



Ordinal regression neural networks based on concentric hyperspheres



Pedro Antonio Gutiérrez^{a,*}, Peter Tiño^b, César Hervás-Martínez^a

^a University of Córdoba, Department of Computer Science and Numerical Analysis, Rabanales Campus, Albert Einstein building, 14071 - Córdoba, Spain

^b School of Computer Science, The University of Birmingham, Birmingham B15 2TT, United Kingdom

ARTICLE INFO

Article history:

Received 7 February 2014

Received in revised form 15 May 2014

Accepted 7 July 2014

Available online 11 July 2014

Keywords:

Ordinal regression
Ordinal classification
Neural networks
Latent variable

ABSTRACT

Threshold models are one of the most common approaches for ordinal regression, based on projecting patterns to the real line and dividing this real line in consecutive intervals, one interval for each class. However, finding such one-dimensional projection can be too harsh an imposition for some datasets. This paper proposes a multidimensional latent space representation with the purpose of relaxing this projection, where the different classes are arranged based on concentric hyperspheres, each class containing the previous classes in the ordinal scale. The proposal is implemented through a neural network model, each dimension being a linear combination of a common set of basis functions. The model is compared to a nominal neural network, a neural network based on the proportional odds model and to other state-of-the-art ordinal regression methods for a total of 12 datasets. The proposed latent space shows an improvement on the two performance metrics considered, and the model based on the three-dimensional latent space obtains competitive performance when compared to the other methods.

© 2014 Elsevier Ltd. All rights reserved.

1. Introduction

When we face an ordinal regression (OR) problem, the objective is to predict the label y_i of an input vector \mathbf{x}_i , where $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^k$ and $y_i \in \mathcal{Y} \in \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_Q\}$. This is done by estimating a classification rule or function $F: \mathcal{X} \rightarrow \mathcal{Y}$ to predict the labels of new samples. In a supervised setting, we are given a training set of N points, $D = \{(\mathbf{x}_i, y_i), 1 \leq i \leq N\}$. All these considerations can be also found in standard nominal classification, but, for OR, a natural label ordering is included, which is given by $\mathcal{C}_1 < \mathcal{C}_2 < \dots < \mathcal{C}_Q$. The symbol $<$ is an order relation representing the nature of the classification problem and expressing that a label is before another in the ordinal scale.

OR problems are very common in real settings, although the machine learning community has often treated them from a standard (nominal) perspective, ignoring the order relationship, $<$, between classes. Some examples of application fields where OR is found are credit rating (Dikkers & Rothkrantz, 2005), econometric modelling (Mathieson, 1996), medical research (Cardoso, da Costa, & Cardoso, 2005) or face recognition (Kim & Pavlovic, 2010), to name a few. Considering the order relationship between classes can result in two significant benefits: (1) minimisation of specific

classification errors, and (2) incorporation of the ordering into the classifier. With respect to the first benefit, it is clear that one should focus on predicting categories as close as possible to the real one when tackling an OR problem. Hence, OR methods are aimed to minimise those errors that involve large category gaps in the ordinal scale. As an example, consider a tumour classification problem where the categories are $\{benign, dangerous, malignant\}$. Misclassification of *malign* tumours as *dangerous* is preferred to assign the label *benign* to a *malign* tumour and OR methods will generally minimise this second type of errors. The second benefit comes from the fact that label order is usually present, in a direct way in the input space or through a latent space representation (Sánchez-Monedero, Gutiérrez, Tino, & Hervás-Martínez, 2013). Imbuing a classifier with this ordering will generally improve generalisation performance, as the classifier is better representing the nature of the task.

The field of OR has experienced significant development in the last decade, with many new methods adapted from traditional machine learning methodologies, from support vector machine (SVM) formulations (Chu & Keerthi, 2007) to Gaussian processes (Chu & Ghahramani, 2005) or discriminant learning (Sun, Li, Wu, Zhang, & Li, 2010). For all these methods, although classifier construction is motivated and undertaken from different points of view, the final models share a common structure or nature. They exploit the fact that it is natural to assume that an unobserved continuous variable underlies the ordinal response variable (e.g. the actual age of the person appearing in a picture

* Corresponding author. Tel.: +34 9572 18349; fax: +34 9572 18630.

E-mail addresses: pagutierrez@uco.es, zamarck@gmail.com (P.A. Gutiérrez).

for an age classification problem). This variable is called latent variable, and methods based on that assumption are known as threshold models (Verwaeren, Waegeman, & De Baets, 2012). Indeed, this structure can be found in one of the first models for OR, the proportional odds model (POM) (McCullagh, 1980), which is a probabilistic model estimating the cumulative probabilities of the different ordered categories and leading to linear decision boundaries. Threshold models methodologies estimate:

- A function $f(\mathbf{x})$ that tries to predict the values of the latent variable.
- A set of thresholds $\mathbf{b} = (b_1, b_2, \dots, b_{Q-1}) \in \mathbb{R}^{Q-1}$ to represent intervals in the range of $f(\mathbf{x})$, which must satisfy the constraints $b_1 \leq b_2 \leq \dots \leq b_{Q-1}$.

From a practical perspective, threshold models are basically trying to find a one-dimensional projection ($f(\mathbf{x})$) where patterns are ordered according to the class labels. Finding this projection can be a problem for real world datasets. If we consider linear models for $f(\mathbf{x})$, the chances that the patterns exhibit a linear ordering relationship are certainly very low. If we consider nonlinear models, the pressure to find this nonlinear projection can result in unnatural or too complex projections leading to poorer generalisation performance. This paper proposes to relax this pressure by allowing a higher dimensional representation of the latent space. This is done by ordering patterns in an L -dimensional space, where each class region is limited by concentric hyperspheres (centred in the origin). The ordering of the classes is imposed by assuring that the radii of the hyperspheres are also ordered.

Another popular way to tackle OR problems is to decompose the original task into several binary tasks, where each binary task consists of predicting if the patterns belong to a category higher (in the ordinal scale) than a given label \mathcal{C}_q . One model is estimated for each class in the ordinal scale. The approach is presented in the work of Frank and Hall (2001), also proposing a way to fuse probabilities given for all binary tasks. Later on, there have been two different lines of research where binary classification and OR were linked in a more direct way (Cardoso & da Costa, 2007; Li & Lin, 2007; Lin & Li, 2012). Instead of learning Q different binary classifiers, a single binary classifier is learnt, where the category examined is included as an additional feature and training patterns are replicated and weighted. The framework in Li and Lin (2007) and Lin and Li (2012) is more generic, in the sense that it can be applied to different cost matrices.

The Error-Correcting Output Codes (ECOC) methodology is a popular and effective coding method to learn complex class targets, which can be used also for OR. The main idea is to associate each class \mathcal{C}_q with a column of a binary coding matrix $\mathbf{M}_{R \times Q}$, where each entry of the matrix $\mathbf{M}(j, q) \in \{-1, +1\}$, Q is the number of classes, R is the number of binary classifiers, $1 \leq q \leq Q$ and $1 \leq j \leq R$. After training the binary classifiers, prediction is then accomplished by choosing the column of \mathbf{M} closest to the set of decision values, where the distance function should be selected according to the error function minimised during learning (Allwein, Schapire, & Singer, 2001; Dietterich & Bakiri, 1995).

In the field of neural networks, there have been some proposals for OR problems. The first one dates back to 1996, when Mathieson proposed a non-linear version of the POM (Mathieson, 1996, 1999) by setting the projection $f(\mathbf{x})$ to be the output of a neural network. Although the results were quite promising, the method was evaluated for a very specific dataset. A more extensive battery of experiments should be done to further validate the proposal.

Costa (1996) derived another neural network architecture to exploit the ordinal nature of the data. It was based on a “partitive approach”, where probabilities are assigned to the joint prediction of constrained concurrent events.

Other approach (Cheng, Wang, & Pollastri, 2008) applies the coding scheme of Frank and Hall and a decision rule based on

examining output nodes with an order and selecting the first one whose output is higher than a predefined threshold T . The problem of this method is that inconsistencies can be found in the predictions (i.e. a sigmoid with value higher than T after the index selected).

The ordinal neural network (oNN) of Cardoso and da Costa (2007) adapts the previously discussed data replication method to neural networks (a single model for binary decomposition using an extended and replicated version of the dataset), allowing the derivation of nonlinear decision boundaries.

Additionally, extreme learning machines (ELMs) have been used as a very fast method to fit single layer neural networks, where the hidden neurons weights are random, and the output weights are analytically obtained (Huang, Zhou, Ding, & Zhang, 2012). They have been adapted to OR (Deng, Zheng, Lian, Chen, & Wang, 2010), considering again the Frank and Hall coding scheme and a prediction based on the ECOC loss-based decoding approach (Allwein et al., 2001), i.e. the chosen label is that which minimises the exponential loss. Another recent paper by Riccardi, Fernandez-Navarro, and Carloni (in press) introduces a cost-sensitive approach for adapting the stagewise additive modelling using a multiclass exponential boosting algorithm (SAMME, which is the multiclass version of the well-known AdaBoost) to OR problems. They consider ELMs as the base classifier and they introduce three different loss functions, affecting the update rule of the error estimation and/or of the pattern weights (Riccardi et al., in press). From the three variants introduced in the paper, the third one (which adapts the update rule of both the error estimation and the pattern weights) obtains the best results. The OR model proposed in Fernandez-Navarro, Riccardi, and Carloni (in press) adapts ELM to OR problems by imposing monotonicity constraints in the weights connecting the hidden layer with the output layer. The optimum of the inequality constrained least squares problem is determined analytically according to the closed-form solution estimated from the Karush–Kuhn–Tucker conditions.

A conceptually different methodology is proposed by da Costa, Alonso, and Cardoso (2008) and da Costa and Cardoso (2005) for training OR models, with a special attention to neural networks. They assume that the random variable class of a pattern should follow a unimodal distribution. Two possible implementations are considered: a parametric one, where a specific discrete distribution is assumed and the associated free parameters are estimated by a neural network; and a non-parametric one, where no distribution is assumed but the error function is modified to avoid errors from distant classes. Finally, the approach in Dobriska, Wang, and Blackburn (2012) is a distribution-independent methodology for OR based on pairwise preferences. The strength of dependency between two data instances (continuous preferences) is shown to improve algorithmic performance, obtaining competitive results.

In this paper, we extend the proposal of Mathieson (1996, 1999), deriving a nonlinear version of the POM based on neural networks. We present a common learning framework to fit the parameters of a nominal neural network (NNN) and the neural network based on the POM (POMNN). The framework is then used to fit an extended version of the POMNN, where, as previously discussed, the latent space is assumed to be L -dimensional and the patterns are ordered by considering $Q - 1$ concentric hyperspheres. The underlying motivation is to relax the imposition of projecting all patterns in a real line.

With regard to the relationship between ECOC and the proposal of this paper, one single model is used for learning the ordinal target, and the problem is not decomposed into several binary ones. In this way, the latent space structure relates each pattern to the posterior probabilities without learning multiple binary classifiers.

This paper is organised as follows: Section 2 is devoted to a brief analysis of the POM model, closely related to the models proposed;

the description of the different ordinal neural network models is carried out in Section 3; Section 4 contains the experimental results; and finally, Section 5 summarises the conclusions of our work.

2. Proportional odds model (POM)

This model is a direct extension of binary logistic regression for the case of OR. It was first presented by McCullagh (1980) and dates back to 1980. POM can be grouped under a wider family of models, the cumulative link models (CLMs) (Agresti, 2002), which predict probabilities of adjacent categories, taking the ordinal scale into account. The key idea of CLMs is to estimate cumulative probabilities as follows:

$$P(y \leq c_q | \mathbf{x}) = P(y = c_1 | \mathbf{x}) + \dots + P(y = c_q | \mathbf{x}),$$

$$P(y = c_q | \mathbf{x}) = P(y \leq c_q | \mathbf{x}) - P(y \leq c_{q-1} | \mathbf{x}),$$

for $1 < q \leq Q$, considering by definition that $P(y \leq c_Q | \mathbf{x}) = 1$ and $P(y = c_1 | \mathbf{x}) = P(y \leq c_1 | \mathbf{x})$. CLMs relate a linear model of the input variables to these cumulative probabilities:

$$f(\mathbf{x}) = g^{-1}(P(y \leq c_q | \mathbf{x})) = b_q - \mathbf{w}^T \mathbf{x},$$

where $g^{-1} : [0, 1] \rightarrow (-\infty, +\infty)$ is a monotonic transformation (the inverse link function), b_q is the threshold defined for class c_q , and \mathbf{w} is the coefficient vector of the linear model. The most common choice for the link function is the logistic function (which is indeed the one selected for the POM McCullagh, 1980), although probit, complementary log–log, negative log–log or cauchit functions could also be used (Agresti, 2002). The logit link function is the inverse of the standard logistic cumulative distribution function (cdf), with the following expression:

$$g^{-1}(P(y \leq c_q | \mathbf{x})) = \ln \left(\frac{P(y \leq c_q | \mathbf{x})}{1 - P(y \leq c_q | \mathbf{x})} \right), \quad (1)$$

while the probit link function is the inverse of the standard normal cdf:

$$g^{-1}(P(y \leq c_q | \mathbf{x})) = \Phi^{-1}(P(y \leq c_q | \mathbf{x})), \quad (2)$$

where $\Phi(x)$ is the normal distribution, and the cauchit link function is the inverse of the Cauchy cdf (characterised for having long thick tails):

$$g^{-1}(P(y \leq c_q | \mathbf{x})) = \tan\{\pi \cdot [P(y \leq c_q | \mathbf{x}) - 0.5]\}. \quad (3)$$

Under the assumption that $f(\mathbf{x})$ follows a logistic cdf and following the idea of the POM model (McCullagh, 1980), the cumulative likelihood of a pattern being associated with a class less than or equal to class c_q is defined as:

$$P(y \leq c_q | \mathbf{x}) = \frac{1}{1 + \exp(f(\mathbf{x}) - b_q)}, \quad (4)$$

where $q = 1, \dots, Q$, and, by definition, $P(y \leq c_Q | \mathbf{x}) = 1$. Therefore, this model approximates the posterior probability of a class q as:

$$\begin{aligned} P(y = c_q | \mathbf{x}) &= P(y \leq c_q | \mathbf{x}) - P(y \leq c_{q-1} | \mathbf{x}) \\ &= \frac{1}{1 + \exp(f(\mathbf{x}) - b_q)} - \frac{1}{1 + \exp(f(\mathbf{x}) - b_{q-1})}. \end{aligned} \quad (5)$$

Thresholds must satisfy the constraint $b_1 < b_2 < \dots < b_{Q-1}$ and their role is to divide the real line into Q contiguous intervals; these intervals map the function value $f(\mathbf{x})$ into the discrete variable, while forcing a proper probability interpretation (increasing probability when a higher class in the ordinal scale is examined).

3. Neural network classification algorithms

This section explains the characteristics of the different neural network models considered for the present work, including the ones proposed. The first subsection will introduce the learning algorithm which will be based on a set of estimated probabilities $\mathbf{p}(\mathbf{x}) = \{P(y = c_1 | \mathbf{x}), P(y = c_2 | \mathbf{x}), \dots, P(y = c_Q | \mathbf{x})\}$. The different ways of obtaining this probability vector $\mathbf{p}(\mathbf{x})$ constitute the different neural network models considered.

3.1. Learning algorithm

As previously stated, our aim is to estimate a classification rule F based on a training set D . If we consider a “1-of- Q ” encoding vector for representing the class labels, then we define a vector $\mathbf{y}_i = (y_i^{(1)}, y_i^{(2)}, \dots, y_i^{(Q)})$ for each training label y_i , such that $y_i^{(q)} = 1$ if \mathbf{x}_i corresponds to an example belonging to class c_q and $y_i^{(q)} = 0$ otherwise. We will denote θ to the vector of free parameters of the model to be learnt, which will be specified in following subsections. To perform the maximum likelihood estimation of the parameter vector θ , we choose to minimise the cross-entropy error function:

$$L(\theta) = -\frac{1}{N} \sum_{n=1}^N \sum_{q=1}^Q y_n^{(q)} \log P(y = c_q | \mathbf{x}_n, \theta). \quad (6)$$

An individual should be assigned to the class which has the maximum probability, given the measurement vector \mathbf{x} :

$$F(\mathbf{x}) = \hat{q}, \quad \text{where } \hat{q} = \arg \max_q P(y = c_q | \mathbf{x}_n, \theta), \quad q = 1, \dots, Q.$$

Considering the $L(\theta)$ error function, we optimise the model parameters by gradient descent using the *iRprop+* local improvement procedure (Igel & Hüsken, 2003). This training scheme combines the local information (i.e. the sign of the partial derivative of the error with respect to a weight like *Rprop*) with more global information (i.e. the error value at each iteration) in order to decide whether to revert an update step for each weight individually, resulting in very robust performance (Igel & Hüsken, 2003).

For the sake of simplicity, let S be the total number of parameters of the model. The gradient vector is given by:

$$\nabla L(\theta) = \left(\frac{\partial L(\theta)}{\partial \theta_1}, \frac{\partial L(\theta)}{\partial \theta_2}, \dots, \frac{\partial L(\theta)}{\partial \theta_S} \right). \quad (7)$$

Considering (6), each of its component will be defined as:

$$\frac{\partial L}{\partial \theta_s} = -\frac{1}{N} \sum_{n=1}^N \sum_{q=1}^Q \frac{y_n^{(q)}}{P(y = c_q | \mathbf{x}_n, \theta)} \cdot \frac{\partial P(y = c_q | \mathbf{x}_n, \theta)}{\partial \theta_s},$$

where $s = 1, \dots, S$. These derivatives will depend on the actual model form and they will be specified in the following subsections.

3.2. Nominal neural network (NNN)

Even when an ordinal classification problem is considered, one could use as a baseline a nominal neural network (NNN) to estimate $\mathbf{p}(\mathbf{x})$. A feed-forward multilayer perceptron can be configured with Q output nodes and one hidden layer. For a more robust model and learning process, the output of this NNN should be transformed by considering the softmax transformation:

$$P(y = c_q | \mathbf{x}_n, \theta) = \frac{\exp(g_q(\mathbf{x}, \theta_q))}{\sum_{j=1}^Q \exp(g_j(\mathbf{x}, \theta_j))}, \quad (8)$$

where $1 \leq q \leq Q$, $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_Q)$ and $g_q(\mathbf{x}, \boldsymbol{\theta}_q)$ is the output of the q th node of the output layer. A proper probability distribution should assure that $\sum_{i=1}^Q P(y = C_i | \mathbf{x}_n, \boldsymbol{\theta}) = 1$, which implies that the probability for one of the classes could be expressed as a function of the others, reducing the degrees of freedom of the model. This can be done by setting one class as the reference class (in our case, the last class, C_Q), and dividing the numerator and the denominator by $\exp(g_Q(\mathbf{x}, \boldsymbol{\theta}_Q))$:

$$P(y = C_q | \mathbf{x}_n, \boldsymbol{\theta}) = \frac{\exp(g_q(\mathbf{x}, \boldsymbol{\theta}_q) - g_Q(\mathbf{x}, \boldsymbol{\theta}_Q))}{1 + \sum_{j=1}^{Q-1} \exp(g_j(\mathbf{x}, \boldsymbol{\theta}_j) - g_Q(\mathbf{x}, \boldsymbol{\theta}_Q))}$$

for $q = 1, \dots, Q-1$. Now we set $f_q(\mathbf{x}, \boldsymbol{\theta}_q) = g_q(\mathbf{x}, \boldsymbol{\theta}_q) - g_Q(\mathbf{x}, \boldsymbol{\theta}_Q)$ and the final model reduces to:

$$P(y = C_q | \mathbf{x}_n, \boldsymbol{\theta}) = \frac{\exp(f_q(\mathbf{x}, \boldsymbol{\theta}_q))}{1 + \sum_{j=1}^{Q-1} \exp(f_j(\mathbf{x}, \boldsymbol{\theta}_j))}, \quad 1 \leq q \leq Q-1,$$

$$P(y = C_Q | \mathbf{x}_n, \boldsymbol{\theta}) = \frac{1}{1 + \sum_{j=1}^{Q-1} \exp(f_j(\mathbf{x}, \boldsymbol{\theta}_j))},$$

which is equivalent to Eq. (8) when the last output is set to zero, $f_Q(\mathbf{x}, \boldsymbol{\theta}_Q) = 0$. This way the number of model parameters is reduced. For the rest of classes, the outputs are defined by a linear combination of the hidden nodes, in the following way:

$$f_q(\mathbf{x}, \boldsymbol{\theta}_q) = f(\mathbf{x}, \boldsymbol{\beta}_q, \mathbf{W}) = \beta_0^q + \sum_{j=1}^M \beta_j^q B_j(\mathbf{x}, \mathbf{w}_j),$$

where $1 \leq q \leq Q-1$, $\boldsymbol{\beta}_q = \{\beta_0^q, \beta_1^q, \dots, \beta_M^q\}$, $\mathbf{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_M\}$, $\mathbf{w}_j = \{w_{j0}, w_{j1}, \dots, w_{jk}\}$, and $B_j(\mathbf{x}, \mathbf{w}_j)$ can be any kind of basis function, in our case, sigmoidal units:

$$B_j(\mathbf{x}, \mathbf{w}_j) = \left(1 + \exp\left(-w_{j0} - \sum_{i=1}^k w_{ji} \cdot x_i\right) \right)^{-1}. \quad (9)$$

With this configuration, the derivatives are given in the following way. For the sake of simplicity, let θ_s be any of the parameters of $\boldsymbol{\beta}_q$ or \mathbf{W} , $P(y = C_q | \mathbf{x}_n, \boldsymbol{\theta}) = p_{nq}$ and $f_q(\mathbf{x}_n, \boldsymbol{\theta}_q) = f_{nq}$:

$$\frac{\partial p_{nq}}{\partial \theta_s} = \sum_{j=1}^Q p_{nj} \cdot (I(j=q) - p_{nj}) \cdot \frac{\partial f_{nq}}{\partial \theta_s},$$

where $I(\cdot)$ is the indicator function. The derivatives of the parameters of the model output functions f_{nq} can be expressed as:

$$\frac{\partial f_{nq}}{\partial \beta_0^k} = \begin{cases} 0 & \text{if } q \neq k, \\ 1 & \text{if } q = k. \end{cases}, \quad \frac{\partial f_{nq}}{\partial \beta_j^k} = \begin{cases} 0 & \text{if } q \neq k, \\ B_j(\mathbf{x}_n, \mathbf{w}_j) & \text{if } q = k \end{cases}$$

and the gradient for the hidden layer depends on the kind of basis function used. For sigmoidal nodes:

$$\frac{\partial f_{nq}}{\partial w_{jt}} = \beta_j^q B_j(\mathbf{x}_n, \mathbf{w}_j) (1 - B_j(\mathbf{x}_n, \mathbf{w}_j)) x_{nt}, \quad 1 \leq t \leq k,$$

$$\frac{\partial f_{nq}}{\partial w_{j0}} = \beta_j^q B_j(\mathbf{x}_n, \mathbf{w}_j) (1 - B_j(\mathbf{x}_n, \mathbf{w}_j)).$$

3.3. Proportional odds model neural network (POMNN)

The fact that the POM is linear limits its applicability in real world datasets, given that the parallel linear decision boundaries are often unrealistic. A non-linear version of the POM model can be

formed by setting the projection $f(\mathbf{x})$ to be the output of a neural network. While the POM model approximates $f(\mathbf{x})$ by a simple linear combination of the input variables, the POMNN considers a non-linear basis transformation of the inputs. For each class we will have:

$$f_q(\mathbf{x}, \boldsymbol{\theta}_q) = f_q(\mathbf{x}, \boldsymbol{\beta}, \mathbf{W}, b_q) = b_q - f(\mathbf{x}, \boldsymbol{\beta}, \mathbf{W})$$

where $q = 1 \dots Q-1$, $\boldsymbol{\theta}_q = \{\boldsymbol{\beta}, \mathbf{W}, b_q\}$ and the projection function $f(\mathbf{x}, \boldsymbol{\beta}, \mathbf{W})$ is estimated by:

$$f(\mathbf{x}, \boldsymbol{\beta}, \mathbf{W}) = \sum_{j=1}^M \beta_j B_j(\mathbf{x}, \mathbf{w}_j),$$

where $\boldsymbol{\beta} = \{\beta_1, \dots, \beta_M\}$, $\mathbf{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_M\}$, and $B_j(\mathbf{x}, \mathbf{w}_j)$ can be any kind of basis functions, in our case, sigmoidal units. Note that in this case the neural network model will have only one output node. As compared with the model of Mathieson (1999), that model had skip-layer connections and the error function was endowed with a regularisation term. Moreover, in order to ensure the constraints in the biases, $b_1 \leq \dots \leq b_{Q-1}$, we propose the following definition: $b_q = b_{q-1} + \Delta_q^2$, $q = 1 \dots Q-1$, with padding variables Δ_q , which are squared to make them positive, and $b_1, \Delta_q \in \mathbb{R}$. Consequently, the parameter vector is defined as: $\boldsymbol{\theta} = \{\boldsymbol{\beta}, \mathbf{W}, b_1, \Delta_2, \dots, \Delta_{Q-1}\}$. In this way, we convert the constrained optimisation problem into an unconstrained one, and gradient descent can be applied to all the parameters in $\boldsymbol{\theta}$ without considering the constraints.

With this model structure, the derivatives can now be reformulated to perform gradient-descent optimisation. Let θ_s be any of the parameters of $\boldsymbol{\beta}$ or \mathbf{W} , $P(y = C_q | \mathbf{x}_n, \boldsymbol{\theta}) = p_{nq}$ and $f(\mathbf{x}_n, \boldsymbol{\beta}, \mathbf{W}) = f_n$:

$$p_{nq|q>1} = \frac{1}{1 + \exp(b_q - f_n)} - \frac{1}{1 + \exp(b_{q-1} - f_n)},$$

$$p_{n1} = \frac{1}{1 + \exp(b_1 - f_n)},$$

$$\frac{\partial p_{nq|q>1}}{\partial \theta_s} = \frac{\exp(b_q - f_n)}{(1 + \exp(b_q - f_n))^2} \frac{\partial f_n}{\partial \theta_s} - \frac{\exp(b_{q-1} - f_n)}{(1 + \exp(b_{q-1} - f_n))^2} \frac{\partial f_n}{\partial \theta_s},$$

$$\frac{\partial p_{n1}}{\partial \theta_s} = \frac{\exp(b_1 - f_n)}{(1 + \exp(b_1 - f_n))^2} \frac{\partial f_n}{\partial \theta_s},$$

where the derivatives for the projection parameters θ_s are:

$$\frac{\partial f_n}{\partial \beta_j} = B_j(\mathbf{x}_n, \mathbf{w}_j),$$

$$\frac{\partial f_n}{\partial w_{jt}} = \beta_j B_j(\mathbf{x}_n, \mathbf{w}_j) (1 - B_j(\mathbf{x}_n, \mathbf{w}_j)) x_{nt}, \quad 1 \leq t \leq k,$$

$$\frac{\partial f_n}{\partial w_{j0}} = \beta_j B_j(\mathbf{x}_n, \mathbf{w}_j) (1 - B_j(\mathbf{x}_n, \mathbf{w}_j)).$$

The derivatives for the biases and the padding variables can be formulated in the following way:

$$\frac{\partial p_{nq|q>1}}{\partial b_1} = \frac{\exp(f_n - b_q)}{(1 + \exp(f_n - b_q))^2} - \frac{\exp(f_n - b_{q-1})}{(1 + \exp(f_n - b_{q-1}))^2},$$

$$\frac{\partial p_{n1}}{\partial b_1} = \frac{\exp(f_n - b_1)}{(1 + \exp(f_n - b_1))^2},$$

$$\frac{\partial p_{nq|q>1}}{\partial \Delta_j} = 2\Delta_j \cdot \frac{\partial p_{nq|q>1}}{\partial b_1}, \quad j = 2 \dots Q-1,$$

$$\frac{\partial p_{n1}}{\partial \Delta_j} = 2\Delta_j \cdot \frac{\partial p_{n1}}{\partial b_1}, \quad j = 2 \dots Q-1.$$

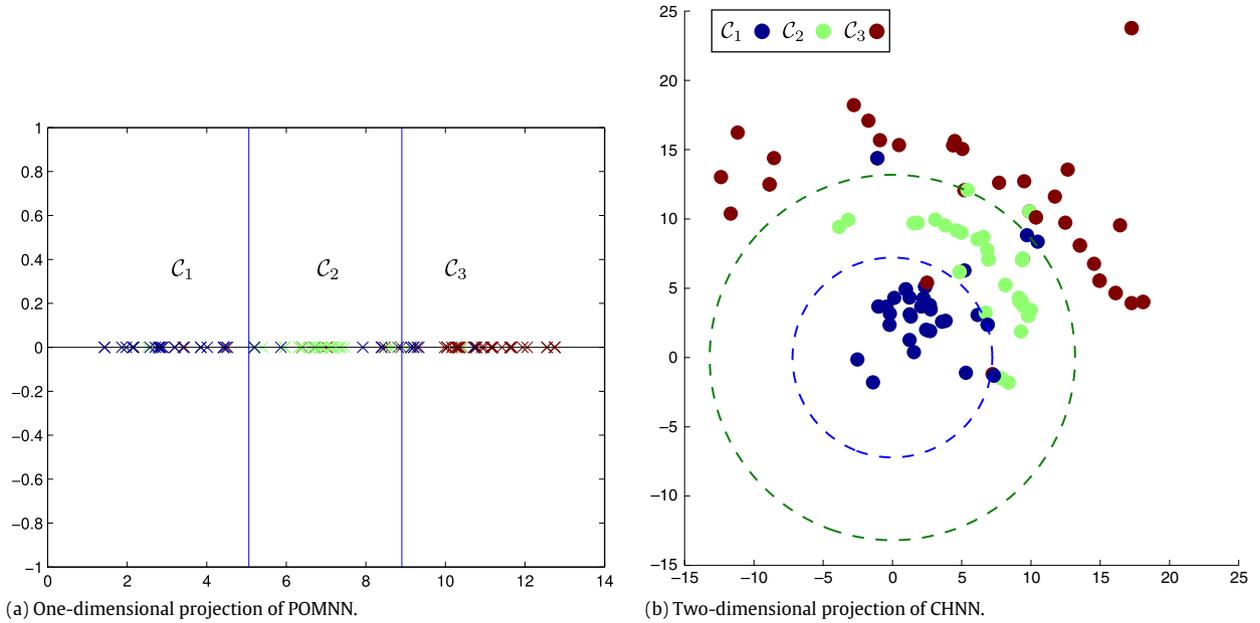


Fig. 1. POMNN projection for the tae dataset and the proposed two-dimensional CHNN projection for the same dataset. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

3.4. Concentric hypersphere neural network (CHNN)

Both POM and POMNN project input patterns to the real line with the difference that POM impose a linear projection and POMNN a nonlinear one. After this, the projection is divided into intervals, each representing a different class. Finding such one-dimensional projection where patterns are ordered according to the class labels can be too harsh an imposition for some datasets. We propose to relax this requirement by allowing a multivariate representation. The class order will be represented by a natural order of nested concentric hyper-spheres (centred at the origin). The proposal tries to find a space made up by several different projections (coordinates) where the order of the class labels is presented in the form of concentric hyperspheres. The smallest hypersphere will contain the patterns of the first class in the ordinal scale.

An example showing the proposed idea is presented in Fig. 1, where the one-dimensional projection of POMNN is compared against the two-dimensional one of the proposed CHNN approach. As can be seen, the patterns can be more easily positioned in their correct region, because they are projected into a two-dimensional space, where each coordinate is approximated separately.

Non-linearity is achieved by letting these projections be the output of different linear combinations of the basis functions. Input conditional class probabilities are related to the distance of the pattern to the centre of the hyperspheres. Specifically, the latent space is defined in the following way:

$$f_l(\mathbf{x}, \theta_l) = f(\mathbf{x}, \beta_l, \mathbf{W}) = \beta_0^l + \sum_{j=1}^M \beta_j^l B_j(\mathbf{x}, \mathbf{w}_j), \quad 1 \leq l \leq L,$$

where L is the latent space dimensionality (hence dimensionality of the hyperspheres), $\beta_l = \{\beta_0^l, \beta_1^l, \dots, \beta_M^l\}$ is the vector of coefficients of the linear combination for the l th projection, and $\mathbf{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_M\}$ is the matrix containing the parameters of the basis functions. The matrix \mathbf{W} is common for all the projections, while the corresponding β_l vectors are specific for each one. In this way, the coordinates of the pattern in the latent space are decided by each of the $f_l(\mathbf{x}, \theta_l)$ functions. Consequently, each input pattern \mathbf{x} is projected using the following mapping $\phi : \mathbb{R}^k \rightarrow \mathbb{R}^L$:

$$\phi(\mathbf{x}) = \{f_1(\mathbf{x}, \theta_1), f_2(\mathbf{x}, \theta_2), \dots, f_L(\mathbf{x}, \theta_L)\}.$$

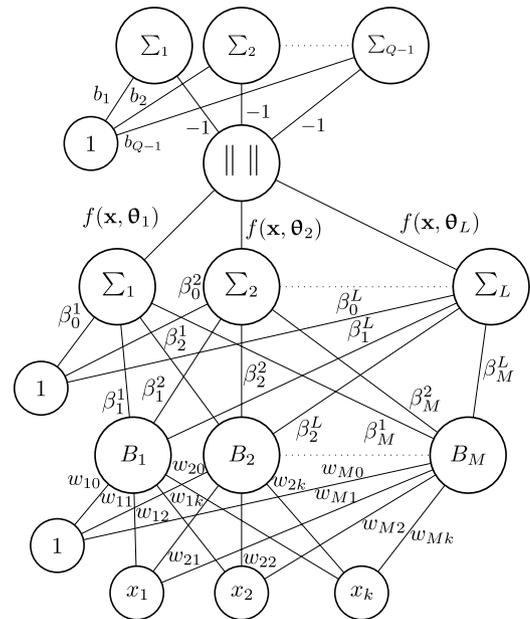


Fig. 2. CHNN structure and topology.

The dimensionality of the latent space is an additional parameter and different values will be considered in the experimental section.

The input conditional class probability is related to the distance to origin of the derived space. In order to simplify the calculus, we obtain the norm of the projection:

$$\|\phi(\mathbf{x})\| = \sqrt{f_1(\mathbf{x}, \theta_1)^2 + f_2(\mathbf{x}, \theta_2)^2 + \dots + f_L(\mathbf{x}, \theta_L)^2},$$

and then the different classes are defined by a set of thresholds $\{b_1, \dots, b_{Q-1}\}$, as in POM:

$$P(y \leq c_q | \mathbf{x}) = \frac{1}{(1 + \exp(b_q - \|\phi(\mathbf{x})\|))}, \quad q = 1, \dots, Q - 1.$$

The scheme of the proposed neural network can be seen in Fig. 2. Again, we need a set of thresholds—radii of the concentric hyperspheres in the latent space. These thresholds have to respect

Table 1

Characteristics of the 12 datasets used for the experiments: number of instances (Size), total number of inputs (#In.), number of classes (#Out.), and number of patterns per-class (NPPC).

Dataset	Size	#In.	#Out.	NPPC
Balance	625	4	3	(288, 49, 288)
Car	1728	21	4	(1210, 384, 69, 65)
ERA	1000	4	9	(92, 142, 181, 172, 158, 118, 88, 31, 18)
ESL	488	4	9	(2, 12, 38, 100, 116, 135, 62, 19, 4)
LEV	1000	4	5	(93, 280, 403, 197, 27)
Newth.	215	5	3	(30, 150, 35)
Pasture	36	25	3	(12, 12, 12)
Sq.-st.	52	51	3	(23, 21, 8)
Sq.-unst.	52	52	3	(24, 24, 4)
SWD	1000	10	4	(32, 352, 399, 217)
Tae	151	54	3	(49, 50, 52)
Toy	300	2	5	(35, 87, 79, 68, 31)

the class order, $b_1 \leq b_2 \leq \dots \leq b_{Q-1}$, so we again introduce padding variables (see previous subsection), $b_q = b_{q-1} + \Delta_q^2$, $1 < q < Q$, $b_1, \Delta_q \in \mathbb{R}$. Consequently, the parameter vector is now $\theta = \{\beta_1, \dots, \beta_L, \mathbf{W}, b_1, \Delta_2, \dots, \Delta_{Q-1}\}$.

With respect to the derivatives, they are the same as for POMNN model, if we replace f_n by $\|\phi(\mathbf{x})\|$. Then:

$$\frac{\partial f_n}{\partial \theta_s} = \frac{2 \cdot f(\mathbf{x}, \theta_l)}{\|\phi(\mathbf{x})\|^4} \cdot f_{nl} \cdot \frac{\partial f_{nl}}{\partial \theta_s},$$

where f_{nl} is $f(\mathbf{x}_n, \theta_l)$ and the derivatives of the parameters of each latent dimension f_{nl} can be expressed as:

$$\frac{\partial f_{nl}}{\partial \beta_0^k} = \begin{cases} 0 & \text{if } l \neq k, \\ 1 & \text{if } l = k. \end{cases}, \quad \frac{\partial f_{nl}}{\partial \beta_j^k} = \begin{cases} 0 & \text{if } l \neq k, \\ B_j(\mathbf{x}_n, \mathbf{w}_j) & \text{if } l = k. \end{cases}$$

The gradient for the hidden layer depends on the basis function. For sigmoidal nodes:

$$\frac{\partial f_{nl}}{\partial w_{jt}} = \beta_j^l B_j(\mathbf{x}_n, \mathbf{w}_j) (1 - B_j(\mathbf{x}_n, \mathbf{w}_j)) x_{nt}, \quad 1 \leq t \leq k,$$

$$\frac{\partial f_{nl}}{\partial w_{j0}} = \beta_j^l B_j(\mathbf{x}_n, \mathbf{w}_j) (1 - B_j(\mathbf{x}_n, \mathbf{w}_j)).$$

4. Experiments

4.1. Experimental design

The proposed methodology was applied to 12 different datasets taken from the UCI repository (Asuncion & Newman, 2007). As can be seen in Table 1, the characteristics vary noticeably among the datasets. We included a more controlled dataset (the toy one), which is a 2D problem synthetically generated, using the instructions given in da Costa et al. (2008). The performance of the different methods was evaluated by the following two metrics (Cruz-Ramírez, Hervás-Martínez, Sánchez-Monedero, & Gutiérrez, 2014):

- The Correctly Classified Ratio (CCR) is the error rate of the classifier, expressed as a percentage:

$$\text{CCR} = \frac{100}{N} \sum_{i=1}^N I(y_i^* = y_i),$$

where y_i is the true label and y_i^* is the predicted label. This metric evaluates the global classification performance task without taking into account the different kinds of errors with regard to the category order.

- The Mean Absolute Error (MAE) is the average deviation in absolute value of the predicted rank from the true one (Baccianella, Esuli, & Sebastiani, 2009):

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |\mathcal{O}(y_i) - \mathcal{O}(y_i^*)|,$$

where the position of a label in the ordinal scale is expressed by the function \mathcal{O} in the form $\mathcal{O}(c_q) = q$, $1 \leq q \leq Q$. $|\mathcal{O}(y_i) - \mathcal{O}(y_i^*)|$ is the distance between the true and predicted ranks. MAE values range from 0 to $Q - 1$ (maximum deviation in number of categories).

All the compared algorithms have been run 30 times for each dataset, considering a random holdout partition with a 75% of patterns for the training set and the remaining 25% for testing generalisation performance of the obtained classifier. The partitions were stratified in the sense that both partitions approximately presented the same class distribution of the complete dataset, and the same random partitions were considered for all the methods.

Apart from the methods previously presented in the paper (POM, NNN, POMNN and CHNN), different state-of-the-art methods were included in the experimentation for comparison purposes.

- Support vector ordinal regression (SVOR) by Chu and Ghahramani (2005) and Chu and Keerthi (2007), who optimised multiple thresholds in order to define parallel discriminant hyperplanes for the ordinal scales. In a first approach with explicit inequality constraints on the thresholds, they derive the optimality conditions for the dual problem, and adapt the SMO algorithm for the solution. We will refer to it as SVOREX. There is a second approach based on implicit constraints (SVORIM), but we consider SVOREX as we have found slightly better results for it in a previous work (Gutiérrez, Pérez-Ortiz, Fernández-Navarro, Sánchez-Monedero, & Hervás-Martínez, 2012).
- Kernel discriminant learning for ordinal regression (KDLOR) (Sun et al., 2010), which is an adaption of kernel discriminant analysis to the field of ordinal regression. Order on classes is imposed by constraining the projection to be obtained from the optimisation problem.
- Extreme learning machine with ordered partitions (ELMOP) (Deng et al., 2010), which is the ordinal regression version of ELMs. They are adapted by relabelling the dataset using the binary coding proposed in Frank and Hall (2001) and then fitting one multi-output model or several binary models based on the Error-Correcting Output Codes (ECOC) framework. The authors present the single model ELM as the one with the best performance (Deng et al., 2010), so this is the configuration chosen for ELMOP.
- AdaBoost for ordinal regression is based on an ordinal cost model for both the error estimation and pattern weights (ABORC3). As proposed in Riccardi et al. (in press), the base classifier is an ELM with Gaussian kernel and a regularisation parameter and the weighted least squares closed-form solution of the error function was considered for estimating the linear parameters of the individuals in the final ensemble model.

Table 2

Statistical results obtained when comparing the different neural network algorithms considered for this work.

CCR (Mean _{SD})						
Dataset	NNN	POMNN	CHNN1D	CHNN2D	CHNN3D	CHNN4D
Balance	95.75 _{1.73}	97.64 _{1.51}	96.79 _{1.45}	97.24 _{1.30}	97.28 _{1.19}	97.11 _{1.52}
Car	98.94 _{0.53}	98.98 _{0.59}	98.91 _{0.71}	99.04 _{0.35}	99.10 _{0.54}	98.99 _{0.57}
ERA	26.49 _{2.78}	26.96 _{2.47}	26.73 _{2.17}	26.96 _{2.51}	26.87 _{2.46}	26.37 _{2.57}
ESL	68.06 _{3.00}	70.08 _{3.12}	70.49 _{2.54}	70.77 _{3.11}	71.37 _{2.88}	69.89 _{3.51}
LEV	62.85 _{2.58}	62.43 _{2.97}	62.67 _{2.21}	62.93 _{2.20}	62.61 _{2.71}	62.85 _{3.15}
Newth.	96.98 _{2.15}	97.04 _{2.21}	96.79 _{2.48}	97.28 _{2.47}	97.22 _{2.32}	96.85 _{2.24}
Pasture	72.59 _{14.51}	68.89 _{13.81}	64.07 _{14.50}	71.85 _{14.22}	74.07 _{12.83}	69.63 _{14.27}
Sq. st.	63.08 _{12.35}	64.62 _{10.99}	63.85 _{10.91}	63.85 _{11.63}	66.92 _{12.48}	64.36 _{13.17}
Sq. unst.	75.13 _{13.20}	76.92 _{10.88}	66.15 _{15.33}	75.64 _{10.71}	76.41 _{13.69}	71.03 _{12.40}
SWD	56.81 _{2.69}	57.36 _{3.23}	57.84 _{3.54}	57.92 _{3.19}	57.96 _{3.25}	57.55 _{3.45}
Tae	57.19 _{8.09}	57.63 _{7.29}	55.61 _{8.92}	59.04 _{5.97}	59.91 _{7.30}	56.93 _{8.78}
Toy	92.93 _{3.40}	93.78 _{2.63}	92.76 _{3.26}	93.11 _{2.74}	93.42 _{2.55}	92.44 _{2.83}
\bar{R}_{CCR}	4.54	3.04	4.88	2.33	1.83	4.38
MAE (Mean _{SD})						
Dataset	NNN	POMNN	CHNN1D	CHNN2D	CHNN3D	CHNN4D
Balance	0.055 _{0.022}	0.025 _{0.018}	0.035 _{0.016}	0.030 _{0.016}	0.029 _{0.014}	0.031 _{0.016}
Car	0.013 _{0.007}	0.010 _{0.006}	0.011 _{0.007}	0.010 _{0.004}	0.009 _{0.005}	0.010 _{0.006}
ERA	1.301 _{0.073}	1.264 _{0.049}	1.263 _{0.051}	1.258 _{0.051}	1.257 _{0.051}	1.262 _{0.056}
ESL	0.346 _{0.035}	0.315 _{0.033}	0.310 _{0.029}	0.305 _{0.034}	0.302 _{0.032}	0.316 _{0.037}
LEV	0.406 _{0.027}	0.412 _{0.031}	0.410 _{0.023}	0.407 _{0.024}	0.409 _{0.029}	0.407 _{0.033}
Newth.	0.030 _{0.021}	0.030 _{0.022}	0.032 _{0.025}	0.027 _{0.025}	0.028 _{0.023}	0.032 _{0.022}
Pasture	0.274 _{0.145}	0.315 _{0.143}	0.367 _{0.163}	0.285 _{0.151}	0.259 _{0.128}	0.304 _{0.143}
Sq. st.	0.408 _{0.143}	0.372 _{0.116}	0.374 _{0.112}	0.372 _{0.120}	0.341 _{0.129}	0.372 _{0.143}
Sq. unst.	0.274 _{0.156}	0.233 _{0.112}	0.344 _{0.159}	0.251 _{0.118}	0.241 _{0.140}	0.292 _{0.125}
SWD	0.464 _{0.029}	0.448 _{0.035}	0.443 _{0.038}	0.443 _{0.034}	0.441 _{0.035}	0.447 _{0.036}
Tae	0.550 _{0.124}	0.539 _{0.096}	0.541 _{0.119}	0.514 _{0.081}	0.497 _{0.106}	0.517 _{0.109}
Toy	0.072 _{0.034}	0.062 _{0.026}	0.072 _{0.033}	0.069 _{0.027}	0.066 _{0.025}	0.076 _{0.028}
\bar{R}_{MAE}	4.75	3.50	4.83	2.41	1.58	3.91

The best result is shown in bold face and the second one in italics.

Regarding hyper-parameter tuning, the following procedure has been applied. For kernel algorithms, i.e. SVOREX, KDLOR and ABORC3, the width of the Gaussian kernel, γ , was adjusted using a grid search with a 5-fold cross-validation, considering the following range: $\gamma \in \{10^{-3}, 10^{-2}, \dots, 10^3\}$. An additional parameter u was also used for KDLOR in order to avoid singularities in the covariance matrices, $u \in \{10^{-6}, 10^{-5}, \dots, 10^{-2}\}$. For SVOREX and ABORC3, the additional cost parameter was adjusted by using the range $C \in \{10^{-3}, 10^{-2}, \dots, 10^3\}$. Finally, the results of the ABORC3 algorithm were directly taken from Riccardi et al. (in press), given that the authors used the same partitions as in this paper (as presented in Gutiérrez et al., 2012), and the same configuration for γ and C parameters. The number of members of the ensemble model was $M = 25$.

For the neural network algorithms, the hyper-parameters (number of hidden neuron, H , and number of iterations of the local search procedure, N), were adjusted using a grid search with a 5-fold cross-validation, considering the following ranges: $H \in \{1, 5, 15, 20, 30, 40\}$ and $N \in \{100, 200, 300, 400, 500\}$. For EL-MOP, a higher number of hidden neurons are considered, $H \in \{5, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100\}$, given that it relies on sufficiently informative random projections (Huang et al., 2012). With respect to $iRprop^+$, the parameter values are set to those specified in Igel and Hüsken (2003). Although more advanced methods could have been considered for adjusting the structure of the neural networks (Bueno-Crespo, García-Laencina, & Sancho-Gómez, 2013; Xi-zhao, Qing-yan, Qing, & Jun-hai, 2013), we have selected cross-validation to ease the comparisons against other state-of-the-art methods.

The CHNN algorithm has been run with different options for the dimensionality of the latent space ($L \in \{1, 2, 3, 4\}$), leading to four versions of the algorithm: CHNN1D, CHNN2D, CHNN3D and CHNN4D. Indeed, the CHNN1D is very similar to POMNN (they both project the patterns into a one dimensional space), but the classes

are arranged in nested intervals for CHNN1D, while the intervals are consecutive in the case of POMNN.

4.2. Comparison of the different alternatives for the neural network algorithms

In the first part of this section, we compare the different neural network algorithms in order to check (1) whether considering an ordinal regression configuration is able to improve the performance of the NNN classifier, and (2) if the additional dimensions for latent space help to better separate the classes. The results are shown in Table 2. This table includes the average ranking (\bar{R}) obtained for each metric and all datasets ($R = 1$ for the best method, $R = 6$ for the worst one).

First of all, attention should be drawn to the results of the nominal version (NNN) of the model presented in this paper. As can be seen, there are only two datasets (LEV and pasture), where the nominal version of the algorithm obtains a first or a second position (for CCR and MAE) when compared to the ordinal versions. A further analysis should be done on those two datasets to check if the order scale of the target variable can be found (in a linear or non-linear form) in the input space distribution of the training patterns. Regarding the ordinal methods, the one dimensional projection of CHNN1D (where the intervals are nested) seems to be too restrictive, and the results are a bit worse than those of POMNN. However, when the dimensionality is increased (CHNN2D and CHNN3D), the results are considerably better, although a too high dimensionality seems to compromise the learning capability (CHNN4D).

The projections of the CHNN2D model for the tae dataset are included in Fig. 1(b), showing a higher degree of flexibility than the POMNN model (Fig. 1(a)). The projection of the patterns in this space can be used to have better understanding about how they are organised with regard to the classification task, serving as a

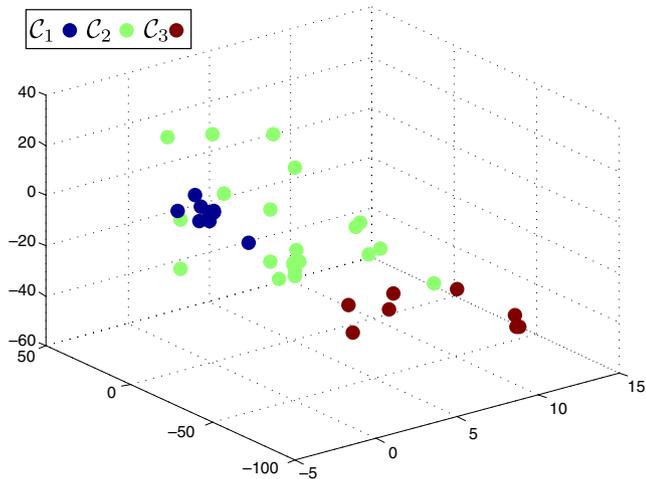


Fig. 3. Three-dimensional CHNN projection for the sq.-st. dataset. Hyperspheres have been omitted for clarity. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

visualisation tool. In Fig. 3, the result of CHNN3D has been included for one of the runs of sq.-st. dataset, where the model only confuses one green pattern, labelled as red.

An analysis of the significance of the differences found was also performed. We consider a procedure for comparing multiple classifiers over multiple datasets, following the guidelines of Demšar (2006). It begins with the Friedman test (Friedman, 1940; Zhai, Xu, & Wang, 2012), using the CCR and MAE ranking as the test variable. This test is a non-parametric equivalent to the repeated-measures ANOVA test. We apply it because a previous evaluation of the CCR and MAE ranking values results in rejecting the normality and the equality of variances' hypothesis. Applying this test to the average ranks at Table 2, the test shows that the effect of the method used for classification is statistically significant at a significance level of 5%, as the confidence interval is $C_0 = (0, F_{(\alpha=0.05)} = 2.38)$ and the F-distribution statistical values are: (1) for CCR, $F^* = 9.48 \notin C_0$; and (2) for MAE, $F^* = 10.06 \notin C_0$. Consequently, we reject the null-hypothesis stating that all algorithms perform equally in mean ranking and a post-hoc test is needed. For the post-hoc test, the best performing method (CHNN3D) was considered as the control method, and it was compared to the remaining ones according to their rankings. It has been noted that the approach of comparing all classifiers to each other in a post-hoc test is not as sensitive as the approach of comparing all classifiers to a given classifier (a control method). One approach to this latter type of comparison is the Holm test. The test statistics for comparing the i th and j th methods using this procedure is:

$$z = \frac{R_i - R_j}{\sqrt{\frac{k(k+1)}{6N}}}, \quad (10)$$

where k is the number of algorithms, N is the number of datasets and R_i is the mean ranking of the i th method. The z value is used to find the corresponding probability from the table of normal distribution, which is then compared with an appropriate level of significance α . Holm test adjusts the value for α in order to compensate for multiple comparisons. This is done in a step-up procedure that sequentially tests the hypotheses ordered by their significance. We will denote the ordered p -values by p_1, p_2, \dots, p_k so that $p_1 \leq p_2 \leq \dots \leq p_k$. Holm test compares each p_i with $\alpha'_{\text{Holm}} = \alpha/(k-i)$, starting from the most significant p value. If p_1 is below $\alpha/(k-1)$, the corresponding hypothesis is rejected and we allow to compare p_2 with $\alpha/(k-2)$. If the second hypothesis is rejected, the test proceeds with the third, and so on. As soon

Table 3

Holm test for the comparison of the different neural network algorithms: adjusted p -values using CCR and MAE as the test variables (CHNN3D is the control method).

Variable test: CCR			Variable test: MAE		
Algorithm	p -value	α'_{Holm}	Algorithm	p -value	α'_{Holm}
CHNN1D	0.000*	0.010	CHNN1D	0.000*	0.010
NNN	0.000*	0.013	NNN	0.000*	0.013
CHNN4D	0.001*	0.017	CHNN4D	0.002*	0.017
POMNN	0.114	0.025	POMNN	0.012*	0.025
CHNN2D	0.513	0.050	CHNN2D	0.275	0.050

* Statistically significant differences for $\alpha = 0.05$.

as a certain null hypothesis cannot be rejected, all the remaining hypotheses are retained as well.

The results of the Holm test can be seen in Table 3, using the corresponding p and α'_{Holm} values. From the results of this test, it can be concluded that the CHNN3D methodology obtains a significantly higher ranking when compared to all methods except CHNN2D (for CCR and MAE) and POMNN (for CCR).

The reason the proposed latent space helps to improve the classification for ordinal regression problems is that the ordering imposed by the projection helps to better locate new patterns in the ordinal scale, which is also true for the POMNN model. However, as opposed to PONNN, the CHNN model allows certain flexibility when constructing the multidimensional projection, which encourages more parsimonious models less prone to overfitting.

4.3. Comparison of CHNN3D against other state-of-the-art algorithms

This second part of the experiments compares the results obtained from CHNN3D to other state-of-the-art algorithms in order to check if the method can be considered competitive. The results of this comparison are included in Tables 4 and 5 and follow the same format as in the previous comparison. Again, the Friedman test shows that the effect of the method used for classification is statistically significant at a significance level of 5%, as the confidence interval is $C_0 = (0, F_{(\alpha=0.05)} = 2.38)$ and the F-distribution statistical values are: (1) for CCR, $F^* = 11.23 \notin C_0$; and (2) for MAE, $F^* = 6.72 \notin C_0$. From Table 4, it can be checked that the best ranking is obtained by CHNN3D for MAE and CCR. The second methods are ABORC3 for CCR and SVOREX for MAE. Recall that ABORC3 is an ensemble of kernel models, while CHNN3D is a single model. SVOREX is known to be one of the most competitive methods for ordinal regression (Chu & Keerthi, 2007; Gutiérrez et al., 2012), because it inherits the good properties of binary SVM. The statistical tests of Table 5 conclude that the differences are significant for all methods, except SVOREX and ABORC3. These results confirm that CHNN3D method is able to achieve a very competitive performance when compared to the state-of-the-art algorithms.

5. Discussion and conclusion

This paper is motivated by the fact that the one dimensional projection of threshold models in ordinal regression can be too restrictive, resulting in too complex nonlinear models or unrealistic assumptions for linear models. We propose to relax this projection by extending the latent space and allowing multiple dimensions. The order of the classes is organised in this space by using concentric hyperspheres centred in the origin, in such a way that intermediate classes are bounded by consecutive hyperspheres.

In analogy with kernel methods, increased dimensionality of internal "feature space" representations allows one to use simpler, less complex functions to accomplish the task. In classification

Table 4
Statistical results obtained when comparing the CHNN3D approach against other state-of-the-art algorithms.

CCR (Mean _{SD})						
Dataset	CHNN3D	ELMOP	POM	SVOREX	KDLOR	ABORC3
Balance	97.28 _{1.19}	91.85 _{2.30}	90.55 _{1.86}	99.79 _{0.56}	83.69 _{2.40}	95.12
Car	99.10 _{0.54}	84.85 _{1.10}	15.75 _{30.63}	98.74 _{0.53}	95.35 _{0.83}	98.30
ERA	26.87 _{2.46}	24.55 _{1.69}	25.61 _{2.11}	28.60 _{2.62}	20.51 _{3.46}	27.24
ESL	71.37 _{2.88}	69.51 _{3.30}	70.55 _{3.36}	71.23 _{3.36}	65.38 _{3.68}	71.88
LEV	62.61 _{2.71}	62.85 _{2.58}	62.33 _{2.80}	62.55 _{1.97}	54.93 _{3.46}	65.01
Newth.	97.22 _{2.32}	94.88 _{2.51}	97.22 _{2.22}	96.85 _{2.39}	97.53 _{2.02}	96.96
Pasture	74.07 _{12.83}	61.11 _{13.29}	49.63 _{15.37}	66.67 _{11.67}	66.30 _{12.54}	74.44
Sq. st.	66.92 _{12.48}	56.41 _{15.16}	38.21 _{15.18}	64.36 _{14.21}	64.87 _{11.90}	59.49
Sq. unst.	76.41 _{13.69}	61.28 _{12.37}	34.87 _{14.25}	73.85 _{11.88}	75.13 _{12.73}	65.64
SWD	57.96 _{3.25}	56.53 _{3.00}	56.79 _{2.96}	57.09 _{3.14}	49.08 _{2.93}	59.07
Tae	59.91 _{7.30}	54.39 _{7.52}	50.44 _{7.73}	57.89 _{5.82}	57.63 _{5.76}	57.89
Toy	93.42 _{2.55}	91.73 _{2.87}	28.93 _{2.55}	98.27 _{1.27}	88.58 _{3.05}	95.49
\bar{R}_{CCR}	1.96	4.67	5.04	2.63	4.33	2.38
MAE (Mean _{SD})						
Dataset	CHNN3D	ELMOP	POM	SVOREX	KDLOR	ABORC3
Balance	0.029 _{0.014}	0.089 _{0.026}	0.107 _{0.021}	0.002 _{0.006}	0.164 _{0.025}	0.050
Car	0.009 _{0.005}	0.175 _{0.012}	1.451 _{0.548}	0.013 _{0.005}	0.046 _{0.008}	0.018
ERA	1.257 _{0.051}	1.243 _{0.044}	1.218 _{0.050}	1.209 _{0.057}	1.783 _{0.099}	1.213
ESL	0.302 _{0.032}	0.322 _{0.035}	0.310 _{0.038}	0.302 _{0.035}	0.365 _{0.043}	0.301
LEV	0.409 _{0.029}	0.406 _{0.027}	0.409 _{0.030}	0.411 _{0.021}	0.512 _{0.039}	0.385
Newth.	0.028 _{0.023}	0.052 _{0.025}	0.028 _{0.022}	0.031 _{0.024}	0.025 _{0.020}	0.036
Pasture	0.259 _{0.128}	0.404 _{0.144}	0.585 _{0.204}	0.333 _{0.117}	0.341 _{0.130}	0.242
Sq. st.	0.341 _{0.129}	0.485 _{0.181}	0.813 _{0.248}	0.367 _{0.147}	0.374 _{0.145}	0.419
Sq. unst.	0.241 _{0.140}	0.423 _{0.137}	0.826 _{0.230}	0.262 _{0.119}	0.251 _{0.132}	0.354
SWD	0.441 _{0.035}	0.452 _{0.031}	0.450 _{0.030}	0.446 _{0.031}	0.579 _{0.035}	0.443
Tae	0.497 _{0.106}	0.625 _{0.115}	0.628 _{0.116}	0.466 _{0.061}	0.461 _{0.065}	0.500
Toy	0.066 _{0.025}	0.083 _{0.029}	0.981 _{0.039}	0.017 _{0.013}	0.114 _{0.030}	0.043
\bar{R}_{MAE}	2.29	4.58	4.83	2.45	4.17	2.67

The best result is shown in bold face and the second one in italics.

Table 5
Holm test for the comparison of the different neural network algorithms: adjusted p-values using CCR and MAE as the test variables (CHNN3D is the control method).

Variable test: CCR			Variable test: MAE		
Algorithm	p-value	α'_{Holm}	Algorithm	p-value	α'_{Holm}
POM	0.000*	0.010	POM	0.000*	0.010
ELMOP	0.000*	0.013	ELMOP	0.003*	0.013
KDLOR	0.002*	0.017	KDLOR	0.014*	0.017
SVOREX	0.383	0.025	ABORC3	0.623	0.025
ABORC3	0.546	0.050	SVOREX	0.827	0.050

* Statistically significant differences for $\alpha = 0.05$.

problems, one can always attempt to design potentially complex non-linear decision boundaries in the data space without the need to map to a higher dimensional feature space. However, the advantage of the feature space is that by (implicitly) representing important class structure properties in the higher dimensional space, usually much simpler (e.g. linear) functions are needed to accomplish the final task. These functions are easier to fit, the fit is more robust, and crucially, less prone to over-fitting.

By the very nature of binary classification, provided linear decision boundary is needed, we only need to specify its normal vector (and bias). In that case, the input points are eventually projected into a single dimension where the class membership is decided. The situation is different for ordinal regression. In one dimension, the class order can be naturally represented (as done in this paper) in a nested way, e.g. class C_1 would be represented by interval $[-a_1, a_1]$, class C_2 by $[-a_2, -a_1] \cup (a_1, a_2]$, class C_3 by $[-a_3, -a_2] \cup (a_2, a_3]$, etc. for some positive constants $a_1 < a_2 < \dots < a_{Q-1}$.

Assume the input points \mathbf{x} are projected onto the real line through function $\phi(\mathbf{x})$. By composing ϕ with a function $\psi(u) = |u|$, thus obtaining another projection function $v = \psi \circ \phi$, one can always transform the nested interval structure into the typical threshold structure used in ordinal regression, $[0, a_1], (a_1, a_2]$,

$(a_2, a_3]$, etc. However, if we did not fix ψ a-priori and wanted to fit v directly, the additional modelling burden would lead to more complex and hence more difficult to fit projection models.

We proposed to push this idea one step further: the nested interval structure can be naturally generalised to a series of nested L -dimensional hyper-spheres. Such structure still preserves the idea of class order, while eliminating the need to directly formulate and fit constrained 1-dimensional projections, as commonly done in ordinal regression. In our case we need to learn a less restricted projection function $\phi = (f_1, f_2, \dots, f_L)$ taking \mathbf{x} into \mathbb{R}^L . This projection is then combined with a fixed function $\psi(u_1, u_2, \dots, u_L) = \sqrt{u_1^2 + u_2^2 + \dots + u_L^2}$ to form the final projection $v = \psi \circ \phi$. Crucially, increasing the latent space dimension from 1 to L can allow us to put less strain on nonlinear projection function ϕ , taking us from the input space to the latent space that would be needed if we wanted to fit directly more complex projections to 1-dimensional structure of ordered intervals. In analogy with the feature space metaphor mentioned above, we suggest a “feature space” structure in our latent space that, at the price of increasing dimensions of the latent space (thus increasing its representational capacity), allows us to use a simple decision function ψ on top of the latent space. In contrast with kernel methods, we do not fix the “feature mapping” ϕ while learning the decision function ψ - instead, ψ is fixed and ϕ is learnt.

The model is implemented by using a neural network approach, where each of these dimensions is set to a linear combination of basis functions (sigmoidal nodes, in our case). The proposed model shows better performance than a nominal neural network and a neural network based on the proportional odds model, and competitive performance when compared to the state-of-the-art in ordinal regression. A study of the number of dimensions of this extended latent space is also performed, where both two and three dimensional spaces seem to be a good option, at least for the 12

datasets considered. Finally, the projections learnt by the model are shown to be useful for studying additional characteristics of the dataset, acting as a visualisation tool.

As future research lines, the same latent space structure could be tested with other ordinal regression models, e.g. linear models or kernel models.

Acknowledgments

The research work of P.A. Gutiérrez and C. Hervás-Martínez was partially subsidized by the TIN2011-22794 project of the Spanish Ministerial Commission of Science and Technology (MICYT), FEDER funds and the P11-TIC-7508 project of the “Junta de Andalucía” (Spain). P. Tiño was supported by EPSRC grant EP/L000296/1.

References

- Agresti, A. (2002). *Categorical data analysis* (2nd Ed.). John Wiley and Sons.
- Allwein, E. L., Schapire, R. E., & Singer, Y. (2001). Reducing multiclass to binary: a unifying approach for margin classifiers. *Journal of Machine Learning Research*, 1, 113–141.
- Asuncion, A., & Newman, D. (2007). UCI machine learning repository. URL <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
- Baccianella, S., Esuli, A., & Sebastiani, F. (2009). Evaluation measures for ordinal regression. In *Proceedings of the ninth international conference on intelligent systems design and applications (ISDA'09)* (pp. 283–287).
- Bueno-Crespo, A., García-Laencina, P. J., & Sancho-Gómez, J.-L. (2013). Neural architecture design based on extreme learning machine. *Neural Networks*, 48(0), 19–24. <http://dx.doi.org/10.1016/j.neunet.2013.06.010>.
- Cardoso, J. S., & da Costa, J. F. P. (2007). Learning to classify ordinal data: the data replication method. *Journal of Machine Learning Research*, 8, 1393–1429.
- Cardoso, J. S., da Costa, J. F. P., & Cardoso, M. (2005). Modelling ordinal relations with SVMs: an application to objective aesthetic evaluation of breast cancer conservative treatment. *Neural Networks*, 18(5–6), 808–817.
- Cheng, J., Wang, Z., & Pollastri, G. (2008). A neural network approach to ordinal regression. In *Proceedings of the IEEE international joint conference on neural networks (IJCNN2008, IEEE world congress on computational intelligence)* (pp. 1279–1284). IEEE Press.
- Chu, W., & Ghahramani, Z. (2005). Gaussian processes for ordinal regression. *Journal of Machine Learning Research*, 6, 1019–1041.
- Chu, W., & Keerthi, S. S. (2007). Support vector ordinal regression. *Neural Computation*, 19(3), 792–815.
- Costa, M. (1996). Probabilistic interpretation of feedforward network outputs, with relationships to statistical prediction of ordinal quantities. *International Journal of Neural Systems*, 7(5), 627–638.
- Cruz-Ramírez, M., Hervás-Martínez, C., Sánchez-Monedero, J., & Gutiérrez, P. A. (2014). Metrics to guide a multi-objective evolutionary algorithm for ordinal classification. *Neurocomputing*, 135, 21–31. <http://dx.doi.org/10.1016/j.neucom.2013.05.058>.
- da Costa, J. F. P., Alonso, H., & Cardoso, J. S. (2008). The unimodal model for the classification of ordinal data. *Neural Networks*, 21, 78–91.
- da Costa, J. F. P., & Cardoso, J. (2005). Classification of ordinal data using neural networks. In J. Gama, R. Camacho, P. Brazdil, A. Jorge, & L. Torgo (Eds.), *Lecture notes in computer science: Vol. 3720. Proceedings of the 16th European conference on machine learning*, (ECML 2005), (pp. 690–697). Berlin, Heidelberg: Springer.
- Demšar, J. (2006). Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7, 1–30.
- Deng, W.-Y., Zheng, Q.-H., Lian, S., Chen, L., & Wang, X. (2010). Ordinal extreme learning machine. *Neurocomputing*, 74(1–3), 447–456.
- Dietterich, T. G., & Bakiri, G. (1995). Solving multiclass learning problems via error-correcting output codes. *Journal of Artificial Intelligence Research*, 2(1), 263–286. URL <http://dl.acm.org/citation.cfm?id=1622826.1622834>.
- Dikkers, H., & Rothkrantz, L. (2005). Support vector machines in ordinal classification: an application to corporate credit scoring. *Neural Network World*, 15(6), 491–507.
- Dobrska, M., Wang, H., & Blackburn, W. (2012). Ordinal regression with continuous pairwise preferences. *International Journal of Machine Learning and Cybernetics*, 3(1), 59–70.
- Fernandez-Navarro, F., Riccardi, A., & Carloni, S. (2014). Ordinal neural networks without iterative tuning. *IEEE Transactions on Neural Networks and Learning Systems*, 1. <http://dx.doi.org/10.1109/TNNLS.2014.2304976> (in press).
- Frank, E., & Hall, M. (2001). A simple approach to ordinal classification. In *Proceedings of the 12th European conference on machine learning, EMCL'01*. (pp. 145–156). London, UK: Springer-Verlag.
- Friedman, M. (1940). A comparison of alternative tests of significance for the problem of m rankings. *Annals of Mathematical Statistics*, 11(1), 86–92.
- Gutiérrez, P. A., Pérez-Ortiz, M., Fernandez-Navarro, F., Sánchez-Monedero, J., & Hervás-Martínez, C. (2012). An experimental study of different ordinal regression methods and measures. In *7th international conference on hybrid artificial intelligence systems* (pp. 296–307).
- Huang, G.-B., Zhou, H., Ding, X., & Zhang, R. (2012). Extreme learning machine for regression and multiclass classification. *IEEE Transactions on Systems, Man and Cybernetics, Part B: Cybernetics*, 42(2), 513–529.
- Igel, C., & Hüsken, M. (2003). Empirical evaluation of the improved Rprop learning algorithms. *Neurocomputing*, 50(6), 105–123.
- Kim, M., & Pavlovic, V. (2010). Structured output ordinal regression for dynamic facial emotion intensity prediction. In K. Daniilidis, P. Maragos, & N. Paragios (Eds.), *Lecture notes in computer science: Vol. 6313. Proceedings of the 11th European conference on computer vision (ECCV 2010), part III* (pp. 649–662). Berlin, Heidelberg: Springer.
- Li, L., & Lin, H.-T. (2007). Ordinal regression by extended binary classification. *Advances in Neural Information Processing Systems*, 19, 865–872.
- Lin, H.-T., & Li, L. (2012). Reduction from cost-sensitive ordinal ranking to weighted binary classification. *Neural Computation*, 24(5), 1329–1367.
- Mathieson, M. J. (1996). Ordinal models for neural networks. In J. M. A.-P. N. Refenes, Y. Abu-Mostafa, & A. Weigend (Eds.), *Neural networks in financial engineering, Proceedings of the third international conference on neural networks in the capital markets* (pp. 523–536). World Scientific.
- Mathieson, M. J. (1999). Ordered classes and incomplete examples in classification. In T. P. Michael, C. Mozer, & Michael I. Jordan (Eds.), *Advances in neural information processing systems: Vol. 9. Proceedings of the 1996 conference on neural information processing systems (NIPS)* (pp. 550–556).
- McCullagh, P. (1980). Regression models for ordinal data. *Journal of the Royal Statistical Society. Series B. Statistical Methodology*, 42(2), 109–142.
- Riccardi, A., Fernandez-Navarro, F., & Carloni, S. (2014). Cost-sensitive adaboost algorithm for ordinal regression based on extreme learning machine. *IEEE Transactions on Cybernetics*, <http://dx.doi.org/10.1109/TCYB.2014.2299291> (in press).
- Sánchez-Monedero, J., Gutiérrez, P. A., Tino, P., & Hervás-Martínez, C. (2013). Exploitation of pairwise class distances for ordinal classification. *Neural Computation*, 25(9), 2450–2485.
- Sun, B.-Y., Li, J., Wu, D. D., Zhang, X.-M., & Li, W.-B. (2010). Kernel discriminant learning for ordinal regression. *IEEE Transactions on Knowledge and Data Engineering*, 22(6), 906–910.
- Verwaeren, J., Waegeman, W., & De Baets, B. (2012). Learning partial ordinal class memberships with kernel-based proportional odds models. *Computational Statistics & Data Analysis*, 56(4), 928–942.
- Xi-zhao, W., Qing-yan, S., Qing, M., & Jun-hai, Z. (2013). Architecture selection for networks trained with extreme learning machine using localized generalization error model. In *Advances in extreme learning machines (ELM 2011), Neurocomputing*, 102(0), 3–9. <http://dx.doi.org/10.1016/j.neucom.2011.12.053>.
- Zhai, J.-h., Xu, H.-y., & Wang, X.-z. (2012). Dynamic ensemble extreme learning machine based on sample entropy. *Soft Computing*, 16(9), 1493–1502. <http://dx.doi.org/10.1007/s00500-012-0824-6>.