

Dimensionality Reduction via Regression in Hyperspectral Imagery

Valero Laparra, Jesús Malo and Gustau Camps-Valls, *Senior Member, IEEE*

Abstract—This paper introduces a new *unsupervised* method for dimensionality reduction via regression (DRR). The algorithm belongs to the family of *invertible transforms* that generalize Principal Component Analysis (PCA) by using curvilinear instead of linear features. DRR identifies the nonlinear features through multivariate regression to ensure the reduction in redundancy between the PCA coefficients, the reduction of the variance of the scores, and the reduction in the reconstruction error. More importantly, unlike other nonlinear dimensionality reduction methods, the invertibility, volume-preservation, and straightforward out-of-sample extension, makes DRR interpretable and easy to apply. The properties of DRR enable learning a more broader class of data manifolds than the recently proposed Non-linear Principal Components Analysis (NLPCA) and Principal Polynomial Analysis (PPA). We illustrate the performance of the representation in reducing the dimensionality of remote sensing data. In particular, we tackle two common problems: processing very high dimensional spectral information such as in hyperspectral image sounding data, and dealing with spatial-spectral image patches of multispectral images. Both settings pose collinearity and ill-determination problems. Evaluation of the expressive power of the features is assessed in terms of truncation error, estimating atmospheric variables, and surface land cover classification error. Results show that DRR outperforms linear PCA and recently proposed invertible extensions based on neural networks (NLPCA) and univariate regressions (PPA).

Index Terms—Manifold learning, nonlinear dimensionality reduction, Principal Component Analysis (PCA), Dimensionality Reduction via Regression, hyperspectral sounder, IASI, Landsat.

I. INTRODUCTION

In the last decades, the technological evolution of optical sensors has provided remote sensing analysts with rich spatial, spectral, and temporal information. In particular, the increase in spectral resolution of hyperspectral sensors in general, and of infrared sounders in particular, opens the doors to new application domains and poses new methodological challenges in data analysis. The distinct highly-resolved spectra offered by hyperspectral images (HSI) allow us to characterize land-cover classes with unprecedented accuracy. For instance, hyperspectral instruments such as NASA's Airborne Visible Infra-Red Imaging Spectrometer (AVIRIS) covers the wavelength region

from 0.4 to 2.5 μm using more than 200 spectral channels, at a nominal spectral resolution of 10 nm. The MetOp/IASI infrared sounder poses even more complex image processing problems, as it acquires more than 8000 channels per iFOV. Actually, such improvements in spectral resolution have called for advances in signal processing and exploitation algorithms capable of summarizing the information content in as few components as possible [1]–[4].

In addition to its eventual high dimensionality, the complex interaction between radiation, atmosphere, and objects in the surface leads to irradiance manifolds which consist of non-aligned clusters that may change nonlinearly in different acquisition conditions [5], [6]. Fortunately, it has been shown that, given the spatial-spectral smoothness of the signal, the intrinsic dimensionality of the data is small, and this can be used both for efficient signal coding [3], [7], and for knowledge extraction from a reduced set of features [8], [9]. The high dimensionality problem is not only affecting the hyperspectral data: very often, multispectral data processing applications involve using spatial, multi-temporal or multi-angular features that are combined with the spectral features [10], [11]. In such cases, the representation space becomes more redundant and pose challenging problems of collinearity for the algorithms. In both cases, the key in coding, classification, and in bio-geophysical parameter retrieval applications reduces to finding the appropriate set of features, that should be necessarily flexible and nonlinear.

In order to find these features, in recent years a number of feature extraction and dimensionality reduction methods have been presented. Most of them are based on nonlinear functions to allow describing data manifolds that exhibit nonlinear relations (see [12] for a comprehensive review). Approaches range from local methods [13]–[17], kernel-based and spectral decompositions [9], [18]–[20], neural networks [21]–[23], or projection pursuit formulations [24], [25]. Despite the theoretical advantages of nonlinear methods, the fact is that classical principal component analysis (PCA) [26] is still the most widely used dimensionality reduction technique in real remote sensing applications [3], [27]–[29]. This is mainly because PCA has different properties that make it useful in real examples: it is easy to apply since it involves solving a linear and convex problem, and it has a straightforward out-of-sample extension. Moreover, the PCA transformation is invertible and, as a result, the features extracted can be easily interpreted.

The new dimensionality reduction algorithms that involve nonlinearities rarely fulfill the above properties. Nonlinear models usually have complex formulations, which introduce

Manuscript received March 24, 2015.

Copyright (c) 2014 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org

Image Processing Laboratory (IPL), Universitat de València, Catedrático A. Escardino - 46980 Paterna, València (Spain). E-mail: {valero.laparra, jesus.malo, gcamps}@uv.es

This work was partially supported by the Spanish Ministry of Economy and Competitiveness (MINECO) under project TIN2012-38102-C03-01, and under a EUMETSAT contract.

a number of non-intuitive free parameters. Tuning these parameters implies strong assumptions about the manifold characteristics (e.g. local Gaussianity or special symmetries), or a high computational cost training. This complexity reduces the applicability of nonlinear feature extraction to specific data, i.e. the performance of these methods do not significantly improve that of PCA on many remote sensing problems [3], [9], [27]. Moreover, these methods have problems to obtain out-of-sample predictions, which is mandatory in most of the real applications. Another critical point is that the transform involved by the nonlinear models is hard to interpret. This problem could be alleviated if the methods were invertible since then one could get the data back to the input domain (where units are meaningful) and understand the results therein. Invertibility allows to characterize the transformed domain, and to evaluate its quality. However, invertibility is scarcely achieved in the manifold learning literature. For instance, spectral and kernel methods involve *implicit* mappings between the original and the curvilinear coordinates, but these *implicit* features are not easily invertible nor interpretable [30].

The desirable properties of PCA are straightforward in methods that find projections onto *explicit* features in the input domain. These *explicit* features may be either straight lines or curves. This family of projection methods may be understood as a generalization of linear transforms extending linear components to curvilinear components. This family ranges between two extreme cases: (1) **rigid** approaches where features are straight lines in the input space (e.g. conventional PCA, Independent Components Analysis -ICA- [31]), and (2) **flexible** non-parametric techniques that closely follow the data, such as Self-Organizing Maps (SOM) [32], or the related Sequential Principal Curves Analysis (SPCA) [6], [33]. This family is discussed in Section II below. Both extreme cases are undesirable because of different reasons: limited performance (in too rigid methods), and complex tuning of free parameters and/or unaffordable computational cost (in too flexible methods). In this *projection-onto-explicit-features* context, autoencoders such as Nonlinear-PCA (NLPCA) [23], and approaches based on fitting functional curves, such as Principal Polynomial Analysis (PPA) [34], [35], represent convenient intermediate points between the extreme cases in the family. Note that these methods have shown better performance than PCA on a variety of real data [35], [36]. Actually, in the case of PPA, it is theoretically ensured to obtain better results than PCA. The method proposed here, *Dimensionality Reduction via Regression* (DRR), represents a qualitative step towards the flexible end in this family because of the multivariate nature of the regression (as opposed to the univariate regressions done in PPA) while keeping the convenient properties of PPA and PCA which make it suitable for practical high dimensional problems (as opposed to SPCA and SOM). Therefore, it extends the applicability of PPA to more general manifolds, such as those encountered in remote sensing data analysis.

Following the taxonomy in [35] these three methods (NLPCA, PPA and DRR) could be included in the *Principal Curves Analysis* framework [37]. This framework includes both parametric (fitting analytic curves) [26], [38], [39], and

non-parametric [6], [33], [40]–[42] methods. NLPCA, PPA and DRR exploit the idea behind this framework to define generalizations of PCA of *controlled* flexibility.

Preliminary results of DRR were presented in [43]. Here we extend the theoretical analysis of the method and the experimental confirmation of the performance in hyperspectral images. The remainder of the paper is organized as follows. Section II reviews the properties and shortcomings of the *projection-onto-explicit-features* family pointing out the *qualitative advantages of the proposed DRR*. Section III introduces the mathematical details of DRR. We describe the DRR transform and the key differences with PPA. We derive an explicit expression for the inverse and we prove the volume preservation property of DRR. The theoretical properties of DRR are demonstrated and illustrated in controlled toy examples of different complexity. In Section IV, we address two important high dimensional problems in remote sensing: the estimation of atmospheric state vectors from Infrared Atmospheric Sounding Interferometer (IASI) hyperspectral sounding data, and the dimensionality reduction and classification of spatio-spectral Landsat image patches. In the experiments, DRR is compared with conventional PCA [26], and with recent fast nonlinear generalizations that belong to the same class of invertible transforms, PPA [34], [35] and NLPCA [23]. Comparisons are made both in terms of reconstruction error and of expressive power of the extracted features. We end the paper with some concluding remarks in Section V.

II. FROM RIGID TO FLEXIBLE FEATURES

Here we illustrate how DRR represents a step forward with regard to NLPCA and PPA in the family of projections onto explicit curvilinear features ranging from rigid to flexible extremes. First, we review the basic details of previous projection methods, and then we illustrate the advantages of the proposed method in an easy to visualize example.

A. Representative projections onto lines and curves

Classical techniques such as PCA [26] or ICA [31] represent the *rigid* extreme of the family, where, zero-mean samples $x \in \mathbb{R}^d$ are projected onto d rectilinear features through the projection matrix, \mathbf{V} :

$$\alpha = \mathbf{V} \cdot x$$

where α_i are the Principal Components (PC scores for PCA) or the Independent Components (for ICA), and the d linear features in the input space are the column vectors (straight directions) in \mathbf{V}^{-1} . These rigid techniques use a single set of global features regardless of the input.

On the contrary, flexible techniques adapt the set of features to the local properties of the input. Examples include SOM [32] where a flexible grid is adapted to the data and samples can be represented by projections onto the local axes defined by the edges of the parallelepiped corresponding to the closest node. Similarly, local-PCA [44] and local-ICA [45] project the data onto local axes corresponding to the closest code vector. More generally, local-to-global methods integrate these local-linear representations into a single global curvilinear

representation [46]. In particular, using the fact that local eigenvectors are tangent to first and secondary principal curves [47], Sequential Principal Curves Analysis (SPCA) [6], [33] integrates local PCAs, $\mathbf{V}(x')$, along a sequence of d principal curves to get a curvilinear representation

$$r = \int_{x_0}^x D(x') \cdot \mathbf{V}(x') \cdot dx',$$

where the local metric, $D(x')$, sets the line element along the curves. SPCA is inverted by taking the lengths, r_i , along the sequence of principal curves drawn from the origin, x_0 . Similarly to SOM, SPCA assumes a grid of curves adapted to the data. However, as opposed to SOM, SPCA does not learn the whole grid, but only d segments of principal curves per sample.

The above methods identify explicit curves/features that follow the data, but they are hard to train (e.g. parameters to control their flexibility depend on the problem) and require many samples to be reliable, which make them hard to use in high-dimensional scenarios. Other methods have been proposed to generalize the rigid representations by considering curvilinear features instead of straight lines [26]. For instance, in NLPCA [21], [23] an invertible internal representation is computed through a two stage neural network,

$$r = \mathbf{W}_2 \cdot g_1(\mathbf{W}_1 \cdot x)$$

where the matrices \mathbf{W}_i represent sets of linear receptive fields, and g_1 is a set of fixed point-wise nonlinearities. The inverse of this autoencoder [22] can be used to make the curvilinear coordinates explicit.

Fitting general parametric curves in \mathbb{R}^d , as done in [38], [39], is difficult because of the unconstrained nature of the problem [26], [35]. Alternatively, PPA [35] follows a deflationary sequence in which a single polynomial depending on a straight line (univariate fit) is computed at a time. Specifically, the i -th stage of PPA accounts for the i -th curvilinear dimension by using two elements: (1) one-dimensional projection onto the leading vector $e^{(i)}$, and (2) polynomial prediction of the average at the orthogonal subspace,

$$\begin{aligned} \alpha_i &= e^{(i)\top} \cdot x^{(i-1)} \\ x^{(i)} &= \mathbf{E}_{\perp}^{(i)} \cdot x^{(i-1)} - f^{(i)}(\alpha_i) \end{aligned} \quad (1)$$

where the polynomial prediction, $f^{(i)}(\alpha_i)$, is removed from the data in the orthogonal subspace. Superindices in the above formula represent the stage. As a result, data at the i -th stage is represented by α_i and by the $(d-i)$ -dimensional residual that cannot be predicted from that projection. Prediction using this univariate polynomial is a way to remove possible nonlinear dependencies between the linear subspace of $e^{(i)}$ and its orthogonal complement. Despite its convenience, the univariate nature of the fits restricts the kind of dependencies that can be taken into account since more information about the orthogonal subspace (better predictions) could be obtained if more dimensions were used in the prediction. Moreover, using a single parameter, α_i , to build the i -th polynomial

implies that the i -th curvilinear feature has the same shape along the $(i-1)$ -th curve.

DRR addresses these limitations by using multivariate instead of univariate regressions in the nonlinear predictions. As a result, DRR improves energy compaction and extends the applicability of PPA to more general manifolds while keeping its simplicity, which make it suitable in high dimensional problems (as opposed to SPCA and SOM).

B. Qualitative advantages of DRR

The advantages of DRR are illustrated in Fig. 1 where we compare representative invertible representations of this family on two curved and noisy manifolds of the class introduced by Delicado [47] (in red and blue). This class of manifolds, originally presented to illustrate the concept of *secondary principal curves* [47], is convenient since one can easily control the complexity of the problem by introducing tilt (non-stationarity) on the secondary principal curves (dark color) along the first principal curve (light color). This controlled complexity is useful to point out the limitations of previous techniques (e.g. required symmetry in the manifold) and how these limitations are alleviated by using the (more general) DRR model. The performance is compared in the input domain through the dimensionality reduction error and through the accuracy of the identified curvilinear features. These manifolds come from a two-dimensional space of latent variables (positions along the first and secondary curves). As a result, the dimensionality reduction error depends on the unfolding ability of the forward transform: the closer the transformed data fit a flat rectangle, the smaller the error when truncating the representation. On the other hand, the identified features depend on how the inverse transform bends a cartesian grid in the latent space: the better the model represents the curvature of data, the bigger the fidelity of the identified features.

Let us start by considering the performance on the easy case: manifold in red with no tilt along the second principal curve. The previously reported techniques perform as expected: on the one hand, progressively more flexible techniques (from PCA to SPCA) reduce the distortion after dimensionality reduction (in MSE_{DR} terms) because they better unfold test data. As a result, removing the third dimension in the rigid-to-flexible family progressively introduces less error. On the other hand, the identified features in the input domain are progressively more similar to the actual curvilinear latent variables when going from the rigid to the flexible extremes. In this specific *easy* example the proposed DRR outperforms even the flexible SPCA in MSE_{DR} terms. Moreover, since this particular manifold may not require increased flexibility (and hence may be better suited to the PPA model), PPA slightly outperforms DRR in MSE_F terms.

Results for the more complex manifold (tilted secondary curves, in blue) provide more insight into the techniques since it pushes them (specifically PPA) to their limits. Firstly, note that the increase in complexity is illustrated by an increase in the errors in all methods. For instance, linear PCA, that identifies the same features in both cases, doubles the normalized MSEs. While the reduction in performance is not

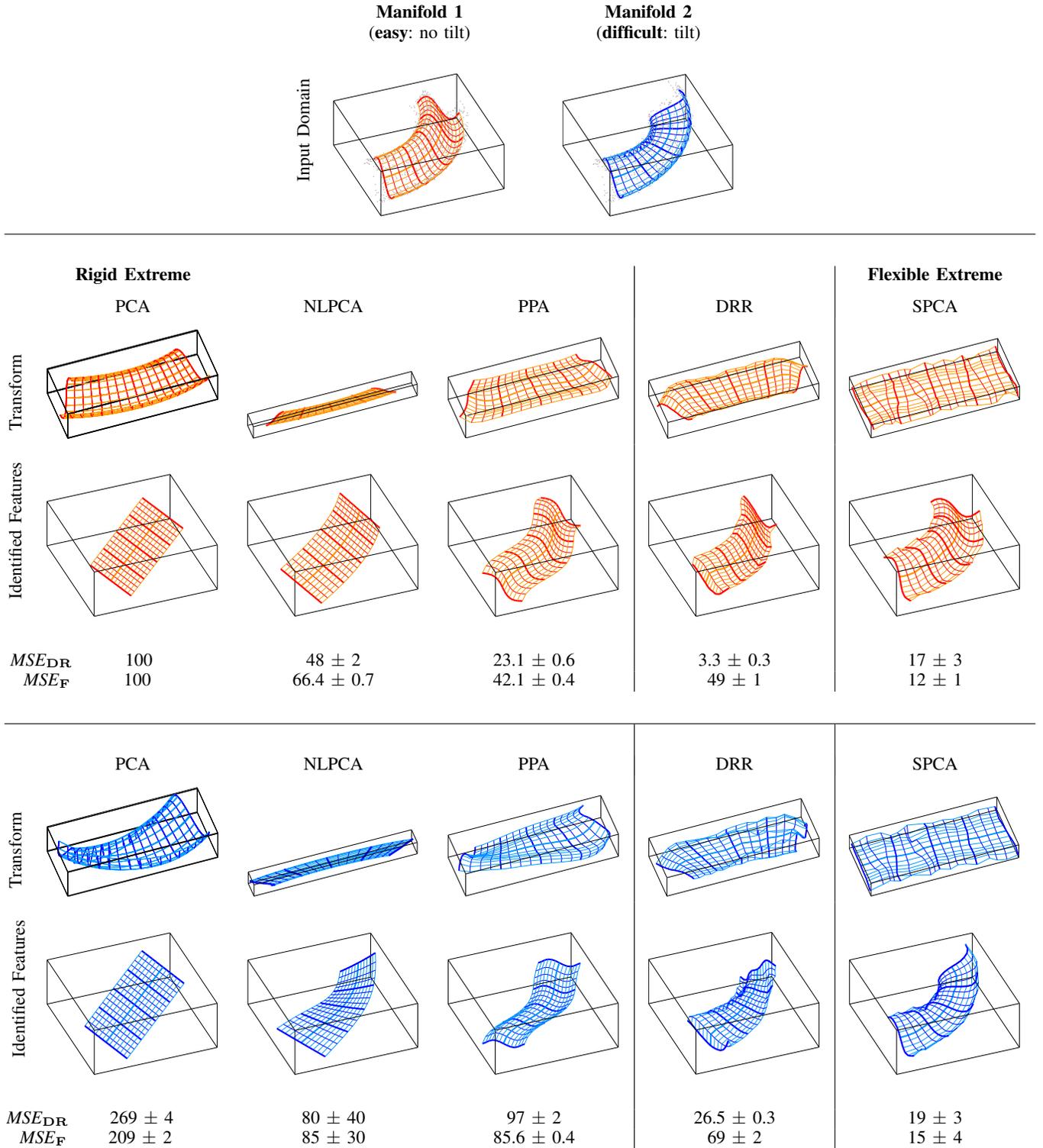


Fig. 1. Performance of the family of invertible representations on illustrative manifolds of different complexity. Complexity of the considered manifolds (top panel) depends on the tilt in secondary principal curves along the first principal curve [47]. Sample data are shown together with the first and secondary principal curves generated by the latent variables (angle and radius) in the input domain. Results of the different techniques for the considered manifolds are depicted in the same color as the input data (red for the no-tilt manifold, and blue for the tilted manifold). Previously reported representations range from rigid schemes such as PCA [26] to flexible schemes such as SPCA [6], [33], including practical nonlinear generalizations of PCA such as NLPCA [23] and PPA [35] which are examples of intermediate flexibility between the extreme cases. Performance is compared in terms of reconstruction error when removing the third dimension (dimensionality reduction MSE_{DR} numbers are relative to the PCA error in the easy case), and in terms of the mean squared distance between the identified and the actual curvilinear features (MSE_F numbers are relative to the PCA error in the easy case). MSE_{DR} is related to the unfolding ability of the model (see the *Transform* rows), and MSE_F is related to its ability to get appropriate curvilinear features from an assumed latent cartesian grid (see the *Identified Features* rows). We used 10^4 training samples and optimal settings in all methods. Dimensionality reduction was tested on the 17×13 highlighted curvilinear grid sampled from the true latent variables. The features in the input space were identified by inverting a 17×13 2-d cartesian grid in the transform domain. This (assumed) latent grid was scaled in every representation to minimize MSE_F . Standard deviations in errors come from models trained on 10 different data set realizations.

that relevant in SPCA (remember these flexible techniques cannot be applied in high dimensional scenarios), this twisted manifold certainly challenges fast generalizations of PCA: the MSEs dramatically increase for NLPCA and PPA. Even though NLPCA identifies certain tilt in the secondary feature along the first curve, it seems too rigid to follow the data structure. PPA displays a different problem: as stated above, by construction, the i -th curvilinear feature in PPA cannot handle relations with the $(i-1)$ -th curve beyond the prediction of the mean. This is because the data in all orthogonal subspaces along the $(i-1)$ -th curve collapse, and are described by a single curve depending on a single parameter (*univariate regression*). This leads to using the same i -th curve all along the $(i-1)$ -th feature (note the repeated secondary curves along the first curve in both, red and blue, cases). This is good enough when data manifolds have the required symmetry (PPA performance is over NLPCA in the first case), but leads to dramatic errors when the method have to deal with relations between three or more variables, as for the manifold in blue, where PPA performance is below NLPCA. This latter effect frequently appears in hyperspectral images, as different (non-stationary) nonlinear relations between spectral channels occur for different objects [3], [48], [49].

Finally, note that DRR clearly improves PPA in the challenging example in blue, mainly because it uses multiple dimensions (instead of a single one) to predict each lower variance dimension in the data. As a result, it can handle non-stationarity along the principal curves leading to better unfolding (lower MSE_{DR}) and tilted secondary features (lower MSE_F). This removes the symmetry requirement in PPA and broadens the class of manifolds suited to DRR.

III. DIMENSIONALITY REDUCTION VIA REGRESSION

PCA removes the second order dependencies between the signal components, i.e. PCA scores are decorrelated [26]. Equivalently, PCA can be casted as the *linear* transform that minimizes reconstruction error when a fixed number of features are neglected. However, for general non-Gaussian sources, and in particular for hyperspectral images, PCA scores still display significant statistical relations, see [3][ch. 2]. The scheme presented here tries to *nonlinearly* remove the information still shared by different PCA components.

A. DRR scheme

It is well known that, even though PCA leads to a domain with decorrelated dimensions, complete independence (or non redundant coefficients) is guaranteed only if the signal has a Gaussian probability density function (PDF). Here, we propose a scheme to remove this redundancy (or uninformative data). The idea is simple: just predict the redundant information in each coefficient that can be extracted from the others. Only the non-predictable information (the residual prediction error) should be retained for data representation. Specifically, we start from the linear PCA representation outlined above: $\alpha = \mathbf{V}x$. Then, we propose to predict each coefficient, α_i , through a multivariate regression function, $f_i(\cdot)$, that takes the

higher variance components as inputs for prediction. The non-predictable information is:

$$y_i = \alpha_i - \hat{\alpha}_i = \alpha_i - f_i(\alpha_1, \alpha_2, \dots, \alpha_{i-1}), \quad (2)$$

and this residual, y_i , is the i -th dimension of the DRR domain. This *prediction+subtraction* is applied $d-1$ times, $\forall i = d, d-1, \dots, 2$, where d is the dimension of the input. As a result, the DRR representation of each input vector x , is:

$$r = \mathbf{R}(x) = (\alpha_1, y_2, y_3, \dots, y_d)^\top.$$

B. Properties of DRR

a) *DRR generalizes PCA*: In the particular case of using linear regressions in $f_i(\cdot)$, i.e. linear-DRR, the prediction $\hat{\alpha}_i$ would be equal to zero. Note that this is the result when using classical (least squares) solution since α_i is uncorrelated with each $\alpha_1, \alpha_2, \dots, \alpha_{i-1}$. Therefore $f_i(\alpha_1, \alpha_2, \dots, \alpha_{i-1}) = 0$, and then $y_i = \alpha_i$, i.e. linear-DRR reduces to PCA.

As a result, if the employed nonlinear functions $f_i(\cdot)$ generalize the linear functions, DRR will obtain at least as good results as PCA. The flexibility of these functions with regard to the linear case will reduce the variance of the residuals, and hence the reconstruction error in dimensionality reduction.

b) *DRR is invertible*: Given the DRR transformed vector, $(\alpha_1, y_2, y_3, \dots, y_d)^\top$, and knowing the functions of model $f_i(\cdot)$, the inverse is straightforward since it reduces to sequentially undo the forward transformation: we first predict coefficient $(i+1)$ -th from previous (known) coefficients using the known regression function, and then, we use the known residual to correct the prediction:

$$\alpha_i = \hat{\alpha}_i + y_i = f_i(\alpha_1, \alpha_2, \dots, \alpha_{i-1}) + y_i \quad (3)$$

c) *DRR has an easy out-of-sample extension*: Note that forward and inverse DRR transforms can be applied to new data (not in the training set) since there is no restriction to apply the prediction functions $f_i(\cdot)$. See Sec. III-C for a discussion on the selected regression functions in this work.

d) *DRR is a volume preserving transform*: A nonlinear transform preserves the volume of the input space if the determinant of its Jacobian is one for all x [50]. Here we prove that the nature of DRR ensures that its Jacobian fulfills this property.

DDR can be thought of as a sequential algorithm in which only one dimension is addressed at a time through the elementary transform \mathbf{R}_i consisting of prediction and subtraction (Eq. (2)). Yet mathematically convenient to formulate the Jacobian, this sequential view is does not prevent the parallelization discussed later. Hence, the (global) Jacobian of DRR, $\nabla \mathbf{R}$, is the product of the Jacobians of the elementary transforms in this sequence times the matrix of the initial PCA as follows:

$$\nabla \mathbf{R}(x) = \left(\prod_{i=2}^d \nabla \mathbf{R}_i \right) \mathbf{V}.$$

The i -th elementary transform leaves all, but the i -th dimension, unaltered. Therefore, each elementary Jacobian is the identity matrix except for the i -th row, which depends

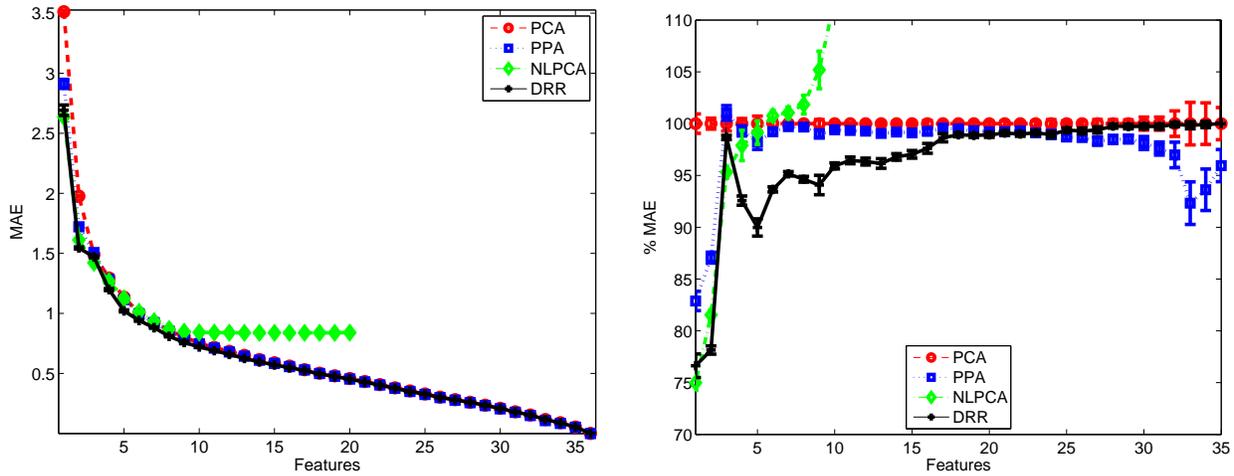


Fig. 2. Reconstruction error results on the contextual multispectral image classification. Comparison between PCA, PPA, NLPCA and DRR for different number of extracted features, in both mean absolute reconstruction error (MAE) (left) and relative MAE with respect to PCA error (right), for which going below the PCA means better results (less error).

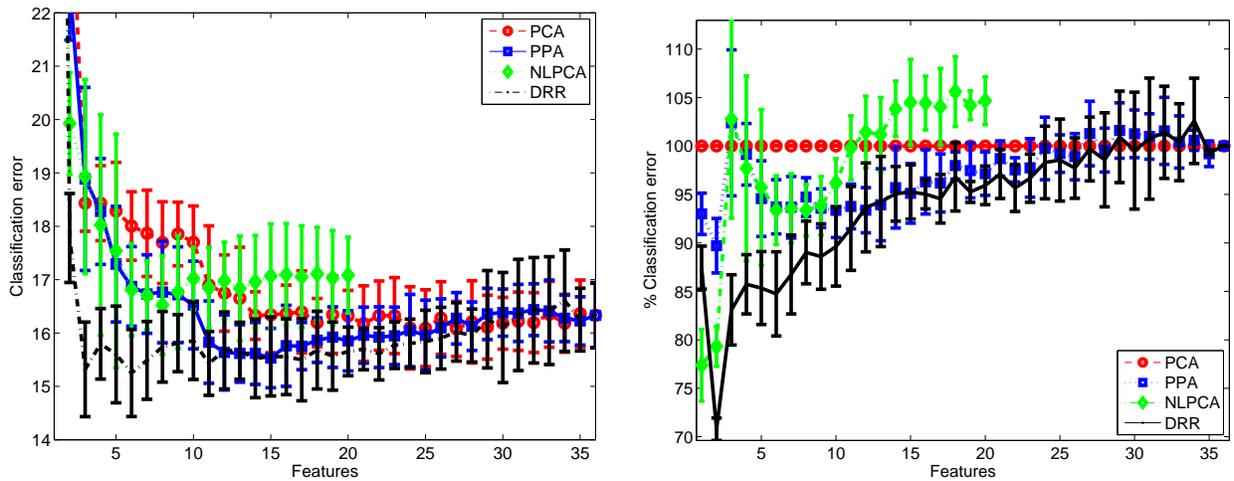


Fig. 3. Classification results on the contextual multispectral image classification. Comparison between PCA, PPA, NLPCA and DRR for different number of extracted features, in both classification error (left) and relative classification error with respect to PCA accuracy (right), for which going below the PCA means better results (less error).

resolution of $80\text{m} \times 80\text{m}$ (all data acquired from a rectangular area approximately 8 km wide)³. Six classes are identified in the image, namely red soil, cotton crop, grey soil, damp grey soil, soil with vegetation stubble and very damp grey soil. A total of 6435 labeled samples are available. Contextual information was included stacking neighboring pixels in 3×3 windows. Therefore, 36-dimensional input samples were generated, with a high degree of redundancy and collinearity. We address two problems with this dataset: a pure spatio-spectral dimensionality reduction problem, and the effect of the reduced dimension in image classification.

1) *Reconstruction accuracy*: In the first problem, we compare the dimensionality reduction performance in terms of Mean Absolute Error (MAE) in the original domain. Note that this kind of evaluation can be used only with invertible methods. For each method, the data are transformed and

then inverted using less dimensions. This is equivalent to truncate dimensions in PCA. In order to illustrate the advantage of using a given method instead of PCA, results are shown in percentage with regard to the PCA performance: $\%MAE_{\text{method}} = 100 \text{ MAE}_{\text{method}} / \text{MAE}_{\text{PCA}}$, where $\text{MAE}_{\text{method}}$ and MAE_{PCA} refer to the MAE obtained with the considered method and PCA, respectively.

Figure 2 shows the results of the experiment. We divided the available labeled data into two sets (training and test) with equal number of samples. The samples of each set have been randomly selected from the original image dataset. The MAE of reconstruction in the test set averaged over ten independent realizations is shown. Several conclusions can be obtained: Specifically, NLPCA obtains good results when a few number of extracted features are obtained, but rapidly degrades its performance with more than 10 extracted features, revealing a clear inability to handle high-dimensional problems. Note

³Image available at <http://www.ics.uci.edu/mllearn/MLRepository.html>

that the available implementation of NLPCA⁴ is restricted to extract at most 20 features. For a given number of extracted features, the reconstruction error increases substantially with regard to PCA (Fig. 2 right). PPA shows better results than NLPCA, and it is better suited than PCA in all the number of extracted features. Nevertheless, it is noticeable that DRR is in all cases better than all the other methods, revealing a maximum gain of +25% over PCA for very few features.

2) *Classification accuracy*: The second problem with this dataset shows the classification results using the inverted data into the original input space of the different methods. We used the standard linear discriminant analysis on top of the inverted data. In all cases, we used 3200 randomly selected examples for training and the same amount for testing. Test results are averaged over five realizations, and are shown in Fig. 3. The performance results indicate similar trends observed in the reconstruction error in Fig. 2. Essentially, DRR outperforms the other methods, especially noticeable when a few number of components are used for reconstruction and classification. As the number of components increase, DRR and PPA show similar results. These results suggest that DRR better compacts the information in a lower number of components, which is useful for both reconstruction and data classification.

3) *Computational load*: Table I shows the computation cost for all considered methods for training and testing⁵. The experiments used 3200 training and 3200 test samples, with $d = 36$. Two main conclusions can be extracted: NLPCA is the most computationally costly algorithm for training and DRR for testing.

TABLE I
COMPUTATIONAL COST LANDSAT DATASET

	PCA	PPA	NLPCA	DRR
Training time (sec)	0.05	0.6	7944	1920
Testing time (sec)	0.007	0.16	0.05	35

B. Experiment 2: Regression from infrared sounding data

We here analyze the benefits of using DRR for the estimation of atmospheric parameters from hyperspectral infrared sounding data with a reduced dimensionality. We first motivate the problem, and then describe the considered dataset. Again, we are interested in analyzing the impact of the reduced dimensionality both in the reconstruction error and in a different task, in this case, the retrieval of geophysical parameters.

Temperature and water vapor are atmospheric parameters of high importance for weather forecast and atmospheric chemistry studies [58], [59]. Observations from spaceborne high spectral resolution infrared sounding instruments can be used to calculate the profiles of such atmospheric parameters

⁴<http://www.nl pca.org/>

⁴While other more sophisticated nonlinear classifiers could be used here, we are actually interested in this setting that allows us to study the expressive power of the extracted features. An homologous setting will be also used in the regression experiments of next subsection.

⁵Experiments were performed using Matlab on an Intel 3.3 GHz processor with 48 GB RAM memory. No parallelization was applied on DRR in this experiment.

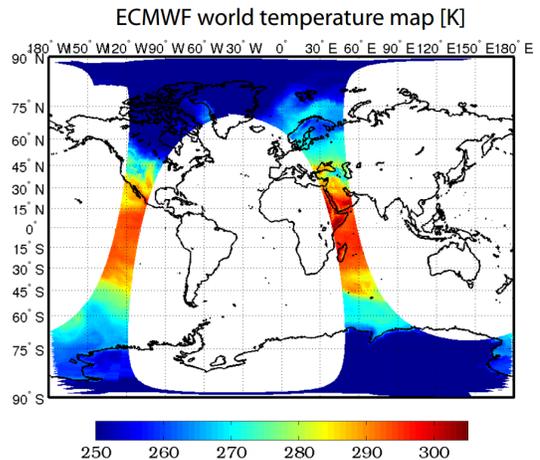


Fig. 4. Surface temperature [in K] world map provided by the official ECMWF model, <http://www.ecmwf.int/>.

with unprecedented accuracy and vertical resolution [60]. In this work we focus on the data coming from the Infrared Atmospheric Sounding Interferometer (IASI), the Microwave Humidity Sensor (MHS) and the Advanced Microwave Sensor Unit (AMSU) onboard of the MetOp-A satellite⁶. The IASI instrument is the one that poses the major dimensionality challenge due to its dense spectrum sampling: while MHS and AMSU spectra consist of about twenty values together, IASI spectra consist of 8461 spectral channels, between 3.62 and 15.5 μm , with a spectral resolution of 0.5 cm^{-1} after apodization [61], [62]. Its spatial resolution is 25 km at nadir with an Instantaneous Field of View (IFOV) size of 12 km at an altitude of 819 km. This huge data dimensionality typically requires simple and computationally efficient processing techniques.

One of the retrieval techniques available in the MetOp-IASI Level 2 Product Processing Facility (L2 PPF) is a computationally inexpensive method based on linear regression from the principal components of the measured brightness spectra and the atmospheric state parameters. We aim to introduce DRR in such scheme as an alternative to PCA. In this application it is important that dimensionality reduction minimizes the reconstruction error and that the identified features are useful in the retrieval stage.

We used a collection of 23 datasets of input data from the different sensors: IASI, MHS and AMSU. The considered output atmospheric variables are diverse, e.g. temperature, moisture, and surface pressure. In each dataset provided by EU-METSAT, the preprocessed input data were 110-dimensional. Each input vector consisted of the following: *one* scalar indicating the secant of satellite zenith angle, 19 radiance values from the AMSU and MHS sensors, and 90 values from the IASI sensor. The data from IASI were actually three separate sets of 30 PC scores each, from three different IASI bands. Note that, despite intra-band decorrelation, the vector elements may still exhibit statistical dependency, which may be significant even at a second order level, among different bands and instruments.

⁶<https://directory.eoportal.org/web/eoportal/satellite-missions/m/metop>

The data to be predicted (or output data) is 277-dimensional. Each output vector consists of the following: 4 data corresponding to the *surface temperature* and *moisture*, the *skin temperature*, and the *surface pressure*; and 273 data corresponding to altitude profiles of *temperature*, *moisture*, and *ozone* at 91 model levels each. An example of surface temperature is shown in Fig. 4. Data was provided by the official European Center for Medium-range Weather Forecasting (ECMWF) model, <http://www.ecmwf.int/>, on March 4th, 2008.

1) *Reconstruction accuracy*: In this experiment, we study the representation power of a small number of features extracted by DRR. The 110 input features are processed with PCA [26], PPA [34], [35], NLPCA [21], [23] and the presented DRR method. Here, the quality of the transformation is evaluated solely with the mean absolute error (MAE) in the input space between the original signal and the reconstructed with the most relevant coefficients retained. Figure 5 illustrates the effect of reconstructing the input data when using PCA, PPA, NLPCA and DRR for different numbers of components. On the one hand, as reported in [35], the performance in PPA is similar or better than in NLPCA in reconstruction error. On the other hand, it is important to note that results in absolute and relative terms show that DRR clearly obtains less reconstruction error than PCA and PPA for an arbitrary number of features.

2) *Retrieval accuracy*: Figure 6 illustrates the effect of using the features either from PCA, PPA or DRR for the retrieval of the physical parameters described before. We used linear regression in the features-to-parameters estimation. We plotted the mean absolute error (MAE) for different number of features. These plots show the effect of using different (linear and non-linear) dimensionality reduction methods for retrieval. Figure 6 shows the results for the first dataset for illustration purposes (similar results were obtained for the remainder datasets). Note that using DRR features to estimate the features has clear benefits. For instance, using just the 20% of the DRR features obtains the same accuracy as PCA when using all the components.

3) *Computational load*: Times for training and testing are shown in Table II (same computer resources as before). In this experiment, we took 10000 training and 10000 test samples, and $d = 110$. As in the previous experiment, NLPCA and DRR are the most expensive in training and test, respectively. In this experiment, however, times for DRR are notably higher due to the increase in dimensionality but mostly to the bigger training set.

TABLE II
COMPUTATIONAL COST IASI DATASET

	PCA	PPA	NLPCA	DRR
Training time (sec)	0.13	16	65389	14424
Testing time (sec)	0.01	0.3	1.3	1112

V. CONCLUSIONS

We introduced a novel unsupervised method for dimensionality reduction via the application of a multivariate nonlinear

regression to approximate each projection from the higher variance scores. The method is shown to generalize PCA and to achieve more data compression (smaller MSE for a fixed number of retained components) and better features for prediction (less error in classification and regression problems) than competitive nonlinear methods like NLPCA and PPA. Besides, unlike other nonlinear dimensionality reduction methods, DRR is easy to apply, it has out-of-sample extension, it is invertible, and the learned transformation is volume-preserving. We focused on the challenging problems of spatial-spectral multispectral land cover classification, and atmospheric parameter retrieval from hyperspectral infrared sounding data. Extension of DRR to cope with multiset/output regression, as well as impact of the data dimensionality and noise sources, will be explored in the future.

VI. ACKNOWLEDGMENTS

The authors wish to thank Tim Hultberg from the European Organisation for the Exploitation of Meteorological Satellites (EUMETSAT) in Darmstadt, Germany, for kindly providing the IASI datasets used in this paper.

REFERENCES

- [1] G. Camps-Valls and L. Bruzzone, "Kernel-based methods for hyperspectral image classification," *Geoscience and Remote Sensing, IEEE Transactions on*, vol. 43, no. 6, pp. 1351–1362, June 2005.
- [2] A. Plaza, J. A. Benediktsson, J. Boardman, J. Brazile, L. Bruzzone, G. Camps-Valls, J. Chanussot, M. Fauvel, P. Gamba, A. Gualtieri, and J. Tilton, "Recent advances in techniques for hyperspectral image processing," *Remote Sensing of Environment*, vol. 113, no. S1, pp. 110–122, Sept 2009.
- [3] G. Camps-Valls, D. Tuia, L. Gómez, S. Jiménez, and J. Malo, *Remote Sensing Image Processing*, ser. Synthesis Lectures on Image, Video and Multimedia Processing, A. Bovik, Ed. Morgan & Claypool Publishers, 2011.
- [4] G. Camps-Valls, D. Tuia, L. Bruzzone, and J. Atli Benediktsson, "Advances in hyperspectral image classification: Earth monitoring with statistical learning methods," *Signal Processing Magazine, IEEE*, vol. 31, no. 1, pp. 45–54, Jan 2014.
- [5] D. Tuia, J. Muñoz Marí, L. Gómez-Chova, and J. Malo, "Graph matching for adaptation in remote sensing," *Geoscience and Remote Sensing, IEEE Transactions on*, vol. 51, no. 1, pp. 329–341, Jan 2013.
- [6] V. Laparra, S. Jiménez, G. Camps-Valls, and J. Malo, "Nonlinearities and adaptation of color vision from sequential principal curves analysis," *Neural Comp.*, vol. 24, no. 10, pp. 2751–2788, 2012.
- [7] B. Penna, T. Tillo, E. Magli, and G. Olmo, "Transform coding techniques for lossy hyperspectral data compression," *IEEE Trans. Geosci. Rem. Sens.*, vol. 45, no. 5, pp. 1408–1421, 2007.
- [8] S. Jiménez and J. Malo, "The role of spatial information in disentangling the irradiance-reflectance-transmittance ambiguity," *IEEE Trans. Geosci. Rem. Sens.*, vol. 52, no. 8, pp. 4881–4894, 2014.
- [9] J. Arenas-García, K. Petersen, G. Camps-Valls, and L. Hansen, "Kernel multivariate analysis framework for supervised subspace learning: A tutorial on linear and kernel multivariate methods," *Signal Processing Magazine, IEEE*, vol. 30, no. 4, pp. 16–29, 2013.
- [10] M. Fauvel, J. A. Benediktsson, J. Chanussot, and J. R. Sveinsson, "Spectral and spatial classification of hyperspectral data using SVMs and morphological profiles," *IEEE Trans. Geosci. Remote Sens.*, vol. 46, no. 11, pp. 3804 – 3814, 2008.
- [11] D. Tuia, F. Pacifici, M. Kanevski, and W. Emery, "Classification of very high spatial resolution imagery using mathematical morphology and support vector machines," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 11, pp. 3866–3879, 2009.
- [12] J. A. Lee and M. Verleysen, *Nonlinear dimensionality reduction*. Springer, 2007.
- [13] J. B. Tenenbaum, V. Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, vol. 290, no. 5500, pp. 2319–2323, December 2000.

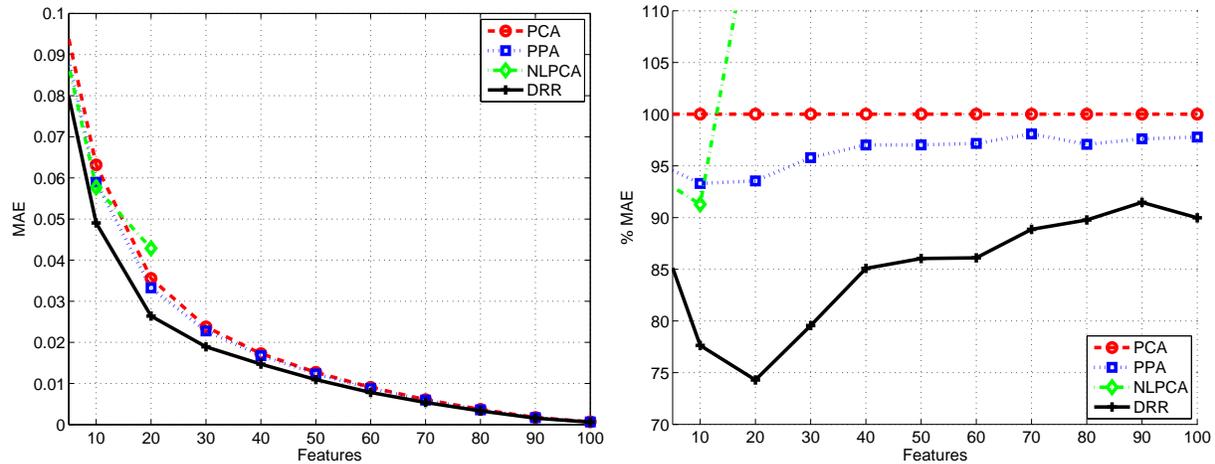


Fig. 5. Reconstruction error. Left: Absolute reconstruction error for different number of retained features obtained when using different DR methods on the first (just one) dataset. Right: Relative error (percentage) with regard to the error in PCA, mean and standard deviation have been obtained over the 23 (all) datasets.

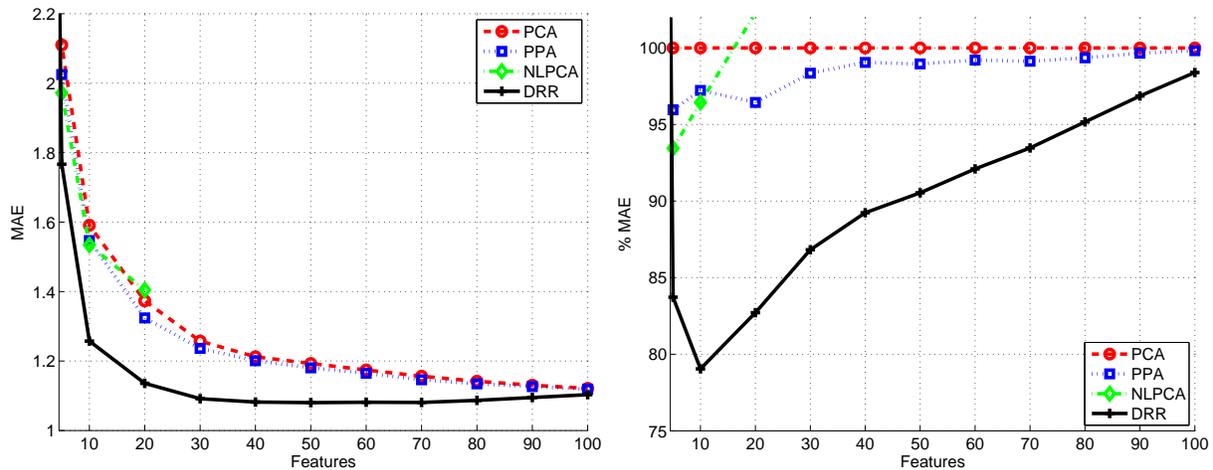


Fig. 6. Retrieval performance. Accuracy of the parameter retrieval (MAE) with regard to the number of retained features. Results are given for different feature extraction (PCA, PPA, NLPCA, DRR) methods. Left: Absolute MAE for the first dataset. Right: Relative (to the PCA MAE in each dimension) results. Results for the remainder 23 are similar.

- [14] S. T. Roweis, L. K. Saul, and G. E. Hinton, "Global coordination of local linear models," in *Advances in Neural Information Processing Systems 14*. MIT Press, 2002, pp. 889–896.
- [15] J. J. Verbeek, N. Vlassis, and B. Krose, "Coordinating principal component analyzers," in *In Proc. International Conference on Artificial Neural Networks*. Springer, 2002, pp. 914–919.
- [16] Y. W. Teh and S. Roweis, "Automatic alignment of local representations," in *NIPS 15*. MIT Press, 2003, pp. 841–848.
- [17] M. Brand, "Charting a manifold," in *NIPS 15*. MIT Press, 2003, pp. 961–968.
- [18] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, vol. 290, no. 5500, pp. 2323–2326, December 2000.
- [19] B. Schölkopf, A. J. Smola, and K.-R. Müller, "Nonlinear component analysis as a kernel eigenvalue problem," *Neural Comp.*, vol. 10, no. 5, pp. 1299–1319, 1998.
- [20] K. Q. Weinberger and L. K. Saul, "Unsupervised learning of image manifolds by semidefinite programming," in *Proc. IEEE CVPR*, 2004, pp. 988–995.
- [21] M. A. Kramer, "Nonlinear principal component analysis using autoassociative neural networks," *AICHe Journal*, vol. 37, no. 2, pp. 233–243, 1991.
- [22] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, July 2006.
- [23] M. Scholz, M. Fraunholz, and J. Selbig, *Nonlinear principal component analysis: neural networks models and applications*. Springer, 2007, ch. 2, pp. 44–67.
- [24] P. Huber, "Projection pursuit," *Annals of Statistics*, vol. 13, no. 2, pp. 435–475, 1985.
- [25] V. Laparra, G. Camps-Valls, and J. Malo, "Iterative gaussianization: from ICA to random rotations," *IEEE Trans. Neur. Net.*, vol. 22, 2011.
- [26] I. Jolliffe, *Principal component analysis*. Springer, 2002.
- [27] G. Camps-Valls, J. Muñoz and, L. Gómez, L. Guanter, and X. Calbet, "Nonlinear statistical retrieval of atmospheric profiles from MetOp-IASI and MTG-IRS infrared sounding data," *IEEE Trans. Geosci. Rem. Sens.*, vol. 50, no. 5, pp. 1759–1769, 2012.
- [28] T. M. Lillesand, R. W. Kiefer, and J. Chipman, *Remote Sensing and Image Interpretation*. New York: John Wiley & Sons, 2008.
- [29] C.-I. Chang, *Hyperspectral Data Exploitation: Theory and Applications*. NJ, USA: Wiley-Interscience, 2008.
- [30] P. Honeine and C. Richard, "The pre-image problem in kernel-based machine learning," *IEEE Signal Processing Magazine*, vol. 28, no. 2, pp. 77–88, 2011.
- [31] A. Hyvärinen, J. Karhunen, and E. Oja, *Independent Component Analysis*. New York, USA: John Wiley & Sons, 2001.
- [32] T. Kohonen, "Self-organized formation of topologically correct feature maps," *Biological Cybernetics*, vol. 43, no. 1, pp. 59–69, January 1982. [Online]. Available: <http://dx.doi.org/10.1007/BF00337288>
- [33] V. Laparra and J. Malo, "Visual aftereffects and nonlinearities from a single statistical framework," *Submitted to Front. Human Neurosci.*, 2014.

- [34] V. Laparra, D. Tuia, S. Jiménez, G. Camps-Valls, and J. Malo, "Non-linear data description with principal polynomial analysis," in *IEEE Workshop on Machine Learning for Signal Processing*, 2012, pp. 1–6.
- [35] V. Laparra, S. Jiménez, D. Tuia, G. Camps-Valls, and J. Malo, "Principal polynomial analysis," *Int. J. Neur. Syst.*, vol. 26, no. 7, 2014.
- [36] M. Scholz, "Validation of nonlinear PCA," *Neural processing letters*, pp. 1–10, 2012.
- [37] T. Hastie, "Principal curves and surfaces," Ph.D. dissertation, Stanford University, 1984.
- [38] D. Donnell, A. Buja, and W. Stuetzle, "Analysis of additive dependencies and concavities using smallest additive principal components," *The Annals of Statistics*, vol. 22, no. 4, pp. 1635–1668, 1994.
- [39] P. C. Besse and F. Ferraty, "Curvilinear fixed effect model," *Computational Statistics*, vol. 10, pp. 339–351, 1995.
- [40] J. Einbeck, G. Tutz, and L. Evers, "Local principal curves," *Statistics and Computing*, vol. 15, pp. 301–313, 2005.
- [41] J. Einbeck, L. Evers, and B. Powell, "Data compression and regression through local principal curves and surfaces," *Int. J. Neur. Syst.*, vol. 20, no. 03, pp. 177–192, 2010.
- [42] U. Ozertem and D. Erdogmus, "Locally defined principal curves and surfaces," *JMLR*, vol. 12, pp. 1249–1286, 2011.
- [43] V. Laparra, J. Malo, and G. Camps-Valls, "Dimensionality reduction via regression on hyperspectral infrared sounding data," in *IEEE Workshop on Hyperspectral Image and Signal Processing*, 2014.
- [44] N. Kambhatla and T. Leen, "Dimension reduction by local PCA," *Neural Computation*, vol. 9, no. 7, pp. 1493–1500, 1997.
- [45] J. Karhunen, S. Malaroui, and M. Ilmoniemi, "Local linear independent component analysis based on clustering," *Intl. J. Neur. Syst.*, vol. 10, no. 6, December 2000.
- [46] J. Malo and J. Gutiérrez, "V1 non-linear properties emerge from local-to-global non-linear ICA," *Network: Comp. Neur. Syst.*, vol. 17, no. 1, pp. 85–102, 2006.
- [47] P. Delicado, "Another look at principal curves and surfaces," *J. Multivar. Anal.*, vol. 77, pp. 84–116, 2001.
- [48] C. Bachmann, T. Ainsworth, and R. Fusina, "Exploiting manifold geometry in hyperspectral imagery," *IEEE Trans. Geosc. Rem. Sens.*, vol. 43, no. 3, pp. 441–454, Mar 2005.
- [49] —, "Improved manifold coordinate representations of large-scale hyperspectral scenes," *IEEE Trans. Geosc. Rem. Sens.*, vol. 44, no. 10, pp. 2786–2803, Oct 2006.
- [50] B. Dubrovin, S. Novikov, and A. Fomenko, *Modern Geometry: Methods and Applications*. New York: Springer Verlag, 1982.
- [51] J. Shawe-Taylor and N. Cristianini, *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- [52] M. Lázaro-Gredilla, J. Q. Candela, C. E. Rasmussen, and A. R. Figueiras-Vidal, "Sparse spectrum gaussian process regression," *Journal of Machine Learning Research*, vol. 11, pp. 1865–1881, 2010.
- [53] J. Arenas-García, K. B. Petersen, G. Camps-Valls, and L. K. Hansen, "Kernel multivariate analysis framework for supervised subspace learning," *IEEE Signal Processing Magazine*, p. 1, 2013.
- [54] Y. Zhang, J. C. Duchi, and M. J. Wainwright, "Divide and conquer kernel ridge regression," in *COLT*, 2013, pp. 592–617.
- [55] J. A. K. Suykens, T. Van Gestel, J. De Brabanter, B. De Moor, and J. Vandewalle, Eds., *Least Squares Support Vector Machines*. Singapore: World Scientific Pub. Co., 2002.
- [56] G. Camps-Valls, L. Guanter, J. Muñoz, L. Gómez, and X. Calbet, "Nonlinear retrieval of atmospheric profiles from MetOp-IASI and MTG-IRS data," in *Proc. Im. Sig. Proc. Rem. Sens. XVI*, vol. 7830. SPIE, 2010, p. 78300Z.
- [57] G. Camps-Valls, V. Laparra, J. Muñoz, L. Gómez, and X. Calbet, "Kernel-based retrieval of atmospheric profiles from IASI data," in *IEEE Proc. IGARSS 11*, Jul 2011, pp. 2813–2816.
- [58] K. N. Liou, *An Introduction to Atmospheric Radiation*, 2nd ed. Hampton, USA: Academic Press, 2002.
- [59] F. Hilton, N. C. Atkinson, S. J. English, and J. R. Eyre, "Assimilation of IASI at the Met Office and assessment of its impact through observing system experiments," *Q. J. R. Meteorol. Soc.*, vol. 135, pp. 495–505, 2009.
- [60] H. L. Huang, W. L. Smith, and H. M. Woolf, "Vertical resolution and accuracy of atmospheric infrared sounding spectrometers," *J. Appl. Meteor.*, vol. 31, pp. 265–274, 1992.
- [61] G. Chalon, F. Cayla, and D. Diebel, "IASI: an advanced sounder for operational meteorology," in *Proceedings of the 52nd Congress of IAF*, Toulouse, France, 2001.
- [62] C. G. Siméoni D., Singer C., "Infrared atmospheric sounding interferometer," *Acta Astronautica*, vol. 40, pp. 113–118, 1997.



Valero Laparra was born in València (Spain) in 1983, and received a B.Sc. degree in Telecommunications Engineering (2005), a B.Sc. degree in Electronics Engineering (2007), a B.Sc. degree in Mathematics degree (2010), and a PhD degree in Computer Science and Mathematics (2011). He is a postdoc in the Image Processing Laboratory (IPL) at Universitat de València, and currently doing a stay in the Laboratory for Computer Vision at the NYU, USA. More details in <http://www.uv.es/lapeva>.



Jesús Malo (1970) received the M.Sc. degree in Physics in 1995 and the Ph.D. degree in Physics in 1999 both from the Universitat de València (Spain). He was the recipient of the Vistakon European Research Award in 1994 for his work in Physiological Optics. In 2000 and 2001 he worked as Fulbright Postdoc at the Vision Group of the NASA Ames Research Center, and at the Lab of Computational Vision of the Center for Neural Science, New York University. He came back to the NYU as visiting Research Specialist in 2013. He served as Associate

Editor of the IEEE Transactions on Image Processing, and currently he is Academic Editor of PLoS ONE, dealing with manuscripts in the intersection between vision science and machine learning. He is with the Image and Signal Processing Group at the Universitat de València (<http://isp.uv.es/>). He is member of the Asociación de Mujeres Investigadoras y Tecnólogas (AMIT). His scientific interests include low-level models of human vision, their relations with statistics and information theory (e.g. feature extraction and sensory organization), and their applications to image processing and vision science experimentation.



Gustau Camps-Valls (M'04, SM'07) received a B.Sc. degree in Physics (1996), in Electronics Engineering (1998), and a Ph.D. degree in Physics (2002) all from the Universitat de València. He is currently an associate professor (hab. Full professor) in the Department of Electronics Engineering. His research is conducted in the Image and Signal Processing (ISP) group, <http://isp.uv.es>. He has been Visiting Researcher at the Remote Sensing Laboratory (Univ. Trento, Italy) in 2002, the Max Planck Institute for Biological Cybernetics (Tübingen, Germany) in

2009, and as Invited Professor at the Laboratory of Geographic Information Systems of the École Polytechnique Fédérale de Lausanne (Lausanne, Switzerland) in 2013. He is interested in the development of machine learning algorithms for geoscience and remote sensing data analysis. He is an author of 120 journal papers, more than 150 conference papers, 20 international book chapters, and editor of the books "Kernel methods in bioengineering, signal and image processing" (IGI, 2007), "Kernel methods for remote sensing data analysis" (Wiley & Sons, 2009), and "Remote Sensing Image Processing" (MC, 2011). He's a co-editor of the forthcoming book "Digital Signal Processing with Kernel Methods" (Wiley & sons, 2015). He holds a Hirsch's h index $h = 38$, entered the ISI list of Highly Cited Researchers in 2011, and Thomson Reuters ScienceWatch[®] identified one of his papers on kernel-based analysis of hyperspectral images as a Fast Moving Front research. In 2015, he got an ERC consolidator grant on statistical learning for Earth observation data analysis. He is a referee of many international journals and conferences. Since 2009 he is Associate Editor of the "IEEE Transactions on Signal Processing", "IEEE Signal Processing Letters", "IEEE Geoscience and Remote Sensing Letters" and acted as Guest Editor of "IEEE Journal of Selected Topics in Signal Processing". Visit <http://www.uv.es/gcamps> for more information.