

Article

Tensor-based Factorization Algorithms for Pixel-wise Classification of Hyperspectral Data Using Deep Convolutional Networks

Josué López ^{1,*}, Deni Torres ¹ and Clement Atzberger ²

¹ Center for Research and Advanced Studies of the National Polytechnic Institute, Telecommunications Group, Av del Bosque 1145, Zapopan 45017, Mexico; dtorres@gdl.cinvestav.mx

³ University of Natural Resources and Life Science, Institute of Geomatics, Peter Jordan 82, Vienna 1180, Austria; clement.atzberger@boku.ac.at

* Correspondence: josue.lopez@cinvestav.mx

Version November 9, 2020 submitted to Remote Sens.

Abstract: Tensor-based algorithms for data compression have evolved in recent years according to the needs of several research areas. Tucker Decomposition (TKD) is one of the most popular factorization methods based on tensor algebra, but it is clear that it is not the only algorithm which can produce a factorization for a given input data set. Besides, this decomposition does not have singular solutions, i.e., it converges to local minima. Hence, depending on the input data, tensor-based decompositions can achieve better solution to a specific input. The phenomenology of Remote Sensing (RS) Hyperspectral Images (HSI) belongs to the set of natural numbers, i.e., the set of positive integers. Then, a non-negative tensor factorization suggest a more suitable decomposition for positive data by nature. The main purpose of this work is to prove the benefits in processing time, as well as in accuracy, of using a well-posed factorization algorithm. Specifically, this paper performs a quantitative analysis of tensor-based factorization algorithms applied to semantic segmentation of HSI using Deep Convolutional Networks (DCN).

Keywords: deep convolutional networks; hyperspectral imagery; tensor decomposition

1. Introduction

Reducing the dimensionality of input data for machine learning algorithms has been one of the most active research areas in ~~recente~~ recent years [1]. The insertion of tensor-based algorithms for ~~this sort of tasks drove a revolution~~ these types of tasks inspired a change in several areas, such as image processing [2].

Most of the researches in the image processing area require data acquired by multiple sensors. Even in the simplest case, color images, data acquired by three sensors that perceive reflectance of an object are processed. Each sensor receives reflectance in different wavelength ranges, and by merging that data, a color image is produced. Thus, the images can be represented in a form of three-dimensional arrays or third-order tensors. Two dimensions represent the spatial properties and the third dimension denotes the depth, i.e., the spectral bands. Medical analysis [3], minorology [4], agriculture [5], radar images [6] and, above all, remote sensing multi- and hyper-spectral images [7] can be represented, by nature, as third-order tensors.

Spectral ~~imagery remarkably aids images make~~ certain image processing tasks much easier [8]. Recently, the use of this type of data has grown exponentially in various areas such as agriculture [9], medical analysis [10], biomedical [11], natural disaster prediction [12], security affairs [13], among others. The ability to obtain information about a target not only by its reflectance in the spatial domain, but also

31 by response at different wavelengths, has driven a growth in accuracy and precision in tasks such as
32 classification and segmentation [].

33 Few years ago, several unsupervised classification and segmentation algorithms [] were
34 developed, taking advantage of the properties that spectral data produce. Subsequently, with the
35 introduction of supervised machine learning algorithms such as SVM [], kNN [] and ANN [], it was
36 found that, under certain conditions, there is a direct relationship between the number of bands used
37 and the performance of these algorithms []. However, with the aim of improving results, neural
38 network models evolved into deep neural networks []. This caused the computational complexity to
39 rise considerably and spectral image processing was not easily achievable. The foregoing requires
40 having robust computer equipment to achieve competitive results in time.

41 Several works opted for matrix factorization algorithm to reduce the high-dimensionality of
42 spectral images []. More recently, with the development of tensor factorization algorithms [], it has
43 been found that some algorithms based on tensor algebra produce advantages over those based on
44 matrices []. Nevertheless, the data produced by both of them are hard to understand for supervised
45 classification algorithms that need spatial relation between pixels to produce a wise prediction [].

46 In this work, ~~we propose the proposal is to find~~ an alternative solution to the problem
47 described previously. ~~To reduce processing times in supervised classification algorithms, such as deep~~
~~neural networks, we~~ We propose a model that ~~, from the benefits offered by tensor decomposition~~
~~algorithms, aids compression of spectral images~~ reduces the computational load of hyperspectral
50 imagery classification supervised algorithms through tensor decomposition models. This produces a
51 lower dimensional tensor while preserving the structural and numerical nature of the original data.

52 1.1. *State of the art* Previous work

53 There are several works focused on the development of frameworks that reduce computational
54 complexity of machine learning algorithms for semantic segmentation of hyperspectral datasets []. The
55 crucial factor, which is addressed in this work, is to achieve compression of the input data to reduce
56 the high number of computations, but without sacrificing pixel accuracy, overall accuracy, precision and
57 recall in the classification task.

58 Before the introduction of tensor decomposition algorithms, the way to use hyperspectral images
59 as input for supervised classification algorithms was by band selection [23] and [22]. Later, matrix
60 decomposition algorithms were used, such as PCA in [?], and even non-negative matrix decomposition
61 methods [?]. In 2015 Zhang et al. [24] were pioneers in experimenting with multilinear algebra-based
62 decompositions on hyperspectral images.

63 On the other hand, there was also the possibility of using multispectral images due to the small
64 number of spectral bands, which still made efficient results in classification without dimensionality
65 reduction achievable, as done in [11], [18], [21] and [?]. However, the need to increase classification
66 performance forces researchers to use data with more features that favor and aid the classification of
67 various classes, which are difficult to differentiate with little spectral data. Thus, more recent researches
68 have decided to use hyperspectral images with tensor decompositions, which has increased the results
69 in classification accuracy [26], [27], [29], [?] and [?].

70 Recently, Sayeh et al. [?] published a work close to our research. They proposed a non-negative
71 tensor decomposition of hyperspectral images but, different to our research, they try to preserve certain
72 spatial-spectral features into the so called abundance maps, i.e. the projection matrices, while this work
73 pursues to preserve the nature of the image just compressing the main information in the positive core
74 tensor.

75 Table 1 summarizes some of the most cited related papers, which deal with the
76 compression-classification issue.

Table 1. Related work in spectral imagery semantic segmentation.

Reference	Input	Decomposition	Reduction	Classifier
Li, S. et al. [23] (2014)	HSI	-	Band selection	SVM
Zhang, L. et al. [24] (2015)	HSI	TKD	Spatial-Spectral	-
Wan, Q. et al. [22] (2016)	HSI	-	Band selection	SVM/kNN/CART
Kemker, R. et al. [11] (2017)	MSI	-	-	CNN
Tong L. et al. [] (2017)	HSI	NMF	Unmixing	-
Hamida, A. et al. [21] (2017)	MSI	-	-	CNN
Chien, J. et al. [] (2017)	RGB	TFNN	Spatial-Spectral	TFNN
Dewa, M. et al. [] (2018)	HSI	PCA	Spectral	PCA
Xu, Y. et al. [] (2018)	HSI	-	-	CNN
Li, J. et al. [28] (2019)	MSI	NTD-CNN	Spatial-spectral	-
An, J. et al. [27] (2019)	HSI	T-MLRD	Spatial-spectral	SVM/1NN
An, J. et al. [29] (2019)	HSI	TDA	Spatial-spectral	SVM/1NN
Lopez, J. et al. [] (2019)	MSI	TKD	Spectral	FCN
Sayeh, M. et al. [] (2019)	HSI	NTD	Spatial-Spectral	3D-CNN
Our framework	MSI/HSI	NTKD<i>i</i>NTD/NTD	Spectral	CNN

77 1.2. *ContributionMotivation*

78 Nowadays, RS image processing is applied in several areas related to caring of the planet.
 79 Nevertheless, task such as classification becomes more complex due to low spatial resolutions, this is
 80 offset by the use of devices with other features such as spectral sensors.

81 In the state of the art we can find a variety of frameworks looking for dimensionality reduction
 82 of images of any type, that is , RGB [?], medical [?], SAR [?], multispectral [?], among others.
 83 In particular, the large number of spectral bands in hyperspectral images has focused efforts in
 84 this category [?]. Spatial [?], spectral [?] and spatial-spectral [?] compression have been achieved
 85 with On the other hand, CNNs have been widely used in recent years in the area of RGB image
 86 semantic classification and segmentation. Its performance is highly competitive and the development
 87 of various improvement strategies have considerably reduced its computational cost []. However,
 88 the computational complexity of the algorithm means that the increase in the dimensionality of the
 89 input data produces a significant increase in the computational load []. This is why the processing of
 90 high-dimensional images such as multi and hyperspectral images becomes unfeasible.

91 Some data compression strategies have favorably reduced the dimensionality of the data.
 92 Decomposition methods based on the matrix and tensor decompositions. It is important to highlight
 93 that, a decomposition as approaches have been applied as pre-processing for a classification algorithm,
 94 instead of favoring, could be counterproductive and greatly affect its performance. It is therefore
 95 important to make a proper selection of the type of decomposition that would produce a suitable data
 96 set for post-processing of input data of neural networks. In tensor decompositions, the processing of
 97 the data in its natural format, i.e., as N-order tensors, improves the decomposition process because
 98 it considers the dependence of the data in its different modes. Although a decomposition can
 99 compress the data, it is also important to note that the inappropriate selection of some decomposition
 100 parameters could lead to information losses, which would penalize the performance of a CNN.

101 Under these considerations, this work is motivated to develop a low computational complexity
 102 and competitive in performance framework that helps various fields of application of remote sensing
 103 image processing to solve classification tasks.

104 1.3. *Contribution*

105 Unlike previous works, this work seeks to adapt the data in the best way to be a more efficient
 106 way to the input of deep convolutional networks. Convolutional network models are designed to
 107 extract and interpret all the spatial properties of an image by moving the kernels over the input data [].
 108 Therefore, producing uncorrelated data in space and spectrum, would make harder the interpretation

of the data in the convolutional network [?]. Thus, the proposed framework maintains the **positive integer and nonnegativity** tensor nature of the spectral images and the spatial dimensionality to preserver spatial-spectral correlation of the data while reducing the spectral dimensionality, in order to decrease computational load in the pixel-wise classification process. **Besides, we approach the problem of setting the compression range selection hyperparameters by the measurement of divergences, which aids to estimate a best rank approximation.**

We can summarize the contribution of this work with the following two points:

1. The framework **NTKD3-CNN-INTD1-CNN** proposed in this work, develops a new strategy to improve **accuracy results performance** of semantic segmentation convolutional neural networks by finding suitable tensor data, preserving spatial correlation and values in the set of the natural numbers while compressing the spectral domain and in turn decreasing computational load.
2. **Furthermore, this work proposes a strategy for defining the range in mode 3 of the compression models based on the Tucker decomposition from the information theory point of view.**
3. This work also presents an exhaustive performance analysis measuring and comparing its efficiency with the most popular metrics, i.e., as pixel accuracy (PA), also PA in function of the number of new tensor bands, precision, recall, F1, orthogonality degree of the factor matrices and the core tensor, reconstruction error of the original tensor, and execution time.

The remainder of this work is organized as follows. Section ?? introduces tensor algebra notation and basic concepts to familiarize the reader with the symbology used in this paper. Section ?? presents the problem statement of this work and the mathematical definition. In Section ??, CNN theory is described for classification and semantic segmentation. Section ?? presents the framework proposed for compression and semantic segmentation of spectral images. Experimental results are presented in Section ???. Finally, Sections ?? and 8 present a discussion and conclusions based on the results obtained in the experiments.

2. Tensor-Based Factorizations

Matrix-based factorizations, such as PCA [] and SVD [] have been significant and useful tools for dimensionality reduction and other approaches. Nevertheless, they are limited **by representations of data in two modes to data representations in 2-dimensional spaces**. Most of current applications have data structures often as higher-order arrays, e.g. dimensions of space, time, and frequency. This 2-way view in matrix factorizations may be inadequate and it is natural to use tensor decomposition approaches [?].

We can define a tensor \mathbf{A} **tensor can be defined** as a multi-way or multidimensional array. The order of a tensor is the number of dimensions, also known as modes, i.e., an N -order tensor $\mathbf{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is an N -dimesional array, which elements $y_{i_1, i_2, \dots, i_n} x_{i_1, i_2, \dots, i_n}$ are indexed by $i_n \in 1, 2, \dots, I_n$ for $1 \leq n \leq N$.

Throughout this paper, the mathematical notation used by Kolda et al. [17] has been adopted. Table 2 summarize this notation. **Tensor algebra notation summary** $\mathbf{A}, \mathbf{A}, \mathbf{a}, \mathbf{a}$ **Tensor, matrix, vector and scalar respectively** $\mathbf{A} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ **N -order tensor of size** $I_1 \times \dots \times I_N$. $a_{i_1 \dots i_N}$ **An element of a tensor** $\mathbf{a}_{i_2 i_3}, \mathbf{a}_{i_1 i_3}, \mathbf{a}_{i_1 i_2}$: **Column, row and tube fibers of a third order tensor** $\mathbf{A}_{:, :, :}, \mathbf{A}_{:, :, :}, \mathbf{A}_{:, :, :}$ **Horizontal, lateral and frontal slices for a third order tensor** $\mathbf{A}^{(n)}, \mathbf{a}^{(n)}$ **A matrix/vector element from a sequence of matrices/vectors** $\mathbf{A}^{(n)}$ **Mode- n matricization of a tensor.** $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times \prod_{m \neq n} I_m}$ $\mathbf{X} = \mathbf{a}^{(1)} \circ \dots \circ \mathbf{a}^{(N)}$ **Outer product of N vectors, where** $x_{i_1 i_2 \dots i_N} = a_{i_1}^{(1)} \dots a_{i_N}^{(N)}$ $\langle \mathbf{A}, \mathbf{B} \rangle$ **Inner product of two tensors.** $\mathbf{B} = \mathbf{A} \times_n \mathbf{U}$ **n -mode product of tensor** $\mathbf{A} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ **by a matrix** $\mathbf{U} \in \mathbb{R}^{J \times I_n}$ **along axis n .**

2.1. Basic concepts

It is also necessary to introduce some tensor algebra operations and basic concepts used in later explanations. **These notations were taken textually from [17].**

¹⁵⁶ 2.0.1. **Matricization**

¹⁵⁷ 2.1. Matricization

¹⁵⁸ The mode- n matricization is the process of reordering the elements of a tensor into a matrix along
¹⁵⁹ axis n and it is denoted as $\mathbf{A}_{(n)} \in \mathbb{R}^{I_n \times \prod_{m \neq n} I_m}$.

¹⁶⁰ 2.1.1. **Outer Product**

¹⁶¹ The outer product of N vectors $\mathbf{x} = \mathbf{a}^{(1)} \circ \cdots \circ \mathbf{a}^{(N)}$ produces a tensor $\mathbf{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ where \circ
¹⁶² denotes the outer product and $\mathbf{a}^{(n)}$ denotes a vector in a sequence of N vectors and each element of
¹⁶³ the tensor is the product of the corresponding vector elements; i.e., $x_{i_1 i_2 \dots i_N} = a_{i_1}^{(1)} \dots a_{i_N}^{(N)}$.

¹⁶⁴ 2.1.1. **Inner Product**

¹⁶⁵ The

¹⁶⁶ 2.2. Inner Product

¹⁶⁷ The inner product of two tensors $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ is the sum of the products of their entries;
¹⁶⁸ i.e., $\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{i_1=1}^{I_1} \cdots \sum_{i_N=1}^{I_N} a_{i_1 \dots i_N} b_{i_1 \dots i_N}$.

¹⁶⁹ 2.2.1. **N -Mode Product**

¹⁷⁰ 2.3. N -Mode Product

¹⁷¹ It means the multiplication of a tensor $\mathbf{A} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ by a matrix $\mathbf{U} \in \mathbb{R}^{J \times I_n}$ or vector $\mathbf{u} \in \mathbb{R}^{I_n}$ in
¹⁷² mode n ; i.e., along axis n . It is represented by $\mathbf{B} = \mathbf{A} \times_n \mathbf{U}$, where $\mathbf{B} \in \mathbb{R}^{I_1 \times \cdots \times I_{n-1} \times J \times I_{n+1} \times \cdots \times I_N}$ [17].

¹⁷³ 2.3.1. **Rank-One Tensor**

¹⁷⁴ 2.4. Rank-One Tensor

¹⁷⁵ A tensor $\mathbf{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ is rank one if it can be written as the outer product of N vectors, i.e.,
¹⁷⁶ $\mathbf{X} = \mathbf{a}^{(1)} \circ \cdots \circ \mathbf{a}^{(N)}$.

¹⁷⁷ 2.4.1. **Rank-R Tensor**

¹⁷⁸ The rank of a tensor $\text{rank}(\mathbf{X})$ is the smallest number of components in a CPD

$$\mathbf{X} = \mathbf{a}^{(1)} \circ \cdots \circ \mathbf{a}^{(N)} \quad (1)$$

¹⁷⁹ where \circ denotes the outer product and $\mathbf{a}^{(n)}$ denotes a vector in a sequence of N
¹⁸⁰ vectors. Each element of the tensor is the product of the corresponding vector elements;
¹⁸¹ i.e., the smallest number of rank-one tensors that generate \mathbf{X} as their sum [17].
¹⁸² $x_{i_1 i_2 \dots i_N} = a_{i_1}^{(1)} \dots a_{i_N}^{(N)}$.

¹⁸³ 2.4.1. **N -Rank**

¹⁸⁴ The n -rank of a tensor $\mathbf{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ denoted $\text{rank}_n(\mathbf{X})$, is the column rank of $\mathbf{X}_{(n)}$; i.e., the
¹⁸⁵ dimension of the vector space spanned by the mode- n fibers. Hence, if $R_n \equiv \text{rank}_n(\mathbf{X})$ for $n = 1, \dots, N$,
¹⁸⁶ we can say that \mathbf{X} has a rank – (R_1, \dots, R_N) tensor [17].

Table 2. Tensor algebra notation summary

$\mathcal{A}, \mathbf{A}, \mathbf{a}, a$	Tensor, matrix, vector and scalar respectively
$\mathcal{A} \in \mathbb{R}^{I_1 \times \dots \times I_N}$	N -order tensor of size $I_1 \times \dots \times I_N$.
$a_{i_1 \dots i_N}$	An element of a tensor
$\mathbf{a}_{i_1 i_2}, \mathbf{a}_{i_1 i_3}, \text{ and } \mathbf{a}_{i_2 i_3}$	Column, row and tube fibers of the third order tensor \mathcal{A}
$\mathbf{A}_{i_1 \dots}, \mathbf{A}_{i_2 \dots}, \mathbf{A}_{i_3 \dots}$	Horizontal, lateral and frontal slices of the third order tensor \mathcal{A}
$\mathbf{A}^{(n)}, \mathbf{a}^{(n)}$	A matrix/vector element from a sequence of matrices/vectors
$\mathbf{A}_{(n)}$	Mode- n matricization of a tensor. $\mathbf{A}_{(n)} \in \mathbb{R}^{I_n \times \prod_{m \neq n} I_m}$
$\mathbf{a}^{(1)} \circ \dots \circ \mathbf{a}^{(N)}$	Outer product of N vectors
$(\mathcal{A}, \mathcal{B})$	Inner product of two tensors.
$\mathcal{B} = \mathcal{A} \times_{\eta} \mathbf{U}$	n -mode product of tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ by a matrix $\mathbf{U} \in \mathbb{R}^{J \times I_n}$ along axis n .

3. Tensor decompositions (TDs)

As an extension of the matrix-based singular value decomposition, two main specific tensor decompositions can be considered; Tucker Decomposition (TKD) [1] and CANDECOMP/PARAFAC (CP) [1]. There are many other tensor decompositions; INDSCAL, PARAFAC2, CANDELINC, DEDICOM, PARATUCK2, among others [17]. Furthermore, there are also nonnegative variants of all of the above. With the aim of preserving particular characteristics of hyperspectral images for pixel-wise classification, this study is limited to use decompositions based on the Tucker model.

3.1. Tucker Decomposition (TKD)

The TKD [17], for the particular case of third-order tensors, can be formally formulated as follows [?]. Given a third-order data tensor $\mathfrak{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ and three positive indices J_1, J_2 and J_3 , find a core tensor $\mathfrak{G} \in \mathbb{R}^{J_1 \times J_2 \times J_3}$ and three component matrices called factor matrices $\mathbf{U}_1 \in \mathbb{R}^{I_1 \times J_1}, \mathbf{U}_2 \in \mathbb{R}^{I_2 \times J_2}$ and $\mathbf{U}_3 \in \mathbb{R}^{I_3 \times J_3}, \mathbf{U}^1 \in \mathbb{R}^{I_1 \times J_1}, \mathbf{U}^2 \in \mathbb{R}^{I_2 \times J_2}$ and $\mathbf{U}^3 \in \mathbb{R}^{I_3 \times J_3}$ which perform the following approximate decomposition:

$$\mathfrak{X} = \mathfrak{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{U}^{(3)} + \mathcal{E} \quad (2)$$

where \mathcal{E} denotes the approximation error tensor. The core tensor \mathfrak{G} preserves the level of interaction for each factor or projection matrix $\mathbf{U}^{(n)}$. The factor matrices are commonly considered orthogonal, but in Tucker models with non-negativity constraints, that is not necessarily imposed [?]. These matrices can be seen as the principal components in each mode [17] (see Figure 1). J_n represents the number of components in the decomposition; i.e., the rank – (R_1, R_2, R_3) .

We can also denote the TKD using the matricization approach and express it by

$$\mathbf{X}_{(1)} = \mathbf{U}^{(1)} \mathbf{G}_{(1)} (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)})^T \quad (3a)$$

$$\mathbf{X}_{(2)} = \mathbf{U}^{(2)} \mathbf{G}_{(2)} (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(1)})^T \quad (3b)$$

$$\mathbf{X}_{(3)} = \mathbf{U}^{(3)} \mathbf{G}_{(3)} (\mathbf{U}^{(2)} \otimes \mathbf{U}^{(1)})^T \quad (3c)$$

where \otimes denotes the Kronecker product and $\mathbf{X}_{(n)}$ and $\mathbf{G}_{(n)}$ are the n -mode matricized versions of tensor \mathfrak{X} and \mathfrak{G} respectively.

Starting from (2), the reconstruction of an approximated tensor can be given by

$$\hat{\mathfrak{X}} = \mathfrak{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{U}^{(3)} \quad (4)$$

where $\hat{\mathfrak{X}}$ is the reconstructed tensor.

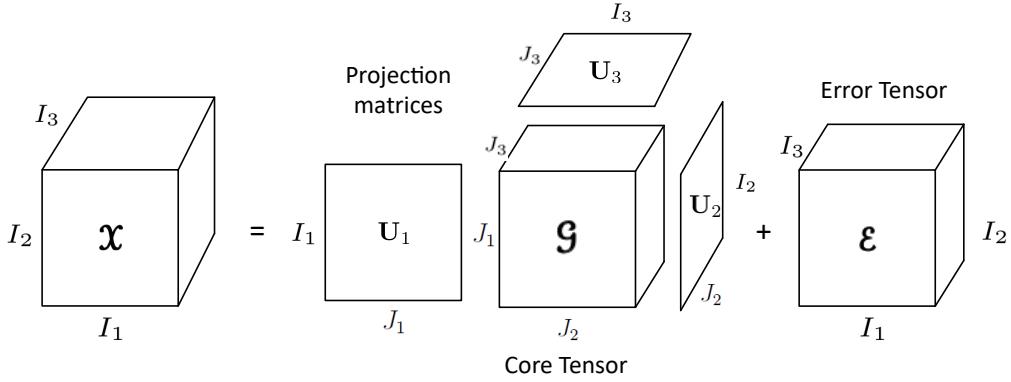


Figure 1. Tucker decomposition for a third-order tensor.

Then, we can acquire the core tensor \mathfrak{G} by the multilinear projection

$$\mathfrak{G} = \mathfrak{X} \times_1 \mathbf{U}^{(1)\top} \times_2 \mathbf{U}^{(2)\top} \times_3 \mathbf{U}^{(3)\top} \quad (5)$$

where $\mathbf{U}^{(n)\top}$ denotes the transpose matrix of $\mathbf{U}^{(n)}$ for $n = 1, \dots, N$. The reconstruction error ξ can be computed as

$$\xi(\hat{\mathfrak{X}}) = \|\mathfrak{X} - \hat{\mathfrak{X}}\|_F^2 \quad (6)$$

and $\|\cdot\|_F$ represents the Frobenius norm. To compute the best rank approximation of a tensor, it is required to use an iterative algorithm. HOSVD, ALS and HOOI are algorithms commonly used in TKD as ALS, HALS, HOOI after a HOSVD initialization [?].

HOOI initializes the factors matrices using HOSVD and assumes that orthogonal matrices are known, so that the core tensor is obtained with (5). Then, it maximizes the cost function

$$\max_{\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}} \|\mathfrak{X} \times_1 \mathbf{U}^{(1)\top} \times_2 \mathbf{U}^{(2)\top} \times_3 \mathbf{U}^{(3)\top}\|_F^2 \quad (7)$$

with $\mathbf{U}^{(n)}$ unknown. Fixing all factor matrices but one, tensor \mathfrak{X} can be projected onto the $\{R_1, \dots, R_{n-1}, R_{n+1}, \dots, R_N\}$ -dimensional space as

$$\mathfrak{W}^{(-n)} = \mathfrak{X} \times_1 \mathbf{U}^{(1)\top} \dots \times_{n-1} \mathbf{U}^{(n-1)\top} \times_{n+1} \mathbf{U}^{(n+1)\top} \dots \times_N \mathbf{U}^{(N)\top} \quad (8)$$

and the orthogonal matrices can be estimated as an orthonormal basis for the dominant subspace of the projection by applying the standard matrix SVD for n -mode unfolded matrix $\mathfrak{W}_{(n)}^{(-n)}$ for $n = 1, 2, 3$ [?].

3.1.1. Non-negative Tucker Decomposition (NTD)

The NTD is a Tucker decomposition with nonnegativity constraints, which decomposes based on the Tucker model. It develops a new tensor factorization method. For the third-order case, the NTD, as defined by Cichocky [15], can be formulated as follows. Given a third-order tensor $\mathfrak{X} \in \mathbb{R}_+^{I_1 \times I_2 \times I_3}$ find a core tensor $\mathfrak{G} \in \mathbb{R}_+^{I_1 \times I_2 \times I_3}$ and the factor matrices $\mathbf{U}_1 \in \mathbb{R}_+^{I_1 \times J_1}$, $\mathbf{U}_2 \in \mathbb{R}_+^{I_2 \times J_2}$ and $\mathbf{U}_3 \in \mathbb{R}_+^{I_3 \times J_3}$ which performs the approximation given in Eq. (2). The TKD method is similarly followed for the NTD.

The HOOI algorithm can be applied for the particular case of third-order tensor. Furthermore, it can be slightly adjusted for preserving any input dimension and develop the decomposition in a specific order. As well as for the TKD model, the best rank approximation of a nonnegative tensor can be computed by an iterative algorithm as HOOI, maximizing the cost function given in equation 7.

²³⁴ Algorithm 1 shows the HOOI algorithm for a NTD.

²³⁵

Algorithm 1: HOOI algorithm to compute a rank- (R_1, \dots, R_N) NTD for an N th-order tensor $\mathbf{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$.

Function HOOI($\mathbf{X}, J_1, \dots, J_N$):

 initialize $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times J_n}$ for $n = 1, \dots, N$ using HOSVD or random

repeat

for $n = 1, \dots, N$ **do**

$\mathbf{W}^{(-n)} \leftarrow \mathbf{X} \times_{-n} \{\mathbf{U}^T\}$

$[\mathbf{U}^{(n)}, \Sigma^{(n)}, \mathbf{V}^{(n)}] \leftarrow \text{svds}(\mathbf{W}_{(n)}^{(-n)}, J_n, 'LM')$

$\mathbf{U}^{(n)} \leftarrow [\mathbf{U}^{(n)}]_+$

end

until fit ceases to improve or maximum iterations exhausted;

$\mathbf{G} \leftarrow \mathbf{W}^{(-N)} \times_N \mathbf{U}^{(N)T}$

Output: $\mathbf{U}^{(n)} \in \mathbb{R}_+^{I_n \times J_n}, \mathbf{G} \in \mathbb{R}_+^{J_1 \times J_2 \times \dots \times J_N}$

²³⁷ 4. Problem phenomenology

²³⁸ 4.1. Spectral Imagery

²³⁹ Multi- or Hyper-spectral images are by nature multidimensional integer nonnegative arrays.
²⁴⁰ A spectral image can be sorted and represented as a third-order tensor $\mathbf{X} \in \mathbb{R}_+^{I_1 \times I_2 \times I_3}$, where
²⁴¹ $\mathbf{X} \in \mathbb{N}^{I_1 \times I_2 \times I_3}$, where N denotes the space of natural numbers. I_1, I_2 and I_3 represent its the height,
²⁴² width and spectral bands respectively. In RS image processing, spectral images are frequently used for
²⁴³ classification of different material in a scene of interest. However, due to the low spatial resolution
²⁴⁴ produced by the distance between the sensor an the target, spatial features are not sufficient to discern
²⁴⁵ certain classes. That is why spectral resolution plays an important role in this type of task.

²⁴⁶ The separation into spectral bands allows perception of reflectance at different wavelengths. This
²⁴⁷ helps to better characterize various materials, in order to simplify the process of discernment between
²⁴⁸ classes. The effort to obtain these spectral features generates a greater amount of data, which increases
²⁴⁹ the processing complexity. This is where the spectral decomposition task becomes relevant.

²⁵⁰ 4.2. Problem Statement

²⁵¹ Given a Spectral Image Let $\mathbf{X} \in \mathbb{N}^{I_1 \times I_2 \times I_3}$ be a spectral image represented as a third-order
²⁵² tensor $\mathbf{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, and its corresponding pixelwise classification, and $\mathbf{Y} \in \mathbb{N}^{I_1 \times I_2}$ its corresponding
²⁵³ ground truth matrix $\mathbf{Y} \in \mathbb{R}^{I_1 \times I_2}$ for a specific number of classes C, find, through Non-negative
²⁵⁴ Tensor Decomposition. Find the best rank- (R_1, R_2, R_n) approximation and its respective core
²⁵⁵ tensor $\mathbf{G} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, through Non-negative Tensor Decompositions. The rank of the decomposition
²⁵⁶ rank_u(\mathbf{X}) is set as rank- (I_1, I_2, I_3) , where $R_1 = I_1, R_2 = I_2$ and $R_3 = I_3 \ll I_3$, to be the input $I_3 \leq I_3$.
²⁵⁷ This built the input space of a pixel-wise classification Convolutional Neural Network using CNNs
²⁵⁸ and produce an output matrix $\hat{\mathbf{Y}}$ of predicted classes, achieving competitive performance metrics for
²⁵⁹ pixel-wise classification while decreasing computational load in the classification process.

²⁶⁰ 4.3. Mathematical Definition

²⁶¹ We can mathematically define the problem statement described above as an optimization problem.

$$\begin{aligned}
& \min_{\mathbf{G}, \mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}} \|\mathbf{X} - \mathbf{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{U}^{(3)}\|_F^2 \\
\text{subject to} \quad & \mathbf{U}^{(n)} \in \mathbb{R}_+^{I_n \times J_n} \quad \text{for } n = 1, 2, 3 \quad \text{and} \quad \mathbf{G} \in \mathbb{R}_+^{J_1 \times J_2 \times J_3} \\
& J_1 = I_1, J_2 = I_2 \quad \text{no compression in the spatial domain,} \\
& J_3 < I_3 \quad \text{reduced spectral domain at the core tensor,} \\
& D(\mathbf{G}_{j_3}) - D(\mathbf{G}_{j_3+1}) < D_s \quad \text{rank searching stop criterion}
\end{aligned} \tag{9}$$

262 5. Deep Convolutional Neural Networks (DCNNs)

CNNs are supervised feed-forward DL-ANNs for computer vision. The idea of applying a sort of convolution of the synaptic weights of a neural network through the input data yields to a preservation of spatial features, which alleviates the hard task of classification and in turn semantic segmentation. This type of ANN works under the same linear regression model as every machine learning (ML) algorithm. Since images are three-dimensional arrays, we can use tensor algebra notation to describe the input of CNNs as a tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, where I_1, I_2 , and I_3 represent height, width, and depth of the third-order array respectively; i.e., the spatial and spectral domain of an image. We can write generally the linear regression model used for ANNs as

$$\hat{\mathbf{y}} = \sigma(\mathbf{W}\mathbf{g} + \mathbf{b})$$

where $\hat{\mathbf{y}}$ represents the output prediction of the network; σ denotes an activation function; \mathbf{g} is the input dataset; \mathbf{W} and \mathbf{b} are the matrix of synaptic weights and the bias vector, respectively. These parameters are adjustable; i.e., their values are modified every iteration looking for convergence to minimize the loss in the prediction through optimization algorithms [32]. For simplicity, the bias vector can be ignored, assuming that matrix \mathbf{W} will update until convergence independently of another parameter [32]. Considering that the input dataset to a CNN is a multidimensional array, we can represent and use tensor algebra notation as

$$\hat{\mathbf{y}} = \sigma(\mathbf{W}\mathbf{g})$$

where $\hat{\mathbf{y}}$ represents the prediction output tensor of the ANN (in our case, a second-order tensor or matrix $\hat{\mathbf{Y}}$), \mathbf{g} is the input dataset, and \mathbf{W} is a $K_1 \times K_2 \times F_1$ tensor called filter or kernel with the adaptable synaptic weights. Different to conventional ANN, in CNNs, \mathbf{W} is a shiftable square tensor is much smaller in height and width than the input data, i.e., $K_1 = K_2$ and $K_s \ll I_s$ for $s = 1, 2$; F_1 denotes the number of input channels; i.e., $F_1 = I_3$. For hidden layers, instead of the prediction tensor $\hat{\mathbf{y}}$, the output is a matrix called activation map $\mathbf{M} \in \mathbb{R}^{I_1 \times I_2}$, which preserves features from the original data in each domain. Actually, it is necessary to use much kernels $\mathbf{W}^{(f_2)}$ as activation maps, with different initialization values to preserve diverse features of the image. Hence, we can also define activation maps as a tensor $\mathbf{M} \in \mathbb{R}^{I_1 \times I_2 \times F_2}$ where F_2 denotes the number of activation maps produced by each filter (see Figure ??). Kernels are displaced through the whole input image as a discrete convolution operation. Then, each element of the output activation map $m_{i_1 i_2 f_2}$ is computed by the summary of the Hadamard product of kernel $\mathbf{W}^{(f_2)}$ and a subtensor from the input tensor \mathbf{g} centered in position (i, j) and with same dimensions of \mathbf{W} , as follows-

$$m_{i_1 i_2 f_2} = \sigma \left[\sum_{k_1=1}^{K_1} \sum_{k_2=1}^{K_2} \sum_{f_1=1}^{F_1} w_{k_1, k_2, f_1} g_{i_1+k_1-o_1, i_2+k_2-o_2, f_1} \right]$$

263 where $m_{i_1 i_2 j_2}$ denotes the value of the output activation map f_2 at position i_1, i_2 ; σ represents the
 264 activation function; and o_1 and o_2 are offsets in spatial dimensions which depend on the kernel size,
 265 and equal $\frac{K_1+1}{2}$ and $\frac{K_2+1}{2}$ respectively (see Figure ??). Convolutional layer with a $K_1 \times K_2 \times F_1 \times F_2$
 266 kernel. Input channels F_1 must equal the spectral bands I_3 . To preserve original dimensions at the
 267 output, zero padding is needed [18]. Output dimensions also depend on stride $S = 1$ to consider
 268 every piece of pixel information and to preserve original dimensions.

An ANN is trained by using iterative gradient-based optimizers, such as Stochastic gradient descent, Momentum, RMSprop, and Adam [32]. This drive the cost function $L(\mathbf{W})$ to a very low value by updating the synaptic weights \mathbf{W} . We can compute the cost function by any function that measures the difference between the training data and the prediction, such as Euclidean distance or cross-entropy [10]. Besides, the same function is used to measure the performance of the model during testing and validation. In order to avoid overfitting [32], the total cost function used to train an ANN combines one of the cost functions mentioned before, plus a regularization term.

$$J(\mathbf{W}) = L(\mathbf{W}) + R(\mathbf{W}),$$

where $J(\mathbf{W})$ denotes the total cost function and $R(\mathbf{W})$ represents a regularization function. Then, we can decrease $J(\mathbf{W})$ by updating the synaptic weights in the direction of the negative gradient. This is known as the method of steepest descent or gradient descent.

$$\mathbf{W}' = \mathbf{W} - \phi \nabla_{\mathbf{W}} J(\mathbf{W}),$$

269 where \mathbf{W}' represents the synaptic weights tensor in next iteration during training, ϕ denotes the
 270 learning rate parameter, and $\nabla_{\mathbf{W}} J(\mathbf{W})$ the cost function gradient. Gradient descent converges when
 271 every element of the gradient is zero, or in practice, very close to zero [10].

272 CNNs has been successfully used in many image classification frameworks. This variation
 273 in architecture from other typical ANN models yields the network to learn spatial and spectral
 274 features, which are highly profitable for image classification. Besides, FCNs, constructed with only
 275 convolutional layers are able to classify each element of the input image; i.e., they yield pixel-wise
 276 classification, or in other words, semantic segmentation.

277 5. Methodology

278 The following subsections described the methodology followed for the framework propose in this
 279 work. We can summarize the big picture in three steps: the HSI modeling, the tensor decomposition,
 280 the classifier and the decision making.

281 5.1. Tensor modeling

282 Consider an input dataset $\mathcal{X} \in \mathbb{N}^{I_1 \times I_2 \times I_3}$ with $I_1 \times I_2 \times I_3$ samples in the space of the natural
 283 numbers \mathbb{N} , where a fiber $\mathbf{x}_{i_1 i_2}$ represents the spectra or endmember of pixel $i_1 i_2$ and can be represented
 284 by the Linear Mixing Model (LMM) as follows

$$\mathbf{x}_{i_1 i_2} = \sum_{c=1}^C (\alpha_{i_1 i_2 c} \mathbf{m}_c + \boldsymbol{\eta}) \quad (10)$$

285 where α_c is the contribution of material c at pixel $i_1 i_2$, \mathbf{m}_c denotes the endmember of a specific material
 286 c , and $\boldsymbol{\eta}$ represents an additive noise vector. The abundance vectors $\alpha_{i_1 i_2}$ must always satisfy two
 287 constraints, i) the non-negativity, $\alpha_{i_1 i_2 c} \geq 0$ for all $c = 1, \dots, C$, and ii) the sum-to-one restriction,
 288 $\sum_{c=1}^C \alpha_{i_1 i_2 c} = 1$. Figure 2a shows the spectral signatures for the Indian Pines dataset.

²⁸⁹ 5.2. *Tensor factorization*

²⁹⁰ Consider $\mathbf{Y} \in \mathbb{C}^{I_1 \times I_2}$ as the matrix of actual classes corresponding to our dataset \mathbf{X} , and $\hat{\mathbf{Y}} \in$
²⁹¹ $\mathbb{C}^{I_1 \times I_2}$ as the prediction matrix, where \mathbb{C} defines the set of C different classes. In order to reduce
²⁹² data dimensionality of the input dataset \mathbf{X} while keeping classifier performance, we propose to
²⁹³ use the restricted NTD denoted as \mathcal{T} , producing a core tensor $\mathbf{G} \in \mathbb{R}^{I_1 \times I_2 \times J_3}$ and n factors matrices
²⁹⁴ $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times J_n}$, expressed as

$$\mathbf{X} \xrightarrow{\mathcal{T}} (\mathbf{G}, \mathbf{U}^{(n)}) \quad (11)$$

²⁹⁵ where the decomposition is restricted to preserve the spatial domain and to be only in the 3rd-mode
²⁹⁶ by the Tucker1 model

$$\mathbf{X} = \mathbf{G} \times_1 \mathbf{I} \times_2 \mathbf{I} \times_3 \mathbf{U}^{(3)} \quad (12)$$

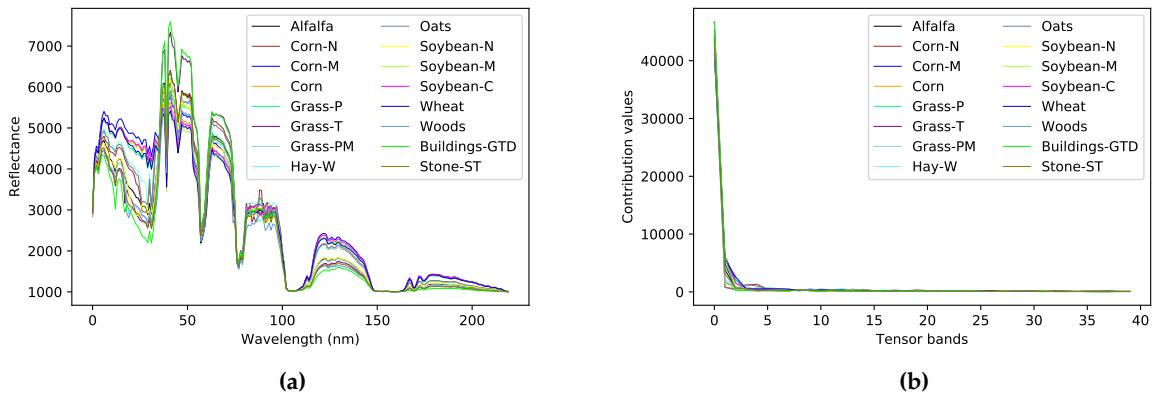


Figure 2. Behavior of the 16 classes of the Indian Pines dataset, a) in the spectral domain (spectral signatures) and b) in the tensor bands domain.

Hence, each fiber $\mathbf{x}_{i_1 i_2}$ of the core tensor takes a new representation in the tensor bands domain and can be mathematically defined as follows

$$\mathbf{g}_{i_1 i_2} = \sum_{c=1}^C (\beta_{i_1 i_2 c} \mathbf{s}_c + \boldsymbol{\eta}) \quad (13)$$

²⁹⁷ where β_c is the contribution of material c at pixel $i_1 i_2$ and \mathbf{s}_c denotes the endmember of a specific
²⁹⁸ material c . We can see in Figure 2b the new tensor bands values for each class.

²⁹⁹ We also propose a variant of the NTD that does an integer decomposition, i. e., the Integer
³⁰⁰ Non-negative Tucker decomposition (INTD). The INTD follows the same Tucker model described in
³⁰¹ Section 3.1. It considers the additional restriction of decomposing a tensor in the set of the natural
³⁰² numbers.

³⁰³ 5.3. *Classifier*

³⁰⁴ The tensor decompositions generates a core tensor based on the Tucker1 model produce a core
³⁰⁵ tensor, where the first tensor bands provide the most of the information a signature enough to
³⁰⁶ differentiate the classes of interest from of the input dataset. Then, the core tensor $\mathbf{G} \in \mathbb{N}^{I_1 \times I_2 \times J_3}$,
³⁰⁷ with $J_3 < I_3$, and its corresponding ground truth \mathbf{Y} form the input tuple of the classifier Θ , which
³⁰⁸ produce a predicted label for each element of the input, i.e.,

$$(\mathbf{G}, \mathbf{Y}) \xrightarrow{\Theta} \hat{\mathbf{Y}} \quad (14)$$

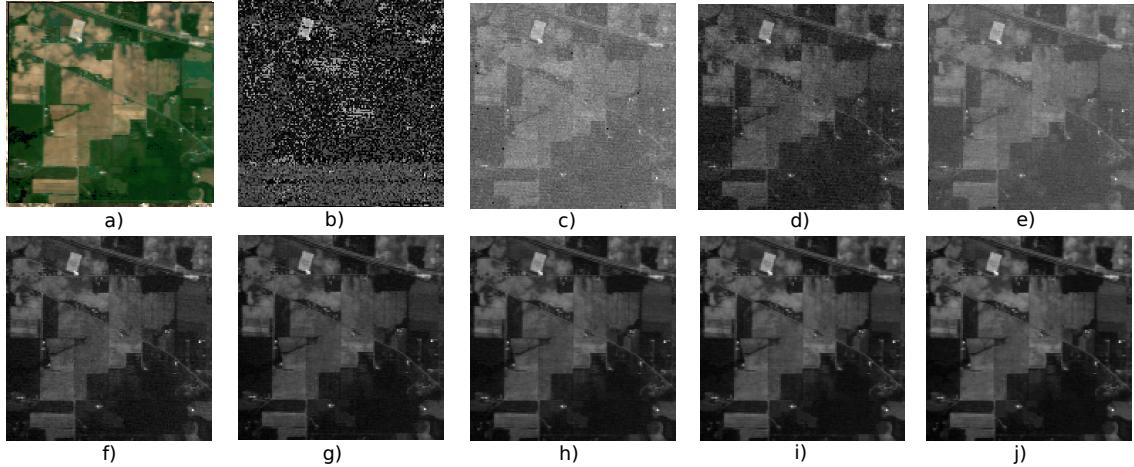


Figure 3. Indian Pines dataset a) True color image b) - j) 1st to 9th spectral band.

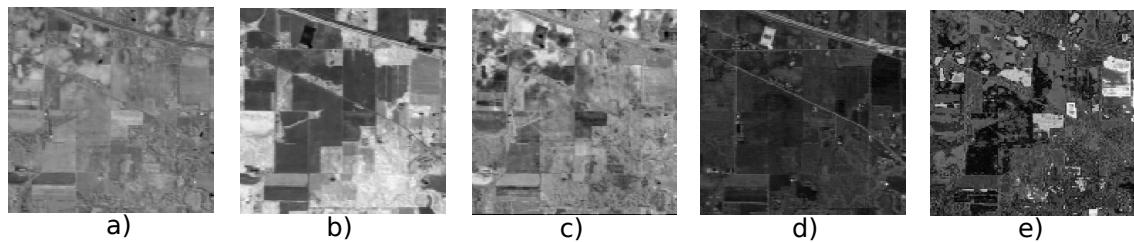


Figure 4. TKD tensor bands of the Indian Pines dataset a) - e) 1st to 5th tensor band.

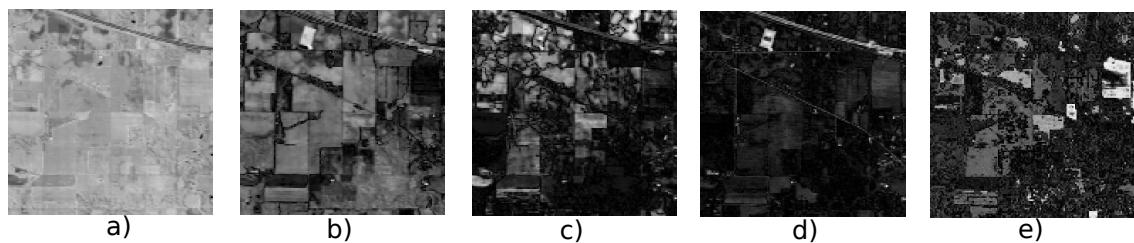


Figure 5. NTD tensor bands of the Indian Pines dataset a) - e) 1st to 5th tensor band.

The performance of our classification model can be measured by the cross-entropy loss, whose output is a probability value. The cross-entropy loss increases as the predicted probability diverges from the actual label and it is computed as

$$J(\mathbf{W}) = -\mathbb{E}_{\mathbf{G}, \mathbf{Y} \sim \hat{p}} \log p(\mathbf{Y} | \mathbf{G}) \quad (15)$$

where $J(\mathbf{G})$ represents the loss function. For a multiclass probability distribution, the cross entropy cost function can be written as

$$H(y, p) = - \sum_{c=1}^C y_c \log(p_c) \quad (16)$$

where $H(y, p)$ denotes the cross entropy of targets y with a probability p .

We use the softmax function as the output of our classifier, to represent the probability distribution over C different classes. Formally, the softmax function is given by

$$\delta(\mathbf{z})_c = \frac{e^{z_c}}{\sum_{l=1}^L e^{z_l}} \quad (17)$$

where $\delta(\mathbf{z})_c$ denotes the softmax function of vector \mathbf{z} , which is each 3rd-mode fiber of the activation maps at the last convolutional layer. Hence, the softmax function produces a normalized probability distribution for every input pixel, which can be seen as the contribution parameter in the LMM Eq. 10.

In this paper, we aim to feed supervised classifiers based on 3D-CNN, with a lower dimensionality tensor than the original dataset. This has three particular motivations: 1) to avoid overfitting the DCNN, 2) to reduce the computational complexity, and 3) keep the classifier performance while reducing the execution time. Figure 6 shows the big picture of the framework proposed.

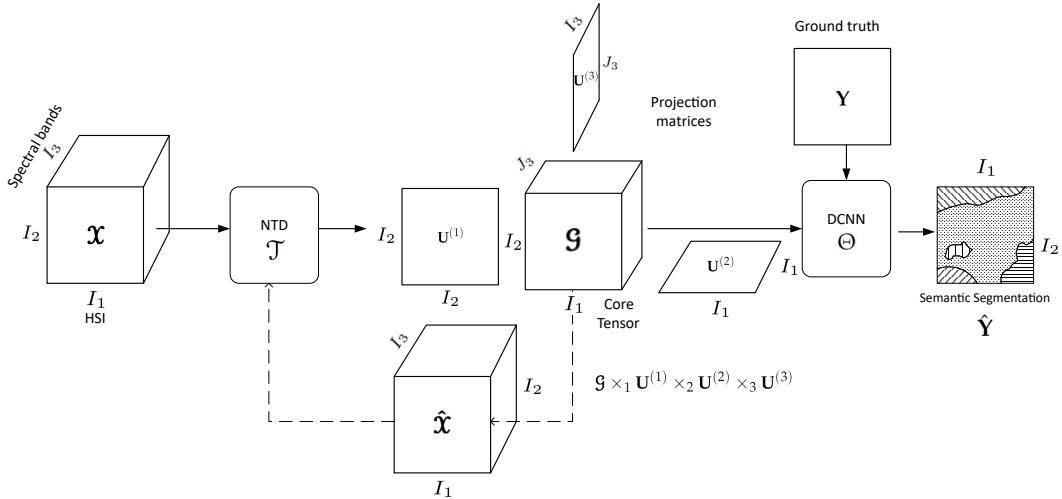


Figure 6. Big picture of the framework proposed.

5.4. Decomposition analysis

Our proposed framework faces two main challenges. The first is the selection of the decomposition method. We are looking for the best representation of the input dataset for a 3D-CNN, so, we limit the set of methods to those that produce a decomposition tensor with the same structure as the input tensor. TKD, NTD and the INTD proposed generate a core tensor with the desired properties.

The second challenge is the search for the rank – (J_1, J_2, J_3). As we want to preserve the spatial domain $J_1 = I_1$ and $J_2 = I_2$, but J_3 has to be selected so that the performance of the classifier does not decrease considerably. Both decisions are made under a criterion under the probabilistic point of view.

The Kullback-Leibler divergence is a metric for quantifying the difference between the probability distributions of the original MSI X and the core tensor G as

$$D_{KL}(X||G) = \sum_{i=1}^I p_i(x) \log \frac{p_i(x)}{g_i(x)} \quad (18)$$

where $D_{KL}(X||G)$ represents the KL divergence of the two probability distributions. Figure 7a shows the results of this metric for the three decomposition used in this work.

Besides, a symmetric version of the KL divergence is used, the Jensen Shanon divergence. This is another method of measuring the similarity between two probability distributions. Defining $M = \frac{X+G}{2}$, we can write the JS divergence as

$$D_{JS}(X||G) = \frac{1}{2}D_{KL}(X||M) + \frac{1}{2}D_{KL}(G||M) \quad (19)$$

where $D_{JS}(X||G)$ represents the JS divergence of the probability distributions X and G . Figure 7b shows the JS divergence for the TKD, NTD and INTD.

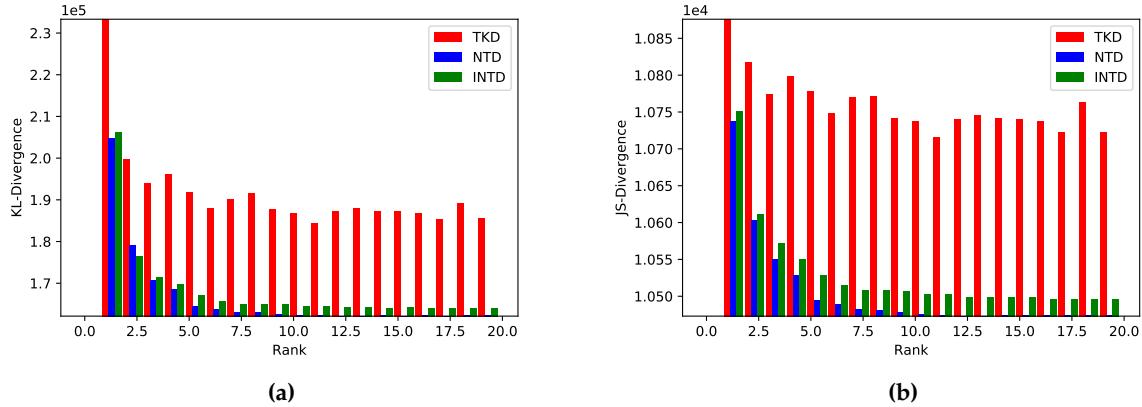


Figure 7. Divergences between the input dataset and the core tensor generated by TKD, NTD and INTD, a) KL-divergence and b) JS-divergence.

341 6. Experimental Results

342 6.1. Input Data

343 For this work, we chose three of the most popular multi- and hyperspectral dataset for
344 classification.

345 6.1.1. Sentinel-2

346 This dataset [developed propose](#) by Lopez et al. [?] is composed of RS Sentinel-2 scenarios from
347 central Europe. It has 100 scenarios for the training space and 10 scenarios for testing, all of them
348 with 128×128 pixels with spatial resolution of $20m^2$ and 9 spectral bands in the range $490 - 2190nm$.
349 The labels are semi-manually assigned for five classes of interest: vegetation, soil, water, clouds and
350 shadows. Data are available in the link [Sentinel-2 Dataset](#).

Table 3. Average of contribution per class in Sentinel-2 dataset.

Class	Dataset (%)	Train (%)	Test (%)
Soil	23.58 425,147	22.87 374,792	30.73 50,355
Shadow	2.52 45,418	2.58 42,281	1.91 3,137
Cloud	13.87 250,020	14.37 235,546	8.83 14,474
Water	12.13 218,620	11.93 195,512	14.10 23,108
Vegetation	47.88 863,035	48.23 790,269	44.41 72,766

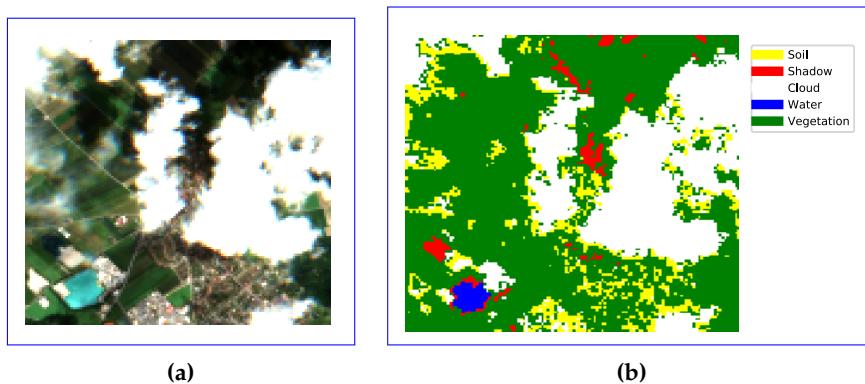


Figure 8. [Sentinel dataset](#), a) True color and b) Ground truth.

351 6.1.2. Indian Pines

352 This dataset is a scene produced by AVIRIS in North-western Indiana and consists of 145×145
 353 pixels and 224 spectral bands in the wavelength range $0.4\text{--}2.5\mu\text{m}$. The Indian Pines scene contains
 354 two-thirds agriculture, and one-third forest or other natural perennial vegetation. There are two major
 355 dual lane highways, a rail line, as well as some low density housing, other built structures, and smaller
 356 roads. Since the scene is taken in June some of the crops present, corn, soybeans, are in early stages of
 357 growth with less than 5% coverage. The ground truth available is designated into sixteen classes and is
 358 not all mutually exclusive. Indian Pines data are available at [Indian Pines dataset](#). [Figure 9 shows the](#)
 359 [true color image of Salinas dataset, as well as the ground truth with each of its corresponding classes](#)
 360 [and Table 4 the number of samples for each class.](#)

Table 4. Average of contribution per class in Indian Pines dataset.

Class	Samples
Alfalfa	46
Corn-notill	1428
Corn-mintill	830
Corn	237
Grass-pasture	483
Grass-tress	730
Grass-pasture-mowed	28
Hay-windrowed	478
Oats	20
Soybean-notill	972
Soybean-mintill	2455
Soybean-clean	593
Wheat	205
Woods	1265
Building-Grass-Trees-Drives	386
Stone-Steel-Towers	93



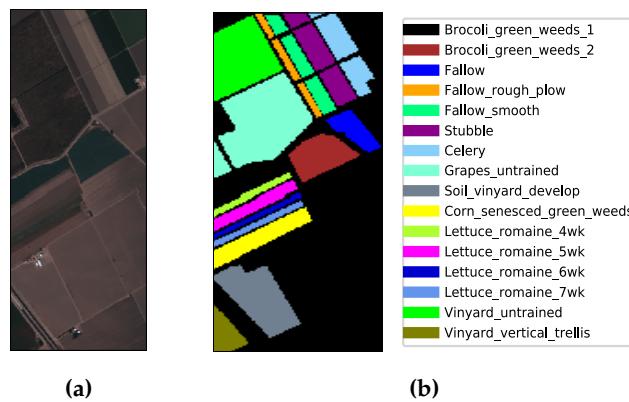
Figure 9. Indian Pines dataset, a) True color and b) Ground truth.

361 6.1.3. Salinas

362 This scene was collected by the AVIRIS sensor over Salinas Valley, California. It has 512×217
 363 pixels with spatial resolution 3.7m , and 224 spectral bands. It includes vegetables, bare soils, and
 364 vineyard fields. Salinas groundtruth contains 16 classes [shows in table 5. Salinas data are available](#)
 365 [at Salinas dataset. Figure 10 shows the true color image of the Salinas dataset, as well as the ground](#)
 366 [truth labeled with each of its corresponding classes and Table 5 the number of samples for each class.](#)

Table 5. Average of contribution per class in Salinas dataset.

Class	Samples
Brocoli-green-weeds-1	2009
Brocoli-green-weeds-2	3726
Fallow	1976
Fallow-rough-plow	1394
Fallow-smooth	2678
Stubble	3959
Celery	3579
Grapes-untrained	11271
Soil-vinyard-develop	6203
Corn-senesced-green-weeds	3278
Lettuce-romaine-4wk	1068
Lettuce-romaine-5wk	1927
Lettuce-romaine-6wk	916
Lettuce-romaine-7wk	1070
Vinyard-untrained	7268
Vinyard-vertical-trellis	1807

**Figure 10.** Salinas dataset, a) True color and b) Ground truth.

368 6.1.4. Pavia University

369 This scene was collected by the ROSIS sensor over Pavia University, nothern Italy. It has
 370 610 × 340 pixels with spatial resolution 1.3m. and 103 spectral bands. Pavia groundtruth contains 9
 371 classes, some of them described in table 6. Pavia University data are available at Pavia University
 372 dataset. Figure 11 shows the true color image of the Pavia dataset, as well as the ground truth labeled
 373 with each of its corresponding classes.

374

Table 6. Average of contribution per class in Pavia dataset.

Class	Samples
Asphalt	6631
Meadows	18649
Gravel	2099
Trees	3064
Painted metal sheets	1345
Bare Soil	5029
Bitumen	1330
Self-Blocking Bricks	3682
Shadows	947

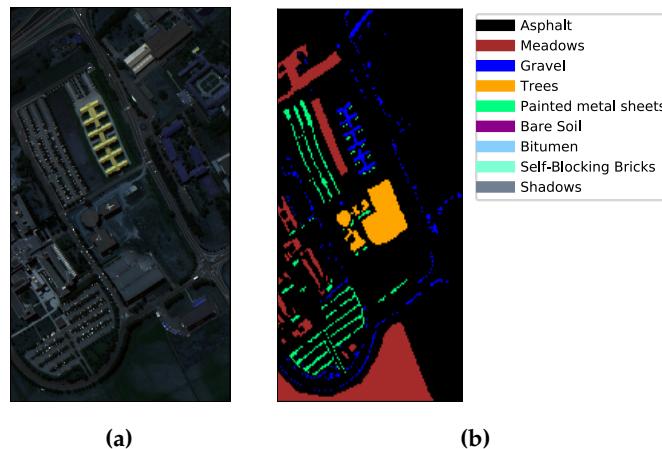


Figure 11. Pavia dataset, a) True color and b) Ground truth.

In the following table 7 it is summarized the datasets used in this work as well as their spatial and spectral characteristics, the number of classes and their samples.

Table 7. Summary of the different dataset used for experiments in this work.

Dataset	Spatial dimensions	Bands	Classes	Samples
Sentinel-2 CNNMSI	128 × 128	9	5	16,220,160 384
Indian Pines	145 × 145	224 220	16	4,709,600 21,025
Salinas	512 × 217	224	16	24,887,296 111,104
Pavia University	610 × 340	103	9	207,400

6.2. MetricsCNN Specifications

The model used to evaluate the framework proposed in this work is Segnet [1]. The strategie of cross-validation was used to set the hyperparameters of the CNN getting the following

- learning rate: 1×10^{-3} 383
- epochs: 100 384
- optimizer: Adam [1] 385
- initialization: Xavier [1]
- kernel dimensions: 3×3
- Activation Function: ReLU / Softmax

6.2.1. Computational Specifications

The software and hardware tools used to run the framework proposed are summarized as follows

- Platform: Python 3.7 392
- AI Framework: Tensorflow 1.13 393
- GPU: NVIDIA GeForce GTX 1050 Ti 394
- Processor: Intel core i7
- RAM: 8GB
- SSD: 128GB / HDD: 1TB

6.3. Algorithms metrics

6.3.1. Relative Mean Square Error (rMSE)

In order to To compute the reconstruction error of any decomposition, it can be used the relative Mean Square Error, given by

$$rMSE(\hat{\mathbf{x}}) = \frac{1}{Q} \sum_{q=1}^Q \frac{\|\hat{\mathbf{x}}_q - \mathbf{x}_q\|_F^2}{\|\mathbf{x}_q\|_F^2}, \quad (20)$$

³⁹⁹ where \mathbf{x}_q represents the q -th MSI from a dataset with Q MSIs and $\hat{\mathbf{x}}_q$ its corresponding reconstruction
⁴⁰⁰ computed by (4).

⁴⁰¹ 6.3.2. Loss

⁴⁰² To quantify the convergence of the tensor \mathbf{x} for the implementation of HOOI, the rMSE was
⁴⁰³ used: iterative process in a supervised classification algorithm, it can be used a metric of difference
⁴⁰⁴ between the lables \mathbf{Y} and its corresponding prediction $\hat{\mathbf{Y}}$ as

$$\text{rMSE} = \frac{1}{Q} \sum_{q=1}^Q \frac{\|\hat{\mathbf{x}}_q - \mathbf{x}_q\|_F^2}{\|\mathbf{x}_q\|_F^2}, \quad (21)$$

⁴⁰⁵ where \mathbf{x}_q represents the q -th CNNMSI from our dataset with Q MSIs and $\hat{\mathbf{x}}_q$ its corresponding
⁴⁰⁶ reconstruction computed by (4).

⁴⁰⁷ 6.4. Evaluation metrics

⁴⁰⁸ 6.4.1. Cohen's Kappa Coefficient

⁴⁰⁹ Cohen's kappa coefficient is a very robust metric used to measure reliability of multi-class and
⁴¹⁰ imbalanced class classification algorithms [1]. It is computed by

$$\kappa = \frac{\rho_o - \rho_e}{1 - \rho_e} \quad (22)$$

⁴¹¹ where ρ_o is the observed agreement, and ρ_e is the expected agreement. This metric will always produce
⁴¹² values less than or equal to 1, where 1 means perfect agreement. Negative values indicate no agreement.
⁴¹³ i.e., futile classification.

⁴¹⁴ 6.4.2. Pixel Accuracy (PA)

⁴¹⁵ We used the PA metric to compute a ratio between the amount of correctly classified pixels and
⁴¹⁶ the total number of pixels as follows. Given a confusion matrix relating the True Positive (TP), True
⁴¹⁷ Negatives (TN), False Positives (FP) and False Negatives (FN), the PA is computed by

$$PA = \frac{\sum_{c=1}^C \tau_{cc}}{\sum_{c=1}^C \sum_{d=1}^C \tau_{cd}} = \frac{TP + TN}{TP + TN + FP + FN} \quad (23)$$

⁴¹⁸ where c is the number of class for $c = 1, \dots, C$ and τ_{cc} is the amount of pixels of class c correctly
⁴¹⁹ assigned to class c i.e., TP and TN, and τ_{cd} is the amount of pixels of class c inferred to belong to class
⁴²⁰ d . Despite this metric is wide used, it is not a totally fair metric for imbalanced classes. In this work we
⁴²¹ used this metric with comparison purposes. However, we also used metrics more in line with multi
⁴²² class classification tasks.

⁴²³ 6.4.3. Precision

⁴²⁴ Another metric used in this work as a performance evaluation metric in multiclass classification
⁴²⁵ is precision. This is a metric that measures the percentage of pixels from class c correctly classified. It
⁴²⁶ is computed with the following equation

$$\Psi_c = \frac{\tau_{cc}}{\sum_{d=1}^C \tau_{cd}} = \frac{TP}{TP + FP} \quad (24)$$

⁴²⁷ where Ψ_c denotes the precision of class c , which is the number of TP divided by the TP plus FP.

428 6.4.4. Recall or Sensitivity

429 Recall, or also known as sensitivity is a metric that indicates the proportion of pixels classified as
430 class c that actually belong to class c . It is computed with the following equation

$$v_c = \frac{\tau_{cc}}{\sum_{d=1}^C \tau_{dc}} = \frac{TP}{TP + FN} \quad (25)$$

431 where v_c denotes the recall of class c for $c = 1, \dots, C$.

432 6.4.5. F1 Score

433 In order to summarize precision and recall in one only metric, we use the F1 score, which is
434 computed by

$$F1 = \frac{2\Psi v}{\Psi + Y} \frac{2\Psi v}{\Psi + v} \quad (26)$$

435 This metric provides a very appropriate measure of multiclass classification for imbalanced dataset.

436 7. Discussion and Comparison

437 8. Conclusions

438 7.1. Figures, Tables and Schemes

439 Original Sentinel-2 spectral bands, a) True color image, b) to j) 1st to 9th spectral band respectively.
440 Tensor bands from the Tucker Decomposition, a) to i) 1st to 9th band respectively. Tensor bands from
441 the Non-negatice Tucker Decomposition, a) to i) 1st to 9th band respectively. Original Sentinel-2 spectral
442 bands. Tucker Decomposition Tensor bands. Nonnegative Tucker Decomposition Tensor bands. Original
443 Sentinel-2 spectral bands. Tucker Decomposition Tensor bands. Nonnegative Tucker Decomposition
444 Tensor bands.

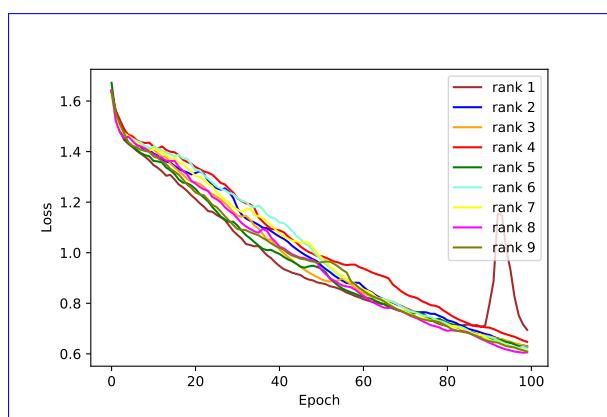


Figure 12. Indian Pines Spectral Signatures a) Original, b) Loss function of the CNN after NTD and c) TKD for rank $J_3 = 1$ to 9.

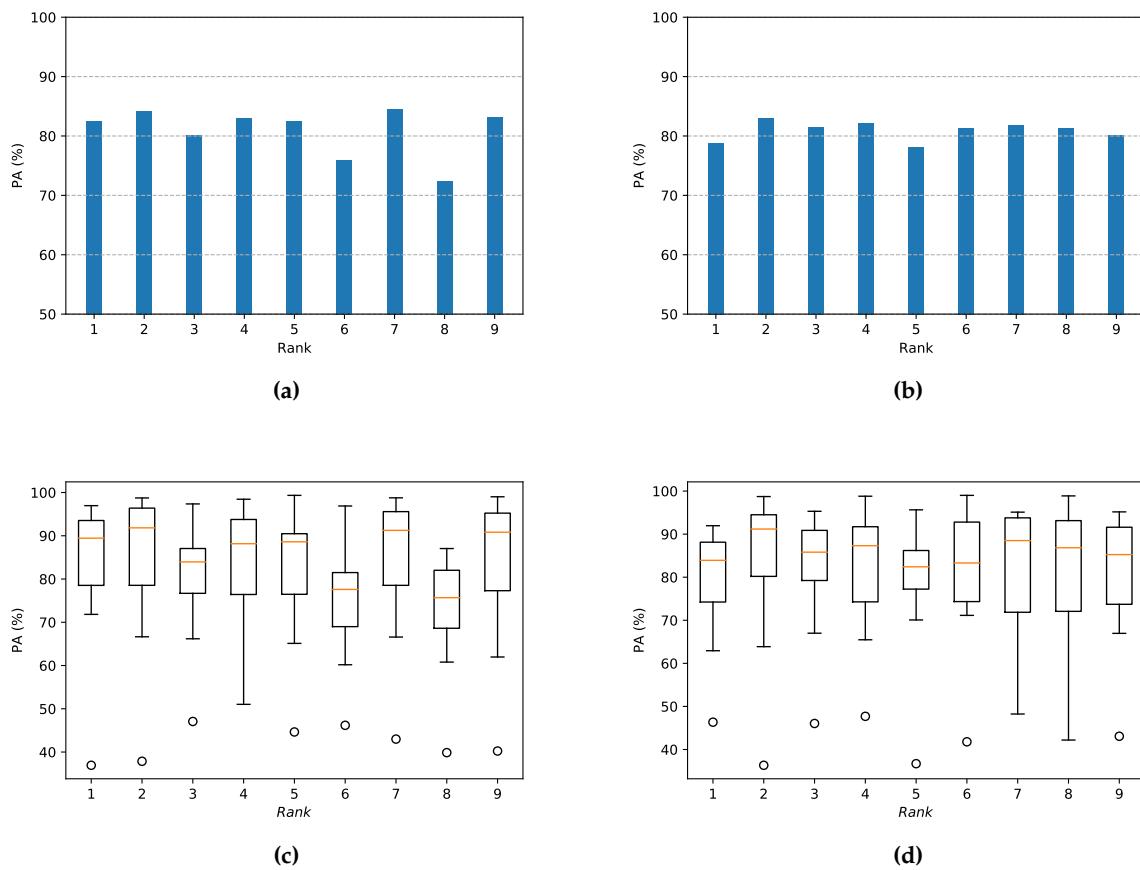


Figure 13. Pixel accuracy vs Rank results a) Comparative bar plot NTKD and TKD, b) Comparative bar plot NTKD and TKD c) Box and whiskers plot for NTKD, and d) Box and whiskers plot for TKD

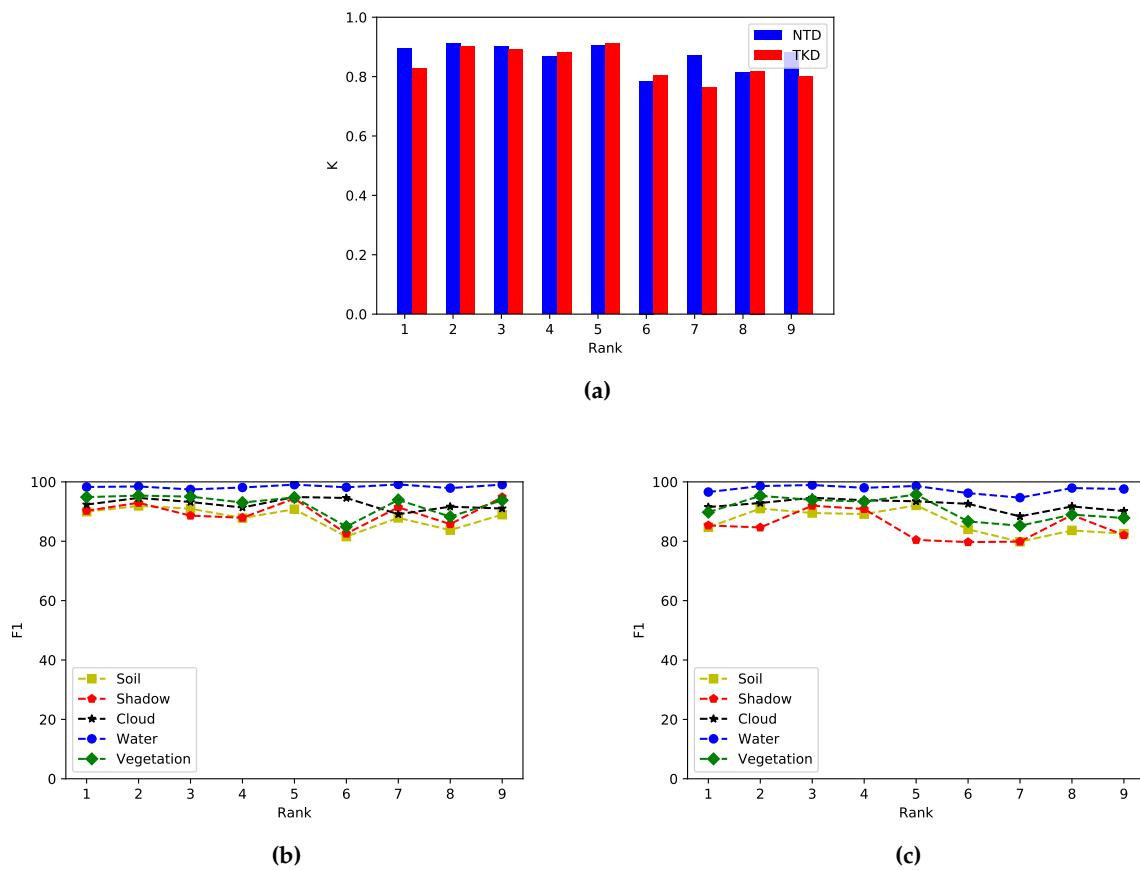


Figure 14. Performance evaluation metrics. a) Kappa score comparison NTD vs TKD, b) F1 for NTD and c) F1 for TKD

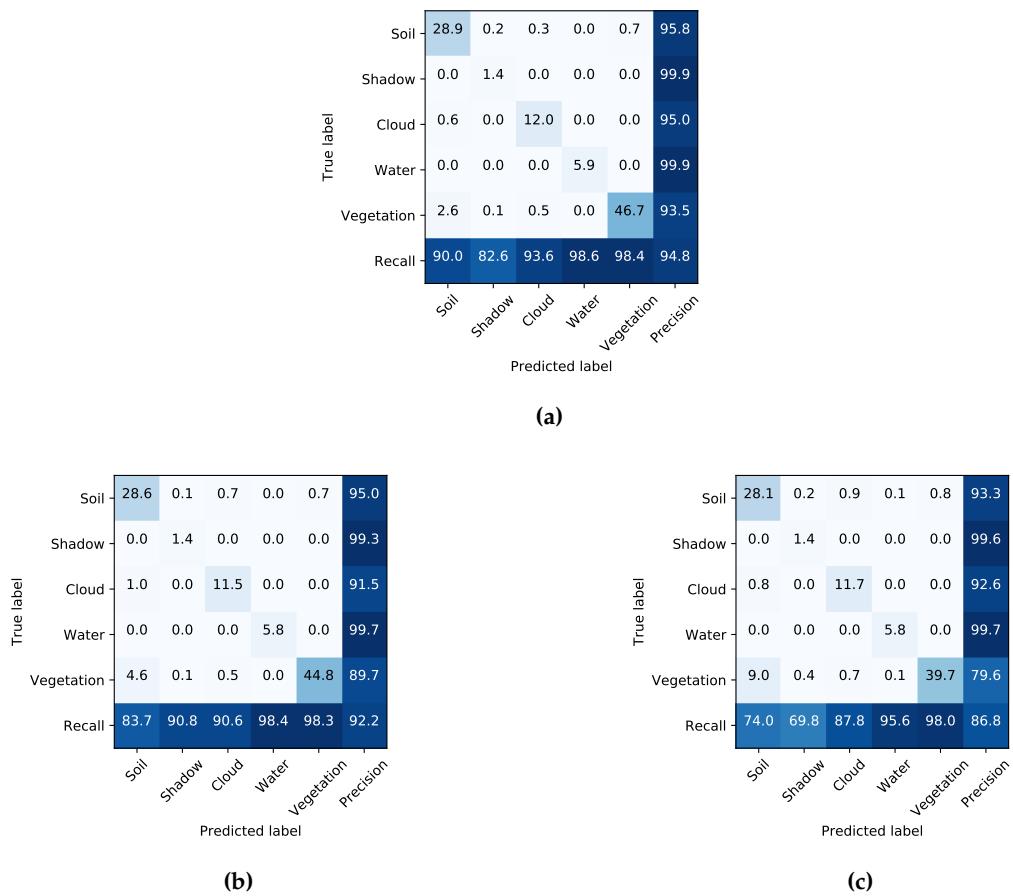


Figure 15. Confusion matrix for a) Sentinel-2 original dataset, b) NTD, c) TKD.

8. Conclusions

Author Contributions: Conceptualization, J.L.; formal analysis, D.T.; investigation, J.L.; methodology, J.L., D.T., and C.A.; resources, C.A.; software, J.L.; supervision, D.T. and C.A.; validation, D.T. and C.A.; writing—original draft, J.L. and D.T.

Funding: This work was supported by the National Council of Science and Technology CONACYT of Mexico under grant XXXXXXXX.

Acknowledgments:

Conflicts of Interest: The authors declare no conflict of interest.

Abbreviations

The following abbreviations are used in this manuscript:

ANN	Artificial Neural Network
CNN	Convolutional neural network
CPD	Canonical Polyadic Decomposition
DL	Deep Learning
FCN	Fully Convolutional Network
HOOI	Higher-Order Orthogonal Iteration
HOSVD	Higher-Order Singular Value Decomposition

457 References

- 458 1. Tempfli, K.; Huurneman, G.; Bakker, W.; Janssen, L.; Feringa, W.; Gieske, A.; Grabmaier, K.; Hecker, C.; Horn, J.; Kerle, N.; et al. *Principles of Remote Sensing: An Introductory Textbook*, 4th ed.; ITC: Geneva, Switzerland, 2009.
- 461 2. He, Z.; Hu, J.; Wang, Y. Low-rank tensor learning for classification of hyperspectral image with limited labeled sample. *IEEE Signal Process.* **2017**, *145*, 12–25.
- 463 3. Richards, A.; Xiuping, J.J. Band selection in sentinel-2 satellite for agriculture applications. In *Remote Sensing Digital Image Analysis*, 4th ed.; Springer-Verlag: Berlin, Germany, 2006.
- 465 4. Zhang, T.; Su, J.; Liu, C.; Chen, W.; Liu, H.; Liu, G. Band selection in sentinel-2 satellite for agriculture applications. In Proceedings of the 23rd International Conference on Automation & Computing, University of Huddersfield, Huddersfield, UK, 7–8 September 2017.
- 468 5. Xie, Y.; Zhao, X.; Li, L.; Wang, H. Calculating NDVI for Landsat7-ETM data after atmospheric correction using 6S model: A case study in Zhangye city, China. In Proceedings of the 18th International Conference on Geoinformatics, Beijing, China, 18–20 June 2010.
- 471 6. Gao, B. NDWI—A normalized difference water index for remote sensing of vegetation liquid water from space. *Remote Sens. Environ.* **1996**, *58*, 1–6.
- 473 7. Ham, J.; Chen, Y.; Crawford, M.; Ghosh, J. Investigation of the random forest framework for classification of hyperspectral data. *IEEE Trans. Geosci. Remote Sens.* **2005**, *43*, 492–501.
- 475 8. Hearst, Marti A. Support Vector Machines. *IEEE Intell. Syst. J.* **1998**, *13*, 18–28.
- 476 9. Huang, X.; Zhang, L. An SVM Ensemble Approach Combining Spectral, Structural, and Semantic Features for the Classification of High-Resolution Remotely Sensed Imagery. *IEEE Trans. Geosci. Remote Sens.* **2013**, *51*, 257–272.
- 479 10. Delalieux, S.; Somers, B.; Haest, B.; Spanhove, T.; Vanden Borre, J.; Mucher, S. Heathland conservation status mapping through integration of hyperspectral mixture analysis and decision tree classifiers. *Remote Sens. Environ.* **2012**, *126*, 222–231.
- 482 11. Kemker, R.; Salvaggio, C.; Kanan, C. Algorithms for semantic segmentation of multispectral remote sensing imagery using deep learning. *ISPRS J. Photogramm. Remote Sens.* **2018**, *145*, 60–77.
- 484 12. Pirotti, F.; Sunar, F.; Piragnolo, M. Benchmark of machine learning methods for classification of a sentinel-2 image. In Proceedings of the XXIII ISPRS Congress, Prague, Czech Republic, 12–19 July 2016.
- 486 13. Mateo-García, G.; Gómez-Chova, L.; Camps-Valls, G. Convolutional neural networks for multispectral image cloud masking. In Proceedings of the IGARSS, Fort Worth, TX, USA, 23–28 July 2017.
- 488 14. Guo, X.; Huang, X.; Zhang, L.; Zhang, L.; Plaza, A.; Benediktsson, J. A. Support Tensor Machines for Classification of Hyperspectral Remote Sensing Imagery. *IEEE Trans. Geosci. Remote Sens.* **2016**, *54*, 3248–3264.
- 490 15. Cichocki, A.; Mandic, D.; De Lathauwer, L.; Zhou, G.; Zhao, Q.; Caiafa, C.; Phan, H. Tensor Decompositions for Signal Processing Applications: From two-way to multiway component analysis. *IEEE Signal Process. Mag.* **2015**, *32*, 145–163.
- 493 16. Jolliffe, I.T. *Principal Component Analysis*, 2nd ed.; Springer Verlag: New York, NY, USA, 2002.
- 494 17. Kolda, T.; Bader, B. Tensor Decompositions and Applications. *SIAM Rev.* **2009**, *51*, 455–500.
- 495 18. Lopez, J.; Santos, S.; Torres, D.; Atzberger, C. Convolutional Neural Networks for Semantic Segmentation of Multispectral Remote Sensing Images. In Proceedings of the LATINCOM, Guadalajara, Mexico, 14–16 November 2018.
- 498 19. European Space Agency. Available online: <https://sentinel.esa.int/web/sentinel/missions/sentinel-2> (accessed on 15 July 2019).
- 500 20. Kemker, R.; Kanan, C. Deep Neural Networks for Semantic Segmentation of Multispectral Remote Sensing Imagery. *arXiv* **2017**, arXiv:abs/1703.06452.
- 502 21. Hamida, A.; Benoît, A.; Lambert, P.; Klein, L.; Amar, C.; Audebert, N.; Lefèvre, S. Deep learning for semantic segmentation of remote sensing images with rich spectral content. In Proceedings of the IGARSS, Fort Worth, TX, USA, 23–28 July 2017.
- 505 22. Wang, Q.; Lin, J.; Yuan, Y. Salient Band Selection for Hyperspectral Image Classification via Manifold Ranking. *IEEE Trans. Neural Netw. Learn. Syst.* **2016**, *27*, 1279–1289.
- 507 23. Li, S.; Qiu, J.; Yang, X.; Liu, H.; Wan, D.; Zhu, Y. A novel approach to hyperspectral band selection based on spectral shape similarity analysis and fast branch and bound search. *Eng. Appl. Artif. Intell.* **2014**, *27*, 241–250.

- 509 24. Zhang, L.; Zhang, L.; Tao, D.; Huang, X.; Du, B. Compression of hyperspectral remote sensing images by
510 tensor approach. *Neurocomputing* **2015**, *147*, 358–363.
- 511 25. Astrid, M.; Lee, Seung-Ik. CP-decomposition with Tensor Power Method for Convolutional Neural Networks
512 compression. In Proceedings of the BigComp, Jeju, Korea, 13–16 February 2017.
- 513 26. Chien, J.; Bao, Y. Tensor-factorized neural networks. *IEEE Trans. Neural Networks Learn. Syst.* **2018**, *29*, 1998–2011.
- 514 27. An, J.; Lei, J.; Song, Y.; Zhang, X.; Guo J. Tensor Based Multiscale Low Rank Decomposition for Hyperspectral
515 Images Dimensionality Reductio. *Remote Sens.* **2019**, *11*, 1485.
- 516 28. Li, J.; Liu, Z. Multispectral Transforms Using Convolution Neural Networks for Remote Sensing Multispectral
517 Image Compression. *Remote Sens.* **2019**, *11*, 759.
- 518 29. An, J.; Song, Y.; Guo, Y.; Ma, X.; Zhang, X. Tensor Discriminant Analysis via Compact Feature Representation
519 for Hyperspectral Images Dimensionality Reduction. *Remote Sens.* **2019**, *11*, 1822.
- 520 30. Absil, P.-A.; Mahony, R.; Sepulchre, R. *Optimization Algorithms on Matrix Manifolds*, 1st ed.; Princeton
521 University Press: Princeton, NJ, USA, 2007.
- 522 31. De Lathauwer, L.; De Moor, B.; Vandewalle, J. On the best rank-1 and rank-(R 1 ,R 2 ,…,R N) approximation
523 of higher-order tensors. *SIAM J. Matrix Anal. Appl.* **2000**, *21*, 1324–1342.
- 524 32. Goodfellow, I.; Bengio, Y.; Courville, A. *Deep Learning*, 1st ed.; MIT Press, 2016.
- 525 33. Sheehan, B. N.; Saad, Y. Higher Order Orthogonal Iteration of Tensors (HOOI) and its Relation to PCA and
526 GLRAM. In Proceedings of the 7th SIAM International Conference on Data Mining, Minneapolis, MN, USA,
527 26–28 April 2007.
- 528 34. Badrinarayanan, V.; Kendall, A.; Cipolla, R. SegNet: A Deep Convolutional Encoder-Decoder Architecture
529 for Image Segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **2017**, *39*, 2481–2495.
- 530 35. De Lathauwer, L.; De Moor, B.; Vandewalle, J. A Multilinear Singular Value Decomposition. *SIAM J. Matrix
531 Anal. Appl.* **2000**, *21*, 1253–1278.
- 532 36. Rodes, I.; Inglada, J.; Hagolle, O.; Dejoux, J.; Dedieu, G. Sampling strategies for unsupervised classification
533 of multitemporal high resolution optical images over very large areas. In Proceedings of the 2012 IEEE
534 International Geoscience and Remote Sensing Symposium, Munich, Germany, 22–27 July 2012.

535 **Sample Availability:** Samples of the compounds are available from the authors.

536 © 2020 by the authors. Submitted to *Remote Sens.* for possible open access publication
537 under the terms and conditions of the Creative Commons Attribution (CC BY) license
538 (<http://creativecommons.org/licenses/by/4.0/>).