

Compressing arrays of classifiers using Volterra-Neural Network: application to face recognition

M. Rubiolo · G. Stegmayer · D. Milone

Received: date / Accepted: date

Abstract Model compression is required when large models are used, for example, for a classification task, but there are transmission, space, time or computing constraints that have to be fulfilled. Multilayer Perceptron (MLP) models have been traditionally used as classifiers. Depending on the problem, they may need a large number of parameters (neuron functions, weights and bias) to obtain an acceptable performance. This work proposes a technique to compress an array of MLPs, through the weights of a Volterra-Neural Network (Volterra-NN), maintaining its classification performance. It will be shown that several MLP topologies can be well-compressed into the first, second and third order (Volterra-NN) outputs. The obtained results show that these outputs can be used to build an array of (Volterra-NN) that needs significantly less parameters than the original array of MLPs, furthermore having the same high accuracy. The Volterra-NN compression capabilities were tested for solving a face recognition problem. Experimental results are presented on two well-known face databases: ORL and FERET.

Keywords Model Compression · Array of Neural Networks · Volterra neural network · Face Recognition

1 Introduction

The purpose of model compression is to find a fast and compact model to approximate a function learned by, for example, a classifier. Moreover, it is desirable to achieve this without significant loss in performance [Buciluă et al, 2006]. Often the best performing models that use supervised learning are combinations of complex and large classifiers. However, in some situations, it is not enough for a classifier to be highly accurate, it also has to meet some requirements regarding on-line execution time, storage space and limited computational power [Zhang and Wangmeng, 2007].

It is well-known that Multilayer Perceptron (MLP) models are usually considered as a powerful classification model. However, they easily become large models just by the addition of neurons, for example in the hidden layer, which has a direct effect over the number of model weights. Similarly, the introduction of more information to the model, that is to say more input variables in order to better learn the training data and to improve its classification ability, can cause an increment of the model complexity.

M. Rubiolo, G. Stegmayer
CONICET, CIDISI-UTN-FRSF, Lavaise 610, (3000) Santa Fe, Argentina
Tel.: +54-342-4601579/2390 (Int: 258)
E-mail: mrubiolo@santafe-conicet.gov.ar

D. Milone
CONICET, SINC(i)-FICH-UNL, Ciudad Universitaria, RN 168 Km. 472.4, (3000) Santa Fe, Argentina

It has been recently shown that a Volterra model can be extracted from the parameters of a trained Neural Network (NN) [Stegmayer and Chiotti, 2009]. The Volterra model is formed by Volterra kernels, which are associated with the parameters of the trained NN. This has proven to be particularly useful for reproducing the nonlinear and dynamic behavior of new wireless communications devices [Orengo et al, 2007]. Moreover, in [Rubiolo et al, 2010] the Volterra kernels extraction procedure has been used to build a Volterra-Neural Network (Volterra-NN) model, which is a compressed version of a trained MLP model over a very simple classification task, maintaining the same recognition rate than the original MLP model but with fewer parameters. In this work we propose that, since it is possible to obtain a Volterra-NN model from a single trained MLP model, it is also possible to compress an array of MLPs (*a*MLP) using Volterra-NNs. That is to say, the same methodology applied to build a Volterra-NN model from a single MLP can be used to obtain an array of Volterra-NN models from an *a*MLP.

In fact, arrays and ensembles of neural classifiers have proven to reach significant better results in classification problems than single models [Dzeroski and Zenko, 2004][Rahman and Verma, 2011], in particular for the task of Face Recognition (FR) [Zhang and Wangmeng, 2007]. In these systems, the first step consists of face detection through image processing techniques. Secondly, a feature extraction method is applied to extract useful information of the face. Finally, this information is used in a classifier for recognizing faces [Zhao et al, 2003]. In [Capello et al, 2009] an *a*MLP model was proposed for the classification task within a FR system. Specifically, an array of neural networks have been used for classification, consisting of one MLP for each subject (valid or authorized person), with a final decision made over the network outputs of the complete array. The classification was performed by the maximum output calculation among all the networks outputs. This configuration has achieved significant improvements over the performance of a classic MLP. However, the use of this new neural configuration implies much more parameters, and therefore, a larger and more complex FR system. Due to the fact that FR models should be able to run on small processing devices such as mobile phones, ipods and security cams, or to be transmitted online, the model must have an appropriate size in order to adjust to these requirements.

This work presents a novel approach for compressing a face recognition model based on a novel application of the Volterra-NN method to an array of multilayer perceptrons. It will be shown how an *a*MLP model that has learnt a classification problem with a certain (high) accuracy can be compressed into a more compact representation using a Volterra-NN model. This novel representation involves less parameters, maintaining however a high recognition accuracy. Two different face recognition databases have been used to show the effectiveness of the proposed method.

The paper is organized as follows. Section 2 explains in detail the proposed Volterra-NN model and its use as a classifier in a face recognition system. The materials and methods used in the study are presented on Section 3. Results of the model evaluation through two public face databases as well as a discussion of the experiments are shown in Section 4. Finally, Section 5 presents the conclusions and future work.

2 Volterra-Neural Network for neural networks compression

The Volterra series and Volterra theorem was developed in 1887 by Vito Volterra. It is a model for representing nonlinear dynamic behavior frequently used in system identification [Volterra, 1959]. In [Stegmayer and Chiotti, 2009] several formulas for the extraction of Volterra weights, independently of the neural model topology, number of variables involved in the problem and nonlinearity of the system have been presented. The equations are based on architectures having an hyperbolic tangent activation functions in the hidden nodes, trained with a classical backpropagation algorithm [Marquardt, 1963], for multi-input, multi-output systems. The MLP model is trained using the available training data and, after that, the Volterra weights are obtained from the trained network parameters.

This section presents the Volterra-NN model for *a*MLP model compression, extending the simple algorithms for classification proposed in [Rubiolo et al, 2010]. The following subsection presents the

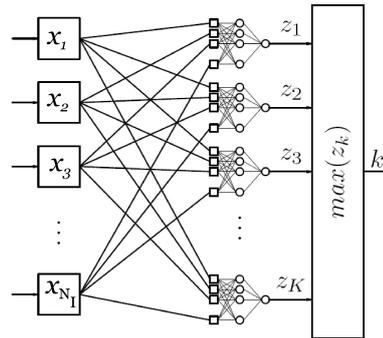


Fig. 1 Array of MLPs (*aMLP*) model for classification.

neural model used for face recognition; after that, some basic concepts on Volterra models necessary to understand the proposed approach are explained. Finally, it is shown how the Volterra-NN can perform as a classifier when an *aMLP* is used.

2.1 Neural models for classification

Arrays and ensembles of single multilayer perceptron networks have proven to reach significant better results in classification problems than single models [Aitkenhead and McDonald, 2003][Bianchini et al, 2005][Capello et al, 2009]. The classifier model is an array of MLPs where there is one MLP model for each class k to be identified, with $k=1..K$, being K the total number of classes. Therefore, the *aMLP* model is formed by K networks like the one shown in Figure 1. Each network output takes a value of 1 if the class is identified or 0 otherwise. The first layer of each MLP in the *aMLP* is a set of N_I input neurons, where each input is an eigenspace vector and there are N_H hidden neurons. When a picture has to be classified, its projected eigenspace vector is used as input for the classifier, that is to say, it is presented to all the k networks defined for the *aMLP* model, and the maximum output obtained among all network outputs is assigned as class label. If a pattern of the k th-class has been presented to the model, a value of (near) 1 is expected at the k th network.

The application of the Volterra weights extraction method for model compression of a classical multilayer perceptron classifier was presented in [Rubiolo et al, 2010]. It was shown how a MLP model which has learnt a classification problem with a certain (high) accuracy, can be compressed into a more compact representation using a Volterra model and its parameters, named Volterra weights. Several MLP topologies can be well-compressed into the first, second and third order Volterra weights, which can be used to build a Volterra model that needs less parameters than the MLP model and, at the same time, has a similar high accuracy. This work proposes to apply and extend the compression procedure based on Volterra models to an array of MLPs to solve a Face Recognition (FR) problem.

The first stage of a complete automatic FR system is face detection. Once detected, the face must be represented through an appropriate feature extraction method. A global representation can be done through a well-known technique such as the eigenfaces [Turk and Pentland, 1991] by applying Principal Component Analysis (PCA) for dimensionality reduction [Kirby and Sirovich, 1990], which is the most widely used feature extraction method for face recognition [Li and Jain, 2004]. The final stage consists of the classification carried out by using an appropriate classifier, which be a very simple method, such as the Euclidean distance or k-nearest-neighbors, or an array or ensemble solution based on hundred of weak classifiers, being MLP is one of the most popular models [Martinez and Kak, 2001] [Kong et al, 2005]. The next subsection presents the details of the new algorithm proposed for extraction of the Volterra

Algorithm 1: Volterra weights extraction.

Data:
 D : training data

Results:
 $v^{(0)}$: 0-order Volterra weight
 $\mathbf{v}^{(1)}$: 1st-order Volterra weights
 $\mathbf{v}^{(2)}$: 2nd-order Volterra weights
 $\mathbf{v}^{(3)}$: 3rd-order Volterra weights

```

1 begin
2    $M \leftarrow$  train MLP model with  $D$ 
3    $v^{(0)} \leftarrow$  calculate zero-order Volterra weight from  $M$  using (1)
4   for  $1 \leq i \leq N_I$  do
5      $v_i^{(1)} \leftarrow$  calculate first-order Volterra weight from  $M$  using (2)
6   Assign each  $v_i^{(1)}$  extracted to  $\mathbf{v}^{(1)}$ 
7   for  $1 \leq i \leq N_I$  do
8     for  $1 \leq j \leq N_I$  do
9        $v_{i,j}^{(2)} \leftarrow$  calculate second-order Volterra weight from  $M$  using (3)
10  Assign each  $v_{i,j}^{(2)}$  extracted to  $\mathbf{v}^{(2)}$ 
11  for  $1 \leq i \leq N_I$  do
12    for  $1 \leq j \leq N_I$  do
13      for  $1 \leq k \leq N_I$  do
14         $v_{i,j,k}^{(3)} \leftarrow$  calculate third-order Volterra weight from  $M$  using (4)
15  Assign each  $v_{i,j,k}^{(3)}$  extracted to  $\mathbf{v}^{(3)}$ 
16 end

```

weights from an array of MLPs used for classification, as well as the procedure used to obtain different order Volterra-NN outputs.

2.2 Volterra-NN model forward computation

[Stegmayer and Chiotti, 2009] have derived equations that allow the calculation of any Volterra kernel order using the weights and hidden neurons bias values of a neural network that has been trained on a problem using hyperbolic tangent hidden functions. The following formulas, instead, allow the calculus of zero, first, second and third order Volterra kernels from a trained MLP having sigmoidal hidden units:

$$h_0 = b_o + \sum_{h=1}^{N_H} w_h^2 \frac{1}{(1 + e^{-b_h})} \quad (1)$$

$$h_1(\cdot) = \sum_{h=1}^{N_H} w_h^2 w_{h,i}^1 \frac{e^{-b_h}}{(1 + e^{-b_h})^2} \quad (2)$$

$$h_2(\cdot) = \sum_{h=1}^{N_H} w_h^2 w_{h,i}^1 w_{h,j}^1 \frac{e^{-b_h} (e^{-b_h} - 1)}{(1 + e^{-b_h})^3} \frac{1}{2!} \quad (3)$$

$$h_3(\cdot) = \sum_{h=1}^{N_H} w_h^2 w_{h,i}^1 w_{h,j}^1 w_{h,k}^1 \frac{-e^{-b_h} (-e^{-2b_h} + 4e^{-1} - 1)}{(1 + e^{-b_h})^4} \frac{1}{3!} \quad (4)$$

where N_H is the number of hidden neurons, N_I is the number of input neurons, w_h and b_h are the weight and the bias associated with a hidden sigmoidal neuron, respectively; for $i, j, k = [1, \dots, N_I]$. These

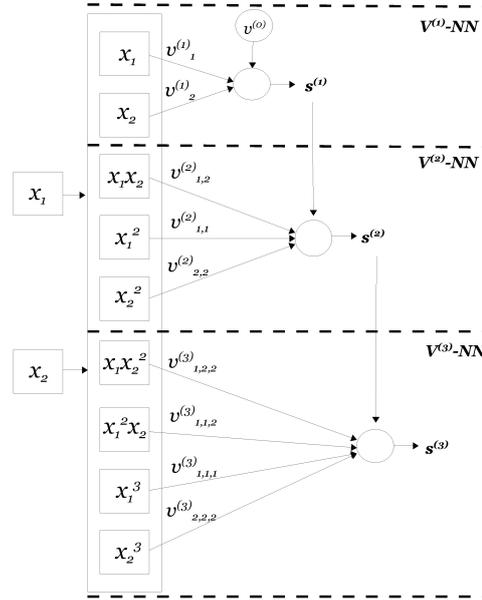


Fig. 2 Example of the topology of a Volterra-NN model for 2 input variables and three V-NN outputs.

formulas are easily extended to any kernel order. Due to space restrictions, only up to third order is shown. To simplify the notation, the Volterra weights will be defined from now on as follows: $v^{(0)} = h_0$, $v_i^{(1)} = h_1(\cdot)$, $v_{i,j}^{(2)} = h_2(\cdot)$ and $v_{i,j,k}^{(3)} = h_3(\cdot)$. Algorithm 1 shows in detail the Volterra weights extraction procedure, in which it is possible to compute the zero, first, second and third-order Volterra weights. The input is the training data and the outputs are the Volterra weights. The first step consists of training a MLP classifier model with the training data D as it is possible to see in line 2. From the trained neural model M , the Volterra weights can be calculated. According to (1), the zero-order Volterra weight is obtained (line 3). Similarly, by applying (2), (3) and (4), it is possible to compute the first (line 9), second (line 8) and third-order (line 7) Volterra weights respectively.

The different order Volterra-NN (V-NN) models that can compress a MLP are depicted graphically in Figure 2. The boxes represent the input variables (in the example, x_1 and x_2) and the arrows that join them symbolize their product with their corresponding Volterra weights. The white circles act as summarizing all the products between input variables and Volterra weights, plus the $(n - 1)$ -order V-NN model ($V^{(n-1)} - NN$). As a result, a new n^{th} -order Volterra model ($V^{(n)} - NN$) is obtained. It is important to highlight that the different cross-products combinations between the input variables are shown in the figure inside a rectangular box, in an effort to clarify the process for obtaining the different Volterra outputs. This V-NN model can be similarly applied to any number of inputs.

The output of the 1^{st} -order V-NN model ($V^{(1)} - NN$) is obtained analytically by adding the 0-order Volterra weight $v^{(0)}$ to the product between each input variable (x_1 and x_2) and their corresponding 1^{st} -order Volterra weights,

$$s^{(1)} = v^{(0)} + v_1^{(1)} x_1 + v_2^{(1)} x_2. \quad (5)$$

Algorithm 2: 3^{rd} -order Volterra-NN outputs forward computation.

Data:
 \mathbf{x} : input point to classify
 K : number of elements in the array
 $\mathbf{v}^{(0)}$: 0-order Volterra weights for the array
 $\mathbf{v}^{(1)}$: array 1^{st} -order Volterra weights
 $\mathbf{v}^{(2)}$: array 2^{nd} -order Volterra weights
 $\mathbf{v}^{(3)}$: array 3^{rd} -order Volterra weights

Results:
 $\mathbf{s}^{(1)}$: $V^{(1)} - NN$ outputs for the array
 $\mathbf{s}^{(2)}$: $V^{(2)} - NN$ outputs for the array
 $\mathbf{s}^{(3)}$: $V^{(3)} - NN$ outputs for the array

```

1 begin
2   for each element in the array do
3      $s^{(1)} \leftarrow v^{(0)}$ 
4     for  $1 \leq i \leq N_I$  do
5        $s^{(1)} = s^{(1)} + v_i^{(1)} x_i$ 
6      $s^{(2)} \leftarrow s^{(1)}$ 
7     for  $1 \leq i \leq N_I$  do
8       for  $1 \leq j \leq N_I$  do
9          $s^{(2)} = s^{(2)} + v_{i,j}^{(2)} x_i x_j$ 
10     $s^{(3)} \leftarrow s^{(2)}$ 
11    for  $1 \leq i \leq N_I$  do
12      for  $1 \leq j \leq N_I$  do
13        for  $1 \leq k \leq N_I$  do
14           $s^{(3)} = s^{(3)} + v_{i,j,k}^{(3)} x_i x_j x_k$ 
15    Assign each  $s^{(1)}$  calculated to  $\mathbf{s}^{(1)}$ 
16    Assign each  $s^{(2)}$  calculated to  $\mathbf{s}^{(2)}$ 
17    Assign each  $s^{(3)}$  calculated to  $\mathbf{s}^{(3)}$ 
18    Determine upper and lower thresholds for each array element (class)
19 end

```

The second order V-NN model ($V^{(2)} - NN$) output is obtained by adding $s^{(1)}$ together with the products among x_1^2, x_2^2 , the cross-product $x_1 x_2$ and their corresponding 2^{nd} -order Volterra weights, resulting in:

$$s^{(2)} = s^{(1)} + v_{1,1}^{(2)} x_1^2 + v_{2,2}^{(2)} x_2^2 + v_{1,2}^{(2)} x_1 x_2. \quad (6)$$

The $V^{(3)} - NN$ model output is obtained similarly:

$$s^{(3)} = s^{(2)} + v_{1,1,1}^{(3)} x_1^3 + v_{2,2,2}^{(3)} x_2^3 + v_{1,1,2}^{(3)} x_1^2 x_2 + v_{1,2,2}^{(3)} x_1 x_2^2. \quad (7)$$

Similarly, higher order V-NN outputs can be calculated. As can be seen, each high order Volterra model includes its own parameters (Volterra weights) plus the lower order ones. The new algorithm for an array of V-NN models forward computation is presented in detail in Algorithm 2. It receives a point from where the number of input variables N_I is determined, the number of elements in the array and the Volterra weights for each element in the array, obtained by using Algorithm 1. The output of this algorithm is an array of third-order ($V^{(3)} - NN$) Volterra outputs, but it is also possible to obtain only the first-order ($V^{(1)} - NN$) (lines 3 to 5) and second-order ($V^{(2)} - NN$) outputs (lines 7 to 10), thus providing different compressed versions of the original aMLP classifier. The next subsection shows how the compressed model can be used as a classifier.

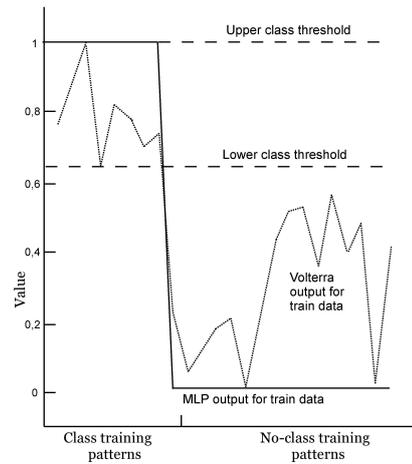


Fig. 3 Model output signal analysis for classification. Full and dotted lines: output signal corresponding to the MLP and V-NN models, respectively, for training data. Dashed lines: class thresholds.

2.3 Volterra-NN as a classifier

The Volterra-NN model can be applied to a classification problem in which different classes can be recognized from the output signal. When an array of MLP models is compressed by using V-NN, it is possible to identify only a class from each of the output signals. That is to say, each V-NN model is specialized on a class and the output signal analysis determines if the pattern is part of the class. During the training phase, data are shown in an ordered way, showing first those pattern belonging to the class associated to the model. In this way, the model output can be considered as a signal having two levels, with a bound between levels corresponding to the class limit [Korenberg et al, 2001]. For instance, if we have a three classes problem, the first model is trained with data where the patterns corresponding to the first-class are activated (they have a 1 as target) and other patterns, which do not represent the first class, are not activated (they have a 0 as target). Similarly, training data for learning the second and third classes are presented to their corresponding models.

Figure 3 shows the output signal analysis that has to be performed for determining upper and lower thresholds for each class. It is possible to see the output signal corresponding to the MLP (full line) and V-NN (dotted line) models, both for the training data. The lower and the upper threshold (dashed lines) of the class are measured in order to determine the output-signal level change for the V-NN model, considering the values obtained for the training set. These changes (thresholds) will allow us to identify a new data point received for classification as belonging (or not) to the class. The membership of a new pattern (test point) to the class is identified by analyzing if the output associated to this pattern has a value between the lower and upper class thresholds.

Since we are analyzing the output of an *a*MLP model, *a* different signals, each one corresponding to a class, must be considered in order to identify which class is activated as a result. First of all, each individual signal must be evaluated as it was presented in the above paragraph. Secondly, it is necessary to discover if more than one class model was activated for each pattern. If this situation occurs, a *max* criteria is applied. That is to say, the higher value of the activated classes for each pattern will be stated as the winner class.



Fig. 4 Samples face images from the AT&T Laboratories Cambridge ORL Database of faces [Li and Jain, 2004].

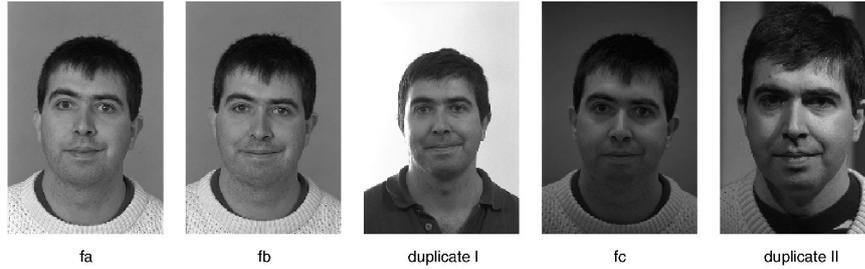


Fig. 5 The Facial Recognition Technology (FERET) database [Phillips et al., 2000]. Examples of different categories of probes (images). The duplicate I image was taken within one year of the fa image and the duplicate II and fa images were taken at least one year apart.

3 Materials and Methods

3.1 Face Databases

Two independent face databases have been used in this study. They provide typical experimental setups for face recognition. They are explained in detail in the following paragraphs.

3.1.1 AT&T Laboratories - ORL

The first experiments were done by using the AT&T Laboratories Cambridge ORL Database of Faces¹ since it is widely used in the face recognition literature [Li and Jain, 2004]. In this database, there are 10 different images of each of 40 persons of different gender, ethnic background and age (see Figure 4). For some subjects, the images had been taken at different times, varying the lighting, facial expressions (open/closed eyes, smiling/not smiling) and facial details (glasses/no glasses). All the images were taken against a dark homogeneous background with the subjects in an upright, frontal position (with tolerance for some side movement). The complete database contains 400 grayscale face images of size 92 x 112 pixels. From this dataset, three subject were used in this study having ten pictures associated with each one.

Training a distinct classifier for each class (in this case, subject) requires sufficient training data per class. However, in face recognition tasks it is common to have a small number of pictures per person. In fact, since different data partitions are used to train and validate the classifier, two pictures for subject are available for testing. Hence, noise addition to the test dataset has been performed in order to enlarge it and to obtain more test samples to prove the performance of the Volterra-NN model. The noise used was *Gaussian noise* with a mean of 0 and a variance between 0.01 and 0.1. After this procedure, a total of 66 patterns per class have been obtained for the testing dataset.

¹ www.cl.cam.ac.uk/research/dtg/attarchive/face/database.html

3.1.2 FERET

The Facial Recognition Technology (FERET) database [Phillips et al, 2000] ran from 1993 through 1997², sponsored by the Department of Defense's Counterdrug Technology Development Program through the Defense Advanced Research Products Agency (DARPA). The aim was to develop automatic face recognition capabilities that could be employed to assist security, intelligence and law enforcement personnel in the performance of their duties. The FERET image corpus was assembled to support government monitored testing and evaluation of face recognition algorithms using standardized tests and procedures. The final corpus consists of 14051 eight-bit grayscale images of human heads with views ranging from frontal to left and right profiles.

The facial images were collected in 15 sessions between August 1993 and July 1996. Collection sessions lasted one or two days. To maintain a degree of consistency throughout the database, the same physical setup and location was used in each photography session. Images of an individual were acquired in sets of 5 to 11 images. Two frontal views were taken (fa and fb); a different facial expression was requested for the second frontal image. For 200 sets of images, a third frontal image was taken with a different camera and different lighting (fc image). The remaining images were collected at various angles between right and left profile (see Figure 5). To add variations to the database, a second set of images was taken, for which the subjects were asked to put on their glasses and/or pull their hair back. The set of images referred to as a duplicate indicates a second set of images of a person that was taken on a later date, resulting in variations in scale, pose, expression, and illumination of the face. Similarly to the previous dataset, three subject were used in the study having ten pictures associated with each one, adding gaussian noise to the pictures to increase the number of training/testing patterns.

3.2 Classifier architecture

Figure 6 shows the neural network topology used in this study, an array of MLP models, each one associated with a subject to recognize. Several MLP topologies have been evaluated from where the corresponding Volterra weights have been extracted after training. For simplicity in the analysis of the results, in particular which respects the performance of the V-NN compression capabilities, only three classes ($k=3$) of each database are presented in the tables. The full results obtained on both face databases can be found as supplementary material.

In face recognition problems, the training and test datasets generally have high dimensionality due to the pictures size. Therefore, an appropriate feature extraction method is needed. A global representation can be done using a widely used technique such as the eigenfaces [Turk and Pentland, 1991] by applying Principal Component Analysis (PCA) for dimensionality reduction [Kirby and Sirovich, 1990]. Although PCA can highly reduces the feature space, the application of a technique called *Scree Test* [Jackson, 1991] further reduces the dimensionality of the features space focusing only on those most representative eigenfaces. Scree Test is applied for determining the number of principal componenets of the training and test datasets, according to the percentage of variability of the data to be shown. This technique allows obtaining the different components in an orderly way acording to the variability of the data, so the first principal components represent the data of more variability.

Regardless of the data that are modeled with PCA, it is common to use a value of 85% of the total variance of the space of characteristics for identification [Jackson, 1991]. Applying Scree Test in the ORL dataset, the 85% of the variance of each training set is represented by 11 eigenfaces, so each MLP model consists of 11 input neurons. As regards hidden neurons number, a simple heuristic will be adopted in which the number of neurons N_H are 11,22 and 33; and an output that takes a value of 1 if the subject is recognized. Considering the case of the FERET dataset, the 85% of the variance of each training set is

² www.itl.nist.gov/iad/humanid/feret

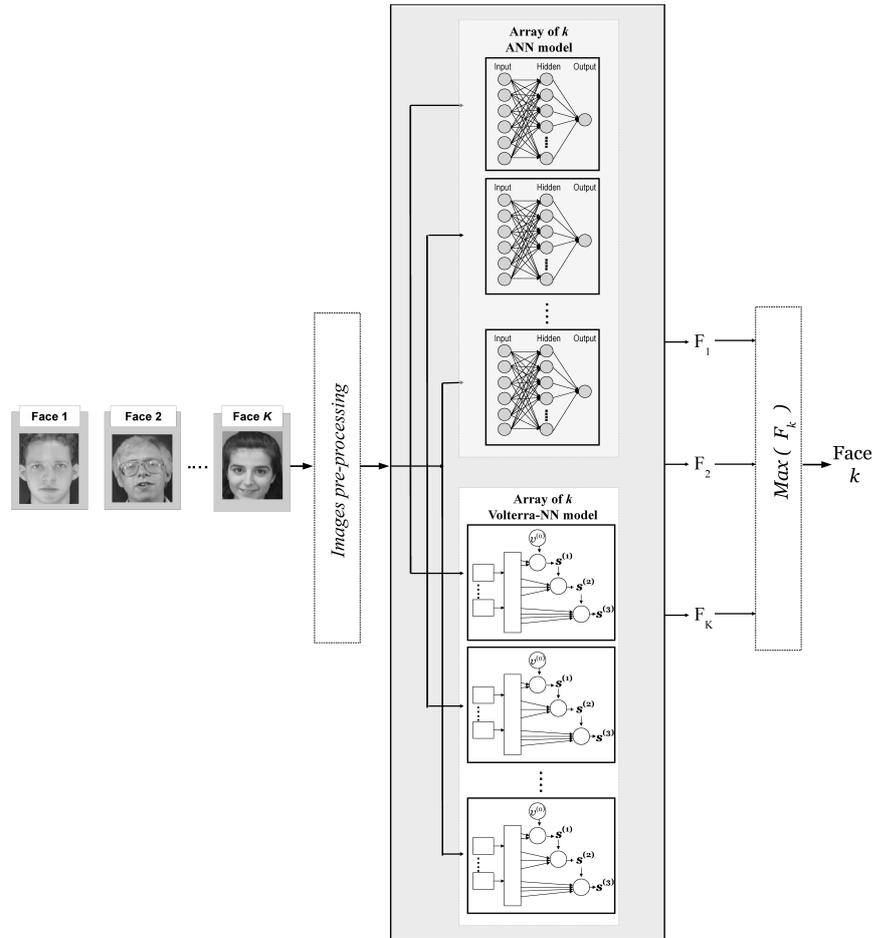


Fig. 6 aMLP model (upper part of the classifier) for Face Recognition problem, which can be compressed into an array of Volterra-NN models (lower part of the classifier)

represented by 9 eigenfaces. Therefore, each MLP model consist of 9 input neurons, a number of hidden neuron that can be 9,18 or 27, and also only one output. For each input signal arriving at the array, a V-NN output of a certain order can be obtained. For example, to obtain a first order $V^{(1)} - NN$ output, the maximum $s^{(1)}$ from the array is selected.

The model parameters (weights and biases) are initialized with random values uniformly distributed between 0 and 1. Neurons in the hidden and output layers have sigmoid activation functions. All the MLPs are trained with the Levenberg-Marquardt algorithm [Madsen et al, 2004] in order to guarantee a fast convergence. To avoid overfitting, a k -fold cross-validation procedure [Haykin, 1999] has been used, using the standard setup of splitting the available images of each person into 80% of each class data for training and 20% for testing. The complete dataset has been randomly split into $k = 3$ mutually exclusive subsets of equal size, repeating the experiments three times in each fold. The cross validation estimate of the overall accuracy of a model has been calculated by simply averaging the accuracy measures over the test datasets. In each experiment and repetition, MLP and V-NN models having less than 100% classification rate for each class of the training dataset have not been considered in this study. This restriction was imposed to the classifier performance in order to be able to measure precisely any loss of accuracy originated by the proposed Volterra-NN compression method.

3.3 Performance measures

This subsection presents two classical measures for comparing models performance. Besides, a new measure for trade-off analysis and final model selection is proposed.

3.3.1 Recognition rate and space saving rate

For comparing each output from the proposed V-NN models against the corresponding MLP classifier, two performance measures are used: recognition rate (RR) and data space savings (SS), adapted here for an *a*MLP classifier. In classification problems, the primary source of performance measurements is the overall accuracy of a classifier estimated through the classification or recognition rate [Duda and Hart, 2003]. For measuring compression, data compression ratio can be used to quantify the reduction in data representation size produced by a compression algorithm. However, the space saving measure is given here instead, defined as the reduction in size relative to the uncompressed space [Salomon, 2007], often reported as a percentage, which gives a better idea of compression power. For both cases, the greater the rate, the better the result.

To estimate how much is the compression level obtained, SS is calculated as the relation between the number of parameters needed for a Volterra-NN output (the Volterra weights) and the number of parameters of each corresponding MLP architecture (the weights and biases). It is well-known that for a MLP model, the number of model parameters can be calculated as follows:

$$P_{MLP} = N_I \times N_H + N_{BH} + N_H \times N_O + N_{BO}, \quad (8)$$

where N_I , N_H and N_O are the number of neurons at input, hidden and output layers, respectively; and N_{BH} and N_{BO} are the number of bias for each neuron in the hidden and output layers, respectively. In an array of MLP models, these measures are affected by the array length a . Hence, the number of model parameters for an *a*MLP can be calculated as

$$P_{aMLP} = a \times P_{MLP}. \quad (9)$$

From each MLP that is part of the array of MLPs, a Volterra-NN output (of a particular order) can be extracted: $s^{(1)}$ that includes only the zero and first-order Volterra weights; $s^{(2)}$ includes up to the second-order Volterra weights; and $s^{(3)}$ that corresponds to a third-order model output. The following equations show the number of parameters necessary to build each of these outputs, for a given training set.

For the $V^{(1)} - NN$ model output we have

$$P_{s^{(1)}} = N_I. \quad (10)$$

In this case, the number of parameters necessary to build $s^{(1)}$ are each first-order Volterra weight that multiplies each input variable.

For the second-order model $V^{(2)} - NN$ we have the following output

$$P_{s^{(2)}} = P_{s^{(1)}} + N_I^2 - \frac{N_I^2 - N_I}{2}. \quad (11)$$

The number of parameters necessary to build $s^{(1)}$ are summed up together with the number of parameters necessary for $s^{(2)}$ (the second-order Volterra weights). There is a second-order weight for each input variable squared, plus the cross-products between each input variable. With respect to these last ones, since the symmetrical weights are equivalent, they are considered only once. That is why half of the cross-product weights are counted.

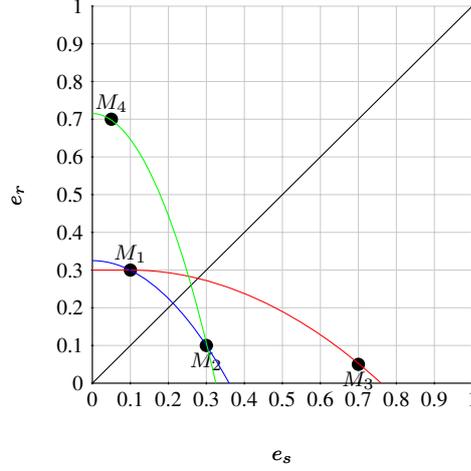


Fig. 7 Trade-off measure determination by analyzing the possible solutions of compressing a model.

The number of parameters necessary to build $V^{(3)} - NN$ are

$$P_{s^{(3)}} = P_{s^{(2)}} + N_I^2, \quad (12)$$

that is to say, the number of weights associated with $s^{(2)}$ plus the product of each third-order weight corresponding to each input variable, counting only once the symmetrical weights.

As Volterra-NN can be applied to the compression of an $aMLP$ model, the output now will be an array of V-NN outputs of different order. That is to say, for example, for each MLP inside the array three different order V-NN outputs can be obtained. In this case, the number of V-NN models parameters has to be redefined as

$$P_{as^{(i)}} = a \times P_{s^{(i)}}, \quad (13)$$

where $i = 1, 2, 3$. Therefore, the space saving measure SS is calculated as

$$SS = 1 - \frac{P_{as^{(i)}}}{P_{aMLP}}. \quad (14)$$

3.3.2 Model selection: measure for trade-off solution indication

Different solutions can be evaluated by calculating the performance measures presented above. A V-NN model output can have a value associated to the recognition rate or to the space saving measures. However, sometimes these values can differ from each other having even completely opposite meaning and the problem of how determining the best set of parameters arises. For instance, one solution can have a high performance as regards RR but a very low SS rate, or vice versa.

It is possible to consider the error for both measures as a point in a coordinate axis graph, where the X axis represents $e_s = 1 - SS$, and the Y axes is associated to $e_r = 1 - RR$. Therefore, a point can be defined as a pair $[e_s, e_r]$, being $[0, 0]$ the optimum. Figure 7 shows examples of models and their associated errors in this new space. Let us suppose a model 1 that has $SS=0.9$ and $RR=0.7$ and an error to the optimum represented in the point $M_1 = [0.1, 0.3]$. For model 2, where $SS=0.7$ and $RR=0.9$, the error is $M_2 = [0.3, 0.1]$. Model 3, in which $SS=0.3$ and $RR=0.95$, has an error $M_3 = [0.7, 0.05]$; and a model 4, with $SS=0.95$ and $RR=0.3$, has an error $M_4 = [0.05, 0.7]$.

Considering these errors, a trade-off solution could be found as the minimum Euclidean distance from the point $[e_s, e_r]$ to the optimum $[0,0]$. However, when calculating these distances, it is possible to have cases with the same result, but opposite values. For example, points M_1 and M_2 are equidistant from the main diagonal and have the same Euclidean distance to the optimum; similarly to any other models that could be located on the blue parabola depicted in Figure 7. However, M_2 is a better classifier than M_1 , while M_1 is better compressed than the first one. Hence, it is necessary to discriminate which of the two points is the best solution for the problem under study.

Since FR is a classification problem, where it is reasonable to think that a better RR would be considered more important than SS, it is possible to infer that the trade-off solution must always be above the diagonal of the coordinate axes graph, because in a classification problem it could be more important to have a high recognition rate than a high space savings rate. But, in some applications, it can be needed a better compressed model than an excellent classifier; in that case SS would be more important than RR.

In order to evaluate which is the trade-off solution (more adequate model) for a problem, in the family of solutions that arose from the experiences, we define a new measure ∂ which represents a trade-off between RR and SS as follows:

$$\partial = \sqrt{(\gamma e_r)^2 + ((1 - \gamma)e_s)^2}, \quad (15)$$

where γ is a regularization parameter. Precisely, to be able to perform the discrimination needed above, we propose to modify the distance calculation including a regularization parameter γ . In order to consider γ as a weight that allows us to select between two models, its value can be modified according to whether a better compressor or a better classifier is needed, with $\gamma \in [0, \dots, 1]$.

When $\gamma = 0.5$, both e_r and e_s are considered with the same importance; the preference for *RR* is emphasized with a $\gamma > 0.5$ (for example, red line in the figure), while the preference for *SS* is obtained choosing a $\gamma < 0.5$ (for example, green line in the figure). Applying (15) to each model M_k , presented in Figure 7 and considering RR as a preference over SS with, for example $\gamma = 0.8$, we obtain $\partial_{M_1} = 0.241$, $\partial_{M_2} = 0.1$, $\partial_{M_3} = 0.146$, and $\partial_{M_4} = 0.56$. Now it is possible to conclude, without any doubt, that the best compromise solution is M_2 for the presented example, since it has the minimum distance to the optimum, according to the new distance calculation proposed.

4 Results and discussion

The experimental results obtained on two well-known face recognition problems are shown in this section. First of all, the results on class by class recognition rates are presented. Then, global RR and SS values for all the models obtained are shown. At last, global results using the new measure for model selection are presented.

4.1 Recognition performance for each class

From each $aMLP_{I,H,O}$ model considered in this study, their corresponding zero-order, first-order, second-order and third-order Volterra weights have been extracted according to Algorithm 1 and their corresponding array of Volterra-NN outputs $s^{(1)}$, $s^{(2)}$ and $s^{(3)}$ have been obtained using Algorithm 2. Table 1 presents the results obtained for the models recognition rate (RR_i) for each class $i = 1, 2, 3$ over the face recognition task, calculated over the ORL Database test set; whereas Table 2 shows similar results but calculated over the FERET database.

In Table 1 it is possible to see that, when the *aMLP* classifier has 11 hidden neurons ($3MLP_{11,11,1}$), the RR values regarding class 1 (RR_1) are between 91% and 94% for all $3V-NN_{s^{(1)}}$, $3V-NN_{s^{(2)}}$ and $3V-NN_{s^{(3)}}$; for the second class (RR_2), these values vary from 90% to 94%; and they are between 91%

Table 1 Recognition rates (RR_i) for each class $i = 1, 2, 3$ for three aV -NN ($a=3$) classifiers and their corresponding first-order $s^{(1)}$, second-order $s^{(2)}$ and third-order $s^{(3)}$ outputs, for the AT&T Laboratories Cambridge ORL Database of Faces. Bold numbers highlight the best value of RR for each order V-NN output.

RR_i [%]	$3MLP_{11,11,1}$			$3MLP_{11,22,1}$			$3MLP_{11,33,1}$		
	RR_1	RR_2	RR_3	RR_1	RR_2	RR_3	RR_1	RR_2	RR_3
$aMLP_{I,H,O}$	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
$3V-NN_{s^{(1)}}$	94.28	93.94	97.47	92.76	92.25	91.92	97.14	91.92	94.28
$3V-NN_{s^{(2)}}$	91.92	92.09	91.25	92.76	91.24	94.27	95.29	94.44	90.57
$3V-NN_{s^{(3)}}$	91.58	90.23	91.58	90.23	89.06	88.89	86.19	90.74	93.26

Table 2 Recognition rates (RR_i) for each class $i = 1, 2, 3$ for three aV -NN ($a=3$) classifiers and their corresponding first-order $s^{(1)}$, second-order $s^{(2)}$ and third-order $s^{(3)}$ outputs, for the Facial Recognition Technology (FERET) Database. Bold numbers highlight the best value of RR for each order V-NN output.

RR_i [%]	$3MLP_{9,9,1}$			$3MLP_{9,18,1}$			$3MLP_{9,27,1}$		
	RR_1	RR_2	RR_3	RR_1	RR_2	RR_3	RR_1	RR_2	RR_3
$aMLP_{I,H,O}$	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
$3V-NN_{s^{(1)}}$	94.11	94.27	95.62	86.70	95.62	93.94	88.38	91.41	96.63
$3V-NN_{s^{(2)}}$	91.08	88.38	91.58	88.89	87.88	94.11	92.25	90.07	95.29
$3V-NN_{s^{(3)}}$	90.40	85.35	93.10	89.56	85.52	85.86	91.92	84.34	85.01

and approximately 97% for the third class (RR_3). As regards the second FR problem (Table 2) when the $aMLP$ classifier has 9 hidden neurons ($3MLP_{9,9,1}$), the RR values are similarly high for class 1, 2 and 3.

Focusing on Table 1, it can be seen that a worse RR rate is obtained for the classification problem if $3V-NN_{s^{(3)}}$ is used instead of the original $aMLP$, specially for class 2. This lower RR is, still, of 90%. Furthermore, the $3V-NN_{s^{(1)}}$ output requires a significant lower number of parameters than the original $3MLP_{11,11,1}$ classifier (a relation of 1 : 12). Similar conclusions can be drawn from Table 2, in which the FERET Database is used: the simplest V-NN output is, at the same time, the best one. But if we take into account that the three $3MLP_{I,H,1}$ -model topologies in both cases, with $N_I = 11, N_H = 11, 22, 33$ for ORL and $N_I = 9, N_H = 9, 18, 27$ for FERET, are solutions to the same problem, actually it is not necessary to consider just only one topology as the best solution. That is to say, there are many possible solutions to the same problem and it is possible to obtain the best according to the goals that are followed. For instance, the best recognition rate in ORL Database case is obtained by the first-order Volterra-NN output for $3MLP_{11,33,1}$, which is approximately as high as 97.50%; while the best recognition rate in FERET Database case is obtained by the first-order Volterra-NN output for $3MLP_{9,27,1}$, achieving an RR of 96.63%.

Tables 3 and 4 show the space savings rate for the models discussed in this paper and the mean values of the recognition rate for these models. The upper part of each table shows the number of parameters and global RR measure for the three $aMLP$ topologies. The second part of the table shows parameters and performance measure values related to the Volterra-NN outputs. First of all, the number of parameters needed for each neural classifier architecture (P_{aMLP}) involved in this study is shown at the first row, while the number of Volterra weights needed for each Volterra-NN model ($P_{as^{(i)}}$) is shown in the first column. The values which fill the RR and SS columns are the average recognition rate and space saving capabilities, respectively, for each combination between a MLP classifier and the corresponding Volterra-NN output.

As stated before (equations (8) and (9)), the number of parameters required for the neural classifier depends not only on N_I but also on the number of neurons at the hidden layer N_H , as well as in the number of elements in the array. Therefore, considering the ORL Database, the number of parameters for $3MLP_{11,11,1}$ is 432, for $3MLP_{11,22,1}$ it is 861, and the $3MLP_{11,33,1}$ model has a total of 1290 parame-

Table 3 Global recognition rate (RR) and space saving(SS) comparison for three aV -NN ($a=3$) classifiers and their corresponding first-order $s^{(1)}$, second-order $s^{(2)}$ and third-order $s^{(3)}$ outputs, for the the AT&T Laboratories Cambridge ORL Database of Faces.

		$3MLP_{11,11,1}$		$3MLP_{11,22,1}$		$3MLP_{11,33,1}$	
$P_{aMLP_{I,H,O}} \rightarrow$		432		861		1290	
RR[%] \rightarrow		100.00		100.00		100.00	
$3V$ -NN	$P_{as(i)}$	RR[%]	SS[%]	RR[%]	SS[%]	RR[%]	SS[%]
$3V$ -NN $_{s^{(1)}}$	33	95.23	92.36	92.31	96.16	94.44	97.44
$3V$ -NN $_{s^{(2)}}$	231	91.75	46.53	92.76	73.17	93.43	82.09
$3V$ -NN $_{s^{(3)}}$	594	91.13	-37.50	89.39	31.01	90.07	53.95

Table 4 Global recognition rate (RR) and space saving(SS) comparison for three aV -NN ($a=3$) classifiers and their corresponding first-order $s^{(1)}$, second-order $s^{(2)}$ and third-order $s^{(3)}$ outputs, for the Facial Recognition Technology (FERET) Database.

		$3MLP_{9,9,1}$		$3MLP_{9,18,1}$		$3MLP_{9,27,1}$	
$P_{aMLP_{I,H,O}} \rightarrow$		300		597		894	
RR[%] \rightarrow		100.00		100.00		100.00	
$3V$ -NN	$P_{as(i)}$	RR[%]	SS[%]	RR[%]	SS[%]	RR[%]	SS[%]
$3V$ -NN $_{s^{(1)}}$	27	94.67	91.00	92.09	95.48	92.14	96.98
$3V$ -NN $_{s^{(2)}}$	162	90.35	46.00	90.29	72.86	92.54	81.88
$3V$ -NN $_{s^{(3)}}$	405	89.62	-35.00	86.98	32.16	87.09	54.70

ters (see Table 3). Regarding the FERET Database, the number of parameters for $3MLP_{9,9,1}$ is 300, for $3MLP_{9,18,1}$ it is 597, and the $3MLP_{9,27,1}$ model has a total of 894 parameters (see Table 4).

4.2 Recognition rate and space saving rate

In the case of the different order Volterra-NN outputs, the number of parameters only depends on the number of inputs, as it is possible to see in (10), (11) and (12). However, for an array, this number must be multiplied by the array size. Therefore, the number of parameters for $3V$ -NN $_{s^{(1)}}$ is 33, and the number of parameters needed for $3V$ -NN $_{s^{(2)}}$ and $3V$ -NN $_{s^{(3)}}$ are 231 and 594, respectively. In the FERET Database, the number of parameters for $3V$ -NN $_{s^{(1)}}$ is 27, and the number of parameters needed for $3V$ -NN $_{s^{(2)}}$ and $3V$ -NN $_{s^{(3)}}$ are 162 and 405, respectively.

Focusing on Table 3, it is possible to see that when using $3V$ -NN $_{s^{(1)}}$ instead of the $3MLP_{11,11,1}$ classifier, a global recognition rate for the face recognition problem of 95.23% can be obtained and a compression or space saving rate of 92.36% can be achieved. Almost half of this space saving rate is obtained if the $3V$ -NN $_{s^{(2)}}$ output is used, and there is no compression when using $3V$ -NN $_{s^{(3)}}$. For the $3MLP_{11,22,1}$ model, the compression achieved by the Volterra-NN models is higher because of the amount of parameters associated with this model; in this case, the SS value achieves a maximum of 96.16% when the smallest possible number of parameters is considered (a first-order Volterra-NN output). A similar case is related to $3MLP_{11,33,1}$ model, where this trend is also verified. The results obtained with the other dataset are quite similar (Table 4). For instance, a global recognition rate of 94.67% can be obtained and a compression or space saving rate of 91% can be achieved when the $3V$ -NN $_{s^{(1)}}$ is used instead of the $3MLP_{9,9,1}$ classifier.

From both Tables 3 and 4 it is possible to conclude that an array of MLPs for classification can be well-compressed by using Volterra-NN model; and this compression rate can be, in some cases, even higher than 90%. Moreover, very high recognition rates are achieved. In fact, the $aMLP$ model architecture can be more and more complex, can have many more hidden units, but the number of weights needed to build each Volterra-NN output will remain the same while the number of input variables of the problems are the same, and this will certainly be reflected in even higher space saving rates.

Table 5 RR with the same SS in aV -NN, OBD, and OBS (ORL Database).

$3MLP_{11,11,1}$			
$P_{aMLP_{I,H,O}} \rightarrow$		432	
RR[%] \rightarrow		100.00	
	$P_{as(i)}$	RR[%]	SS[%]
$3V$ -NN ⁽¹⁾		95.23	92.36
OBD	33	40.57	
OBS		61.11	

Table 6 SS in aV -NN, OBD, and OBS, with the same RR (ORL Database).

$3MLP_{11,11,1}$			
$P_{aMLP_{I,H,O}} \rightarrow$		432	
RR[%] \rightarrow		100.00	
	$P_{as(i)}$	RR[%]	SS[%]
$3V$ -NN ⁽¹⁾	33	≈ 95	92.36
OBD	192		58.44
OBS	285		34.03

Table 7 RR with the same SS in aV -NN, OBD, and OBS (FERET Database).

$3MLP_{9,9,1}$			
$P_{aMLP_{I,H,O}} \rightarrow$		300	
RR[%] \rightarrow		100.00	
	$P_{as(i)}$	RR[%]	SS[%]
$3V$ -NN ⁽¹⁾		94.67	91.00
OBD	27	38.38	
OBS		40.24	

Table 8 SS in aV -NN, OBD, and OBS, with the same RR (FERET Database).

$3MLP_{11,11,1}$			
$P_{aMLP_{I,H,O}} \rightarrow$		300	
RR[%] \rightarrow		100.00	
	$P_{as(i)}$	RR[%]	SS[%]
$3V$ -NN ⁽¹⁾	27	≈ 95	91.00
OBD	204		32.00
OBS	198		34.00

Furthermore, in some cases, these high compression rates have not to be payed by lower model classification rates. For example, for the $3MLP_{11,11,1}$ model in the ORL Database, the best RR value is related to the $3V$ -NN_{s(1)} output, achieving 95.23%, and having a space saving rate of more than 92%. The exactly same case can be seen in the FERET Database, in which the $3MLP_{9,9,1}$ model achieves an RR of 94.67%, having a space saving rate of 91%.

Considering the $3MLP_{11,22,1}$ case, in Table 3 the best RR value is associated with the $3V$ -NN_{s(2)} output, which has only a SS of 73.17%. But if a better compression is necessary for this topology, it is possible to choose the $3V$ -NN_{s(1)} which has a SS of 96.16% and preserves a very high RR value of 92.31%. Even though the differences are minimum between $3V$ -NN_{s(1)} and $3V$ -NN_{s(2)} outputs regarding the second model in both tables, in Table 4 the results are slightly different if it is considered the model $3MLP_{9,18,1}$. Here, the best RR and SS are obtained in the same case, achieving a space saving of 95.48% preserving the 92.09% of recognition ability if the $3V$ -NN_{s(1)} output is chosen.

For the $3MLP_{11,33,1}$ model, in Table 3, the best RR value is associated with the $3V$ -NN_{s(1)} Volterra-NN model and it is possible to achieve a SS of 97.44%. This is a very interesting result, because with approximately less than 3% of the parameters of the original $aMLP$ classifier, the classification capacity of $3V$ -NN_{s(1)} is very high (94.44%) and very close to the $aMLP$ model. This particular case can be considered as the best overall result obtained in this study, where the $3MLP_{11,33,1}$ classifier can be compressed in almost a 98% using the corresponding $3V$ -NN_{s(1)} Volterra-NN output without significantly losing recognition capability. Similar results can be highlighted in Table 4. Once again, a small difference of 0.4% can be seen as regards the output ($3V$ -NN_{s(1)} or $3V$ -NN_{s(2)}) that is related to the best RR obtained in the $3MLP_{9,27,1}$ classifier. Despite this, it is possible to state that the best overall result obtained in this case is related to the $3V$ -NN_{s(1)} Volterra-NN output, in which the classifier can be compressed in almost a 97%, maintaining a 92% of recognition rate.

We have compared our proposal against two classical weight pruning algorithms (as simpler ways of compression) for MLP models, such as Optimal Brain Damage (OBD) [Cun et al, 1990] and Optimal Brain Surgeon (OBS) [Hassibi et al, 1993]. Two aspects have been compared: i) RR of each method while maintaining the same number of parameters; ii) SS considering approximately the same RR value for all the three methods. In ORL database case, for the 33 parameters used for $3V$ -NN_{s(1)} in order to achieve an RR of 95.23%, OBS and OBD obtained RR values of 40.57% and 61.11%, respectively (see

Table 9 Measure for trade-off solution ∂ comparison for 3Volterra-NN models corresponding to a 3MLP classifier, for the the AT&T Laboratories Cambridge ORL Database of Faces. The best value for each model is highlighted in bold.

∂	3V-NN	$\gamma_1 = 0.25$	$\gamma_2 = 0.5$	$\gamma_3 = 0.75$
3MLP _{11,11,1}	3V-NN _{s(1)}	0.059	0.045	0.041
	3V-NN _{s(2)}	0.402	0.271	0.147
	3V-NN _{s(3)}	1.031	0.689	0.350
3MLP _{11,22,1}	3V-NN _{s(1)}	0.035	0.043	0.058
	3V-NN _{s(2)}	0.202	0.139	0.086
	3V-NN _{s(3)}	0.518	0.349	0.190
3MLP _{11,33,1}	3V-NN _{s(1)}	0.024	0.031	0.042
	3V-NN _{s(2)}	0.135	0.095	0.067
	3V-NN _{s(3)}	0.346	0.236	0.137

Table 10 Measure for trade-off solution ∂ comparison for 3Volterra-NN models corresponding to a 3MLP classifier, for the Facial Recognition Technology (FERET) Database. The best value for each model is highlighted in bold.

∂	3V-NN	$\gamma_1 = 0.25$	$\gamma_2 = 0.5$	$\gamma_3 = 0.75$
3MLP _{9,9,1}	3V-NN _{s(1)}	0.069	0.052	0.046
	3V-NN _{s(2)}	0.406	0.274	0.153
	3V-NN _{s(3)}	1.013	0.677	0.346
3MLP _{9,18,1}	3V-NN _{s(1)}	0.039	0.046	0.060
	3V-NN _{s(2)}	0.205	0.144	0.100
	3V-NN _{s(3)}	0.510	0.345	0.196
3MLP _{9,27,1}	3V-NN _{s(1)}	0.030	0.042	0.059
	3V-NN _{s(2)}	0.137	0.098	0.072
	3V-NN _{s(3)}	0.341	0.236	0.149

Table 5). As regards FERET database, considering that 3V-NN_{s(1)} needs 27 parameters for an RR value of 94.67%, OBD obtained an RR value of 38.38%, and OBS achieved an RR of 40.24% for the same number of parameters (see Table 7). With respect to compression, as it is possible to see in Tables 6 and 8, for maintaining an RR of 95% in both databases, while V-NN needs approximately 30 parameters, OBS and OBD require approximately 200 parameters or more.

As a limitation of the a V-NN model, it can be noted that the model reduces its compression capability as long as the number of units in the hidden layer of the original MLP decreases. That is to say, as smaller the number of hidden neurons needed to solve the problem, less compression will be obtained by the a V-NN model.

From the previous paragraphs, it can be seen that it is necessary to look at both RR and SS tables, for each dataset, to select the best model compressed, which can be confusing. The next subsection will present this analysis in a simplified manner, through the use of the proposed ∂ trade-off measure.

4.3 Global results for model selection

Tables 9 and 10 present the ∂ values for each 3V-NN, considering three different γ values. The rows group the three possible 3V-NN outputs ($s^{(1)}$, $s^{(2)}$ and $s^{(3)}$) for each 3MLP topologies that are shown, which vary in the number of hidden neurons inside the single MLP model. The three main columns represent the application of three different γ values, each one to emphasize the priority of SS over RR ($\gamma = 0.25$), RR over SS ($\gamma = 0.75$), or both of them equally measured ($\gamma = 0.5$).

Taking into account the results at the first column ($\gamma = 0.25$) of Table 9 on the one hand, it is possible to conclude that the best trade-off solution is reached by 3V-NN_{s(1)} when the topology of the MLP model is 3MLP_{11,33,1}, achieving the minimum ∂ value 0.024 (it is the best SS (97.44%) and its RR is up to 94.44%). On the other hand, focusing on the same column of Table 10, similar conclusions can be drawn

due to the best trade-off solution is also reached by $3V\text{-NN}_{s^{(1)}}$, when the MLP model is $3\text{MLP}_{9,27,1}$, achieving the minimum ∂ value 0.030 (SS is 96.98% while its RR is as higher as 92.14%).

For the second column ($\gamma = 0.5$), the trade-off solution is associated with the same solution that in the previous case: $3V\text{-NN}_{s^{(1)}}$ for the MLP model $3\text{MLP}_{11,33,1}$ in Table 9 and for the MLP model $3\text{MLP}_{9,27,1}$ in Table 10. For this case, both measures at both experiences have high values and determine the trade-off solution when there is no priority criteria between RR and SS rates.

In the third column, in which $\gamma = 0.75$, the $3V\text{-NN}_{s^{(1)}}$ model for the $3\text{MLP}_{11,11,1}$ model is the trade-off solution at Table 9, achieving the minimum ∂ value 0.041. Its RR value is the best one (95.23%) and its SS is up to 92.36%. In Table 10, the $3V\text{-NN}_{s^{(1)}}$ model for the $3\text{MLP}_{9,9,1}$ model is also the trade-off solution for the third column, with an RR value of 94.67% and a SS of 91%.

It is necessary to note that the last ∂ in Table 9 has a very close value ($\partial = 0.042$) for $3V\text{-NN}_{s^{(1)}}$ when the topology of the MLP model is $3\text{MLP}_{11,33,1}$. But, in this last case, SS is higher than RR and this situation is penalized by the value of γ , because of the priorities in model selection. It is important to note that for all previously options, $3V\text{-NN}_{s^{(1)}}$ is the best trade-off model, which requires fewer parameters to be represented, and therefore is the smallest one.

Another important point to be highlighted is that the model selection step is a necessary task that has to be performed in order to be certain of which model and parameters are the most suited to a dataset. In this sense, the ∂ measure can help in the comparisons among models. Furthermore, this measure could help finding the best possible classifier for each class by combining different order V-NN outputs, obtained from different neural topologies and configurations.

Finally, from the analysis of Tables 9 and 10, it can be seen that consistent results are obtained with respect to the detailed analysis performed on Tables 3 and 4, achieved however in a more compact and simpler way, thanks to the new proposed trade-off measure.

5 Conclusions and future work

This paper has shown a method to obtain a compact representation of an array of MLPs using the different-order $aV\text{-NN}$ model outputs. Two algorithms that implement the proposed approach have been explained. Algorithm 1 allowed extracting the Volterra kernels from the MLP parameters (after training). Algorithm 2 allowed obtaining the third-order ($V^{(3)} - NN$) Volterra output by using the Volterra weights when a new data point to be classified is received. The $aV\text{-NN}$ has been tested on a face recognition task, obtaining almost the same accuracy than three different configurations of arrays of MLP classifiers. They have been significantly compressed into less parameters. Experimental results have demonstrated the capabilities of the proposed V-NN model to compress a solution to the face recognition problem with very high recognition and space savings rates. Furthermore, a new trade-off measure for model selection was proposed, allowing to consider in a simple manner a priority of one rate over the other one. This new measure ∂ allowed us to evaluate different solutions in a compact way by only considering one value, which arises from the relationship between the recognition rate and the space savings. This measure was useful for indicating one of the obtained outputs as the best one, considering both rates at the same time but with a different weights.

Future work involves further application of the proposed method to more complex problems, involving more classes for classification and data having an evolution on time. Besides, the V-NN compression capabilities could be tested on ensembles of different kinds of neural network models. These different NN classifier topologies should be further studied in order to establish their impact on the compression capabilities offered by the proposed Volterra-NN model approach.

References

- Aitkenhead MJ, McDonald AJS (2003) A neural network face recognition system. *Engineering Applications of Artificial Intelligence* 16(3):167–176
- Bianchini M, Maggini M, Sarti L, Scarselli F (2005) Recursive neural networks learn to localize faces. *Pattern Recognition Letters* 26(12):1885–1895
- Buciluă C, Caruana R, Niculescu-Mizil A (2006) Model compression. In: *KDD '06: Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, pp 535–541
- Capello D, Martinez C, Milone D, Stegmayer G (2009) Array of multilayer perceptrons with no-class resampling training for face recognition. *Revista Iberoamericana de Inteligencia Artificial* 13(44):5–13
- Cun YL, Denker JS, Solla SA (1990) Optimal brain damage. In: *Advances in Neural Information Processing Systems*, Morgan Kaufmann, pp 598–605
- Duda R, Hart P (2003) *Pattern Classification and Scene Analysis*. Wiley
- Dzeroski S, Zenko B (2004) Is combining classifiers with stacking better than selecting the best one? *Machine Learning* 54:255–273
- Hassibi B, Stork DG, Com SCR (1993) Second order derivatives for network pruning: Optimal brain surgeon. In: *Advances in Neural Information Processing Systems 5*, Morgan Kaufmann, pp 164–171
- Haykin S (1999) *Neural Networks: A Comprehensive Foundation*. Prentice-Hall
- Jackson J (1991) *A user's guide to principal components*. Wiley series in probability and mathematical statistics: Applied probability and statistics, Wiley, URL <http://books.google.com.ar/books?id=qhQYvH8CFQQC>
- Kirby M, Sirovich L (1990) Application of the Karhunen-Loeve procedure for the characterization of human faces. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 12(1):103–108
- Kong S, Heo J, Abidi B, Palk J, Abidi M (2005) Recent advances in visual and infrared face recognition—a review. *Computer Vision and Image Understanding* 97(1):103–135
- Korenberg M, David R, Hunter I, Solomon J (2001) Parallel cascade identification and its application to protein family prediction. *Journal of Biotechnology* 91:35–47
- Li S, Jain A (eds) (2004) *Handbook of Face Recognition*. Springer-Verlag
- Madsen K, Nielsen HB, Tingleff O, Modelling M (2004) Imm methods for non-linear least squares problems
- Marquardt D (1963) An algorithm for least-squares estimation of nonlinear parameters. *SIAM Journal on Applied Mathematics* 11(2):431–441
- Martinez A, Kak A (2001) Pca versus lda. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 23(2):228–233
- Orengo G, Colantonio P, Serino A, F, Stegmayer G, Pirola M, Ghione G (2007) Neural networks and Volterra-series for time-domain behavioral models. *International Journal of RF and Microwave CAD Engineering* 17(2):160–168
- Phillips PJ, Moon H, Rizvi SA, Rauss PJ (2000) The feret evaluation methodology for face-recognition algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22:1090–1104
- Rahman A, Verma B (2011) Novel layered clustering-based approach for generating ensemble of classifiers. *IEEE Transactions on Neural Networks* 22(5):781–792
- Rubiolo M, Stegmayer G, Milone D (2010) Compressing a neural network classifier using a volterra-neural network model. In: *IEEE International Joint Conference on Neural Networks (IJCNN)*, Barcelona, Spain, pp 1–7
- Salomon D (2007) *Data Compression: The Complete Reference*. Springer
- Stegmayer G, Chiotti O (2009) Volterra NN-based behavioral model for new wireless communications devices. *Neural Computing and Applications* 18:283–291
- Turk M, Pentland A (1991) Eigenfaces for recognition. *Journal of Cognitive Neuroscience* 3(1)
- Volterra V (1959) *Theory of Functionals and Integral and Integro-Differential Equations*. Dover
- Zhang D, Wangmeng Z (2007) Computational intelligence-based biometric technologies. *IEEE Computational Intelligence Magazine* 2(2):26–36
- Zhao W, Chellappa R, Phillips P, Rosenfeld A (2003) Face recognition: A literature survey. *ACM Computing Surveys* 35(4):399–458

Supplementary Table 1 Recognition rates (RR_i) for each class $i = 1, 2, \dots, 36$ for an aV -NN ($a=36$) classifier and its corresponding first-order $s^{(1)}$, second-order $s^{(2)}$ and third-order $s^{(3)}$ outputs, for the AT&T Laboratories Cambridge ORL Database of Faces. The experiment was performed by using a $36MLP_{40,40,1}$ model, trained with the information of 36 subjects of the ORL Database.

$36MLP_{40,40,1}$												
RR_i [%]	RR_1	RR_2	RR_3	RR_4	RR_5	RR_6	RR_7	RR_8	RR_9	RR_{10}	RR_{11}	RR_{12}
$aMLP_{I,H,O}$	97.22	97.22	97.22	97.22	97.22	97.22	97.22	97.22	97.22	97.22	97.22	90.28
$36V\text{-NN}_{s^{(1)}}$	95.83	94.44	95.83	95.83	95.83	100.00	97.22	97.22	94.44	95.83	95.83	94.44
$36V\text{-NN}_{s^{(2)}}$	95.83	95.83	94.44	98.61	94.44	94.44	100.00	97.22	98.61	98.61	95.83	95.83
$36V\text{-NN}_{s^{(3)}}$	93.06	98.61	94.44	98.61	98.61	95.83	97.22	95.83	95.83	100.00	95.83	97.22
RR_i [%]	RR_{13}	RR_{14}	RR_{15}	RR_{16}	RR_{17}	RR_{18}	RR_{19}	RR_{20}	RR_{21}	RR_{22}	RR_{23}	RR_{24}
$aMLP_{I,H,O}$	97.22	97.22	97.22	97.22	97.22	97.22	97.22	88.89	97.22	97.22	97.22	97.22
$36V\text{-NN}_{s^{(1)}}$	95.83	94.44	94.44	95.83	98.61	97.22	94.44	95.83	98.61	97.22	98.61	98.61
$36V\text{-NN}_{s^{(2)}}$	93.06	97.22	95.83	94.44	95.83	97.22	93.06	95.83	94.44	95.83	94.44	100.00
$36V\text{-NN}_{s^{(3)}}$	93.06	95.83	95.83	98.61	98.61	95.83	94.44	95.83	100.00	98.61	94.44	97.22
RR_i [%]	RR_{25}	RR_{26}	RR_{27}	RR_{28}	RR_{29}	RR_{30}	RR_{31}	RR_{32}	RR_{33}	RR_{34}	RR_{35}	RR_{36}
$aMLP_{I,H,O}$	97.22	97.22	97.22	97.22	97.22	97.22	97.22	97.22	97.22	97.22	97.22	97.22
$36V\text{-NN}_{s^{(1)}}$	98.61	95.83	91.67	94.44	95.83	97.22	98.61	97.22	95.83	98.61	95.83	94.44
$36V\text{-NN}_{s^{(2)}}$	95.83	97.22	94.44	93.06	98.61	98.61	91.67	97.22	95.83	97.22	94.44	94.44
$36V\text{-NN}_{s^{(3)}}$	95.83	98.61	97.22	98.61	94.44	97.22	95.83	93.06	94.44	95.83	95.83	93.06

Supplementary Table 2 Global recognition rate (RR) and space saving(SS) comparison for an aV -NN ($a=36$) classifier and its corresponding first-order $s^{(1)}$, second-order $s^{(2)}$ and third-order $s^{(3)}$ outputs, for the the AT&T Laboratories Cambridge ORL Database of Faces. The experience was performed by using an $36MLP_{40,40,1}$ model, thus representing 36 subject of the ORL Database.

$36MLP_{40,40,1}$			
$P_{aMLP_{I,H,O}} \rightarrow$		60516	
$RR[\%] \rightarrow$		96.79	
$aV\text{-NN}$	$P_{as^{(i)}}$	$RR[\%]$	$SS[\%]$
$36V\text{-NN}_{s^{(1)}}$	1440	96.29	97.62
$36V\text{-NN}_{s^{(2)}}$	30960	95.99	48.84
$36V\text{-NN}_{s^{(3)}}$	88560	96.37	-46.34

Supplementary Table 3 Recognition rates (RR_i) for each class $i = 1, 2, \dots, 40$ for an aV -NN ($a=40$) classifier and its corresponding first-order $s^{(1)}$, second-order $s^{(2)}$ and third-order $s^{(3)}$ outputs, for the Facial Recognition Technology (FERET) Database. The experiment was performed by using a $40MLP_{37,37,1}$ model, trained with the information of 40 subjects of the FERET Database.

$40MLP_{37,37,1}$										
RR_i [%]	RR_1	RR_2	RR_3	RR_4	RR_5	RR_6	RR_7	RR_8	RR_9	RR_{10}
$aMLP_{I,H,O}$	97.16	97.50	100.00	96.25	97.73	97.61	97.50	97.50	98.52	97.95
$40V-NN_{s^{(1)}}$	100.00	95.11	97.50	95.00	97.05	95.00	99.09	94.20	95.45	95.00
$40V-NN_{s^{(2)}}$	97.50	93.64	98.41	97.73	95.45	96.14	97.05	94.55	98.18	95.00
$40V-NN_{s^{(3)}}$	93.07	93.86	100.00	96.59	97.73	96.48	96.02	96.36	97.16	95.00
RR_i [%]	RR_{11}	RR_{12}	RR_{13}	RR_{14}	RR_{15}	RR_{16}	RR_{17}	RR_{18}	RR_{19}	RR_{20}
$aMLP_{I,H,O}$	97.95	97.73	94.66	99.55	98.30	97.50	96.36	97.50	97.50	96.82
$40V-NN_{s^{(1)}}$	95.91	99.89	96.93	99.43	98.86	95.23	97.27	95.00	95.00	95.00
$40V-NN_{s^{(2)}}$	94.43	98.30	95.00	97.05	99.09	97.39	95.00	96.59	95.34	93.98
$40V-NN_{s^{(3)}}$	95.00	99.77	95.00	95.34	98.41	98.98	93.52	95.00	95.00	94.77
RR_i [%]	RR_{21}	RR_{22}	RR_{23}	RR_{24}	RR_{25}	RR_{26}	RR_{27}	RR_{28}	RR_{29}	RR_{30}
$aMLP_{I,H,O}$	97.50	97.50	97.50	97.50	97.50	97.50	97.50	97.50	95.23	97.50
$40V-NN_{s^{(1)}}$	98.98	94.66	97.61	95.91	99.43	96.70	95.00	98.52	99.55	95.00
$40V-NN_{s^{(2)}}$	98.41	95.00	97.05	96.93	99.32	94.43	95.00	95.45	94.89	95.00
$40V-NN_{s^{(3)}}$	94.20	97.05	98.75	96.48	97.39	94.43	95.00	96.14	98.41	95.00
RR_i [%]	RR_{31}	RR_{32}	RR_{33}	RR_{34}	RR_{35}	RR_{36}	RR_{37}	RR_{38}	RR_{39}	RR_{40}
$aMLP_{I,H,O}$	95.11	92.84	91.59	97.50	97.27	97.39	97.50	96.82	98.52	97.50
$40V-NN_{s^{(1)}}$	97.95	97.73	94.32	96.36	95.00	95.23	99.43	98.30	95.00	95.11
$40V-NN_{s^{(2)}}$	94.66	97.73	96.25	97.84	96.59	95.57	99.66	98.30	94.09	95.80
$40V-NN_{s^{(3)}}$	97.95	98.07	95.68	96.36	96.36	97.50	93.64	94.32	95.45	96.25

Supplementary Table 4 Global recognition rate (RR) and space saving(SS) comparison for an aV -NN ($a=40$) classifier and its corresponding first-order $s^{(1)}$, second-order $s^{(2)}$ and third-order $s^{(3)}$ outputs, for the Facial Recognition Technology (FERET) Database. The experience was performed by using an $40MLP_{37,37,1}$ model, thus representing 40 subject of the FERETE Database.

$40MLP_{37,37,1}$			
$P_{aMLP_{I,H,O}} \rightarrow$		59200	
$RR[\%] \rightarrow$		97.16	
$aV-NN$	$P_{as^{(i)}}$	$RR[\%]$	$SS[\%]$
$40V-NN_{s^{(1)}}$	1480	96.69	97.50
$40V-NN_{s^{(2)}}$	29600	96.34	50.00
$40V-NN_{s^{(3)}}$	84360	96.37	-42.50