

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 6, 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

Big-data compression—Reducing the dimensionality of input data for machine learning algorithms has been one of the most active research areas in recent years [1]. The insertion of tensor-based algorithms for this sort of tasks drove a revolution 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], mineralogy [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 certain image processing tasks [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 by response at

31 different wavelengths, has driven a growth in accuracy and precision in tasks such as classification
32 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 an alternative solution to the problem described previously. To reduce
47 processing times in supervised classification algorithms, such as deep neural networks, we propose a
48 model that, from the benefits offered by tensor decomposition algorithms, aids compression of spectral
49 images. This produces a lower dimensional tensor while preserving the structural and numerical
50 nature of the original data.

51 1.1. State of the art

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

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

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

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

74 Table 1 summarizes some of the most cited related papers, which deal with the
75 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	Spectral	CNN

76 1.2. Contribution

77 In the state of the art we can find a variety of frameworks looking for dimensionality reduction
 78 of images of any type, that is, RGB [?], medical [?], SAR [?], multispectral [?], among others.
 79 In particular, the large number of spectral bands in hyperspectral images has focused efforts in this
 80 category [?]. Spatial [?], spectral [?] and spatial-spectral [?] compression have been achieved with
 81 matrix and tensor decompositions. It is important to highlight that, a decomposition as pre-processing
 82 for a classification algorithm, instead of favoring, could be counterproductive and greatly affect its
 83 performance. It is therefore important to make a proper selection of the type of decomposition that
 84 would produce a suitable data set for post-processing.

85 Unlike previous works, this work seeks to adapt the data in the best way to be the input of deep
 86 convolutional networks. Convolutional network models are designed to extract and interpret all the
 87 spatial properties of an image by moving the kernels over the input data []. Therefore, producing
 88 uncorrelated data in space and spectrum, would make harder the interpretation of the data in the
 89 convolutional network [?]. Thus, the proposed framework maintains the positive tensor nature of
 90 the spectral images and the spatial dimensionality to preserver spatial-spectral correlation of the data
 91 while reducing the spectral dimensionality, in order to decrease computational load in the pixel-wise
 92 classification process.

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

- 94 1. The framework NTKD3-CNN proposed in this work, develops a new strategy to improve
 95 accuracy results of semantic segmentation convolutional neural networks by finding suitable
 96 tensor data, preserving spatial correlation and values in the set of the natural numbers while
 97 compressing the spectral domain and in turn decreasing computational load.
- 98 2. This work also presents an exhaustive performance analysis measuring and comparing its
 99 efficiency with the most popular metrics, i.e., as pixel accuracy (PA), also PA in function of the
 100 number of new tensor bands, precision, recall, F1, orthogonality degree of the factor matrices
 101 and the core tensor, reconstruction error of the original tensor, and execution time.

102 The remainder of this work is organized as follows. Section ?? introduces tensor algebra notation
 103 and basic concepts to familiarize the reader with the symbology used in this paper. Section ?? presents
 104 the problem statement of this work and the mathematical definition. In Section 4, CNN theory is
 105 described for classification and semantic segmentation. Section ?? presents the framework proposed
 106 for compression and semantic segmentation of spectral images. Experimental results are presented in
 107 Section ???. Finally, Sections ?? and 8 present a discussion and conclusions based on the results obtained
 108 in the experiments.

109 2. Tensor-Based Factorizations

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

115 We can define a tensor as a multi-way or multidimensional array. The order of a tensor is the
 116 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
 117 array, which elements x_{i_1, i_2, \dots, i_N} are indexed by $i_n \in 1, 2, \dots, I_n$ for $1 \leq n \leq N$.

118 Throughout this paper, the mathematical notation used by Kolda et al. [17] has been adopted.
 119 Table 2 summarize this notation.

Table 2. Tensor algebra notation summary

$\mathbf{A}, \mathbf{A}, \mathbf{a}, a$	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
$\mathbf{a}_{:i_2:i_3}, \mathbf{a}_{i_1::i_3},$ and $\mathbf{a}_{i_1:i_2::}$	Column, row and tube fibers of a third order tensor
$\mathbf{A}_{i_1::}, \mathbf{A}_{:i_2::}, \mathbf{A}_{::i_3}$	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 .

120 2.1. Basic concepts

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

123 2.1.1. Matricization

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

126 2.1.2. Outer Product

127 The outer product of N vectors $\mathbf{X} = \mathbf{a}^{(1)} \circ \dots \circ \mathbf{a}^{(N)}$ produces a tensor $\mathbf{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$
 128 where \circ denotes the outer product and $\mathbf{a}^{(n)}$ denotes a vector in a sequence of N vectors
 129 and each element of the tensor is the product of the corresponding vector elements; i.e.,
 130 $x_{i_1 i_2 \dots i_N} = a_{i_1}^{(1)} \dots a_{i_N}^{(N)}$.

131 2.1.3. Inner Product

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

134 2.1.4. N -Mode Product

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

137 2.1.5. Rank-One Tensor

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

140 2.1.6. Rank-R Tensor

141 The rank of a tensor \mathfrak{X} is the smallest number of components in a CPD; i.e., the smallest
 142 number of rank-one tensors that generate \mathfrak{X} as their sum [17].

143 2.1.7. N-Rank

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

147 2.2. *Nonnegative* Tucker Decomposition (~~NTKD~~~~TKD~~)*The TKD, also called the Tucker3*

148 *The TKD*[17], for the particular case of third-order tensors, *or best rank (J_1, J_2, J_3) approximation*,
 149 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
 150 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
 151 factor *or loading* 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 perform the following
 152 approximate decomposition:

$$\mathfrak{X} = \mathfrak{G} \times_1 \mathbf{U}^{(1)} \times_{\underline{N_2}} \mathbf{U}^{(N)(2)} \times_3 \mathbf{U}^{(3)} + \mathcal{E} \quad (1)$$

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

158 We can also denote the TKD using the matricization approach and *expressed* it by

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

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

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

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

161 Starting from (1), the reconstruction of an approximated tensor can be given by

$$\hat{\mathfrak{X}} = \mathfrak{G} \times_1 \mathbf{U}^{(1)} \times_{\underline{N_2}} \mathbf{U}^{(N)(2)} \times_3 \mathbf{U}^{(3)} \quad (3)$$

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

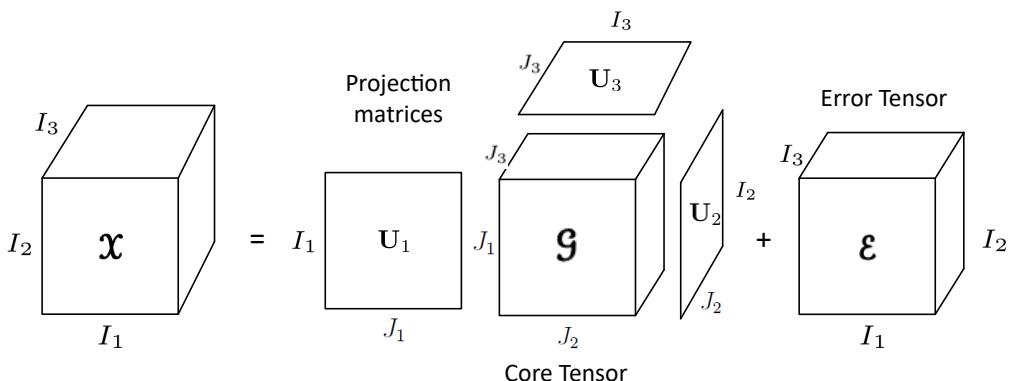


Figure 1. Tucker decomposition for a third-order tensor.

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

$$\mathbf{G} = \mathbf{X} \times_1 \mathbf{U}^{(1)\top} \times_{\textcolor{red}{N}2} \mathbf{U}^{\textcolor{red}{(N)\top}} \times_{\textcolor{blue}{3}} \mathbf{U}^{(3)\top} \quad (4)$$

164 where $\mathbf{U}^{(n)\top}$ denotes the transpose matrix of $\mathbf{U}^{(n)}$ for
 165 $n = 1, \dots, N$. The reconstruction error ξ can be computed as

$$\xi(\hat{\mathbf{X}}) = \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 \quad (5)$$

166 where $\|\cdot\|_F$ represents the Frobenius norm. To compute the best rank approximation of a tensor,
 167 it is required an iterative algorithm. HOSVD, ALS and HOOI are algorithms commonly used in TKD
 168 [?].

169 HOOI initializes the factors matrices using HOSVD and assumes that orthogonal matrices are
 170 known, so that the core tensor is obtained with (4). Then, it maximizes the cost function

$$\max_{\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}} \|\mathbf{X} \times_1 \mathbf{U}^{(1)\top} \times_2 \mathbf{U}^{(2)\top} \cdots \times_{\textcolor{red}{N}3} \mathbf{U}^{\textcolor{red}{(N)\top}}\|_F^2 \quad (6)$$

171 with $\mathbf{U}^{(n)}$ unknown. Fixing all factor matrices but one, tensor \mathbf{X} can be projected onto the
 172 $\{R_1, \dots, R_{n-1}, R_{n+1}, \dots, R_N\}$ -dimensional space as

$$\mathbf{W}^{(-n)} = \mathbf{X} \times_1 \mathbf{U}^{(1)\top} \cdots \times_{n-1} \mathbf{U}^{(n-1)\top} \times_{n+1} \mathbf{U}^{(n+1)\top} \cdots \times_N \mathbf{U}^{(N)\top} \quad (7)$$

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

176 2.2.1. Non-negative Tucker Decomposition (NTD)

177 The NTD is a Tucker decomposition with nonnegativity constraints, which develops a new tensor
 178 factorization method [1]. Thus, for the third-order case, the NTD can be formulated as follows. Given a
 179 third-order tensor $\mathbf{X} \in \mathbb{R}_+^{I_1 \times I_2 \times I_3}$ find a core tensor $\mathbf{G} \in \mathbb{R}_+^{I_1 \times I_2 \times I_3}$ and the factor matrices $\mathbf{U}_1 \in \mathbb{R}_+^{I_1 \times J_1}$,
 180 $\mathbf{U}_2 \in \mathbb{R}_+^{I_2 \times J_2}$ and $\mathbf{U}_3 \in \mathbb{R}_+^{I_3 \times J_3}$ which performs the approximation given in Eq. (1). The TKD method
 181 is similarly followed for the NTD.

182 The HOOI algorithm evidently can be applied for the particular case of third-order tensor.
 183 Furthermore, it can be slightly adjusted for preserving any input dimension and develop the
 184 decomposition in a specific order. This will be explained in next chapters Algorithm 1 shows the

185 HOOI algorithm for a NTD.

186

Algorithm 1: HOOI algorithm to compute a rank- (R_1, \dots, R_N) ~~TKD-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: $\mathbf{U}^{(n)} \in \mathbb{R}_+^{I_n \times J_n}, \mathbf{G} \in \mathbb{R}_+^{J_1 \times J_2 \times \dots \times J_N}$

188 3. Problem phenomenology**189 3.1. Spectral Imagery**

190 Multi- or Hyper-spectral images are by nature multidimensional nonnegative arrays. A spectral
 191 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 I_1, I_2 and I_3 represent
 192 its height, width and spectral bands respectively. In RS image processing, spectral images are frequently
 193 used for classification of different material in a scene of interest. However, due to the low spatial
 194 resolution produced by the distance between the sensor and the target, spatial features are not sufficient
 195 to discern certain classes. That is why spectral resolution plays an important role in this type of task.

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

200 3.2. Problem Statement

201 Given a Spectral Image as a third-order tensor $\mathbf{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, and its corresponding pixelwise
 202 classification ground truth matrix $\mathbf{Y} \in \mathbb{R}^{I_1 \times I_2}$ for a specific number of classes C , find, through
 203 Non-negative Tensor Decomposition, the best rank- (R_1, R_2, R_n) approximation and its respective
 204 core tensor $\mathbf{G} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, where $R_1 = I_1, R_2 = I_2$ and $R_3 = I_3 << I_3$, to be the input of a pixel-wise
 205 classification Convolutional Neural Network and produce an output matrix $\hat{\mathbf{Y}}$ of predicted classes,
 206 achieving competitive performance metrics for pixel-wise classification while decreasing computational
 207 load in the classification process.

208 3.3. Mathematical Definition

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

$$\begin{aligned}
& \min_{\mathbf{g}, \mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}} \|\mathbf{x} - \mathbf{g} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{U}^{(3)}\|_F^2 \\
\text{subject to} \quad & \mathbf{U}^{(n)} \in \mathbb{R}_+^{I_n \times J_n} \quad \text{for } n = 1, 2, 3 \quad \text{and} \quad \mathbf{g} \in \mathbb{R}_+^{J_1 \times J_2 \times J_3} \\
& J_1 = I_1, J_2 = I_2 \quad \text{no compression in the spatial domain,} \\
& J_3 \ll I_3 \quad \text{reduced spectral domain at the core tensor,} \\
& \xi(\hat{\mathbf{x}}) \rightarrow 0 \quad \text{and} \quad \text{PA} \geq \psi \quad \text{competitive pixel accuracy up to a threshold}
\end{aligned} \tag{8}$$

210 4. 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

$$\hat{\mathbf{y}} = \sigma(\mathbf{W}\mathbf{g} + \mathbf{b}) \tag{9}$$

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 (9) using tensor algebra notation as

$$\hat{\mathbf{y}} = \sigma(\mathbf{W}\mathbf{g}) \tag{10}$$

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 2). 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

$$m_{i_1 i_2 f_2} = \sigma \left[\sum_{k_1=1}^{K_1} \sum_{k_2=1}^{K_2} \sum_{f_1=1}^{F_1} w_{k_1, k_2, f_1} g_{i_1+k_1-o_1, i_2+k_2-o_2, f_1} \right] \tag{11}$$

211 where $m_{i_1 i_2 f_2}$ denotes the value of the output activation map f_2 at position i_1, i_2 ; σ represents the
212 activation function; and o_1 and o_2 are offsets in spatial dimensions which depend on the kernel size,

²¹³ and equal $\frac{K_1+1}{2}$ and $\frac{K_2+1}{2}$ respectively (see Figure 2).

²¹⁴

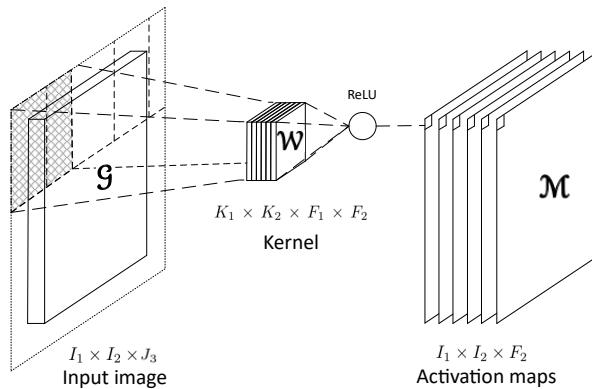


Figure 2. Convolutional layer with a $K_1 \times K_2 \times F_1 \times F_2$ kernel. Input channels F_1 must equal the spectral bands J_3 . To preserve original dimensions at the output, zero padding is needed [18]. Output dimensions also depend on stride $S = 1$ to consider every piece of pixel information 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.

$$J(\mathbf{W}) = L(\mathbf{W}) + R(\mathbf{W}), \quad (12)$$

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.

$$\mathbf{W}' = \mathbf{W} - \phi \nabla_{\mathbf{W}} J(\mathbf{W}), \quad (13)$$

²¹⁵ where \mathbf{W}' represents the synaptic weights tensor in next iteration during training, ϕ denotes the
²¹⁶ learning rate parameter, and $\nabla_{\mathbf{W}} J(\mathbf{W})$ the cost function gradient. Gradient descent converges when
²¹⁷ every element of the gradient is zero, or in practice, very close to zero [10].

²¹⁸ CNNs has been successfully used in many image classification frameworks. This variation in
²¹⁹ architecture from other typical ANN models yields the network to learn spatial and spectral features,
²²⁰ which are highly profitable for image classification. Besides, FCNs, constructed with only convolutional
²²¹ layers are able to classify each element of the input image; i.e., they yield pixel-wise classification, or in
²²² other words, semantic segmentation.

²²³ 5. NTD for DCNNs Methodology

²²⁴ The following subsections described the methodology followed for the framework propose
²²⁵ in this work. We can summarize the big picture in three steps: the HSI modeling, the tensor
²²⁶ decomposition, the classifier and the decision making.

227 5.1. Tensor modeling

228 Consider an input dataset $\mathbf{X} \in \mathbb{N}^{I_1 \times I_2 \times I_3}$ with $I_1 \times I_2 \times I_3$ samples in the space of the natural
 229 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
 230 by the Linear Mixing Model (LMM) as follows

$$\mathbf{x}_{i_1 i_2} = \sum_{c=1}^C (\alpha_{i_1 i_2 c} \mathbf{m}_c + \boldsymbol{\eta}) \quad (14)$$

231 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
 232 c , and $\boldsymbol{\eta}$ represents an additive noise vector. The abundance vectors $\alpha_{i_1 i_2}$ must always satisfy two
 233 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,
 234 $\sum_{c=1}^C \alpha_{i_1 i_2 c} = 1$. Figure 3a shows the spectral signatures for the Indian Pines dataset.

235 Let

236 5.2. Tensor factorization

237 Consider $\mathbf{Y} \in \mathbb{C}^{I_1 \times I_2}$ be as the matrix of actual classes corresponding to our dataset \mathbf{X} , and
 238 $\hat{\mathbf{Y}} \in \mathbb{C}^{I_1 \times I_2}$ as the prediction matrix, where \mathbb{C} defines the set of C different classes, and $\hat{\mathbf{Y}} \in \mathbb{C}^{I_1 \times I_2}$
 239 be the prediction matrix. In order to reduce data dimensionality of the input dataset \mathbf{X} , we while
 240 keeping classifier performance, we propose to use the restricted NTD denoted as \mathcal{T} , producing a core
 241 tensor $\mathbf{G} \in \mathbb{N}^{I_1 \times I_2 \times I_3}$, $\mathbf{G} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ and n factors matrices $\mathbf{U}^{(n)}$, i.e., $\mathbf{X} \xrightarrow{\mathcal{T}} (\mathbf{G}, \mathbf{U}^{(n)})$, where the core
 242 tensor $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times J_n}$ expressed as

$$\mathbf{X} \xrightarrow{\mathcal{T}} (\mathbf{G}, \mathbf{U}^{(n)}) \quad (15)$$

243 where the decomposition is restricted to to preserve the spatial dimensions domain and to be in the
 244 set of the natural numbers only in the 3rd-mode by the Tucker1 model

$$\mathbf{X} = \mathbf{G} \times_1 \mathbf{I} \times_2 \mathbf{I} \times_3 \mathbf{U}^{(3)} \quad (16)$$

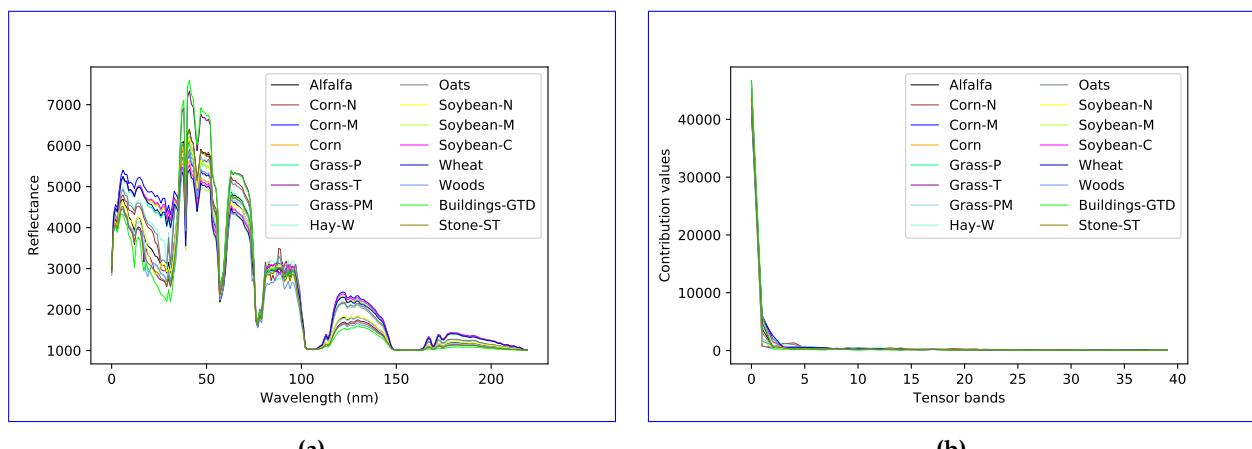


Figure 3. Behavior of the 16 classes of the Indian Pines dataset, a) in the spectral domain (spectral signatures) and b) in the tensor bands domain.

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 [Behavior of the 16 classes of the Indian Pines dataset, a\) in the spectral domain \(spectral signatures\) and b\) in the the tensor bands domain.](#)

$$\mathbf{g}_{i_1 i_2} = \sum_{c=1}^C (\beta_{i_1 i_2 c} \mathbf{s}_c + \boldsymbol{\eta}) \quad (17)$$

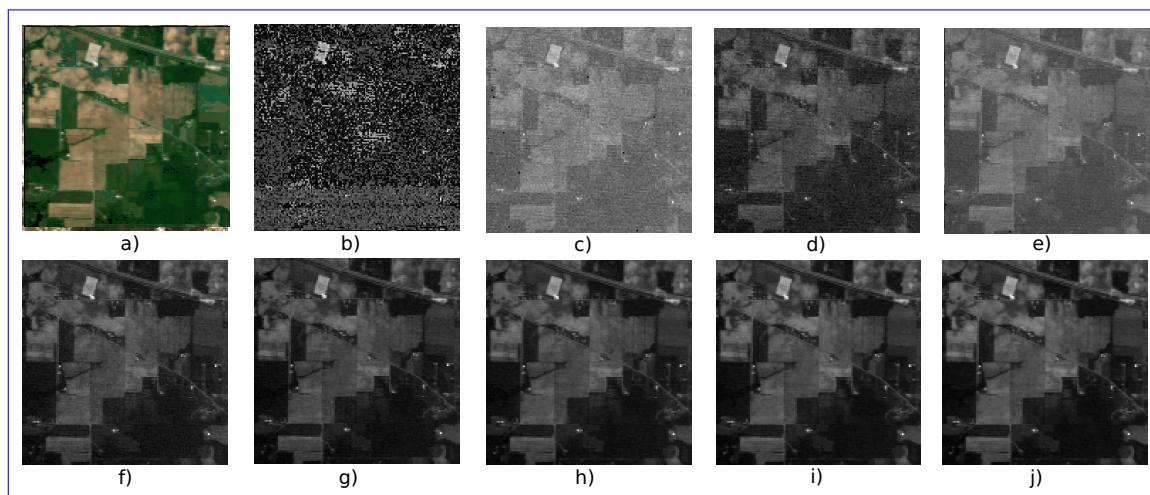
where β_c is the contribution of material c at pixel $i_1 i_2$ and \mathbf{s}_c denotes the endmember of a specific material c . We can see in Figure 3b the new tensor bands values for each class. [The](#)

[We also propose a variant of the NTD that does an integer decomposition, i. e., the Integer Non-negative Tucker decomposition \(INTD\). The INTD follows the same Tucker model described in Section 2.2. It considers the additional restriction of decomposing a tensor in the set of the natural numbers.](#)

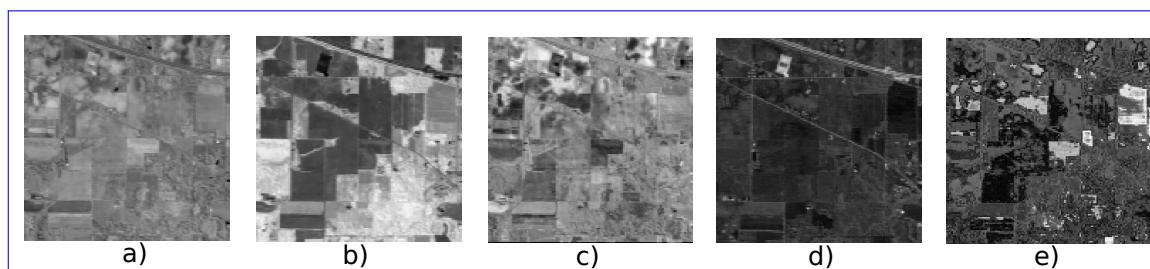
5.3. Classifier

[The tensor decompositions generates a core tensor where the first tensor bands provide the most of the information to differentiate the classes of interest from the input dataset. Then, the core tensor \$\mathbf{G} \in \mathbb{N}^{I_1 \times I_2 \times J_3}\$, with \$J_3 < I_3\$, and its corresponding ground truth \$\mathbf{Y}\$ form the input tuple of the classifier \$\Theta\$, which produce a predicted label for each element of the input, i.e.,](#)

$$(\mathbf{G}, \mathbf{Y}) \xrightarrow{\Theta} \hat{\mathbf{Y}} \quad (18)$$



[Figure 4. Indian Pines dataset a\) True color image b\) - j\) 1st to 9th spectral band.](#)



[Figure 5. TKD tensor bands of the Indian Pines dataset a\) - e\) 1st to 5th tensor band.](#)

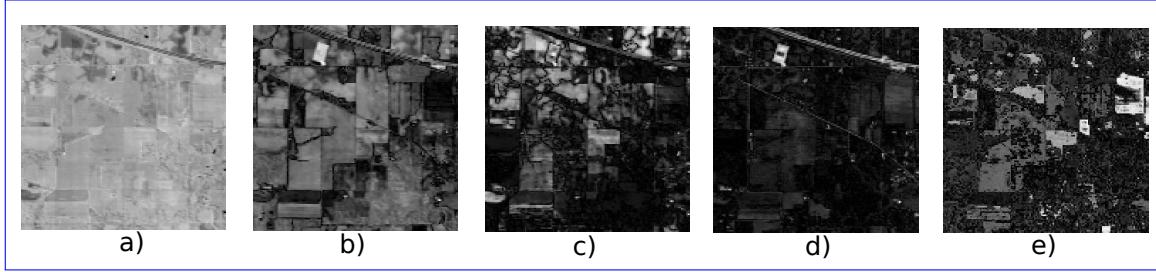


Figure 6. NTD tensor bands of the Indian Pines dataset a) - e) 1st to 5th tensor band.

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

$$J(\mathbf{W}) = -\mathbb{E}_{\mathbf{g}, \mathbf{Y} \sim p} \log p(\mathbf{Y} | \mathbf{g}) \quad (19)$$

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

$$H(y, p) = - \sum_{c=1}^C y_c \log(p_c) \quad (20)$$

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

$$\delta(\mathbf{z})_c = \frac{e^{z_c}}{\sum_{l=1}^L e^{z_l}} \quad (21)$$

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. 14.

In this paper, we aim to feed supervised classifiers based on DCNN3D-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 7 shows the big picture of the framework proposed.

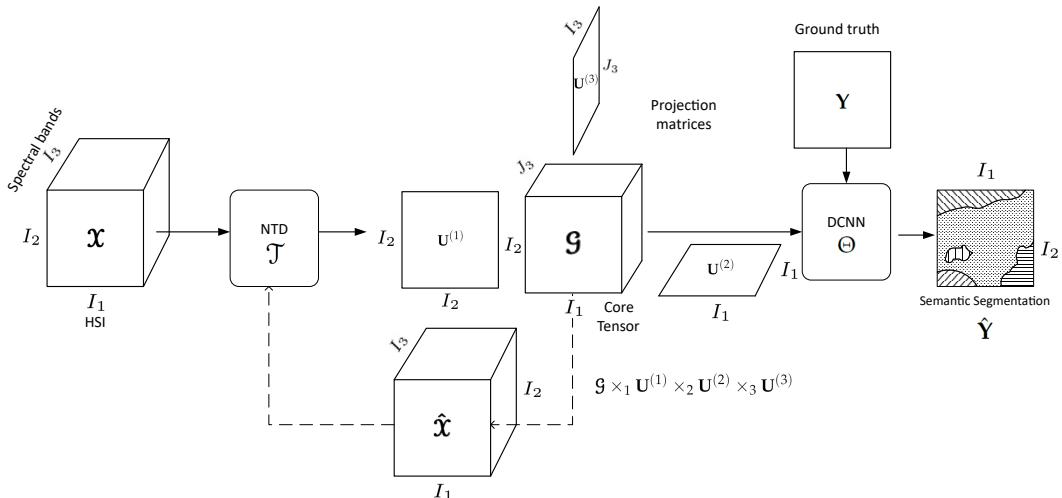


Figure 7. Big picture of the framework proposed.

272 5.4. *Decomposition analysis*

273 Our proposed framework faces two main challenges. The first is the selection of the decomposition
 274 method. Since we are looking for the best representation of the input dataset for a 3D-CNN, so,
 275 we limit the set of methods to those that produce a decomposition tensor with the same structure as
 276 the input tensor. TKD and NTD, NTD and the INTD proposed generate a core tensor with the desired
 277 properties. We also propose a variant of the NTD that does an integer decomposition, i.e., the integer
 278 non-negative Tucker decomposition (INTD).

279 The second challenge is the search for the rank – (J_1, J_2, J_3). As we want to preserve the spatial
 280 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
 281 decrease considerably. Both decisions are made under a criterion from under the probabilistic point of
 282 view.

283 We use the The Kullback-Leibler divergence as a distance metric that quantifies is a metric for
 284 quantifying the difference between two probability distributions as

$$D_{KL}(P\|Q) = \sum_{i=1}^I p_i(x) \log \frac{p_i(x)}{q_i(x)}$$

285 the probability distributions of the original MSI X and G represent the probability distribution of the
 286 input data and the core tensor respectively, as discrete random variables, $KL(X\|G)$ G as

$$D_{KL}(X\|G) = \sum_{i=1}^I p_i(x) \log \frac{p_i(x)}{g_i(x)} \quad (22)$$

287 where $D_{KL}(X\|G)$ represents the KL divergence of the two probability distributions. We also use a the
 288 Figure 8a shows the results of this metric for the three decomposition used in this work.

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

$$D_{JS}(X\|G) = \frac{1}{2} D_{KL}(X\|M) + \frac{1}{2} D_{KL}(G\|M) \quad (23)$$

292 where $D_{JS}(X\|M)$ $D_{JS}(X\|G)$ represents the JS divergence of two-the probability distributions X and
 293 G. Figure 8b shows the JS divergence for the TKD, NTD and INTD.

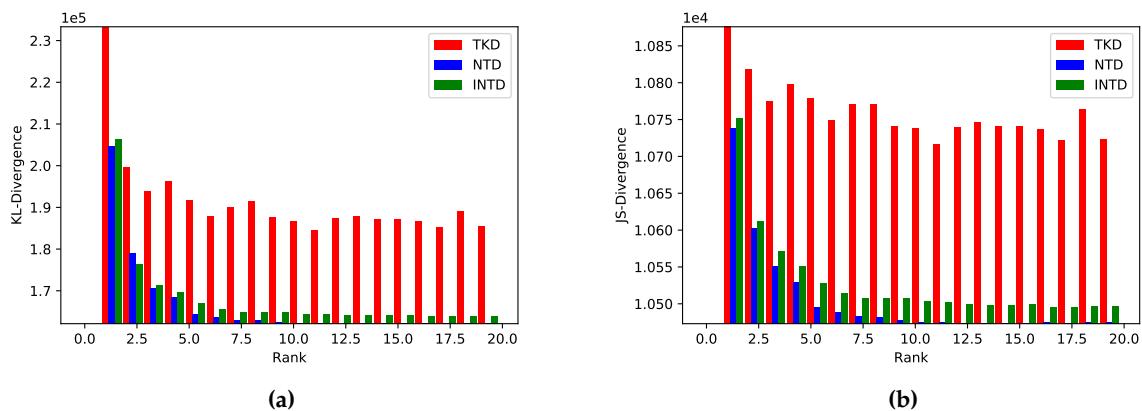


Figure 8. Divergences between the input dataset and the core tensor generated by TKD, NTD and INTD, a) KL-divergence and b) JS-divergence.

294 6. Experimental Results

295 6.1. Input Data

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

298 6.1.1. Sentinel-2 [CNNMSI](#)

299 This dataset developed by Lopez et al. [?] is composed of RS Sentinel-2 scenarios from central
 300 Europe. It has 100 scenarios for the training space and 10 scenarios for testing, all of them with
 301 128×128 pixels with spatial resolution of $20m^2$ and 9 spectral bands in the range $490 - 2190nm$.
 302 The labels are semi-manually assigned for five classes of interest: vegetation, soil, water, clouds and
 303 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	22.87	30.73
Shadow	2.52	2.58	1.91
Cloud	13.87	14.37	8.83
Water	12.13	11.93	14.10
Vegetation	47.88	48.23	44.41

304 6.1.2. Indian Pines

305 This dataset is a scene produced by AVIRIS in North-western Indiana and consists of 145×145
 306 pixels and 224 spectral bands in the wavelength range $0.4\text{--}2.5\mu m$. The Indian Pines scene contains
 307 two-thirds agriculture, and one-third forest or other natural perennial vegetation. There are two major
 308 dual lane highways, a rail line, as well as some low density housing, other built structures, and smaller
 309 roads. Since the scene is taken in June some of the crops present, corn, soybeans, are in early stages of
 310 growth with less than 5% coverage. The ground truth available is designated into sixteen classes and
 311 is not all mutually exclusive. Indian Pines data are available at [Indian Pines dataset](#).

312 6.1.3. Salinas

313 This scene was collected by the AVIRIS sensor over Salinas Valley, California. It has 512×217
 314 pixels with spatial resolution $3.7m$, and 224 spectral bands. It includes vegetables, bare soils, and
 315 vineyard fields. Salinas groundtruth contains 16 classes.

Table 4. 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
Indian Pines	145×145	224	16	4,709,600
Salinas	512×217	224	16	24,887,296

316 6.2. Metrics

317 6.2.1. Relative Mean Square Error (rMSE)

318 In order to compute the reconstruction error of the tensor \mathfrak{X} for the implementation of HOOI,
 319 the rMSE was used:

$$320 rMSE(\hat{\mathfrak{X}}) = \frac{1}{Q} \sum_{q=1}^Q \frac{\|\hat{\mathfrak{X}}_q - \mathfrak{X}_q\|_F^2}{\|\mathfrak{X}_q\|_F^2}, \quad (24)$$

320 where \mathbf{x}_q represents the q -th CNNMSI from our dataset with Q MSIs and $\hat{\mathbf{x}}_q$ its corresponding
321 reconstruction computed by (3).

322 6.2.2. KL Divergence

323 The Kullback–Leibler divergence is a distance metric that quantifies the difference between two
324 probability distributions. If P and Q represent the probability distribution of a discrete random
325 variable, the Kullback–Leibler divergence is calculated as

$$326 \quad D_{KL}(P\|Q) = \sum_{i=1}^I p_i(x) \log \frac{p_i(x)}{q_i(x)}$$

326 where $D_{KL}(P\|Q)$ represents the KL divergence of two probability distribution of discrete random
327 variables P and Q .

328 6.2.2. JS Divergence

329 The Jensen–Shanon divergence is another method of measuring the similarity between two
330 probability distributions. It is a symmetric version of the KL divergence. Defining $M = \frac{P+Q}{2}$, we
331 can write the JS divergence as

$$332 \quad D_{JS}(P\|Q) = \frac{1}{2}D_{KL}(P\|M) + \frac{1}{2}D_{KL}(Q\|M)$$

332 where $D_{JS}(P\|Q)$ represents the JS divergence of two probability distribution of discrete random
333 variables P and Q .

334 6.2.2. Cohen's Kappa Coefficient

335 Cohen's kappa coefficient is a very robust metric used to measure reliability of multi-class and
336 imbalanced class classification algorithms. It is computed by

$$337 \quad \kappa = \frac{\rho_o - \rho_e}{1 - \rho_e} \quad (25)$$

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

340 6.2.3. Pixel Accuracy (PA)

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

$$344 \quad PA = \frac{\sum_{c=1}^C \tau_{cc}}{\sum_{c=1}^C \sum_{d=1}^C \tau_{cd}} = \frac{TP + TN}{TP + TN + FP + FN} \quad (26)$$

344 where c is the number of class for $c = 1, \dots, C$ and τ_{cc} is the amount of pixels of class c correctly
345 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
346 d . Despite this metric is wide used, it is not a totally fair metric for imbalanced classes. In this work we
347 used this metric with comparison purposes. However, we also used metrics more in line with multi
348 class classification tasks.

³⁴⁹ 6.2.4. Precision

³⁵⁰ Another metric used in this work as a performance evaluation metric in multiclass classification
³⁵¹ is precision. This is a metric that measures the percentage of pixels from class c correctly classified. It
³⁵² is computed with the following equation

$$\Psi_c = \frac{\tau_{cc}}{\sum_{d=1}^C \tau_{cd}} = \frac{TP}{TP + FP} \quad (27)$$

³⁵³ where Ψ_c denotes the precision of class c , which is the number of TP divided by the TP plus FP.

³⁵⁴ 6.2.5. Recall or Sensitivity

³⁵⁵ Recall, or also known as sensitivity is a metric that indicates the proportion of pixels classified as
³⁵⁶ class c that actually belong to class c . It is computed with the following equation

$$v_c = \frac{\tau_{cc}}{\sum_{d=1}^C \tau_{dc}} = \frac{TP}{TP + FN} \quad (28)$$

³⁵⁷ where v_c denotes the recall of class c for $c = 1, \dots, C$.

³⁵⁸ 6.2.6. F1 Score

³⁵⁹ In order to summarize precision and recall in one only metric, we use the F1 score, which is
³⁶⁰ computed by

$$F1 = \frac{2\Psi v}{\Psi + v} \quad (29)$$

³⁶¹ This metric provides a very appropriate measure of multiclass classification for imbalanced dataset.

³⁶² 7. Discussion and Comparison

³⁶³ 8. Conclusions

³⁶⁴ 8.1. Figures, Tables and Schemes

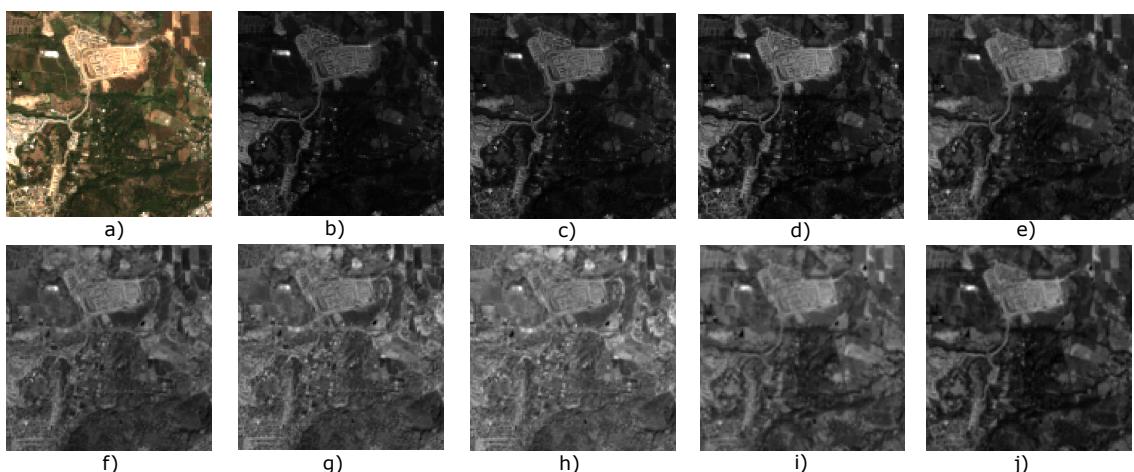


Figure 9. Original Sentinel-2 spectral bands, a) True color image, b) to j) 1st to 9th spectral band respectively.

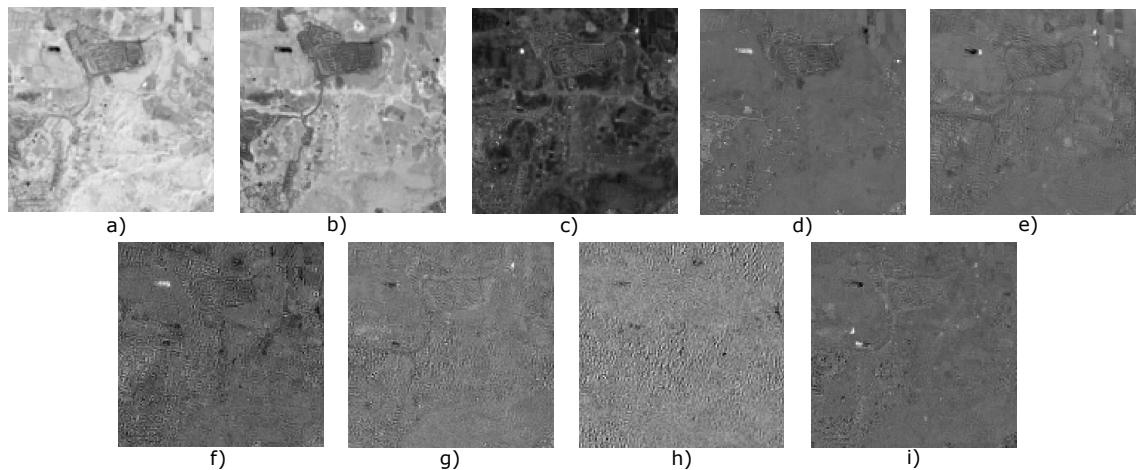


Figure 10. Tensor bands from the Tucker Decomposition, a) to i) 1st to 9th band respectively.

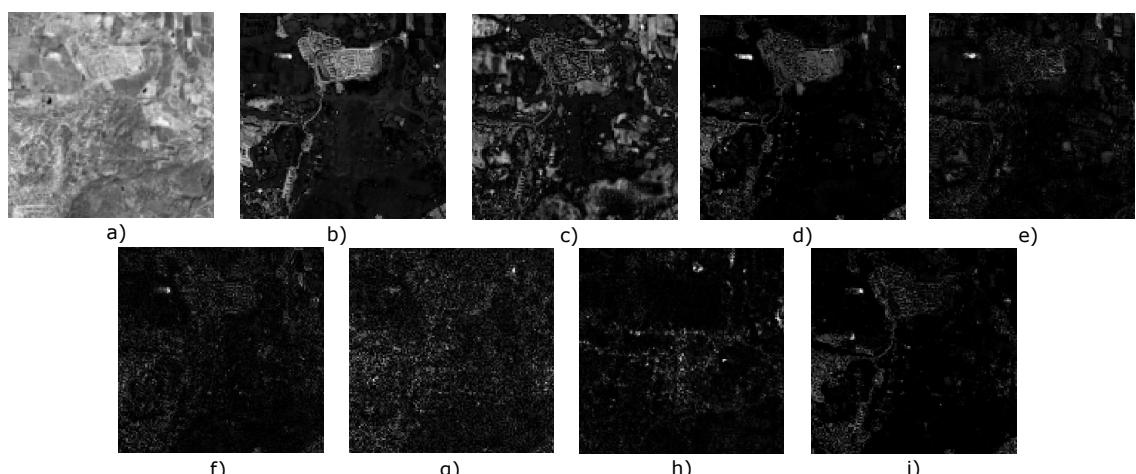


Figure 11. Tensor bands from the Non-negatice Tucker Decomposition, a) to i) 1st to 9th band respectively.

Original Sentinel-2 spectral bands: Tucker Decomposition Tensor bands: Nonnegative Tucker Decomposition Tensor bands:

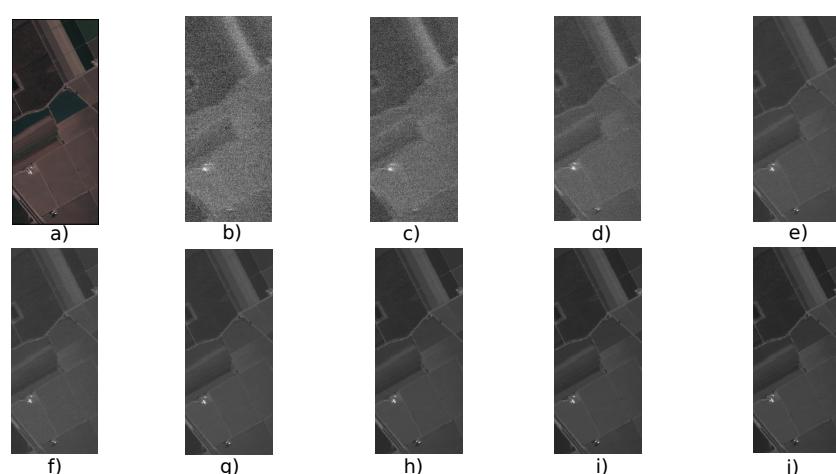


Figure 12. Original Sentinel-2 spectral bands.

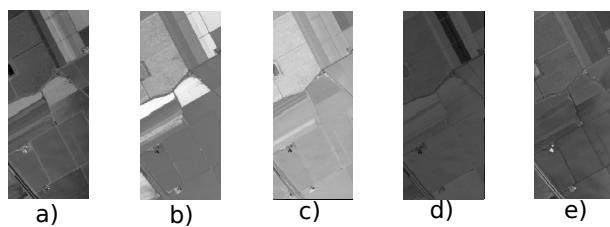


Figure 13. Tucker Decomposition Tensor bands.

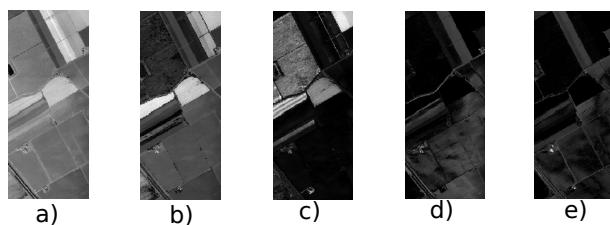


Figure 14. Nonnegative Tucker Decomposition Tensor bands.

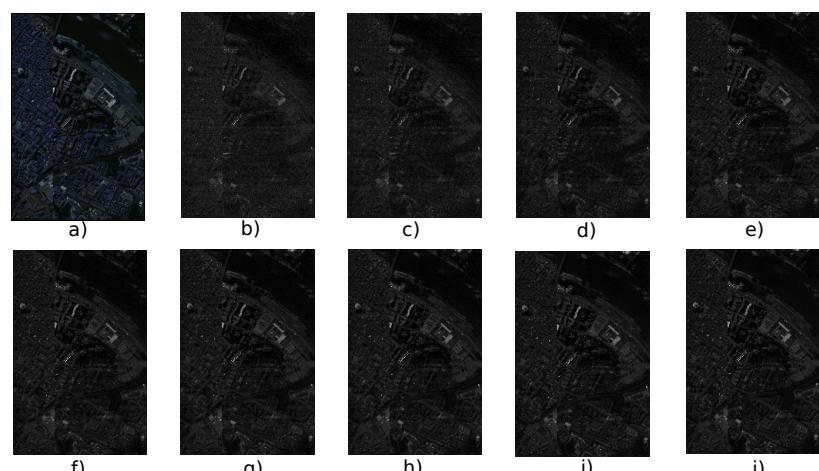


Figure 15. Original Sentinel-2 spectral bands.

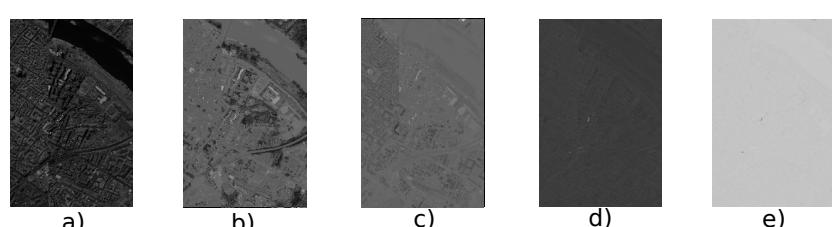


Figure 16. Tucker Decomposition Tensor bands.

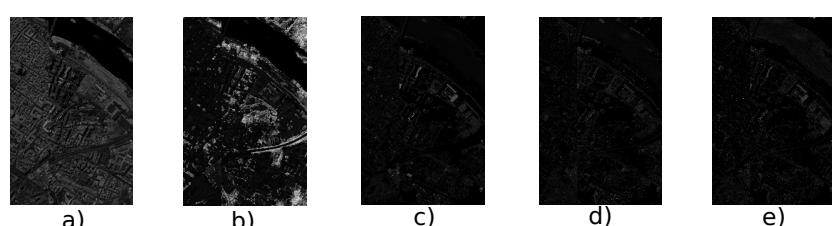
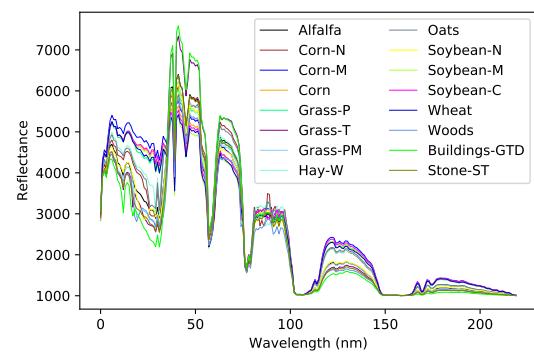
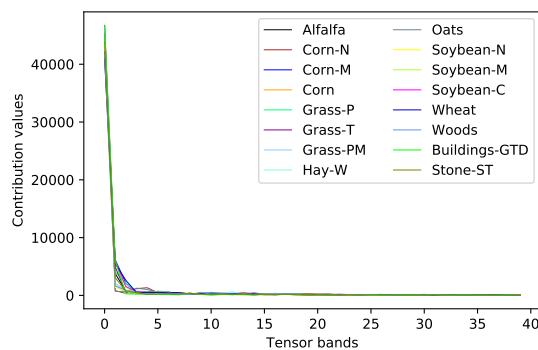


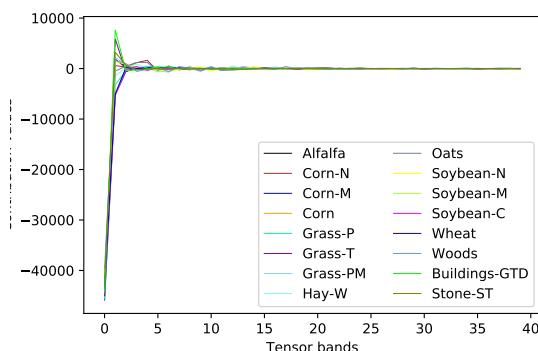
Figure 17. Nonnegative Tucker Decomposition Tensor bands.



(a)



(b)



(c)

Figure 18. Indian Pines Spectral Signatures a) Original, b) NTD and c) TKD

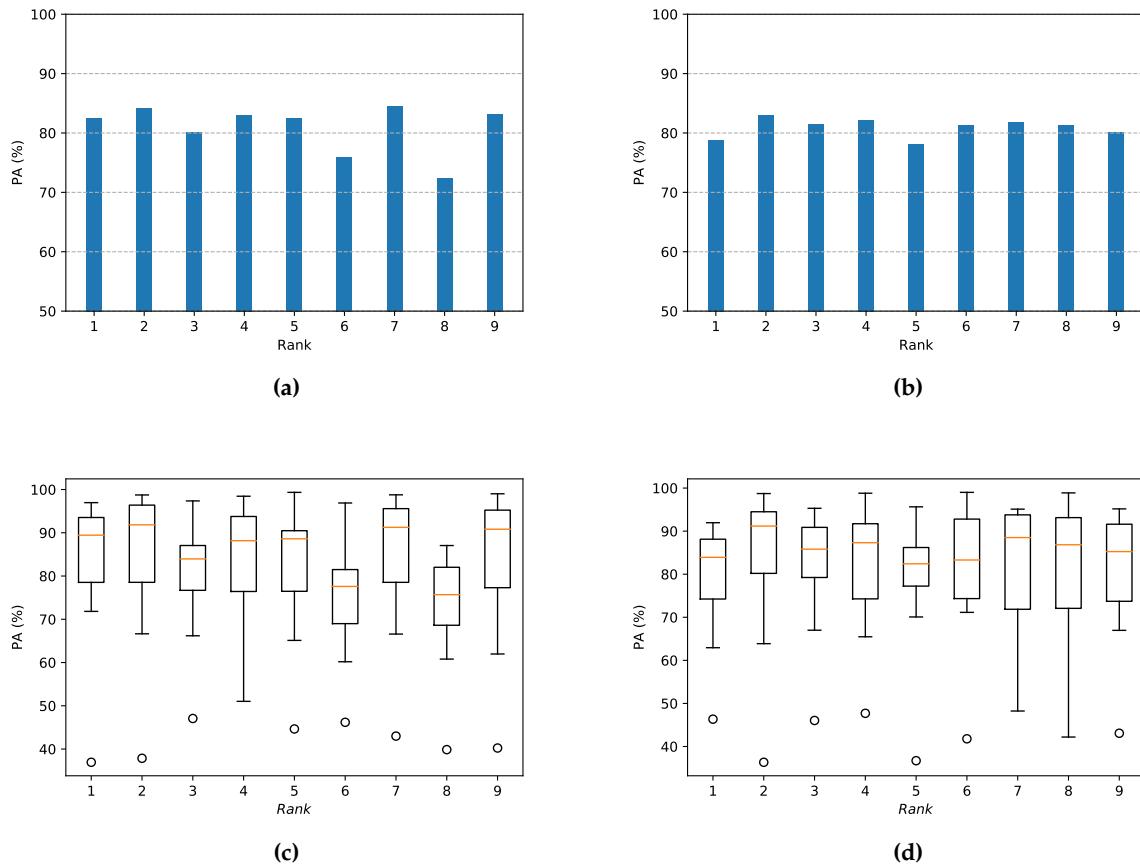


Figure 19. Pixel accuracy vs Rank results a) Comparative bar plot NTKD and TKD, b) Comparative bar plot NTKD and TKD c) Box and whiskers plot for NTKD, and d) Box and whiskers plot for TKD

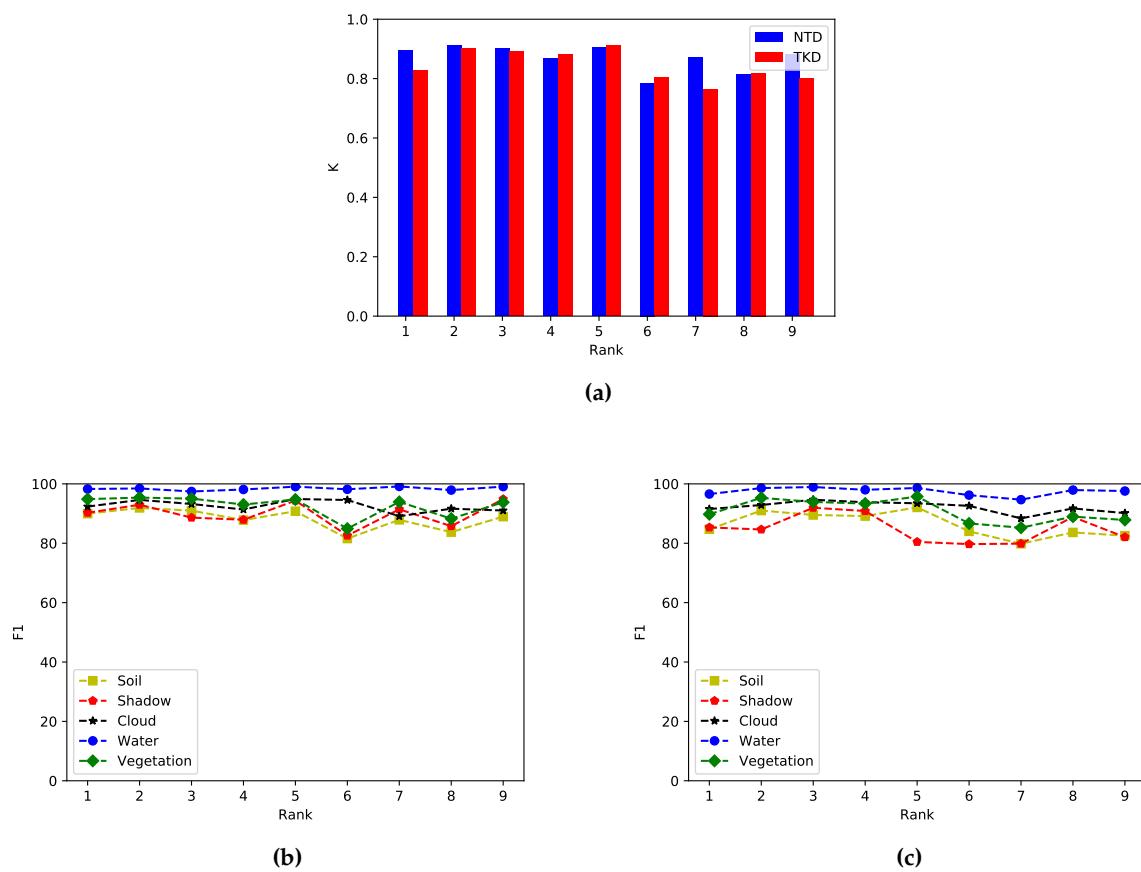
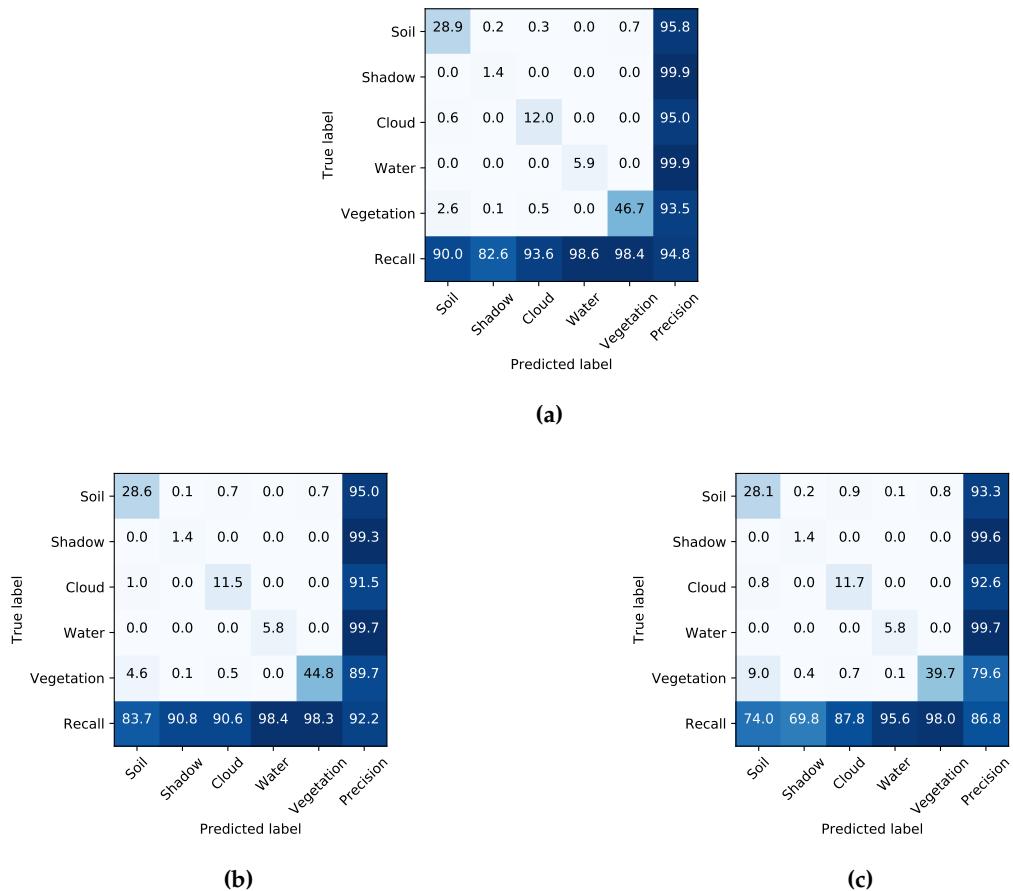


Figure 20. Performance evaluation metrics. a) Kappa score comparison NTD vs TKD, b) F1 for NTD and c) F1 for TKD

**Figure 21.** Confusion matrix for a) Sentinel-2 original dataset, b) NTD, c) TKD.

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

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

372 Acknowledgments:

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

374 Abbreviations

375 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

378 References

- 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.

- 382 2. He, Z.; Hu, J.; Wang, Y. Low-rank tensor learning for classification of hyperspectral image with limited
383 labeled sample. *IEEE Signal Process.* **2017**, *145*, 12–25.
- 384 3. Richards, A.; Xiuping, J.J. Band selection in sentinel-2 satellite for agriculture applications. In *Remote Sensing*
385 *Digital Image Analysis*, 4th ed.; Springer-Verlag: Berlin, Germany, 2006.
- 386 4. Zhang, T.; Su, J.; Liu, C.; Chen, W.; Liu, H.; Liu, G. Band selection in sentinel-2 satellite for agriculture
387 applications. In Proceedings of the 23rd International Conference on Automation & Computing, University
388 of Huddersfield, Huddersfield, UK, 7–8 September 2017.
- 389 5. Xie, Y.; Zhao, X.; Li, L.; Wang, H. Calculating NDVI for Landsat7-ETM data after atmospheric correction
390 using 6S model: A case study in Zhangye city, China. In Proceedings of the 18th International Conference on
391 Geoinformatics, Beijing, China, 18–20 June 2010.
- 392 6. Gao, B. NDWI—A normalized difference water index for remote sensing of vegetation liquid water from
393 space. *Remote Sens. Environ.* **1996**, *58*, 1–6.
- 394 7. Ham, J.; Chen, Y.; Crawford, M.; Ghosh, J. Investigation of the random forest framework for classification of
395 hyperspectral data. *IEEE Trans. Geosci. Remote Sens.* **2005**, *43*, 492–501.
- 396 8. Hearst, Marti A. Support Vector Machines. *IEEE Intell. Syst.* **1998**, *13*, 18–28.
- 397 9. Huang, X.; Zhang, L. An SVM Ensemble Approach Combining Spectral, Structural, and Semantic Features
398 for the Classification of High-Resolution Remotely Sensed Imagery. *IEEE Trans. Geosci. Remote Sens.*
399 **2013**, *51*, 257–272.
- 400 10. Delalieux, S.; Somers, B.; Haest, B.; Spanhove, T.; Vanden Borre, J.; Mucher, S. Heathland conservation
401 status mapping through integration of hyperspectral mixture analysis and decision tree classifiers.
Remote Sens. Environ. **2012**, *126*, 222–231.
- 402 11. Kemker, R.; Salvaggio, C.; Kanan, C. Algorithms for semantic segmentation of multispectral remote sensing
403 imagery using deep learning. *ISPRS J. Photogramm. Remote Sens.* **2018**, *145*, 60–77.
- 404 12. Pirotti, F.; Sunar, F.; Piragnolo, M. Benchmark of machine learning methods for classification of a sentinel-2
405 image. In Proceedings of the XXIII ISPRS Congress, Prague, Czech Republic, 12–19 July 2016.
- 406 13. Mateo-García, G.; Gómez-Chova, L.; Camps-Valls, G. Convolutional neural networks for multispectral image
407 cloud masking. In Proceedings of the IGARSS, Fort Worth, TX, USA, 23–28 July 2017.
- 408 14. Guo, X.; Huang, X.; Zhang, L.; Zhang, L.; Plaza, A.; Benediktsson, J. A. Support Tensor Machines for
409 Classification of Hyperspectral Remote Sensing Imagery. *IEEE Trans. Geosci. Remote Sens.* **2016**, *54*, 3248–3264.
- 410 15. Cichocki, A.; Mandic, D.; De Lathauwer, L.; Zhou, G.; Zhao, Q.; Caiafa, C.; Phan, H. Tensor Decompositions
411 for Signal Processing Applications: From two-way to multiway component analysis. *IEEE Signal Process. Mag.*
412 **2015**, *32*, 145–163.
- 413 16. Jolliffe, I.T. *Principal Component Analysis*, 2nd ed.; Springer Verlag: New York, NY, USA, 2002.
- 414 17. Kolda, T.; Bader, B. Tensor Decompositions and Applications. *SIAM Rev.* **2009**, *51*, 455–500.
- 415 18. Lopez, J.; Santos, S.; Torres, D.; Atzberger, C. Convolutional Neural Networks for Semantic Segmentation
416 of Multispectral Remote Sensing Images. In Proceedings of the LATINCOM, Guadalajara, Mexico,
417 14–16 November 2018.
- 418 19. European Space Agency. Available online: <https://sentinel.esa.int/web/sentinel/missions/sentinel-2>
419 (accessed on 15 July 2019).
- 420 20. Kemker, R.; Kanan, C. Deep Neural Networks for Semantic Segmentation of Multispectral Remote Sensing
421 Imagery. *arXiv* **2017**, arXiv:abs/1703.06452.
- 422 21. Hamida, A.; Benoît, A.; Lambert, P.; Klein, L.; Amar, C.; Audebert, N.; Lefèvre, S. Deep learning for
423 semantic segmentation of remote sensing images with rich spectral content. In Proceedings of the IGARSS,
424 Fort Worth, TX, USA, 23–28 July 2017.
- 425 22. Wang, Q.; Lin, J.; Yuan, Y. Salient Band Selection for Hyperspectral Image Classification via Manifold
426 Ranking. *IEEE Trans. Neural Netw. Learn. Syst.* **2016**, *27*, 1279–1289.
- 427 23. Li, S.; Qiu, J.; Yang, X.; Liu, H.; Wan, D.; Zhu, Y. A novel approach to hyperspectral band selection based on
428 spectral shape similarity analysis and fast branch and bound search. *Eng. Appl. Artif. Intell.* **2014**, *27*, 241–250.
- 429 24. Zhang, L.; Zhang, L.; Tao, D.; Huang, X.; Du, B. Compression of hyperspectral remote sensing images by
430 tensor approach. *Neurocomputing* **2015**, *147*, 358–363.
- 431 25. Astrid, M.; Lee, Seung-Ik. CP-decomposition with Tensor Power Method for Convolutional Neural Networks
432 compression. In Proceedings of the BigComp, Jeju, Korea, 13–16 February 2017.
- 433 26. Chien, J.; Bao, Y. Tensor-factorized neural networks. *IEEE Trans. Neural Networks Learn. Syst.* **2018**, *29*, 1998–2011.
- 434

- 435 27. An, J.; Lei, J.; Song, Y.; Zhang, X.; Guo J. Tensor Based Multiscale Low Rank Decomposition for Hyperspectral
436 Images Dimensionality Reductio. *Remote Sens.* **2019**, *11*, 1485.
- 437 28. Li, J.; Liu, Z. Multispectral Transforms Using Convolution Neural Networks for Remote Sensing Multispectral
438 Image Compression. *Remote Sens.* **2019**, *11*, 759.
- 439 29. An, J.; Song, Y.; Guo, Y.; Ma, X.; Zhang, X. Tensor Discriminant Analysis via Compact Feature Representation
440 for Hyperspectral Images Dimensionality Reduction. *Remote Sens.* **2019**, *11*, 1822.
- 441 30. Absil, P.-A.; Mahony, R.; Sepulchre, R. *Optimization Algorithms on Matrix Manifolds*, 1st ed.; Princeton
442 University Press: Princeton, NJ, USA, 2007.
- 443 31. De Lathauwer, L.; De Moor, B.; Vandewalle, J. On the best rank-1 and rank-(R_1, R_2, \dots, R_N) approximation
444 of higher-order tensors. *SIAM J. Matrix Anal. Appl.* **2000**, *21*, 1324–1342.
- 445 32. Goodfellow, I.; Bengio, Y.; Courville, A. *Deep Learning*, 1st ed.; MIT Press, 2016.
- 446 33. Sheehan, B. N.; Saad, Y. Higher Order Orthogonal Iteration of Tensors (HOOI) and its Relation to PCA and
447 GLRAM. In Proceedings of the 7th SIAM International Conference on Data Mining, Minneapolis, MN, USA,
448 26–28 April 2007.
- 449 34. Badrinarayanan, V.; Kendall, A.; Cipolla, R. SegNet: A Deep Convolutional Encoder-Decoder Architecture
450 for Image Segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **2017**, *39*, 2481–2495.
- 451 35. De Lathauwer, L.; De Moor, B.; Vandewalle, J. A Multilinear Singular Value Decomposition. *SIAM J. Matrix
452 Anal. Appl.* **2000**, *21*, 1253–1278.
- 453 36. Rodes, I.; Inglada, J.; Hagolle, O.; Dejoux, J.; Dedieu, G. Sampling strategies for unsupervised classification
454 of multitemporal high resolution optical images over very large areas. In Proceedings of the 2012 IEEE
455 International Geoscience and Remote Sensing Symposium, Munich, Germany, 22–27 July 2012.

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

457 © 2020 by the authors. Submitted to *Remote Sens.* for possible open access publication
458 under the terms and conditions of the Creative Commons Attribution (CC BY) license
459 (<http://creativecommons.org/licenses/by/4.0/>).