

# Deep ordinal classification based on cumulative link models

Víctor-Manuel Vargas, Pedro-Antonio Gutiérrez, *Senior Member, IEEE*, and  
César Hervás-Martínez, *Senior Member, IEEE*

**Abstract**—This paper proposes a deep convolutional neural network model for ordinal regression by considering a family of probabilistic ordinal link functions in the output layer. The link functions are those used for cumulative link models, which are traditional statistical linear models based on projecting each pattern into a 1-dimensional space. A set of ordered thresholds splits this space into the different classes of the problem. In our case, the projections are estimated by a non-linear deep neural network. To further improve the results, we combine these ordinal models with a loss function that takes into account the distance between the categories, based on the weighted Kappa index. Three different link functions are studied in the experimental study, and the results are contrasted with statistical analysis. The experiments run over two different ordinal classification problems, and the statistical tests confirm that these models improve the results of a nominal model and outperform other proposals considered in the literature.

**Index Terms**—Deep learning, ordinal regression, cumulative link models, Kappa index.

## I. INTRODUCTION

DEEP LEARNING, introduced by Lecun et al. [1], combines multiple machine learning techniques and allows computational models that are composed of numerous processing layers to learn representations of data with various levels of abstraction. These methods have dramatically improved the state-of-the-art in many domains, such as image classification [2], [3], [4], speech recognition [5], control problems [6], object detection [7], [8], privacy and security protection [9], [10], [11], recovery of human pose [12], semantic segmentation [13] and image retrieval [14], [15]. Convolutional Neural Networks (CNN) are one of the types of deep networks that are designed to process data that come in the form of multiple arrays. CNNs are appropriate for images, video, speech and audio processing, and they have been used extensively in the last years for automatic classification tasks [16], [17], [18]. For images, each colour channel is represented by a 2D array, and convolutional layers extract the main features from the pixels, and, after that, a fully connected layer classify every sample based on its extracted features. At the output of the CNN, a softmax function provides the probabilities of the set of classes predefined in the model for classification tasks. However, the softmax may not be the best option depending on the problem considered.

Ordinal classification problems are those classification tasks where labels are ordered, and there are different

inter-classes importances for each pair of categories. This kind of problem can be treated as nominal classification, but this discards the ordinal information. A better approach is to use specific methods that take the ordinality into account to improve the performance of the classification model. The Proportional Odds Model (POM) [19] is an ordinal alternative to the binary logistic regression. It belongs to a wider family of models called Cumulative Link Models (CLMs) [20]. CLMs are inspired in the concept of a latent variable that is projected into a 1-dimensional space and a set of thresholds that divides the projection into the different ordinal levels. A link function needs to be specified, which can be of different types, although the most common option is the `logit`, which is used in POM. In this paper, we explore different existing alternatives, as explained in depth in Section III-A.

In this paper, we propose the use of CLMs for deriving deep learning ordinal classifiers<sup>1</sup>. In the case of CNNs, the model projection used by the threshold model can be obtained from the last layer of the network. Given that we work with a 1-dimensional space, the last layer would have only one neuron (projection of the pattern), and its value could be used to classify the sample into the corresponding class according to the thresholds. Some previous works have used the `logit` in shallow neural networks [21], but this strategy has not been considered for deep learning, and alternative link functions have not been evaluated. To further improve the results, we train these models by minimising an ordinal loss function based on the Weighted Kappa index [22], instead of using the standard cross-entropy.

An experimental study evaluating the three most common link functions is performed. Also, other parameters that can affect the training process and the model performance are studied, such as the learning rate of the optimization algorithm, the batch size, and their interaction. The nominal version of this model is used as a baseline for comparison. We contrast the results obtained with a statistical analysis to provide more robust conclusions. An ANOVA III test [23], followed by a posthoc Tukey's test [24], is performed over 5 runs of the experiments, because of the demands of computational time required to run a higher number of executions. The experiments are run using two different ordinal datasets: Diabetic

<sup>1</sup>The source code is available at <https://github.com/ayrna/deep-ordinal-clm>.

Retinopathy [22], which contains high-resolution fundus images related with diabetes disease, and Adience [25], which includes human faces images associated with an age range.

The paper is organized as follows: in Section II, we analyse previous related works. Section III presents a formal description of the proposal in this paper, which combines a CLM with an ordinal loss function. In Section IV, we describe the experiments and the datasets used, while, in Section V, we present the results obtained and the statistical analysis. Finally, Section VI exposes the conclusions of this work.

## II. RELATED WORKS

There are many works related to the application and development of CNN models [26], but few works focus on ordinal classification problems. The existing deep ordinal approaches are mainly based on simply using an ordinal evaluation metric, on solving the ordinal problem as multiple binary sub-problems, on using an ordinal loss functions or on constraining the probability distribution of the output layer. These works are described in the following subsections.

### A. Simply using an ordinal evaluation metric.

ALALI et al. [27] proposed a complex CNN architecture for solving Twitter Sentiment Classification as an ordinal problem. They checked that using average pooling preserves significant features that provide more expressiveness to ordinal scale. They didn't propose any method to include the ordinal information into the classifier, but they tried to find the best CNN model architecture based on an ordinal metric.

### B. Solving the ordinal problem as multiple binary sub-problems

Niu et al. [28] proposed a learning approach to address ordinal regression problems using CNNs. They divided the problem into a series of binary classification sub-problems and proposed a multiple output CNN optimization algorithm to collectively solve these classification sub-problems, taking into account the correlation between them.

Li et al. [29] applied deep learning techniques for solving the ordinal problem of Alzheimer's diagnosis and detecting the different levels of the disease as multiple binary sub-problems.

Liu et al. [26] proposed a new approach which transforms the ordinal regression problem to binary classification sub-problems and use triplets with instances from different categories to train deep neural networks. In this way, high-level features describing the ordinal relationships are extracted automatically. Given that triplets must be generated, this approach is only recommended for small datasets.

Chen et al. [30] proposed a deep learning method termed Ranking-CNN. This method combines multiple

binary CNNs that are trained with ordinal age labels. The binary outputs are aggregated for the final age prediction. They achieved a tighter error bound for ranking-based age estimation.

In general, all these approaches increment the number of parameters to adjust, as several binary classifiers are simultaneously learnt.

### C. Using an ordinal loss function

de la Torre et al. [22] proposed the use of a continuous version of the quadratic weighted kappa (QWK) metric as loss function for the optimization algorithm. They compared this cost function against the traditional log-loss function using three different datasets, including the Diabetic Retinopathy database as the most complex one. They proved that their function could improve the results as it reduces overfitting and training time. Also, they checked the importance of hyper-parameter tuning.

Rios et al. [31] presented a CNN model designed to handle ordinal regression tasks on psychiatric notes. They combined an ordinal loss function, a CNN model and conventional feature engineering. Also, they applied a technique called Locally Interpretable Model-agnostic Explanation (LIME) to make the non-linear model more interpretable.

Fu et al. [32] applied deep learning techniques to Monocular Depth Estimation. They introduced a spacing-increasing discretization strategy to treat the problem as an ordinal regression problem. They improved the performance when training the network with an ordinal regression loss. Also, they used a multi-scale network structure that avoids unnecessary spatial pooling.

Pal et al. [33] defined a loss function for CNN that is based on the Earth Mover's Distance and takes into account the ordinal class relationships.

Liu et al. [34] proposed a constrained optimization formulation for the ordinal regression problem which minimizes the negative loglikelihood for a multi-class problem constrained by the order relationship between instances.

Although the use of these losses introduces the ordinality in model learning, the nature of the models remain nominal.

### D. Unimodal probability distributions

Beckham and Pal [25] proposed a straightforward technique to constrain discrete ordinal probability distributions to be unimodal, via the use of the Poisson and binomial probability distributions. The parameters of these distributions were learnt by using a deep neural network. They evaluated this approach on two large ordinal image datasets, including the Adience dataset used in this paper, obtaining promising results. They also included a simple squared-error reformulation [35] that was sensitive to class ordering.

This approach is the one most related to the CLMs considered in this paper. However, CLMs indirectly model a latent space together with the set of threshold separating

the ordered classes, which provides a more flexible and interpretable approach to deep ordinal classification.

### III. MODEL PROPOSAL

Based on the previous analysis of the state-of-the-art, our proposal is to combine a flexible threshold model in the output layer (different forms of a CLM) with an ordinal loss function, in order to better introduce ordinal constraints during learning.

#### A. Cumulative Link Model (CLM)

An ordinal classification problem consists in predicting the label  $y$  of an input vector  $\mathbf{x}$ , where  $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^K$  and  $y \in \mathcal{Y} = \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_Q\}$ , i.e.  $\mathbf{x}$  is in a  $K$ -dimensional input space, and  $y$  is in a label space of  $Q$  different labels. The objective in an ordinal problem is to find a function  $r : \mathcal{X} \rightarrow \mathcal{Y}$  to predict the labels or categories of new patterns, given a training set of  $N$  samples,  $D = \{(\mathbf{x}_i, y_i), i = 1, \dots, N\}$ . Labels have a natural ordering in ordinal problems:  $\mathcal{C}_1 \prec \mathcal{C}_2 \prec \dots \prec \mathcal{C}_Q$ . The order between labels gives us the possibility to compare two different elements of  $\mathcal{Y}$  by using the relation  $\prec$ . This is not possible under the nominal classification setting. In regression (where  $y \in \mathbb{R}$ ), real values in  $\mathbb{R}$  can be ordered by the standard  $<$  operator, but labels in ordinal regression ( $y \in \mathcal{Y}$ ) do not carry metric information, i.e. the category serves as a qualitative indication of the pattern rather than a quantitative one.

The Proportional Odds Model (POM) arises from a statistical background and is one of the first models designed explicitly for ordinal regression [19]. It dated back to 1980 and is a member of a wider family of models lately recognised as Cumulative Link Models (CLMs) [20]. CLMs predict probabilities of groups of contiguous categories, taking the ordinal scale into account. In this way, cumulative probabilities  $P(y \prec \mathcal{C}_q | \mathbf{x})$  are estimated, which can be directly related to standard probabilities:

$$\begin{aligned} P(y \preceq \mathcal{C}_q | \mathbf{x}) &= P(y = \mathcal{C}_1 | \mathbf{x}) + \dots + P(y = \mathcal{C}_q | \mathbf{x}), \\ P(y = \mathcal{C}_q | \mathbf{x}) &= P(y \preceq \mathcal{C}_q | \mathbf{x}) - P(y \preceq \mathcal{C}_{q-1} | \mathbf{x}), \end{aligned}$$

with  $q = 2, \dots, Q-1$ , and considering that  $P(y = \mathcal{C}_1 | \mathbf{x}) = P(y \preceq \mathcal{C}_1 | \mathbf{x})$  and  $P(y \preceq \mathcal{C}_Q | \mathbf{x}) = 1$ .

The model is inspired by the notion of a latent variable, where  $f(\mathbf{x})$  represents a one-dimensional mapping. The decision rule  $r : \mathcal{X} \rightarrow \mathcal{Y}$  is not fitted directly, but stochastic ordering of space  $\mathcal{X}$  is satisfied by the following general model form [36]:

$$g^{-1}(P(y \preceq \mathcal{C}_q | \mathbf{x})) = b_q - f(\mathbf{x}), \quad q = 1, \dots, Q-1,$$

where  $g^{-1} : [0, 1] \rightarrow (-\infty, +\infty)$  is a monotonic function often termed as the inverse link function, and  $b_q$  is the threshold defined for class  $\mathcal{C}_q$ . Consider the latent variable  $y^* = f(\mathbf{x})^* = f(\mathbf{x}) + \epsilon$ , where  $\epsilon$  is the random component of the error. The most common choice for the probability distribution of  $\epsilon$  is the logistic function (which is the default function for POM). Label  $\mathcal{C}_q$  is predicted if and only if  $f(\mathbf{x}) \in [b_{q-1}, b_q]$ , where the function  $f$

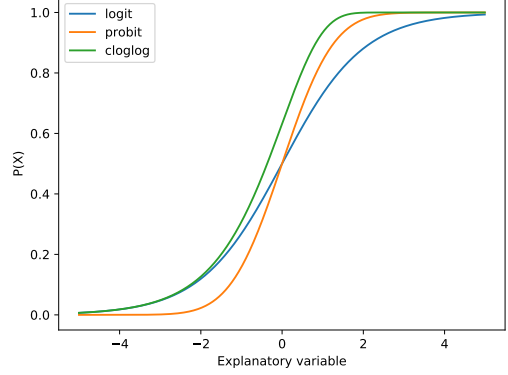


Fig. 1. Different link functions commonly used for CLMs.

and  $\mathbf{b} = (b_0, b_1, \dots, b_{Q-1}, b_Q)$  are to be determined from the data. It is assumed that  $b_0 = -\infty$  and  $b_Q = +\infty$ , so the real line defined by  $f(\mathbf{x}), \mathbf{x} \in \mathcal{X}$ , is divided into  $Q$  consecutive intervals. Each interval corresponds to a category. The constraints  $b_1 \leq b_2 \leq \dots \leq b_{Q-1}$  ensure that  $P(y \preceq \mathcal{C}_q | \mathbf{x})$  increases with  $q$  [19].

In this work, we consider different link functions previously proposed in CLMs for the probability distribution of  $\epsilon$ , including **logit**, **probit** and complementary log-log (**clog-log**). These three types of links are explained below and represented in Figure 1. They all follow the same form  $\text{link}[P(y \preceq \mathcal{C}_q | \mathbf{x})] = b_q - f(\mathbf{x})$ .

- The **logit** link function is the function used for the POM and is defined as:

$$\begin{aligned} \text{logit}[P(y \preceq \mathcal{C}_q | \mathbf{x})] &= \log \frac{P(y \preceq \mathcal{C}_q | \mathbf{x})}{1 - P(y \preceq \mathcal{C}_q | \mathbf{x})} = \\ &= b_q - f(\mathbf{x}), \quad q = 1, \dots, Q-1, \end{aligned}$$

or the equivalent expression:

$$P(y \preceq \mathcal{C}_q | \mathbf{x}) = \frac{1}{1 + e^{-(b_q - f(\mathbf{x}))}}.$$

- The **probit** link function is the inverse of the standard normal cumulative distribution function (cdf)  $\Phi$ . Its expression is:

$$\begin{aligned} \Phi^{-1}[P(y \preceq \mathcal{C}_q | \mathbf{x})] &= b_q - f(\mathbf{x}), \quad q = 1, \dots, Q-1, \\ P(y \preceq \mathcal{C}_q | \mathbf{x}) &= \Phi(b_q - f(\mathbf{x})), \quad q = 1, \dots, Q-1, \end{aligned}$$

which can also be expressed as:

$$P(y \preceq \mathcal{C}_q | \mathbf{x}) = \int_{-\infty}^{b_q - f(\mathbf{x})} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx.$$

- The **clog-log** takes a response that is restricted to the  $(0, 1)$  interval and converts it into a value in the  $(-\infty, +\infty)$  interval (like **logit** and **probit** transformations). The **clog-log** expression is:

$$\log[-\log[1 - P(y \preceq \mathcal{C}_q | \mathbf{x})]] = b_q - f(\mathbf{x}),$$

with  $q = 1, \dots, Q-1$ , that is:

$$P(y \preceq \mathcal{C}_q | \mathbf{x}) = 1 - e^{-e^{b_q - f(\mathbf{x})}}, \quad q = 1, \dots, Q-1.$$

**logit** and **probit** links are symmetric:

$$\text{link}[P(y \preceq C_q | \mathbf{x})] = -\text{link}[1 - P(y \preceq C_q | \mathbf{x})],$$

which means that the response curve for  $P(y \preceq C_q | \mathbf{x})$  is symmetric around the point  $P(y \preceq C_q | \mathbf{x}) = 0.5$ , i.e.  $P(y \preceq C_q | \mathbf{x})$  has the same rate when approaching 0 than when approaching 1. This symmetry property can be demonstrated as follows:

- 1) Let  $P(y \preceq C_q | \mathbf{x}) \equiv p$ . For the **logit** function, we have:

$$\begin{aligned} \text{link}(p) = \text{logit}(p) &= \log\left(\frac{p}{1-p}\right) = \\ &= \log(p) - \log(1-p), \end{aligned}$$

while:

$$\begin{aligned} -\text{link}(1-p) &= -\text{logit}(1-p) = \\ &= -\log\left(\frac{1-p}{p}\right) = -\log(1-p) + \log(p). \end{aligned}$$

- 2) For the **probit**:

$$\begin{aligned} p \equiv P(y \preceq C_q | \mathbf{x}) &= \Phi(b_q - f(\mathbf{x})) = \\ &= \int_{-\infty}^{b_q - f(\mathbf{x})} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx, \end{aligned}$$

which leads to:

$$\begin{aligned} \Phi^{-1}(p) &= \Phi^{-1}(p) = b_q - f(\mathbf{x}), \\ -\Phi^{-1}(1-p) &= \Phi^{-1}(1-p) = -b_q + f(\mathbf{x}), \end{aligned}$$

where:

$$\begin{aligned} 1-p &= \int_{-\infty}^{-b_q + f(\mathbf{x})} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx, \\ p &= 1 - \int_{-\infty}^{-b_q + f(\mathbf{x})} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx. \end{aligned}$$

Unlike **logit** and **probit**, the **clog-log** link is asymmetrical. In this way, when the given data is not symmetric in the  $[0, 1]$  interval and increase slowly at small to moderate value but increases sharply near 1, the **logit** and **probit** models are inappropriate, while **clog-log** can lead to better results.

In this paper, the probabilistic structure of CLMs is proposed as a link function for deep convolutional neural networks. This can be achieved by defining a new type of output layer alternative to the standard softmax layer. In this way, the proposed output layer will transform the one-dimensional projection, previously denoted as  $f(\mathbf{x})$ , into a set of probabilities.  $f(\mathbf{x})$  is estimated from a nonlinear transformation of the set of features learnt by the previous layers,  $l(\mathbf{x}) = f(\mathbf{x})$ , where  $\mathbf{x}$  is the pattern being evaluated and  $l(\mathbf{x})$  is a latent representation of the pattern given by the output of a single neuron. In order to apply unconstrained optimizers while ensuring  $b_1 \leq b_2 \leq \dots \leq b_{Q-1}$ , we can redefine the thresholds. All of them can be derived from the first one in the following form:

$$b_q = b_1 + \sum_{q=1}^{q-1} \alpha_q^2, \quad q = 2, \dots, Q,$$

where  $b_1$  and  $\alpha_q$  are the learnable parameters, and  $Q$  is the number of classes.

### B. Continuous Quadratic Weighted Kappa (QWK) loss function

In order to increase the performance of the deep ordinal model, the CLM structure in the output layer is combined with the continuous version of the QWK loss [22] function. The Kappa index is a well-known metric that measures the agreement between two different raters. The Weighted Kappa (WK) [37] is based on the Kappa index and adds different weights to the different types of disagreements based on a weight matrix. It is useful to evaluate the performance in ordinal problems, as it gives a higher weight to the errors that are further from the correct class. This metric is defined as follows:

$$\text{QWK} = 1 - \frac{\sum_{i,j} \omega_{i,j} O_{i,j}}{\sum_{i,j} \omega_{i,j} E_{i,j}}, \quad (1)$$

where  $N$  is the number of samples rated,  $\omega$  is the penalization matrix (in this case, quadratic weights are considered,  $\omega_{i,j} = \frac{(i-j)^2}{(C-1)^2}$ ,  $\omega_{i,j} \in [0, 1]$ ),  $O$  is the confusion matrix,  $E_{ij} = \frac{O_{i\bullet} O_{\bullet j}}{N}$ ,  $O_{i\bullet}$  is the sum of the  $i$ -th row and  $O_{\bullet j}$  is the sum of the  $j$ -th column.

The WK defined above cannot be used as a loss function for the optimization algorithm as it is not continuous. However, it has been previously redefined [22] in terms of probabilities of the predictions:

$$\text{QWK}_c = \frac{\sum_{k=1}^N \sum_{q=1}^Q \omega_{t_k, q} P(y = C_q | \mathbf{x}_k)}{\sum_{i=1}^Q \frac{N_i}{N} \sum_{j=1}^Q (\omega_{i,j} \sum_{k=1}^N P(y = C_j | \mathbf{x}_k))},$$

where  $\text{QWK}_c \in [0, 2]$ ,  $\mathbf{x}_k$  and  $t_k$  are the input data and the real class of the  $k$ -th sample,  $Q$  is the number of classes,  $N$  is the number of samples,  $N_i$  is the number of samples of the  $i$ -th class,  $P(y = C_q | \mathbf{x}_k)$  is the probability that the  $k$ -th sample belongs to class  $C_q$  (estimated using the CLM structure), and  $\omega_{i,j}$  are the elements of the penalization matrix ( $\omega_{i,j} = \frac{(i-j)^2}{(C-1)^2}$ ). This loss function can be minimized using a gradient descent based algorithm.

## IV. EXPERIMENTS

### A. Data

In order to evaluate the different models, we make use of two ordinal datasets:

- 1) *Diabetic Retinopathy (DR)*: DR<sup>2</sup> is a dataset consisting of extremely high-resolution fundus image data. The training set consists of 17563 pairs of images (where a pair includes a left and right eye image corresponding to a patient). In this dataset, we try to predict the correct category from five levels of diabetic retinopathy: no DR

<sup>2</sup><https://www.kaggle.com/c/diabetic-retinopathy-detection/data>

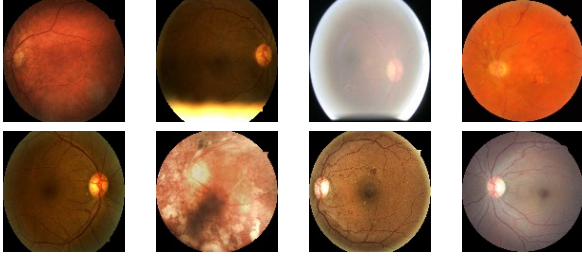


Fig. 2. Examples taken from the Diabetic Retinopathy test set.

(25810 images), mild DR (2443 images), moderate DR (5292 images), severe DR (873 images), or proliferative DR (708 images). The test set contains 26788 pairs of images. These images are taken in variable conditions: by different cameras, conditions of illumination and resolutions. They come from the EyePACS dataset that was used in the DR detection competition hosted on the Kaggle platform. Also, this dataset has been used in different works [22], [38], where the  $QWK_c$  cost function was considered in [22] to achieve better performance. A validation set is set aside, consisting of 10% of the patients in the training set. The images are resized to 128 by 128 pixels and rescaled to  $[0, 1]$  range. Data augmentation techniques, described in Section IV-C, are applied to achieve a higher number of samples. A few test images of this dataset are shown in Figure 2.

2) *Adience*: Adience<sup>3</sup> dataset consists of 26580 faces belonging to 2284 subjects. We use the form of the dataset where faces have been pre-cropped and aligned. The dataset was preprocessed, using the methods described in a previous work [25], so that the images are 256 pixels in width and height, and pixels values follow a  $(0; 1)$  normal distribution. The original dataset was split into five cross-validation folds. The training set consists of merging the first four folds which comprise a total of 15554 images. From this, 10% of the images are held out as part of a validation set. The last fold is used as test set. Some images of this dataset are shown in Figure 3. Adience dataset has been used in other works for human age estimation but most of them solved the problem as a multi-class problem instead of using the ordinal relation between classes. Eidinger et al. [39] presented an approach using support vector machines and neural networks. Chen et al. [40] proposed a coarse-to-fine strategy for deep CNNs. Levi and Hassner [41] presented another convolutional network model for age estimation. As previously discussed, Beckham and Pal [25] proposed a straightforward method to constrain discrete probability distributions to be unimodal.

### B. Model

CNNs have been used for both datasets. The different architectures of CNNs used in these experiments are presented in Table I. The architecture for DR is the same that was used in [22] and the network for Adience is a small



Fig. 3. Examples taken from the Adience test set.

TABLE I  
DESCRIPTION OF THE ARCHITECTURE USED FOR EACH DATASET.

Diabetic Retinopathy		Adience	
Layer	Output	Layer	Output
2 Conv3x3@32s1	252x252x32	Conv7x7@32s2	112x112x32
MaxPool2x2s2	126x126x32	MaxPool3x3s2	55x55x32
2 Conv3x3@64s1	122x122x64	2 ResBlock3x3@64s1	55x55x32
MaxPool2x2s2	61x61x64	ResBlock3x3@128s2	28x28x64
2 Conv3x3@128s1	57x57x128	2 ResBlock3x3@128s1	28x28x64
MaxPool2x2s2	28x28x128	1 ResBlock3x3@256s2	14x14x128
2 Conv3x3@128s1	24x24x128	2 ResBlock3x3@256s1	14x14x128
MaxPool2x2s2	12x12x128	1 ResBlock3x3@512s2	7x7x256
Conv4x4@128s1	9x9x128	2 ResBlock3x3@512s1	7x7x256
		AveragePool7x7s2	1x1x256

Residual Network (ResNet) [3] that was used in [25]. The most important parameters for convolutional layers are the number of filters that are used to make the convolution operation, the size of these filters and the stride, which is the number of pixels that the filter is moved in every operation. Pooling layers have similar parameters: pool size (number of pixels that will be involved in the operation) and stride. For convolutional layers, ConvWxH@FsS stands for filters of size WxH and stride S. For pooling layers, PoolWxHsS corresponds to a pool size of WxH and stride S.

The Exponential Linear Unit (ELU) [42] has been used as the activation function for all the convolutional and dense layers, instead of the ReLU [43] function, as it mitigates the effects of the vanishing gradient problem [44], [45] via the identity for positive values. Also, ELUs lead to faster training and better generalization performance than ReLU and Leaky ReLU (LReLU) [46] functions on networks with more than five layers.

After every ELU activation function of the convolutional layers, Batch Normalization [47] is applied. This method reduces the internal covariate shift by normalizing layer outputs. It allows us to use higher learning rates and be less careful about weight initialization. It also eliminates the need for using regularization techniques like Dropout.

At the output of the network, the CLM is used (see Section III-A). Also, a learnable parameter  $\tau$  has been used to rescale the projections used by the CLM to make it more stable and guarantee the convergence in most cases. The following expression describes the transformation applied to these projections:

$$f(\mathbf{x}) = \frac{l(\mathbf{x})}{\tau}.$$

where  $\tau$  is optimized as a free parameter.

<sup>3</sup><http://www.openu.ac.il/home/hassner/Adience/data.html>

### C. Experimental design

Weights are adjusted using a batch based first-order optimization algorithm called Adam [48]. We study different initial learning rates ( $\eta$ ) in order to find the optimal one for each problem. We apply an exponential decay [49] across training epochs to the initial learning rate ( $\eta_0$ ) following the expression below:

$$\eta = \eta_0 \cdot e^{-0.025 \cdot \text{epoch}}.$$

Both datasets are artificially balanced using data augmentation techniques [50]. However, different transformations are applied to each one. DR dataset augmentation is based on image cropping and zooming, horizontal and vertical flipping, brightness adjustment and random rotations. Horizontal flipping is the only transformation applied to the Adience dataset. These transformations are applied every time a new batch is loaded, and the parameters of each one are randomly chosen from a defined range ([0.8, 1.2] for zooming, [0.5, 1.5] for brightness and [0, 90] degrees for rotation), providing a new set of transformed images for each batch. This technique reduces the overfitting risk and provides an important performance boost as we always work with different but similar images [4].

The epoch size is equal to the number of images in the training set. It could be a higher number as we are using data augmentation, but instead of increasing the epoch size, we rather run the training for more epochs. In this case, we set the maximum number of epochs to 100. However, we always save the best model, that is evaluated when the training finishes.

The models are mainly evaluated using the QWK metric defined in Eq. (1). Also, other evaluation metrics are used to ease the comparison with alternative works:

- Minimum Sensitivity (MS) [51] is the lowest percentage of samples correctly predicted to belong to a class with respect to the number of samples of that class.

$$\text{MS} = \min \left\{ S_q = \frac{O_{qq}}{O_{q\bullet}}, q = 1, \dots, Q \right\},$$

where  $O$  is the confusion matrix and  $Q$  is the number of classes.

- Mean Absolute Error (MAE) [51] is the average absolute deviation of the predicted category from the real one.

$$\text{MAE} = \frac{1}{N} \sum_{i,j=1}^Q |i - j| O_{ij},$$

where  $N$  is the number of samples,  $Q$  is the number of classes and  $O$  is the confusion matrix.

- Accuracy-based metrics. Correct Classification Rate (CCR) or standard accuracy is the most common metric for classification tasks and shows the percentage of correctly classified samples. We also include Top-2 CCR and Top-3 CCR [25], which are similar to CCR, but they take a prediction as correct when the real class is between the two or three classes, respectively, with the highest probability.

- 1-off accuracy [39], [40], [41] marks the prediction as correct when the correct class is at one category of distance (in the ordinal scale) from the predicted one.

QWK, MAE and 1-off accuracy are ordinal evaluation metrics, while MS, CCR, Top-2 CCR and Top-3 CCR do not take the order of categories into consideration.

In order to ease reproducibility, the source code is available in a public repository<sup>4</sup>.

### D. Factors

In our study, three different factors are considered:

- *Learning rate* (LR,  $\eta$ ). LR is one of the most critical hyper-parameters to tune for training deep neural networks. Optimal learning rate can vary depending on the dataset and the CNN architecture. Previous works have presented some techniques that adjust this parameter in order to achieve better performance [52], [53]. In this work, we consider three different values for the initial value of this parameter:  $10^{-2}$ ,  $10^{-3}$  and  $10^{-4}$ .
- *Batch size* (BS). Batch size is also an important parameter as it controls the number of weight updates that are made on every epoch. It can affect the training time and the model performance. In this paper, we try three different batch sizes for each dataset. For the DR dataset, we use 5, 10 and 15, while, for Adience, 64, 128 and 256 images are used. We took the batch sizes that were used in [22] and [25] as a reference, and we expand the range on both sides.
- *Link function* (LF). Different link functions are used for the CLM at the last layer output: **logit**, **probit** and complementary log-log (see Section III-A).

## V. RESULTS

In this Section, we present the results of the experiments. First, in Sections V-A, V-B, and V-C, we perform the study for adjusting the value of the different parameters. Then, Section V-D compares the results against the state of the art.

For each dataset, we show a table with the detailed results of the experiments performed for training the model with each combination of parameters. Every parameter combination was run five times. These tables show the mean value and the standard deviation (SD) of each metric across these five executions for the test set.

### A. Diabetic Retinopathy

Detailed test results for the DR dataset are presented in Table II. The best result for each metric is marked in bold and the second best is in italic font.

The best mean QWK value was obtained with the **clog-log** link function using a BS of 10 and a LR of  $10^{-3}$ . However, the best CCR value was obtained with a BS of 15, the **logit** link and a LR of  $10^{-4}$ . The optimal

<sup>4</sup><https://github.com/ayrna/deep-ordinal-clm>

TABLE II

DR RESULTS. BS STANDS FOR BATCH SIZE, LF FOR LINK FUNCTION AND LR FOR LEARNING RATE. MEAN AND STANDARD DEVIATION ARE REPRESENTED AS MEANS<sub>SD</sub>.

BS	LF	LR	QWK <sub>(SD)</sub>	MS <sub>(SD)</sub>	MAE <sub>(SD)</sub>	CCR <sub>(SD)</sub>	Top-2 <sub>(SD)</sub>	Top-3 <sub>(SD)</sub>	1-off <sub>(SD)</sub>
5	clog-log	10 <sup>-2</sup>	0.414 <sub>(0.057)</sub>	0.075 <sub>(0.042)</sub>	0.177 <sub>(0.023)</sub>	0.556 <sub>(0.057)</sub>	0.833 <sub>(0.042)</sub>	0.968 <sub>(0.011)</sub>	0.816 <sub>(0.021)</sub>
5	clog-log	10 <sup>-3</sup>	0.534 <sub>(0.027)</sub>	0.102 <sub>(0.011)</sub>	0.137 <sub>(0.006)</sub>	0.658 <sub>(0.015)</sub>	0.871 <sub>(0.011)</sub>	0.966 <sub>(0.003)</sub>	0.852 <sub>(0.002)</sub>
5	clog-log	10 <sup>-4</sup>	0.520 <sub>(0.006)</sub>	0.067 <sub>(0.008)</sub>	0.123 <sub>(0.003)</sub>	0.697 <sub>(0.006)</sub>	0.842 <sub>(0.008)</sub>	0.961 <sub>(0.003)</sub>	0.851 <sub>(0.002)</sub>
5	logit	10 <sup>-2</sup>	0.416 <sub>(0.041)</sub>	0.095 <sub>(0.029)</sub>	0.175 <sub>(0.021)</sub>	0.563 <sub>(0.054)</sub>	0.762 <sub>(0.040)</sub>	0.908 <sub>(0.026)</sub>	0.807 <sub>(0.029)</sub>
5	logit	10 <sup>-3</sup>	0.554 <sub>(0.013)</sub>	0.093 <sub>(0.009)</sub>	0.137 <sub>(0.003)</sub>	0.660 <sub>(0.008)</sub>	0.802 <sub>(0.005)</sub>	0.936 <sub>(0.004)</sub>	0.853 <sub>(0.005)</sub>
5	logit	10 <sup>-4</sup>	0.520 <sub>(0.003)</sub>	0.063 <sub>(0.004)</sub>	0.122 <sub>(0.002)</sub>	0.706 <sub>(0.005)</sub>	0.823 <sub>(0.004)</sub>	0.949 <sub>(0.003)</sub>	0.862 <sub>(0.003)</sub>
5	probit	10 <sup>-2</sup>	0.460 <sub>(0.048)</sub>	0.079 <sub>(0.046)</sub>	0.197 <sub>(0.064)</sub>	0.504 <sub>(0.167)</sub>	0.808 <sub>(0.034)</sub>	0.927 <sub>(0.073)</sub>	0.689 <sub>(0.240)</sub>
5	probit	10 <sup>-3</sup>	0.564 <sub>(0.018)</sub>	0.099 <sub>(0.013)</sub>	0.147 <sub>(0.018)</sub>	0.636 <sub>(0.045)</sub>	0.822 <sub>(0.040)</sub>	0.939 <sub>(0.020)</sub>	0.840 <sub>(0.015)</sub>
5	probit	10 <sup>-4</sup>	0.523 <sub>(0.005)</sub>	0.067 <sub>(0.012)</sub>	0.122 <sub>(0.002)</sub>	0.701 <sub>(0.006)</sub>	0.823 <sub>(0.002)</sub>	0.953 <sub>(0.002)</sub>	0.860 <sub>(0.003)</sub>
10	clog-log	10 <sup>-2</sup>	0.423 <sub>(0.239)</sub>	0.062 <sub>(0.051)</sub>	0.127 <sub>(0.017)</sub>	0.684 <sub>(0.046)</sub>	<b>0.894</b> <sub>(0.062)</sub>	<b>0.986</b> <sub>(0.012)</sub>	0.832 <sub>(0.020)</sub>
10	clog-log	10 <sup>-3</sup>	<b>0.582</b> <sub>(0.016)</sub>	0.102 <sub>(0.006)</sub>	0.128 <sub>(0.003)</sub>	0.680 <sub>(0.007)</sub>	<i>0.880</i> <sub>(0.004)</sub>	0.972 <sub>(0.003)</sub>	0.861 <sub>(0.004)</sub>
10	clog-log	10 <sup>-4</sup>	0.537 <sub>(0.010)</sub>	0.064 <sub>(0.004)</sub>	<i>0.116</i> <sub>(0.001)</sub>	0.717 <sub>(0.003)</sub>	0.837 <sub>(0.002)</sub>	0.971 <sub>(0.001)</sub>	0.860 <sub>(0.002)</sub>
10	logit	10 <sup>-2</sup>	0.531 <sub>(0.031)</sub>	0.107 <sub>(0.008)</sub>	0.151 <sub>(0.010)</sub>	0.623 <sub>(0.025)</sub>	0.802 <sub>(0.022)</sub>	0.934 <sub>(0.013)</sub>	0.838 <sub>(0.014)</sub>
10	logit	10 <sup>-3</sup>	0.579 <sub>(0.009)</sub>	0.096 <sub>(0.012)</sub>	0.127 <sub>(0.005)</sub>	0.686 <sub>(0.013)</sub>	0.817 <sub>(0.006)</sub>	0.954 <sub>(0.005)</sub>	0.861 <sub>(0.002)</sub>
10	logit	10 <sup>-4</sup>	0.539 <sub>(0.007)</sub>	0.074 <sub>(0.013)</sub>	0.126 <sub>(0.005)</sub>	0.707 <sub>(0.010)</sub>	0.823 <sub>(0.007)</sub>	0.957 <sub>(0.005)</sub>	0.858 <sub>(0.004)</sub>
10	probit	10 <sup>-2</sup>	0.508 <sub>(0.037)</sub>	0.088 <sub>(0.044)</sub>	0.145 <sub>(0.018)</sub>	0.639 <sub>(0.045)</sub>	0.835 <sub>(0.015)</sub>	0.960 <sub>(0.008)</sub>	0.829 <sub>(0.020)</sub>
10	probit	10 <sup>-3</sup>	0.558 <sub>(0.034)</sub>	<b>0.111</b> <sub>(0.005)</sub>	0.134 <sub>(0.003)</sub>	0.666 <sub>(0.008)</sub>	0.831 <sub>(0.007)</sub>	0.955 <sub>(0.001)</sub>	0.863 <sub>(0.003)</sub>
10	probit	10 <sup>-4</sup>	0.541 <sub>(0.010)</sub>	0.076 <sub>(0.006)</sub>	0.119 <sub>(0.002)</sub>	0.712 <sub>(0.005)</sub>	0.828 <sub>(0.003)</sub>	0.961 <sub>(0.002)</sub>	0.862 <sub>(0.001)</sub>
15	clog-log	10 <sup>-2</sup>	0.564 <sub>(0.016)</sub>	0.108 <sub>(0.014)</sub>	0.143 <sub>(0.006)</sub>	0.640 <sub>(0.015)</sub>	0.879 <sub>(0.011)</sub>	0.972 <sub>(0.005)</sub>	0.851 <sub>(0.006)</sub>
15	clog-log	10 <sup>-3</sup>	0.559 <sub>(0.026)</sub>	<i>0.111</i> <sub>(0.008)</sub>	0.127 <sub>(0.004)</sub>	0.682 <sub>(0.010)</sub>	0.871 <sub>(0.008)</sub>	<i>0.974</i> <sub>(0.002)</sub>	<b>0.868</b> <sub>(0.002)</sub>
15	clog-log	10 <sup>-4</sup>	0.538 <sub>(0.009)</sub>	0.054 <sub>(0.003)</sub>	<b>0.115</b> <sub>(0.002)</sub>	0.720 <sub>(0.006)</sub>	0.835 <sub>(0.007)</sub>	0.970 <sub>(0.003)</sub>	0.860 <sub>(0.006)</sub>
15	logit	10 <sup>-2</sup>	0.551 <sub>(0.020)</sub>	0.104 <sub>(0.008)</sub>	0.139 <sub>(0.011)</sub>	0.654 <sub>(0.027)</sub>	0.815 <sub>(0.017)</sub>	0.948 <sub>(0.016)</sub>	0.856 <sub>(0.015)</sub>
15	logit	10 <sup>-3</sup>	0.551 <sub>(0.010)</sub>	0.106 <sub>(0.016)</sub>	0.129 <sub>(0.008)</sub>	0.680 <sub>(0.019)</sub>	0.818 <sub>(0.008)</sub>	0.952 <sub>(0.007)</sub>	<i>0.866</i> <sub>(0.001)</sub>
15	logit	10 <sup>-4</sup>	0.543 <sub>(0.008)</sub>	0.056 <sub>(0.003)</sub>	0.121 <sub>(0.004)</sub>	<b>0.723</b> <sub>(0.004)</sub>	0.833 <sub>(0.004)</sub>	0.964 <sub>(0.003)</sub>	0.862 <sub>(0.004)</sub>
15	probit	10 <sup>-2</sup>	0.534 <sub>(0.032)</sub>	0.104 <sub>(0.013)</sub>	0.148 <sub>(0.015)</sub>	0.631 <sub>(0.038)</sub>	0.845 <sub>(0.030)</sub>	0.964 <sub>(0.010)</sub>	0.852 <sub>(0.010)</sub>
15	probit	10 <sup>-3</sup>	<i>0.580</i> <sub>(0.021)</sub>	0.104 <sub>(0.016)</sub>	0.129 <sub>(0.008)</sub>	0.680 <sub>(0.018)</sub>	0.832 <sub>(0.010)</sub>	0.959 <sub>(0.007)</sub>	0.866 <sub>(0.003)</sub>
15	probit	10 <sup>-4</sup>	0.533 <sub>(0.004)</sub>	0.065 <sub>(0.005)</sub>	0.117 <sub>(0.002)</sub>	<i>0.721</i> <sub>(0.004)</sub>	0.832 <sub>(0.002)</sub>	0.964 <sub>(0.001)</sub>	0.863 <sub>(0.001)</sub>

configuration depends on the metric we are analysing. In this case, as we are working with an ordinal problem, the most reliable metric is the QWK. However, the rest of the metrics are also included to allow further comparisons with future works.

### B. Adience

Test results for the experiments made with the Adience dataset are shown in Table III. The best result for each metric is marked in bold and the second best is in italic font.

The best mean QWK value was obtained with the **logit** link function using a BS of 64 and a LR of 10<sup>-4</sup>. Also, this configuration obtained the best score for Top-2, Top-3 and 1-off accuracy, and the second best for MS, MAE and CCR. In this case, this configuration can be selected as the optimal for this problem.

### C. Statistical analysis

In this subsection, a statistical analysis will be performed in order to obtain conclusions from the results. The significance and relative importance of the parameters concerning the results obtained, as well as the most suitable values, were obtained using an ANalysis Of the Variance (ANOVA).

The ANOVA test [23] is one of the most widely used statistical techniques. ANOVA is essentially a method of analysing the variance to which a response is subject into

its various components, corresponding to the sources of variation which can be identified. ANOVA, in this case, examines the effects of three quantitative variables (termed factors) on one quantitative response. Considered factors are the LF, the LR for the Adam optimization algorithm, and the BS. We assume that five executions are enough to do the statistical tests because of the computational time limitations.

The ANOVA test results show that there are significant differences in average QWK depending on the LF and also depending on the LR for  $\alpha = 0.05$  (p-value = 0.000). Moreover, an interaction between the LF and the LR can be recognised (p-value = 0.001).

Given that there exist significant differences between the means, we analyse now these differences. A post-hoc multiple comparison test has been performed on the mean QWK obtained. An HSD Tukey's test [24] has been selected under the null hypothesis that the variance of the error of the dependent variable is the same between the groups. The results of this test over the test set are shown in Table IV. They show that the best LF is the **clog-log** but the **probit** link performance is close to it. Also, the best value for the LR parameter is 10<sup>-3</sup>. The BS is not relevant for this dataset with the values considered.

The results of the ANOVA III test for the Adience dataset, first, demonstrate that there exist significant differences in average QWK concerning the three factors (p-value = 0.000). Secondly, we found interactions between all the pairs of factors and between all the three



TABLE III  
ADIENCE TEST RESULTS. BS STANDS FOR BATCH SIZE, LF FOR LINK FUNCTION AND LR FOR LEARNING RATE. MEAN AND STANDARD DEVIATION MEAN<sub>SD</sub>.

BS	LF	LR	QWK <sub>(SD)</sub>	MS <sub>(SD)</sub>	MAE <sub>(SD)</sub>	CCR <sub>(SD)</sub>	Top-2 <sub>(SD)</sub>	Top-3 <sub>(SD)</sub>	1-off <sub>(SD)</sub>
64	clog-log	10 <sup>-2</sup>	0.808 <sub>(0.025)</sub>	0.086 <sub>(0.041)</sub>	0.147 <sub>(0.008)</sub>	0.415 <sub>(0.031)</sub>	0.677 <sub>(0.024)</sub>	0.798 <sub>(0.036)</sub>	0.804 <sub>(0.015)</sub>
64	clog-log	10 <sup>-3</sup>	0.873 <sub>(0.006)</sub>	0.144 <sub>(0.057)</sub>	<b>0.124<sub>(0.003)</sub></b>	<b>0.519<sub>(0.014)</sub></b>	<i>0.764<sub>(0.010)</sub></i>	0.861 <sub>(0.019)</sub>	0.886 <sub>(0.006)</sub>
64	clog-log	10 <sup>-4</sup>	0.799 <sub>(0.010)</sub>	0.000 <sub>(0.000)</sub>	0.174 <sub>(0.001)</sub>	0.324 <sub>(0.015)</sub>	0.616 <sub>(0.020)</sub>	0.795 <sub>(0.012)</sub>	0.771 <sub>(0.014)</sub>
64	logit	10 <sup>-2</sup>	0.778 <sub>(0.019)</sub>	0.074 <sub>(0.041)</sub>	0.159 <sub>(0.006)</sub>	0.366 <sub>(0.025)</sub>	0.636 <sub>(0.015)</sub>	0.785 <sub>(0.010)</sub>	0.775 <sub>(0.015)</sub>
64	logit	10 <sup>-3</sup>	<b>0.881<sub>(0.005)</sub></b>	<i>0.178<sub>(0.023)</sub></i>	<i>0.126<sub>(0.001)</sub></i>	<i>0.518<sub>(0.008)</sub></i>	<b>0.765<sub>(0.015)</sub></b>	<b>0.902<sub>(0.005)</sub></b>	<b>0.894<sub>(0.005)</sub></b>
64	logit	10 <sup>-4</sup>	0.784 <sub>(0.011)</sub>	0.000 <sub>(0.000)</sub>	0.180 <sub>(0.001)</sub>	0.318 <sub>(0.026)</sub>	0.621 <sub>(0.034)</sub>	0.772 <sub>(0.024)</sub>	0.731 <sub>(0.030)</sub>
64	probit	10 <sup>-2</sup>	0.836 <sub>(0.005)</sub>	0.135 <sub>(0.021)</sub>	0.134 <sub>(0.002)</sub>	0.468 <sub>(0.011)</sub>	0.720 <sub>(0.009)</sub>	0.861 <sub>(0.009)</sub>	0.829 <sub>(0.005)</sub>
64	probit	10 <sup>-3</sup>	<i>0.874<sub>(0.004)</sub></i>	0.134 <sub>(0.012)</sub>	0.126 <sub>(0.003)</sub>	0.511 <sub>(0.014)</sub>	0.756 <sub>(0.009)</sub>	<i>0.895<sub>(0.003)</sub></i>	<i>0.889<sub>(0.003)</sub></i>
64	probit	10 <sup>-4</sup>	0.805 <sub>(0.004)</sub>	0.000 <sub>(0.000)</sub>	0.170 <sub>(0.001)</sub>	0.360 <sub>(0.011)</sub>	0.653 <sub>(0.011)</sub>	0.809 <sub>(0.009)</sub>	0.790 <sub>(0.009)</sub>
128	clog-log	10 <sup>-2</sup>	0.832 <sub>(0.013)</sub>	0.123 <sub>(0.031)</sub>	0.135 <sub>(0.004)</sub>	0.463 <sub>(0.013)</sub>	0.705 <sub>(0.019)</sub>	0.813 <sub>(0.025)</sub>	0.832 <sub>(0.006)</sub>
128	clog-log	10 <sup>-3</sup>	0.873 <sub>(0.006)</sub>	<b>0.185<sub>(0.029)</sub></b>	0.128 <sub>(0.002)</sub>	0.513 <sub>(0.007)</sub>	0.758 <sub>(0.008)</sub>	0.870 <sub>(0.011)</sub>	0.880 <sub>(0.009)</sub>
128	clog-log	10 <sup>-4</sup>	0.659 <sub>(0.025)</sub>	0.000 <sub>(0.000)</sub>	0.190 <sub>(0.002)</sub>	0.235 <sub>(0.026)</sub>	0.466 <sub>(0.031)</sub>	0.640 <sub>(0.030)</sub>	0.536 <sub>(0.041)</sub>
128	logit	10 <sup>-2</sup>	0.781 <sub>(0.041)</sub>	0.096 <sub>(0.059)</sub>	0.153 <sub>(0.007)</sub>	0.398 <sub>(0.031)</sub>	0.638 <sub>(0.033)</sub>	0.790 <sub>(0.025)</sub>	0.779 <sub>(0.020)</sub>
128	logit	10 <sup>-3</sup>	0.865 <sub>(0.005)</sub>	0.127 <sub>(0.026)</sub>	0.134 <sub>(0.001)</sub>	0.497 <sub>(0.009)</sub>	0.754 <sub>(0.008)</sub>	0.882 <sub>(0.009)</sub>	0.874 <sub>(0.008)</sub>
128	logit	10 <sup>-4</sup>	0.586 <sub>(0.008)</sub>	0.000 <sub>(0.000)</sub>	0.196 <sub>(0.001)</sub>	0.192 <sub>(0.001)</sub>	0.364 <sub>(0.060)</sub>	0.581 <sub>(0.034)</sub>	0.396 <sub>(0.002)</sub>
128	probit	10 <sup>-2</sup>	0.849 <sub>(0.005)</sub>	0.132 <sub>(0.010)</sub>	0.131 <sub>(0.001)</sub>	0.479 <sub>(0.004)</sub>	0.728 <sub>(0.007)</sub>	0.854 <sub>(0.009)</sub>	0.847 <sub>(0.007)</sub>
128	probit	10 <sup>-3</sup>	0.866 <sub>(0.002)</sub>	0.124 <sub>(0.043)</sub>	0.130 <sub>(0.002)</sub>	0.505 <sub>(0.006)</sub>	0.750 <sub>(0.010)</sub>	0.882 <sub>(0.004)</sub>	0.873 <sub>(0.006)</sub>
128	probit	10 <sup>-4</sup>	0.718 <sub>(0.015)</sub>	0.000 <sub>(0.000)</sub>	0.185 <sub>(0.001)</sub>	0.300 <sub>(0.031)</sub>	0.575 <sub>(0.015)</sub>	0.733 <sub>(0.010)</sub>	0.640 <sub>(0.033)</sub>
256	clog-log	10 <sup>-2</sup>	0.853 <sub>(0.004)</sub>	0.157 <sub>(0.024)</sub>	0.130 <sub>(0.002)</sub>	0.485 <sub>(0.009)</sub>	0.744 <sub>(0.006)</sub>	0.842 <sub>(0.016)</sub>	0.858 <sub>(0.004)</sub>
256	clog-log	10 <sup>-3</sup>	0.840 <sub>(0.017)</sub>	0.095 <sub>(0.017)</sub>	0.144 <sub>(0.005)</sub>	0.456 <sub>(0.021)</sub>	0.720 <sub>(0.022)</sub>	0.840 <sub>(0.018)</sub>	0.842 <sub>(0.018)</sub>
256	clog-log	10 <sup>-4</sup>	0.552 <sub>(0.010)</sub>	0.000 <sub>(0.000)</sub>	0.199 <sub>(0.001)</sub>	0.187 <sub>(0.001)</sub>	0.368 <sub>(0.022)</sub>	0.475 <sub>(0.025)</sub>	0.387 <sub>(0.001)</sub>
256	logit	10 <sup>-2</sup>	0.764 <sub>(0.102)</sub>	0.077 <sub>(0.067)</sub>	0.155 <sub>(0.020)</sub>	0.387 <sub>(0.083)</sub>	0.632 <sub>(0.103)</sub>	0.790 <sub>(0.077)</sub>	0.783 <sub>(0.065)</sub>
256	logit	10 <sup>-3</sup>	0.851 <sub>(0.008)</sub>	0.100 <sub>(0.030)</sub>	0.147 <sub>(0.003)</sub>	0.449 <sub>(0.015)</sub>	0.726 <sub>(0.015)</sub>	0.861 <sub>(0.006)</sub>	0.850 <sub>(0.008)</sub>
256	logit	10 <sup>-4</sup>	0.558 <sub>(0.008)</sub>	0.000 <sub>(0.000)</sub>	0.202 <sub>(0.001)</sub>	0.187 <sub>(0.002)</sub>	0.206 <sub>(0.007)</sub>	0.395 <sub>(0.046)</sub>	0.389 <sub>(0.003)</sub>
256	probit	10 <sup>-2</sup>	0.858 <sub>(0.005)</sub>	0.164 <sub>(0.033)</sub>	0.130 <sub>(0.002)</sub>	0.486 <sub>(0.007)</sub>	0.741 <sub>(0.008)</sub>	0.867 <sub>(0.008)</sub>	0.862 <sub>(0.005)</sub>
256	probit	10 <sup>-3</sup>	0.850 <sub>(0.008)</sub>	0.111 <sub>(0.040)</sub>	0.144 <sub>(0.002)</sub>	0.460 <sub>(0.011)</sub>	0.732 <sub>(0.006)</sub>	0.865 <sub>(0.006)</sub>	0.853 <sub>(0.007)</sub>
256	probit	10 <sup>-4</sup>	0.565 <sub>(0.010)</sub>	0.000 <sub>(0.000)</sub>	0.196 <sub>(0.001)</sub>	0.189 <sub>(0.001)</sub>	0.409 <sub>(0.014)</sub>	0.602 <sub>(0.022)</sub>	0.392 <sub>(0.002)</sub>

factors together (p-values 0.000, 0.000, 0.000 and 0.001, respectively).

As we did for the DR dataset, a post-hoc multiple comparison test has been performed on the average QWK obtained for Adience. Under the null hypothesis that the variance of the error of the dependent variable is the same between the groups, the HSD Tukey's test has been applied. The results of this test over the test set are shown in Table IV.

The results over the test set show that the best LF is the **logit**, the best LR is 10<sup>-3</sup> and the best BS is 128. However, the interactions between these factors made the configuration that uses a **logit** link,  $\eta = 10^{-3}$  and BS of 64, the best configuration. It obtained a mean QWK value of 0.940 for validation and 0.881 for test. The same parameters, but using the **probit** link, achieves the second best result (0.874). The standard deviation is very low for both cases.

To sum up, the results showed that the best parameter configuration depends on the problem that is being solved. The **clog-log** function offers the best results in DR dataset while the **logit** link is the best option for the Adience dataset. However, the best LR for both datasets were 10<sup>-3</sup>. It is recommended to use this value for future datasets. The best BS for DR was 10, while the best value for Adience was 128 (intermediate values considered). Finally, there are more interactions between the three factors for the Adience dataset than for DR. These results highlight the importance of adjusting the

TABLE IV  
TUKEY'S TEST RESULTS FOR BOTH DATASETS.

		DR		Adience	
LF	LF	Mean diff.	Sig.	Mean diff.	Sig.
logit	probit	0.046	0.000	-0.002	0.011
	clog-log	0.084	0.000	0.012	0.000
	logit	-0.046	0.000	0.002	0.011
probit	clog-log	0.038	0.000	0.014	0.248
	logit	-0.084	0.000	-0.012	0.000
clog-log	probit	-0.038	0.000	-0.014	0.248
LR	LR	Mean diff.	Sig.	Mean diff.	Sig.
10 <sup>-2</sup>	10 <sup>-3</sup>	-0.046	0.000	-0.073	0.000
	10 <sup>-4</sup>	0.148	0.000	-0.044	0.000
10 <sup>-3</sup>	10 <sup>-2</sup>	0.046	0.000	0.073	0.000
	10 <sup>-4</sup>	0.194	0.000	0.029	0.023
10 <sup>-4</sup>	10 <sup>-2</sup>	-0.148	0.000	0.044	0.000
	10 <sup>-3</sup>	-0.194	0.000	-0.029	0.023
BS	BS	Mean diff.	Sig.	Mean diff.	Sig.
64	128	-0.041	0.026	-	-
	256	-0.027	0.000	-	-
128	64	0.041	0.026	-	-
	256	0.014	0.000	-	-
256	64	0.027	0.000	-	-
	128	-0.014	0.000	-	-

hyper-parameters for each problem instead of trying to find an optimal configuration for all the datasets.

#### D. Comparison with nominal method and previous works

Once the factor parameters have been studied and selected, experiments are run with the standard cross-entropy loss and the softmax function too in order to prove



TABLE V  
COMPARISON BETWEEN THE BEST RESULTS OF THE CLM NETWORK,  
A NOMINAL NETWORK (USING SOFTMAX AND CROSS-ENTROPY) AND  
PREVIOUS WORKS FOR BOTH DATASETS.

Diabetic Retinopathy			
Method	$\overline{QWK}_{(SD)}$	$\overline{CCR}_{(SD)}$	$\overline{1-off}_{(SD)}$
CLM network (proposal)	0.582 <sub>(0.016)</sub>	0.723 <sub>(0.004)</sub>	0.868 <sub>(0.002)</sub>
Nominal network	0.498 <sub>(0.013)</sub>	0.692 <sub>(0.014)</sub>	0.854 <sub>(0.006)</sub>
de la Torre et al. [22]	0.537 <sub>(-)</sub>	-	-
Nebot et al. [38]	0.555 <sub>(-)</sub>	-	-
Adience			
Method	$\overline{QWK}_{(SD)}$	$\overline{CCR}_{(SD)}$	$\overline{1-off}_{(SD)}$
CLM network (proposal)	0.881 <sub>(0.005)</sub>	0.519 <sub>(0.013)</sub>	0.894 <sub>(0.005)</sub>
Nominal network	0.787 <sub>(0.004)</sub>	0.458 <sub>(0.008)</sub>	0.800 <sub>(0.007)</sub>
Beckham and Pal [25]	0.855 <sub>(0.012)</sub>	0.467 <sub>(0.019)</sub>	0.867 <sub>(0.011)</sub>
Eidinger et al. [39]	-	0.451 <sub>(0.026)</sub>	0.807 <sub>(0.011)</sub>
Chen et al. [40]	-	0.529 <sub>(0.060)</sub>	0.885 <sub>(0.022)</sub>
Levi and Hassner [41]	-	0.507 <sub>(0.051)</sub>	0.847 <sub>(0.022)</sub>

the performance improvement of considering the ordinality of the problem (QWK loss and the CLM). The evaluation metrics remains the same in order to be able to compare. When considering the nominal version, for the DR dataset, the best mean value of QWK was 0.497 and was obtained when using a BS of 10 and a LR of  $10^{-4}$ . In the case of Adience dataset, the highest QWK was 0.787 and was achieved with a BS of 64 and a LR of  $10^{-3}$ . There are some parameter configurations where the training process gets stagnated and a very low QWK is obtained. As we saw in Sections V-A and V-B, this problem is not found when using the ordinal method.

These results are included in Table V, together with the comparison against previous works of the state-of-the-art. All the results are given for the test set, except those from [22] (DR dataset), because the authors only provided validation results for  $128 \times 128$  images (however, validation results are usually better than test results). The results for [25] were obtained by reproducing their experiments because they did not provide the results for test set. As can be checked, the proposed ordinal model outperforms all the other alternatives in terms of QWK.

The performance gain of CLM over the nominal version reaches 16.8% for DR and 11.9% for Adience dataset. The improvement of the ordinal method for DR is higher than that for Adience. It seems that the method proposed in this work offers a more significant improvement as the given problem complexity increases. Moreover, when compared against alternative ordinal methods (many of them with a deep structure), CLMs are very competitive, possibly because of the flexibility provided by the threshold model structure (where the threshold of each class is independently adjusted).

## VI. CONCLUSIONS

This paper introduces a new deep ordinal network based on combining CLM models with a continuous QWK loss function. The proposed model is able to improve the performance of the deep network compared to the equivalent

nominal version and other models proposed in previous works. Also, it is able to reduce the chance that the model gets stuck when training with some parameter configurations. We conclude that the optimal values for the different parameters considered are problem-dependant. The results highlight the importance of making an experimental design where all of these parameters are adjusted for each problem. In summary, the most significant contributions of the model proposal are the performance increase, the reduction of the number of parameters configurations that should be tried to find the best one and the prevention from over-fitting and stagnation.

As future research, it seems that the design of new generalised link functions could be promising, which could be dynamically adapted to any problem based on a learnable parameter.

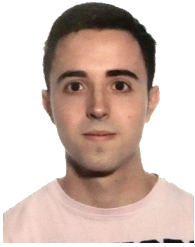
## ACKNOWLEDGMENT

This work has been partially subsidised by the TIN2017-85887-C2-1-P and TIN2017-90567-REDT projects of the Spanish Ministry of Economy and Competitiveness (MINECO), and FEDER funds of the European Union.

## REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] J. Schmidhuber, U. Meier, and D. Ciresan, "Multi-column deep neural networks for image classification," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 3642–3649.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. of the IEEE conf. computer vision and pat. recog.*, 2016, pp. 770–778.
- [4] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Adv. in neural inf. proc. sys.*, 2012, pp. 1097–1105.
- [5] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath *et al.*, "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal processing magazine*, vol. 29, no. 6, pp. 82–97, 2012.
- [6] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, p. 529, 2015.
- [7] X. Jiang, Y. Pang, X. Li, and J. Pan, "Speed up deep neural network based pedestrian detection by sharing features across multi-scale models," *Neurocomputing*, vol. 185, pp. 163–170, 2016.
- [8] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587.
- [9] J. Yu, B. Zhang, Z. Kuang, D. Lin, and J. Fan, "iprivity: image privacy protection by identifying sensitive objects via deep multi-task learning," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 5, pp. 1005–1016, 2017.
- [10] X. Yuan, P. He, Q. Zhu, and X. Li, "Adversarial examples: Attacks and defenses for deep learning," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–20, 2019.
- [11] F. Tan, X. Hou, J. Zhang, Z. Wei, and Z. Yan, "A deep learning approach to competing risks representation in peer-to-peer lending," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 5, pp. 1565–1574, May 2019.
- [12] C. Hong, J. Yu, J. Wan, D. Tao, and M. Wang, "Multimodal deep autoencoder for human pose recovery," *IEEE Transactions on Image Processing*, vol. 24, no. 12, pp. 5659–5670, 2015.

- [13] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
- [14] Z. Li and J. Tang, "Weakly supervised deep metric learning for community-contributed image retrieval," *IEEE Transactions on Multimedia*, vol. 17, no. 11, pp. 1989–1999, 2015.
- [15] L. Jin, K. Li, Z. Li, F. Xiao, G. Qi, and J. Tang, "Deep semantic-preserving ordinal hashing for cross-modal similarity search," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 5, pp. 1429–1440, May 2019.
- [16] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *European conference on computer vision*. Springer, 2014, pp. 184–199.
- [17] Y. Sun, X. Wang, and X. Tang, "Deep convolutional network cascade for facial point detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 3476–3483.
- [18] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [19] P. McCullagh, "Regression models for ordinal data," *Journal of the royal statistical society. Series B (Methodological)*, pp. 109–142, 1980.
- [20] A. Agresti, *Analysis of ordinal categorical data*. J. Wiley & Sons, 2010, vol. 656.
- [21] P. A. Gutierrez, M. Perez-Ortiz, J. Sanchez-Monedero, F. Fernandez-Navarro, and C. Hervás-Martínez, "Ordinal regression methods: survey and experimental study," *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 1, pp. 127–146, 2016.
- [22] J. de la Torre, D. Puig, and A. Valls, "Weighted kappa loss function for multi-class classification of ordinal data in deep learning," *Pattern Recognition Letters*, vol. 105, pp. 144–154, 2018.
- [23] R. G. Miller Jr, *Beyond ANOVA: basics of applied statistics*. Chapman and Hall/CRC, 1997.
- [24] J. W. Tukey, "Comparing individual means in the analysis of variance," *Biometrics*, pp. 99–114, 1949.
- [25] C. Beckham and C. Pal, "Unimodal probability distributions for deep ordinal classification," in *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ser. ICML'17. JMLR.org, 2017, pp. 411–419. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3305381.3305424>
- [26] Y. Liu, A. W.-K. Kong, and C. K. Goh, "Deep ordinal regression based on data relationship for small datasets," in *IJCAI*, 2017, pp. 2372–2378.
- [27] M. ALALI, N. M. Sharef, H. Hamdan, M. A. A. Murad, and N. A. Husin, "Multi-layers convolutional neural network for twitter sentiment ordinal scale classification," in *International Conference on Soft Computing and Data Mining*. Springer, 2018, pp. 446–454.
- [28] Z. Niu, M. Zhou, L. Wang, X. Gao, and G. Hua, "Ordinal regression with multiple output cnn for age estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4920–4928.
- [29] H. Li, M. Habes, and Y. Fan, "Deep ordinal ranking for multi-category diagnosis of alzheimer's disease using hippocampal mri data," *arXiv preprint arXiv:1709.01599*, 2017.
- [30] S. Chen, C. Zhang, M. Dong, J. Le, and M. Rao, "Using ranking-cnn for age estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5183–5192.
- [31] A. Rios and R. Kavuluru, "Ordinal convolutional neural networks for predicting rdcc positive valence psychiatric symptom severity scores," *Journal of biomedical informatics*, vol. 75, pp. S85–S93, 2017.
- [32] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao, "Deep ordinal regression network for monocular depth estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2002–2011.
- [33] A. Pal, A. Chaturvedi, U. Garain, A. Chandra, R. Chatterjee, and S. Senapati, "Severity assessment of psoriatic plaques using deep cnn based ordinal classification," in *OR 2.0 Context-Aware Operating Theaters, Computer Assisted Robotic Endoscopy, Clinical Image-Based Procedures, and Skin Image Analysis*. Springer, 2018, pp. 252–259.
- [34] Y. Liu, A. Wai Kin Kong, and C. Keong Goh, "A constrained deep neural network for ordinal regression," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 831–839.
- [35] C. Beckham and C. Pal, "A simple squared-error reformulation for ordinal classification," *arXiv preprint arXiv:1612.00775*, 2016.
- [36] R. Herbrich, "Large margin rank boundaries for ordinal regression," *Advances in large margin classifiers*, pp. 115–132, 2000.
- [37] A. Ben-David, "Comparison of classification accuracy using cohen's weighted kappa," *Expert Systems with Applications*, vol. 34, no. 2, pp. 825–832, 2008.
- [38] À. Nebot et al., "Diabetic retinopathy detection through image analysis using deep convolutional neural networks," in *A.I. Research and Development: Proc. of the 19th Int. Conf. of the Catalan Association for A.I.* IOS press, 2016, p. 58.
- [39] E. Eiding, R. Enbar, and T. Hassner, "Age and gender estimation of unfiltered faces," *IEEE Transactions on Information Forensics and Security*, vol. 9, no. 12, pp. 2170–2179, 2014.
- [40] J.-C. Chen, A. Kumar, R. Ranjan, V. M. Patel, A. Alavi, and R. Chellappa, "A cascaded convolutional neural network for age estimation of unconstrained faces," in *Proc. of 8th IEEE Conference on Biometrics Theory, Applications and Systems (BTAS)*. IEEE, 2016, pp. 1–8.
- [41] G. Levi and T. Hassner, "Age and gender classification using convolutional neural networks," in *Proc. of the IEEE conf. comp. vision and pat. rec.*, 2015, pp. 34–42.
- [42] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," 11 2015.
- [43] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proceedings of the 27th international conference on machine learning (ICML-10)*, 2010, pp. 807–814.
- [44] Y. Bengio, P. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE Transactions on Neural Networks*, vol. 5, no. 2, pp. 157–166, 1994.
- [45] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in *International Conference on Machine Learning*, 2013, pp. 1310–1318.
- [46] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in *Proc. icml*, vol. 30, no. 1, 2013, p. 3.
- [47] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proceedings of the 32nd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, F. Bach and D. Blei, Eds., vol. 37. Lille, France: PMLR, 07–09 Jul 2015, pp. 448–456. [Online]. Available: <http://proceedings.mlr.press/v37/loff15.html>
- [48] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proceedings of ICLR2015*, ser. ICLR2015, 2014.
- [49] W.-S. Chin, Y. Zhuang, Y.-C. Juan, and C.-J. Lin, "A learning-rate schedule for stochastic gradient methods to matrix factorization," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 2015, pp. 442–455.
- [50] D. A. Van Dyk and X.-L. Meng, "The art of data augmentation," *Journal of Computational and Graphical Statistics*, vol. 10, no. 1, pp. 1–50, 2001.
- [51] M. Cruz-Ramírez, C. Hervás-Martínez, J. Sánchez-Monedero, and P. A. Gutiérrez, "Metrics to guide a multi-objective evolutionary algorithm for ordinal classification," *Neurocomputing*, vol. 135, pp. 21–31, 2014.
- [52] L. N. Smith, "Cyclical learning rates for training neural networks," in *IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2017, pp. 464–472.
- [53] A. Senior, G. Heigold, K. Yang et al., "An empirical study of learning rates in deep neural networks for speech recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2013, pp. 6724–6728.



**Víctor Manuel Vargas** was born in Córdoba, Spain. He received the degree of Computer Engineering from the University of Córdoba, Spain, in 2018. He is currently studying a master's degree at the International University of Menéndez Pelayo, Spain. He is a member of the AYRNA group of the Department of Computer Science and Numerical Analysis of the University of Córdoba. His current research interests include deep neural networks and their

applications to image classification on ordinal problems.



**Pedro Antonio Gutiérrez** (M'07-SM'15) received the B.S. degree in computer science from the University of Sevilla (Spain) in 2006, and the Ph.D. degree in computer science and artificial intelligence from the University of Granada (Spain) in 2009. He is currently an Assistant Professor with the Department of Computer Science and Numerical Analysis, University of Córdoba (Spain). His research interests are in the areas of supervised learning, evolutionary artificial neural networks, ordinal

classification and the application of these techniques to different real world problems, including precision agriculture, renewable energy, climatology and biomedicine, among others. He serves on the Editorial board for the journals *IEEE Transaction on Neural Networks and Learning Systems* and *Progress in Artificial Intelligence* and on the organization/program committees of several computational intelligence conferences.



**César Hervás-Martínez** (M'07-SM'14) was born in Cuenca, Spain. He received the B.S. degree in statistics and operations research from the Universidad Complutense de Madrid, Madrid, Spain, in 1978, and the Ph.D. degree in mathematics from the University of Seville, Seville, Spain, in 1986. He is currently a Professor of Computer Science and Artificial Intelligence with the Department of Computer Science and Numerical Analysis, University of Córdoba, Córdoba, Spain, and an Associate

Professor with the Department of Mathematics and Engineering, Loyola University Andalucía, Spain. His current research interests include neural networks, evolutionary computation, and the modelling of natural systems.