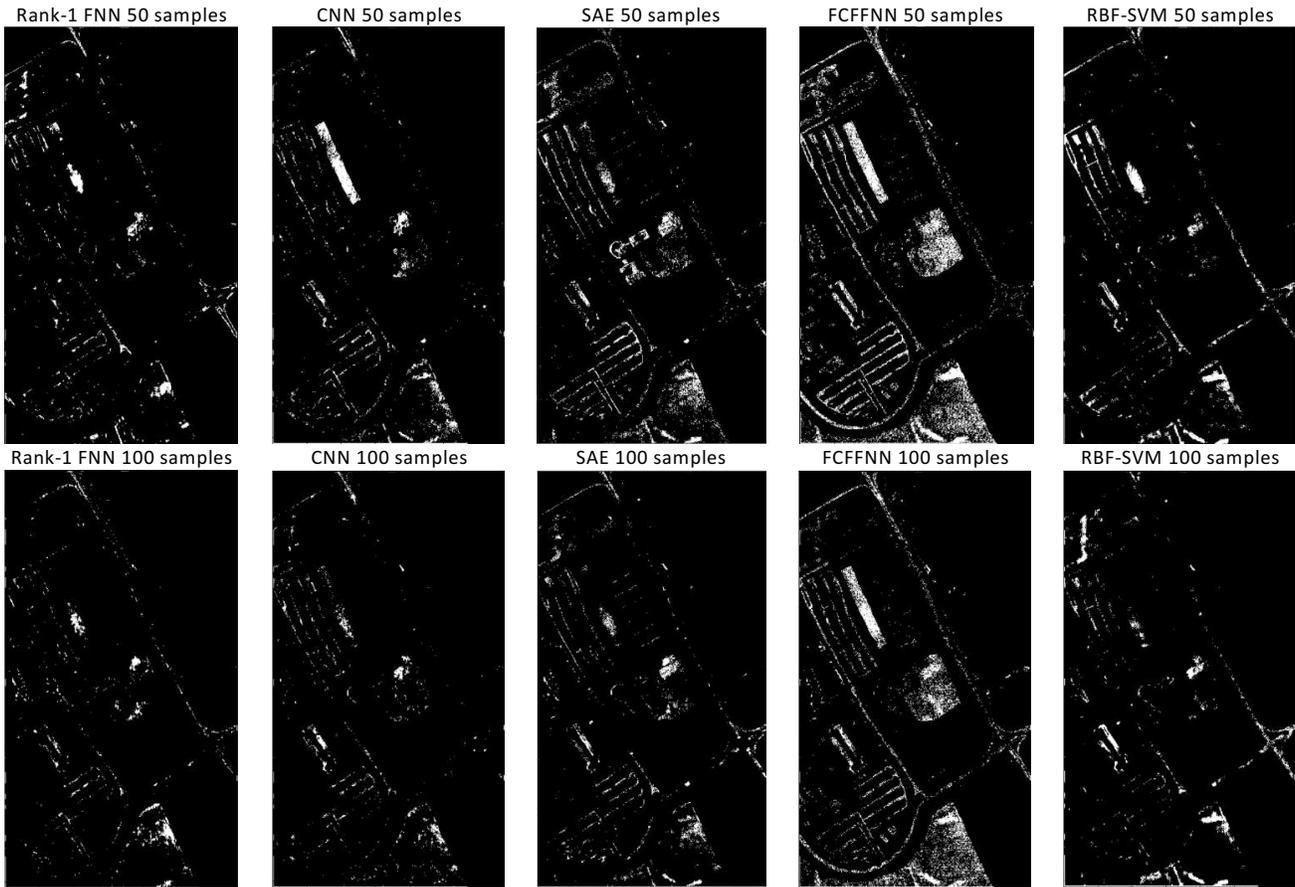


Fig. 10. Visualization of classification accuracy for all tested approaches on Pavia University dataset. White pixels were misclassified.

a modification of a feedforward neural network. Both the proposed models assume a *rank-1* canonical decomposition of the weight parameters. For this reason, we properly modified existing learning schemes to train these models so as to keep the *rank-1* canonical decomposition property. We call the new proposed non-linear classifier as *rank-1* FNN.

We have evaluated the performance of the presented model in terms of classification accuracy and for dimensionality reduction. As far as the dimensionality reduction is concerned, the tensor-based scheme allows selection of the most discriminative spectral bands. It produces an interpretable dimensionality reduction, which can be used with more complicated

classifiers without deteriorating their performance. In terms of classification, the linear classifier outperforms conventional linear models such as logistic regression and SVM.

However, the linear classifier is characterized by low capacity, since it produces classification rules that are linear in the input space. For this reason, in this paper, we have introduced nonlinear classification models the weights of which satisfy the *rank-1* canonical decomposition property. We have also introduced suitable learning algorithms to train the new proposed non-linear model. The performance of the nonlinear model was compared against other nonlinear classifiers, including the state-of-the-art deep learning classifiers. The main results

are that when the size of the training set is small, our newly proposed model presents superior performance against the compared methods.

As future work, one can evaluate the performance of the proposed tensor-based classifier on fused datasets that contain hyperspectral data and LiDAR data [38]. In this case, we need to investigate on the extraction of new features that better model how LiDAR data are related with hyperspectral particularities [39].

## APPENDIX A

### A. Proof of relation (15)

*Proof:* For simplicity, we omit the superscripts of  $\mathbf{w}$ . Let us denote as  $\mathbf{W}$  the tensor

$$\mathbf{W} = \mathbf{w}_1 \circ \cdots \circ \mathbf{w}_D \in \mathbb{R}^{p_1, \dots, p_D}$$

having the same dimensions with the tensor  $\mathbf{X}$ . Then,

$$\text{vec}(\mathbf{W}) = \mathbf{w}_1 \odot \cdots \odot \mathbf{w}_D.$$

We also have that

$$\begin{aligned} \langle \mathbf{w}_1 \odot \cdots \odot \mathbf{w}_D, \mathbf{X} \rangle &= \\ \sum_{i_1}^{p_1} \cdots \sum_{i_D}^{p_D} W_{i_1, \dots, i_D} X_{i_1, \dots, i_D} &= \\ \langle \mathbf{W}_{(d)}, \mathbf{X}_{(d)} \rangle \end{aligned}$$

for any  $d = 1, 2, \dots, D$ . Furthermore, from relation (3) of the main manuscript, we have

$$\begin{aligned} \langle \mathbf{W}_{(d)}, \mathbf{X}_{(d)} \rangle &= \\ \langle \mathbf{w}_d (\mathbf{w}_D \odot \cdots \odot \mathbf{w}_{d+1} \odot \mathbf{w}_{d-1} \odot \cdots \odot \mathbf{w}_1)^T, \mathbf{X}_{(d)} \rangle. \end{aligned} \quad (\text{A.1})$$

Now by using (A.1) and the property of inner and dot products, which states that for any three matrices  $A$ ,  $B$  and  $C$ , the following

$$\langle \mathbf{A}\mathbf{B}^T, \mathbf{C} \rangle = \mathbf{A}\mathbf{B}^T\mathbf{C}^T = \langle \mathbf{A}, \mathbf{C}\mathbf{B} \rangle$$

holds, we conclude that

$$\begin{aligned} \langle \mathbf{w}_1 \odot \cdots \odot \mathbf{w}_D, \mathbf{X} \rangle &= \\ \langle \mathbf{w}_d, \mathbf{X}_{(d)} (\mathbf{w}_D \odot \cdots \odot \mathbf{w}_{d+1} \odot \mathbf{w}_{d-1} \odot \cdots \odot \mathbf{w}_1) \rangle. \end{aligned}$$

The proof is completed.  $\triangleleft$

## REFERENCES

- [1] E. Wycoff, T.-H. Chan, K. Jia, W.-K. Ma, and Y. Ma, "A non-negative sparse promoting algorithm for high resolution hyperspectral imaging," in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. IEEE, 2013, pp. 1409–1413.
- [2] C.-I. Chang, *Hyperspectral data processing: algorithm design and analysis*. John Wiley & Sons, 2013.
- [3] G. Camps-Valls and L. Bruzzone, "Kernel-based methods for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 43, no. 6, pp. 1351–1362, 2005.
- [4] —, *Kernel methods for remote sensing data analysis*. John Wiley & Sons, 2009.
- [5] J. M. Benítez, J. L. Castro, and I. Requena, "Are artificial neural networks black boxes?" *IEEE Transactions on neural networks*, vol. 8, no. 5, pp. 1156–1164, 1997.
- [6] G. Camps-Valls, D. Tuia, L. Bruzzone, and J. Atli Benediktsson, "Advances in hyperspectral image classification: Earth monitoring with statistical learning methods," *Signal Processing Magazine, IEEE*, vol. 31, no. 1, pp. 45–54, 2014.
- [7] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [8] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [9] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural computation*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [10] Y. Bengio, P. Lamblin, D. Popovici, H. Larochelle *et al.*, "Greedy layer-wise training of deep networks," *Advances in neural information processing systems*, vol. 19, p. 153, 2007.
- [11] Y. Chen, Z. Lin, X. Zhao, G. Wang, and Y. Gu, "Deep learning-based classification of hyperspectral data," *IEEE Journal of Selected topics in applied earth observations and remote sensing*, vol. 7, no. 6, pp. 2094–2107, 2014.
- [12] K. Makantasis, K. Karantzalos, A. Doulamis, and N. Doulamis, "Deep Supervised Learning for Hyperspectral Data Classification through Convolutional Neural Networks," in *IEEE International Geoscience and Remote Sensing Symposium (IGARSS 2015)*, July 2015.
- [13] M. Vakalopoulou, K. Karantzalos, N. Komodakis, and N. Paragios, "Building detection in very high resolution multispectral data with deep learning features," in *Geoscience and Remote Sensing Symposium (IGARSS), 2015 IEEE International*. IEEE, 2015, pp. 1873–1876.
- [14] M. Papadomanolaki, M. Vakalopoulou, S. Zagoruyko, and K. Karantzalos, "Benchmarking deep learning frameworks for the classification of very high resolution satellite multispectral data," *ISPRS Annals of Photogrammetry, Remote Sensing & Spatial Information Sciences*, vol. 3, no. 7, 2016.
- [15] V. Mnih and G. E. Hinton, "Learning to label aerial images from noisy data," in *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, 2012, pp. 567–574.
- [16] K. Makantasis, K. Karantzalos, A. Doulamis, and K. Loupos, "Deep learning-based man-made object detection from hyperspectral data," in *International Symposium on Visual Computing*. Springer, 2015, pp. 717–727.
- [17] G. B. Orr and K.-R. Müller, *Neural networks: tricks of the trade*. Springer, 2003.
- [18] Z. Kandylakis, K. Karantzalos, A. Doulamis, and N. Doulamis, "Multiple object tracking with background estimation in hyperspectral video sequences," in *IEEE Workshop on Hyperspectral Image and Signal Processing: Evolution in Remote Sensing (WHISPERS)*, 2015.
- [19] N. Kussul, A. Shelestov, M. Lavreniuk, I. Butko, and S. Skakun, "Deep learning approach for large scale land cover mapping based on remote sensing data fusion," in *Geoscience and Remote Sensing Symposium (IGARSS), 2016 IEEE International*. IEEE, 2016, pp. 198–201.
- [20] K. Makantasis, E. Protopapadakis, A. Doulamis, N. Doulamis, and C. Loupos, "Deep convolutional neural networks for efficient vision based tunnel inspection," in *Intelligent Computer Communication and Processing (ICCP), 2015 IEEE International Conference on*. IEEE, 2015, pp. 335–342.
- [21] H. Zhou, L. Li, and H. Zhu, "Tensor regression with applications in neuroimaging data analysis," *Journal of the American Statistical Association*, vol. 108, no. 502, pp. 540–552, 2013.
- [22] X. Tan, Y. Zhang, S. Tang, J. Shao, F. Wu, and Y. Zhuang, "Logistic tensor regression for classification," in *International Conference on Intelligent Science and Intelligent Data Engineering*. Springer, 2012, pp. 573–581.
- [23] L. De Lathauwer, B. De Moor, and J. Vandewalle, "Computation of the canonical decomposition by means of a simultaneous generalized schur decomposition," *SIAM journal on Matrix Analysis and Applications*, vol. 26, no. 2, pp. 295–327, 2004.
- [24] B. C. Csáji, "Approximation with artificial neural networks," *Faculty of Sciences, Eötvös Loránd University, Hungary*, vol. 24, p. 48, 2001.
- [25] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.
- [26] C. M. Bishop, *Pattern recognition and machine learning*. springer, 2006.
- [27] H. Hung and C.-C. Wang, "Matrix variate logistic regression model with application to eeg data," *Biostatistics*, vol. 14, no. 1, pp. 189–202, 2013.
- [28] R. A. Harshman, "Foundations of the parafac procedure: Models and conditions for an "explanatory" multi-modal factor analysis," 1970.
- [29] J. D. Carroll and J.-J. Chang, "Analysis of individual differences in multidimensional scaling via an n-way generalization of eckart-young decomposition," *Psychometrika*, vol. 35, no. 3, pp. 283–319, 1970.

- [30] A. D. Doulamis, N. D. Doulamis, and S. D. Kollias, "On-line retrainable neural networks: improving the performance of neural networks in image analysis problems," *IEEE Transactions on Neural Networks*, vol. 11, no. 1, pp. 137–155, 2000.
- [31] A. Doulamis, N. Doulamis, K. Ntalianis, and S. Kollias, "An efficient fully unsupervised video object segmentation scheme using an adaptive neural-network classifier architecture," *IEEE Transactions on Neural Networks*, vol. 14, no. 3, pp. 616–630, 2003.
- [32] A. D. Doulamis, N. D. Doulamis, and S. D. Kollias, "An adaptable neural-network model for recursive nonlinear traffic prediction and modeling of mpeg video sources," *IEEE Transactions on Neural Networks*, vol. 14, no. 1, pp. 150–166, 2003.
- [33] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene selection for cancer classification using support vector machines," *Machine learning*, vol. 46, no. 1-3, pp. 389–422, 2002.
- [34] L. Zhang, L. Zhang, D. Tao, and X. Huang, "Tensor discriminative locality alignment for hyperspectral image spectral–spatial feature extraction," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 51, no. 1, pp. 242–256, 2013.
- [35] J. Li, X. Huang, P. Gamba, J. Bioucas-Dias, L. Zhang, J. Benediktsson, and A. Plaza, "Multiple feature learning for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 3, pp. 1592–1606, 2015.
- [36] X. Guo, X. Huang, L. Zhang, L. Zhang, A. Plaza, and J. A. Benediktsson, "Support tensor machines for classification of hyperspectral remote sensing imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 54, no. 6, pp. 3248–3264, 2016.
- [37] H. Lu, K. N. Plataniotis, and A. N. Venetsanopoulos, "MPCA: Multilinear principal component analysis of tensor objects," *IEEE transactions on Neural Networks*, vol. 19, no. 1, pp. 18–39, 2008.
- [38] F. Pacifici, Q. Du, and S. Prasad, "Report on the 2013 IEEE GRSS data fusion contest: Fusion of hyperspectral and lidar data [technical committees]," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 1, no. 3, pp. 136–38, 2013.
- [39] E. Maltezos, N. Doulamis, A. Doulamis, and C. Ioannidis, "Deep convolutional neural networks for building extraction from orthoimages and dense image matching point clouds," *Journal of Applied Remote Sensing*, vol. 11, no. 4, 2017.