

Classification of ordinal data in deep learning: an experimental study

Víctor-Manuel Vargas, Pedro-Antonio Gutiérrez and César Hervás

Abstract—The abstract goes here.

Index Terms—IEEE, IEEEtran, journal, LATEX, paper, template.

I. INTRODUCTION

Habría que reestructurar la introducción: 1) deep learning y CNN, pero corto 2) clasificación ordinal, modelos de umbral, POM y funciones de enlace (más bien breve) 3) clasificación ordinal profunda 4) objetivo del paper (un estudio experimental sobre funciones de enlace y otros parámetros en deep learning ordinal). learning rate y batch normalization, si se mete, debería ser más breve

1) **D**EEP LEARNING, introduced by Yann Lecun [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 protection [9], recovery of human pose [10], semantic segmentation [11] and image retrieval [12]. ~~Deep learning discovers complex structures in large datasets by using the backpropagation algorithm to change the internal parameters of the model and compute better representations in each layer from the features in the previous layer.~~ Convolutional Neural Networks (CNN) are one of the types of deep networks that are designed to process data that comes 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 [13], [14], [15]. On image classification tasks, each colour channel is represented by a 2D array. In this case, convolutional layers extract the main features from the pixels of the images 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. These outputs are compared against the correct values.

2) Ordinal classification problems are those classification tasks where labels are ordered, and there are different inter-classes importances for each pair of classes. This kind of problem can be treated as a nominal classification problem, but this discards ordinal information. A better approach is to use specific methods that take

into account this kind of information to improve the performance of the classification model. ~~One way to use the ordinal information is to evaluate the model using an ordinal metric. Multiple metrics exist in the literature of machine learning and statistics [16], [17]. Kappa index is a well-known statistic coefficient defined by Cohen [18] to measure inter-rater agreement on classifying elements into a set of categories. Later, Weighted Kappa (WK) is a modified version of the Kappa statistic calculated to allow assigning different weights to different levels of aggregation between two variables. An ordinal weighted kappa loss was described in a previous work [19], which is a differentiable cost function based on the WK metric. This previous work proved that using these functions improves the model performance and reduce the overfitting risk. Different weights can be assigned in an ordinal classification problem. In a linear penalization, the weight is proportional to the distance between the predicted and the real class (in number of categories). In the Quadratic Weighted Kappa (QWK) the penalization is proportional to the square of the distance. In this work, we will use the QWK metric and the QWK loss.~~ The Proportional Odds Model (POM) [20] is an ordinal alternative to the softmax function that was designed for ordinal regression. It belongs to a wider family of models called Cumulative Link Models (CLM). It is inspired in the concept of a latent variable that is projected in an n-dimensional space and a set of thresholds that divides this space into the different ordinal levels. This kind of models uses a link function which can be of different types. The most common link function is the *logit* one, that is used in POM. However, there are other functions that can be explored. This model will be described in depth in Section III.

3) Ordinal models can also be applied to deep learning models. In the case of convolutional networks, the model projection used by the threshold model will be obtained from the last layer of the network. When working with a 1-dimensional space, the last layer will have only one neuron, and its value will be used to classify every sample in the corresponding class according to the thresholds. Some previous works have used the POM in traditional neural networks [21], but it has not been applied to convolutional networks yet. Also, there are some link functions that has not been explored.

4) In this paper, an experimental study regarding link functions for Cumulative Link Models will be made. Also, other parameters that can affect the training process and the model performance, like the learning rate of the opti-

mization algorithm and the batch size, and its interaction will be studied. The nominal version of this model will be used as a baseline for comparison. We will contrast the results obtained with statistical analysis to provide more robust conclusions. An approximated ANOVA III [22] test followed by a posthoc Tukey's test [23] will be performed because of the limitations of the computational time required to run a higher number of executions.

The experiments will be run using two different ordinal datasets: Diabetic Retinopathy [19], which contains high-resolution fundus images related with diabetes disease, and Adience [24], which includes human faces images associated with an age range.

This paper is organized as follows: in Section II, we take a look at previous works related to this paper. Section III present a formal description of an ordinal problem and the Cumulative Link Models. In Section IV, we describe the model, 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. However, few works are focused on ordinal classification problems.

J. de la Torre et al. [19] proposed the use of a continuous version of the QWK loss function for the optimization algorithm. They compared this cost function against the traditional log-loss function across three different databases, 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 hyperparameter tuning. First, they defined QWK metric 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 ω is the penalization matrix (in this case, quadratic weights are considered), O is the confusion matrix and E is the normalized outer product between the prediction and the true vector.

Then, they provided the QWK loss (QWK_l) definition below:

$$\text{QWK}_l = \frac{\sum_{k=1}^N \sum_{q=1}^Q \omega_{t_k,q} P(y = C_q | \mathbf{x}_k)}{\sum_{i=1}^Q \frac{N_i}{N} \sum_{j=1}^Q (\omega_{i,j} \sum_{k=1}^N P(y = C_j | \mathbf{x}_k))}, \quad (2)$$

where $\text{QWK}_l \in [0, 2]$, \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 and $\omega_{i,j}$ are the elements of the penalization matrix. Generally, $\omega_{i,j} = \frac{|i-j|^n}{(C-1)^n}$, where

$\omega_{i,j} \in [0, 1]$. In this case, the QWK_l is a function to be minimized while the QWK metric must be maximized.

Z. Niu et al. [25] proposed a learning approach to address ordinal regression problems using convolutional neural networks. They divided the problem into a series of binary classification sub-problems and proposed a multiple output CNN optimization algorithm to collectively resolve these classification sub-problems, taking into account the correlation between them.

Christopher Beckham and Christopher Pal [24] proposed a straightforward technique to constrain discrete ordinal probability distributions to be unimodal, via the use of the Poisson and binomial probability distributions. They evaluated this approach in the context of deep learning on two large ordinal image datasets, including the Adience dataset used in this paper, obtaining promising results. Also, they proposed a simple squared-error reformulation [26] that was sensitive to class ordering.

Adience dataset has been used in other works for human age estimation. E. Eidinger [27] presented an approach using support vector machines and neural networks. J.-C. Chen [28] proposed a coarse-to-fine strategy for deep convolutional networks. G. Levi [29] presented another convolutional network model for age estimation. **Estos son ordinales? No son ordinales pero los cito aqui porque luego los uso para comparar con ellos ya que son los que he encontrado que utilicen las base de datos Adience para predecir edad.**

H. Li et al. [30] applied deep learning techniques for solving the ordinal problem of Alzheimer's diagnosis and detecting the different levels of the disease.

Y. Liu et al. [31] proposed a new approach of which transforms the ordinal regression problem to binary classification problems and use triplets with instances from different categories to train deep neural networks. In this way, high-level features describing the ordinal relationship are extracted automatically.

A. Rios et al. [32] 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.

S. Chen et al. [33] proposed a deep 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 and they achieved a tighter error bound for ranking-based age estimation.

H. Fu et al. [34] 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 ordinary regression loss. Also, they used a multi-scale network structure that avoids unnecessary spatial pooling.

Y. Liu et al. [35] proposed a constrained optimization formulation for the ordinal regression problem which minimizes the negative loglikelihood for multiple categories constrained by the order relationship between instances.

A. Pal et al. [36] defined a loss function for CNN that is based on the Earth Mover's Distance and takes into account the ordinal class relationships.

M. ALALI et al. [37] 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.

Currently, there is a great discussion in terms of finding the best activation function, offering good performance and mitigating the effects of the gradient vanishing problem [38], [39]. Rectified Linear Unit (ReLU) [40] is widely used in most deep learning works, but recently, Clevert et al. proposed the Exponential Linear Unit (ELU) [41]. They proved that ELUs alleviate the vanishing gradient problem via the identity for positive values. In their experiments, ELUs lead not only to faster learning but also to significantly better generalization performance than ReLUs and Leaky ReLU functions (LReLU) on networks with more than five layers. In this way, this paper considers the ELU function for the different experiments.

Sergey Ioffe and Christian Szegedy [42] described Batch Normalization and its benefits. It reduces the internal covariate shift by normalizing layer inputs. Their method draws its strength from making normalization a part of the model architecture and performing the normalization for each training batch. It allows them to use higher learning rates and be less careful about initialization. It also eliminates the need for using regularization techniques like Dropout.

III. CUMULATIVE LINK MODELS (CLM)

An ordinal classification problem consists of 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 of the 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 points, $D = \{(\mathbf{x}_i, y_i), i = 1, \dots, N\}$. 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 makes it possible 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, so 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 [43]. It dated back to 1980 and is a member of a wider family of models recognised as Cumulative Link Models (CLM) [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:

$$P(y \preceq \mathcal{C}_q | \mathbf{x}) = P(y = \mathcal{C}_1 | \mathbf{x}) + \dots + P(y = \mathcal{C}_q | \mathbf{x}), \quad (3)$$

$$P(y = \mathcal{C}_q | \mathbf{x}) = P(y \preceq \mathcal{C}_q | \mathbf{x}) - P(y \preceq \mathcal{C}_{q-1} | \mathbf{x}), \quad (4)$$

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 obtained from the output of the last layer, which has only one neuron. 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 [44]:

$$g^{-1}(P(y \preceq \mathcal{C}_q | \mathbf{x})) = b_q - f(\mathbf{x}), \quad q = 1, \dots, Q - 1, \quad (5)$$

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 ϵ is the random variable of the error. The most common choice for the probability distribution of ϵ 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 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 $b_1 \leq b_2 \leq \dots \leq b_{Q-1}$ ensures that $P(y \preceq \mathcal{C}_q | \mathbf{x})$ increases with q [43]. This order is achieved by defining the first threshold and calculating the rest of thresholds from the first in the following form:

$$b_n = b_1 + \sum_{i=1}^{n-1} \alpha_i^2, \quad n = 2, \dots, N, \quad (6)$$

where b_1 and α_i are learnable parameters, and N is the number of classes.

In this work, we use this ordinal model with different link functions for the probability distribution of ϵ , 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 $P(y \preceq \mathcal{C}_q | \mathbf{x}) = \Phi(b_q - f(\mathbf{x}))$ for a continuous cdf Φ .

- **Logit.** **logit** link function is the function used for the Proportional Odds Model. The **logit** link is:

$$\begin{aligned} \text{logit}[P(y \preceq \mathcal{C}_q | \mathbf{x})] &= \log \frac{P(y \preceq \mathcal{C}_q | \mathbf{x})}{1 - P(y \preceq \mathcal{C}_q | \mathbf{x})} = \\ &= b_q - f(\mathbf{x}), \quad q = 1, \dots, Q - 1, \end{aligned} \quad (7)$$

or the equivalent expression:

$$P(y \preceq \mathcal{C}_q | \mathbf{x}) = \frac{1}{1 + e^{-(b_q - f(\mathbf{x}))}}. \quad (8)$$

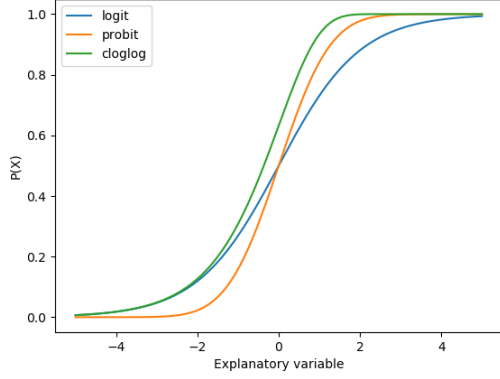


Fig. 1. Representation of link functions.

- **Probit.** **probit** link function is the inverse of the standard normal cumulative distribution function (cdf) Φ . Its expression is:

$$\begin{aligned}\Phi^{-1}[P(y \preceq C_q | \mathbf{x})] &= b_q - f(\mathbf{x}), \quad q = 1, \dots, Q-1, \\ P(y \preceq C_q | \mathbf{x}) &= \Phi(b_q - f(\mathbf{x})), \quad q = 1, \dots, Q-1,\end{aligned}\quad (9)$$

which can also be expressed as:

$$P(y \preceq C_q | \mathbf{x}) = \int_{-\infty}^{b_q - f(\mathbf{x})} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx. \quad (10)$$

- **Complementary log-log.** Like the **logit** and the **probit** transformation, the complementary log-log transformation takes a response that is restricted to the $(0, 1)$ interval and converts it into something in the $(-\infty, +\infty)$ interval. Complementary log-log expression is:

$$\begin{aligned}\log[-\log[1 - P(y \preceq C_q | \mathbf{x})]] &= \\ &= b_q - f(\mathbf{x}), \quad q = 1, \dots, Q-1,\end{aligned}\quad (11)$$

that is:

$$P(y \preceq C_q | \mathbf{x}) = 1 - e^{-e^{b_q - f(\mathbf{x})}}, \quad q = 1, \dots, Q-1. \quad (12)$$

Logit and **probit** links are symmetric:

$$\text{link}[P(y \preceq C_q | \mathbf{x})] = -\text{link}[1 - P(y \preceq C_q | \mathbf{x})], \quad (13)$$

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 symmetric property can be demonstrated as follows:

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

$$\begin{aligned}\text{link}(p) &= \text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \\ &= \log(p) - \log(1-p),\end{aligned}\quad (14)$$

while:

$$\begin{aligned}-\text{link}(1-p) &= -\text{logit}(1-p) = \\ &= -\log\left(\frac{1-p}{p}\right) = -\log(1-p) + \log(p).\end{aligned}\quad (15)$$

- 2) For the **probit**:

$$\begin{aligned}p &\equiv P(y \preceq C_q | \mathbf{x}) = \Phi(b_q - f(\mathbf{x})) = \\ &= \int_{-\infty}^{b_q - f(\mathbf{x})} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx,\end{aligned}\quad (16)$$

which leads to:

$$\text{probit}(p) = \Phi^{-1}(p) = b_q - f(\mathbf{x}), \quad (17)$$

$$-\text{probit}(1-p) = \Phi^{-1}(1-p) = -b_q + f(\mathbf{x}), \quad (18)$$

where:

$$1-p = \int_{-\infty}^{-b_q + f(\mathbf{x})} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx, \quad (19)$$

$$p = 1 - \int_{-\infty}^{-b_q + f(\mathbf{x})} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx. \quad (20)$$

Unlike **logit** and **probit**, the complementary log-log model is asymmetrical. It is frequently used when the probability of an event is very small or very large. 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. However, in this situation, the complementary log-log model might give a satisfying answer.

IV. EXPERIMENTS

A. Data

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

- **Diabetic Retinopathy (DR)**¹. DR is a dataset consisting of extremely high-resolution fundus image data. The training set consists of 17563 pairs of images (where a pair consists of 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 (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, varying conditions of illumination and different resolutions. These images come from the EyePACS dataset that was used in a Diabetic Retinopathy Detection competition hosted on the Kaggle platform. Also, this dataset was used in later works [19], [45], applying an ordinal QWK cost function in [19] 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 images of this dataset are shown in Figure 2.
- **Adience**². This dataset consists of 26580 faces belonging to 2284 subjects. We use the form of the dataset

¹<https://www.kaggle.com/c/diabetic-retinopathy-detection/data>

²<http://www.openul.ac.il/home/hassner/Adience/data.html>

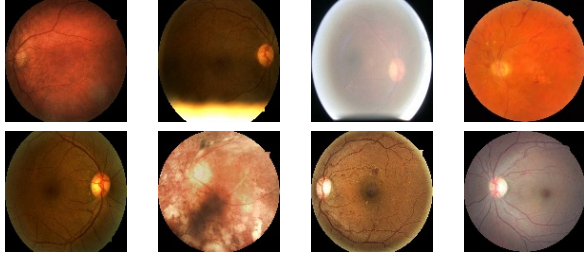


Fig. 2. Examples of the Diabetic Retinopathy test set.



Fig. 3. Examples of the Diabetic Retinopathy test set.

where faces have been pre-cropped and aligned. The dataset was preprocessed, using the methods described in a previous work [24], 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.

B. Model

CNNs have been used for both datasets. The different architectures of CNN used in these experiments are presented in tables I and II. The architecture for DR is the same that was used in [19] and the network for Adience is a small Residual Network (ResNet) [3] that was used in [24]. 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 for obtaining the next pixel. Pooling layers have got similar parameters: the pool size is the number of pixels that will be involved in the pooling operation, and the stride represents the same concept that in convolutional layers. For convolutional layers, $\text{ConvWxH@FsS} = F$ filters of size $W \times H$ and stride S . For pooling layers, $\text{PoolWxHsS} = \text{pool size of } W \times H \text{ and stride } S$.

The Exponential Linear Unit (ELU) [41] has been used for activation function for all the convolutional and dense layers, instead of the ReLU [40] function, as it mitigates the effects of the vanishing gradients problem [38], [39] via the identity for positive values. Also, ELUs lead to faster training and better generalization performance than ReLU

TABLE I
DESCRIPTION OF THE ARCHITECTURE USED IN THE DR
EXPERIMENTS.

Layer	Output shape
2 x Conv3x3@32s1	252x252x32
MaxPool2x2s2	126x126x32
2 x Conv3x3@64s1	122x122x64
MaxPool2x2s2	61x61x64
2 x Conv3x3@128s1	57x57x128
MaxPool2x2s2	28x28x128
2 x Conv3x3@128s1	24x24x128
MaxPool2x2s2	12x12x128
Conv4x4@128s1	9x9x128

TABLE II
DESCRIPTION OF THE ARCHITECTURE USED IN THE ADIANCE
EXPERIMENTS.

Layer	Output shape
Conv7x7@32s2	112x112x32
MaxPool3x3s2	55x55x32
2 x ResBlock3x3@64s1	55x55x32
1 x ResBlock3x3@128s2	28x28x64
2 x ResBlock3x3@128s1	28x28x64
1 x ResBlock3x3@256s2	14x14x128
2 x ResBlock3x3@256s1	14x14x128
1 x ResBlock3x3@512s2	7x7x256
2 x ResBlock3x3@512s1	7x7x256
AveragePool7x7s2	1x1x256

and Leaky ReLU (LReLU) [46] functions on networks with more than five layers.

After every ELU activation function of the convolutional layers, Batch Normalization [42] 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 weights initialization. It also eliminates the need for using regularization techniques like Dropout.

At the output of the network, the CLM is used. Also, a learnable parameter has been used to rescale the projections used by the Cumulative Link Model to make it more stable and guarantee the convergence in most cases.

C. Experimental design

The model is optimized using a batch based first-order optimization algorithm called Adam [47]. We study different initial learning rates in order to find the optimal one for each problem. We apply an exponential decay across training epochs to the initial learning rate.

Both datasets have been artificially equalized using data augmentation techniques [48], [4]. However, different transformations were applied to each one. DR dataset augmentation was based on image cropping and zooming, horizontal and vertical flipping, brightness adjustment and random rotations. Horizontal flipping was the only transformation applied to Adience dataset.

In the case of DR dataset, the epoch size has been fixed to 100000 images per epoch. For the Adience dataset, epoch size is the number of images in the training set.

The model is evaluated using the Quadratic Weighted Kappa metric (QWK) [49]. This evaluation measure gives

a higher weight to the errors that are further from the correct class. Quadratic weighted kappa loss is considered as loss function for the optimizer as it gives better performance for ordinal classification problems [19].

Also, other evaluation metrics have been used to ease the comparison with other works: Minimum Sensitivity (MS) [16], Mean Absolute Error (MAE) [16], Mean Squared Error (MSE) [50], Correct Classification Rate (CCR), Top-2 CCR [24], Top-3 CCR [24] and 1-off accuracy [28], [29], [27].

Experiments were run with the standard cross-entropy loss and the softmax function too in order to prove the performance improvement of considering the ordinality of the problem (QWK loss and the Cumulative Link Model). The results of these experiments are analysed in Section V-D.

D. Factors

Three different factors have been considered: learning rate, batch size and link function for the final output layer.

- *Learning rate* (LR, η). Learning rate 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 [51], [52]. Within this work, we have considered three different values for 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 have tried three separate batch sizes for each dataset. For DR dataset, we have used 5, 10 and 15 while, for the Adience dataset, 64, 128 and 256 images were used. We took the batch sizes that were used in [19] and [24] as a reference, and we expanded the range on both sides.
- *Link function* (LF). Different link functions have been used for the CLM at the last layer output: **logit**, **probit** and complementary log-log.

V. RESULTS

In this section, we present the results of the experiments. For each dataset, we show a table with the detailed experiments performed for training the model with each combination of parameters. Every parameter combination has been 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 Diabetic Retinopathy dataset are presented 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 complementary log-log link function using a batch size of 10 and a learning rate of 10^{-3} . However, the best CCR value was obtained with a batch size of 15, the **logit** link and a learning rate of 10^{-4} . The optimal 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 IV. 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 batch size of 64 and a learning rate of 10^{-4} . Also, this configuration obtained the best rate 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 result.

The significance and relative importance of the parameters concerning the results obtained, as well as suitable values for each of them, were obtained using an ANalysis Of the VAriance (ANOVA).

The ANOVA test [22] 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 link function, ~~which can be logit, probit and complementary log-log~~, the learning rate for the Adam optimization algorithm (~~10^{-2} , 10^{-3} and 10^{-4}~~), and the batch size (~~5, 10 and 15 for DR and 64, 128 and 256 for Adience~~). **Esto está repitiendo mucho**

Following the setup of the previous study, we performed an ANOVA III analysis and multiple comparison tests. We assume that five executions are enough to do the statistical tests because of the computational time limitations.

We denote by $QWK_{i,j,k,l}$ ($i = 1, \dots, 3; j = 1, \dots, 3; k = 1, \dots, 3$) the value observed when the first factor is at the i -th level, the second at the j -th level and the third at the k -th level. We assume that the three factors do not act independently and therefore there exists an interaction between each pair of them and between the three factors. In this case, the observations fit:

$$QWK_{i,j,k,l} = \mu + L_i + P_j + B_k + LP_{i,j} + LB_{i,k} + PB_{j,k} + LPB_{i,j,k} + \epsilon_{i,j,k,l}, \quad (21)$$

TABLE III
DIABETIC RETINOPATHY RESULTS. BS STANDS FOR BATCH SIZE, LF FOR LINK FUNCTION AND LR FOR LEARNING RATE. MEAN AND STANDARD DEVIATION MEANS_{SD}.

BS	LF	LR	QWK _(SD)	MS _(SD)	MAE _(SD)	MSE _(SD)	CCR _(SD)	Top-2 _(SD)	Top-3 _(SD)	1-off _(SD)
5	clog-log	10 ⁻²	0.414 _(0.057)	0.075 _(0.042)	0.177 _(0.023)	0.165 _(0.020)	0.556 _(0.057)	0.833 _(0.042)	0.968 _(0.011)	0.816 _(0.021)
5	clog-log	10 ⁻³	0.534 _(0.027)	0.102 _(0.011)	0.137 _(0.006)	0.123 _(0.004)	0.658 _(0.015)	0.871 _(0.011)	0.966 _(0.003)	0.852 _(0.002)
5	clog-log	10 ⁻⁴	0.520 _(0.006)	0.067 _(0.008)	0.123 _(0.003)	0.104 _(0.001)	0.697 _(0.006)	0.842 _(0.008)	0.961 _(0.003)	0.851 _(0.002)
5	logit	10 ⁻²	0.416 _(0.041)	0.095 _(0.029)	0.175 _(0.021)	0.162 _(0.018)	0.563 _(0.054)	0.762 _(0.040)	0.908 _(0.026)	0.807 _(0.029)
5	logit	10 ⁻³	0.554 _(0.013)	0.093 _(0.009)	0.137 _(0.003)	0.123 _(0.003)	0.660 _(0.008)	0.802 _(0.005)	0.936 _(0.004)	0.853 _(0.005)
5	logit	10 ⁻⁴	0.520 _(0.003)	0.063 _(0.004)	0.122 _(0.002)	0.102 _(0.001)	0.706 _(0.005)	0.823 _(0.004)	0.949 _(0.003)	0.862 _(0.003)
5	probit	10 ⁻²	0.460 _(0.048)	0.079 _(0.046)	0.197 _(0.064)	0.182 _(0.053)	0.504 _(0.167)	0.808 _(0.034)	0.927 _(0.073)	0.689 _(0.240)
5	probit	10 ⁻³	0.564 _(0.018)	0.099 _(0.013)	0.147 _(0.018)	0.132 _(0.014)	0.636 _(0.045)	0.822 _(0.040)	0.939 _(0.020)	0.840 _(0.015)
5	probit	10 ⁻⁴	0.523 _(0.005)	0.067 _(0.012)	0.122 _(0.002)	0.105 _(0.001)	0.701 _(0.006)	0.823 _(0.002)	0.953 _(0.002)	0.860 _(0.003)
10	clog-log	10 ⁻²	0.423 _(0.239)	0.062 _(0.051)	0.127 _(0.017)	0.120 _(0.014)	0.684 _(0.046)	0.894_(0.062)	0.986_(0.012)	0.832 _(0.020)
10	clog-log	10 ⁻³	0.582_(0.016)	0.102 _(0.006)	0.128 _(0.003)	0.115 _(0.002)	0.680 _(0.007)	0.880 _(0.004)	0.972 _(0.003)	0.861 _(0.004)
10	clog-log	10 ⁻⁴	0.537 _(0.010)	0.064 _(0.004)	0.116 _(0.001)	0.096 _(0.001)	0.717 _(0.003)	0.837 _(0.002)	0.971 _(0.001)	0.860 _(0.002)
10	logit	10 ⁻²	0.531 _(0.031)	0.107 _(0.008)	0.151 _(0.010)	0.140 _(0.008)	0.623 _(0.025)	0.802 _(0.022)	0.934 _(0.013)	0.838 _(0.014)
10	logit	10 ⁻³	0.579 _(0.009)	0.096 _(0.012)	0.127 _(0.005)	0.113 _(0.004)	0.686 _(0.013)	0.817 _(0.006)	0.954 _(0.005)	0.861 _(0.002)
10	logit	10 ⁻⁴	0.539 _(0.007)	0.074 _(0.013)	0.126 _(0.005)	0.095 _(0.002)	0.707 _(0.010)	0.823 _(0.007)	0.957 _(0.005)	0.858 _(0.004)
10	probit	10 ⁻²	0.508 _(0.037)	0.088 _(0.044)	0.145 _(0.018)	0.137 _(0.014)	0.639 _(0.045)	0.835 _(0.015)	0.960 _(0.008)	0.829 _(0.020)
10	probit	10 ⁻³	0.558 _(0.034)	0.111_(0.005)	0.134 _(0.003)	0.120 _(0.002)	0.666 _(0.008)	0.831 _(0.007)	0.955 _(0.001)	0.863 _(0.003)
10	probit	10 ⁻⁴	0.541 _(0.010)	0.076 _(0.006)	0.119 _(0.002)	0.098 _(0.001)	0.712 _(0.005)	0.828 _(0.003)	0.961 _(0.002)	0.862 _(0.001)
15	clog-log	10 ⁻²	0.564 _(0.016)	0.108 _(0.014)	0.143 _(0.006)	0.134 _(0.005)	0.640 _(0.015)	0.879 _(0.011)	0.972 _(0.005)	0.851 _(0.006)
15	clog-log	10 ⁻³	0.559 _(0.026)	0.111 _(0.008)	0.127 _(0.004)	0.113 _(0.003)	0.682 _(0.010)	0.871 _(0.008)	0.974 _(0.002)	0.868_(0.002)
15	clog-log	10 ⁻⁴	0.538 _(0.009)	0.054 _(0.003)	0.115_(0.002)	0.093 _(0.001)	0.720 _(0.006)	0.835 _(0.007)	0.970 _(0.003)	0.860 _(0.006)
15	logit	10 ⁻²	0.551 _(0.020)	0.104 _(0.008)	0.139 _(0.011)	0.129 _(0.009)	0.654 _(0.027)	0.815 _(0.017)	0.948 _(0.016)	0.856 _(0.015)
15	logit	10 ⁻³	0.551 _(0.010)	0.106 _(0.016)	0.129 _(0.008)	0.114 _(0.005)	0.680 _(0.019)	0.818 _(0.008)	0.952 _(0.007)	0.866 _(0.001)
15	logit	10 ⁻⁴	0.543 _(0.008)	0.056 _(0.003)	0.121 _(0.004)	0.090_(0.001)	0.723_(0.004)	0.833 _(0.004)	0.964 _(0.003)	0.862 _(0.004)
15	probit	10 ⁻²	0.534 _(0.032)	0.104 _(0.013)	0.148 _(0.015)	0.138 _(0.014)	0.631 _(0.038)	0.845 _(0.030)	0.964 _(0.010)	0.852 _(0.010)
15	probit	10 ⁻³	0.580 _(0.021)	0.104 _(0.016)	0.129 _(0.008)	0.116 _(0.005)	0.680 _(0.018)	0.832 _(0.010)	0.959 _(0.007)	0.866 _(0.003)
15	probit	10 ⁻⁴	0.533 _(0.004)	0.065 _(0.005)	0.117 _(0.002)	0.094 _(0.001)	0.721 _(0.004)	0.832 _(0.002)	0.964 _(0.001)	0.863 _(0.001)

TABLE IV
ADIENCE TEST RESULTS. BS STANDS FOR BATCH SIZE, LF FOR LINK FUNCTION AND LR FOR LEARNING RATE. MEAN AND STANDARD DEVIATION MEANS_{SD}.

BS	LF	LR	QWK _(SD)	MS _(SD)	MAE _(SD)	MSE _(SD)	CCR _(SD)	Top-2 _(SD)	Top-3 _(SD)	1-off _(SD)
64	clog-log	10 ⁻²	0.808 _(0.025)	0.086 _(0.041)	0.147 _(0.008)	0.129 _(0.007)	0.415 _(0.031)	0.677 _(0.024)	0.798 _(0.036)	0.804 _(0.015)
64	clog-log	10 ⁻³	0.873 _(0.006)	0.144 _(0.057)	0.124_(0.003)	0.101 _(0.003)	0.519_(0.014)	0.764 _(0.010)	0.861 _(0.019)	0.886 _(0.006)
64	clog-log	10 ⁻⁴	0.799 _(0.010)	0.000 _(0.000)	0.174 _(0.001)	0.100 _(0.002)	0.324 _(0.015)	0.616 _(0.020)	0.795 _(0.012)	0.771 _(0.014)
64	logit	10 ⁻²	0.778 _(0.019)	0.074 _(0.041)	0.159 _(0.006)	0.137 _(0.007)	0.366 _(0.025)	0.636 _(0.015)	0.785 _(0.010)	0.775 _(0.015)
64	logit	10 ⁻³	0.881_(0.005)	0.178 _(0.023)	0.126 _(0.001)	0.098 _(0.003)	0.518 _(0.008)	0.765_(0.015)	0.902_(0.005)	0.894_(0.005)
64	logit	10 ⁻⁴	0.784 _(0.011)	0.000 _(0.000)	0.180 _(0.001)	0.108 _(0.004)	0.318 _(0.026)	0.621 _(0.034)	0.772 _(0.024)	0.731 _(0.030)
64	probit	10 ⁻²	0.836 _(0.005)	0.135 _(0.021)	0.134 _(0.002)	0.121 _(0.002)	0.468 _(0.011)	0.720 _(0.009)	0.861 _(0.009)	0.829 _(0.005)
64	probit	10 ⁻³	0.874 _(0.004)	0.134 _(0.012)	0.126 _(0.003)	0.105 _(0.003)	0.511 _(0.014)	0.756 _(0.009)	0.895 _(0.003)	0.889 _(0.003)
64	probit	10 ⁻⁴	0.805 _(0.004)	0.000 _(0.000)	0.170 _(0.001)	0.100 _(0.002)	0.360 _(0.011)	0.653 _(0.011)	0.809 _(0.009)	0.790 _(0.009)
128	clog-log	10 ⁻²	0.832 _(0.013)	0.123 _(0.031)	0.135 _(0.004)	0.117 _(0.002)	0.463 _(0.013)	0.705 _(0.019)	0.813 _(0.025)	0.832 _(0.006)
128	clog-log	10 ⁻³	0.873 _(0.006)	0.185_(0.029)	0.128 _(0.002)	0.100 _(0.001)	0.513 _(0.007)	0.758 _(0.008)	0.870 _(0.011)	0.880 _(0.009)
128	clog-log	10 ⁻⁴	0.659 _(0.025)	0.000 _(0.000)	0.190 _(0.002)	0.125 _(0.004)	0.235 _(0.026)	0.466 _(0.031)	0.640 _(0.030)	0.536 _(0.041)
128	logit	10 ⁻²	0.781 _(0.041)	0.096 _(0.059)	0.153 _(0.007)	0.125 _(0.007)	0.398 _(0.031)	0.638 _(0.033)	0.790 _(0.025)	0.779 _(0.020)
128	logit	10 ⁻³	0.865 _(0.005)	0.127 _(0.026)	0.134 _(0.001)	0.099 _(0.003)	0.497 _(0.009)	0.754 _(0.008)	0.882 _(0.009)	0.874 _(0.008)
128	logit	10 ⁻⁴	0.586 _(0.008)	0.000 _(0.000)	0.196 _(0.001)	0.151 _(0.005)	0.192 _(0.001)	0.364 _(0.060)	0.581 _(0.034)	0.396 _(0.002)
128	probit	10 ⁻²	0.849 _(0.005)	0.132 _(0.010)	0.131 _(0.001)	0.115 _(0.001)	0.479 _(0.004)	0.728 _(0.007)	0.854 _(0.009)	0.847 _(0.007)
128	probit	10 ⁻³	0.866 _(0.002)	0.124 _(0.043)	0.130 _(0.002)	0.100 _(0.003)	0.505 _(0.006)	0.750 _(0.010)	0.882 _(0.004)	0.873 _(0.006)
128	probit	10 ⁻⁴	0.718 _(0.015)	0.000 _(0.000)	0.185 _(0.001)	0.110 _(0.002)	0.300 _(0.031)	0.575 _(0.015)	0.733 _(0.010)	0.640 _(0.033)
256	clog-log	10 ⁻²	0.853 _(0.004)	0.157 _(0.024)	0.130 _(0.002)	0.110 _(0.001)	0.485 _(0.009)	0.744 _(0.006)	0.842 _(0.016)	0.858 _(0.004)
256	clog-log	10 ⁻³	0.840 _(0.017)	0.095 _(0.017)	0.144 _(0.005)	0.097 _(0.004)	0.456 _(0.021)	0.720 _(0.022)	0.840 _(0.018)	0.842 _(0.018)
256	clog-log	10 ⁻⁴	0.552 _(0.010)	0.000 _(0.000)	0.199 _(0.001)	0.165 _(0.004)	0.187 _(0.001)	0.368 _(0.022)	0.475 _(0.025)	0.387 _(0.001)
256	logit	10 ⁻²	0.764 _(0.102)	0.077 _(0.067)	0.155 _(0.020)	0.125 _(0.015)	0.387 _(0.083)	0.632 _(0.103)	0.790 _(0.077)	0.783 _(0.065)
256	logit	10 ⁻³	0.851 _(0.008)	0.100 _(0.030)	0.147 _(0.003)	0.094_(0.002)	0.449 _(0.015)	0.726 _(0.015)	0.861 _(0.006)	0.850 _(0.008)
256	logit	10 ⁻⁴	0.558 _(0.008)	0.000 _(0.000)	0.202 _(0.001)	0.191 _(0.002)	0.187 _(0.002)	0.206 _(0.007)	0.395 _(0.046)	0.389 _(0.003)
256	probit	10 ⁻²	0.858 _(0.005)	0.164 _(0.033)	0.130 _(0.002)	0.112 _(0.002)	0.486 _(0.007)	0.741 _(0.008)	0.867 _(0.008)	0.862 _(0.005)
256	probit	10 ⁻³	0.850 _(0.008)	0.111 _(0.040)	0.144 _(0.002)	0.095 _(0.001)	0.460 _(0.011)	0.732 _(0.006)	0.865 _(0.006)	0.853 _(0.007)
256	probit	10 ⁻⁴	0.565 _(0.010)	0.000 _(0.000)	0.196 _(0.001)	0.150 _(0.005)	0.189 _(0.001)	0.409 _(0.014)	0.602 _(0.022)	0.392 _(0.002)

TABLE V

ANOVA III FOR THE ANALYSIS OF THE MAIN FACTORS IN THE DESIGN OF A CONVOLUTIONAL ORDINAL NEURAL NETWORK FOR THE RETINOPATHY DATASET.

Response variable QWK					
Source	S.S.	D.F.	M.S.	F-ratio	Sig.
Model	37.860	9	4.207	1562.840	0.000
L factor	0.057	2	0.029	10.646	0.000
P factor	0.121	2	0.060	22.468	0.000
LP factors	0.057	4	0.014	5.261	0.001
Error	0.339	126	0.003		
Total	38.199	135			

where μ is the fixed effect that is common to all the populations; L_i is the effect associated with the i -th level of the link factor (**logit**, **probit**, complementary log-log); P_j is the effect associated with the j -th level of the learning rate factor and B_k is the effect associated with the k -th level of the batch size factor. The term $LP_{i,j}$ denotes the joint effect of the presence of level i of the first factor and level j of the second one; this, therefore, is denominated the interaction term between L and P factors. The same interaction effect is appreciated on $LB_{i,k}$, $PB_{j,k}$ and $LPB_{i,j,k}$. The term $\epsilon_{i,j,k,l}$ is the influence on the result of everything that could not be assigned or of random factors. $QWK_{i,j,k,l}$ is the quadratic weighted kappa measure, the response variable used to perform the statistical analysis.

We consider some hypotheses testing where the null hypothesis is proposed that each term of the above equation is independent of the levels involved. The hypotheses for the levels of the L factor are $H_0 \equiv L_1 = L_2 = L_3$, and $H_1 \equiv$ some L_i is different.

The same hypotheses are made for the other factors. In this way, we test in the null hypothesis that all of the population means are equal against an alternative hypothesis that there is at least one mean that is not equal to the others.

The hypothesis associated with the interaction between L and P is $H_0 \equiv LP_{i,j} = 0, \forall i, j$, and $H_1 \equiv \exists LP_{i,j} \neq 0$. Similar hypotheses can be assumed for the interaction between the other factors.

The analysis of variance table represents the initial study in a compact form, containing the sum of squares, degrees of freedom, mean square, test statistics and significance level, where non-significative factors and interactions have been removed (p-value > 0.05). These factors and interactions take part of the error component now. In this way, the results of the ANOVA III test for the Diabetic Retinopathy dataset are summarised in Table V. There are significant differences in average QWK depending on the link function and also depending on the learning rate for $\alpha = 0.05$ (p-value = 0.000). Moreover, an interaction between the link function and the learning rate can be recognised (p-value = 0.001). It means that the learning rate and the link functions have a significant impact on the optimization algorithm results.

Given that there exist significant differences between

TABLE VI

TUKEY'S TEST RESULTS FOR THE DIABETIC RETINOPATHY DATASET.

LF	LF	Mean diff.	Sig.
logit	probit	-0.002	0.011
	clog-log	0.012	0.000
probit	logit	0.002	0.011
	clog-log	0.014	0.248
clog-log	logit	-0.012	0.000
	probit	-0.014	0.248
LR	LR	Mean diff.	Sig.
10^{-2}	10^{-3}	-0.073	0.000
	10^{-4}	-0.044	0.000
10^{-3}	10^{-2}	0.073	0.000
	10^{-4}	0.029	0.023
10^{-4}	10^{-2}	0.044	0.000
	10^{-3}	-0.029	0.023

TABLE VII

ANOVA III FOR THE ANALYSIS OF THE MAIN FACTORS IN THE DESIGN OF A CONVOLUTIONAL ORDINAL NEURAL NETWORK FOR THE ADIANCE DATASET.

Response variable QWK					
Source	S.S.	D.F.	M.S.	F-ratio	Sig.
Model	84.372	27	3.125	4414.006	0.000
L factor	0.156	2	0.078	110.103	0.000
P factor	0.925	2	0.462	653.163	0.000
B factor	0.040	2	0.020	28.218	0.000
LP factors	0.284	4	0.071	100.118	0.000
LB factors	0.008	4	0.002	2.837	0.028
PB factors	0.026	4	0.007	9.267	0.000
LPB factors	0.021	8	0.003	3.728	0.001
Error	0.076	108	0.001		
Total	84.449	135			

the means, we analyse now those differences. A post-hoc multiple comparison test has been performed on the mean QWK obtained. An HSD Tukey's test [23] was made 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 VI. They show that the best link function is the complementary log-log but the **probit** link performance is close to it. Also, the best value for the learning rate parameter is 10^{-3} . The batch size is not relevant for this dataset with the values considered.

The results of the ANOVA III test for the Adience dataset are shown in Table VII. First, we observe 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 factors together (p-values 0.000, 0.000, 0.000 and 0.001). It means that the joint action of two or three factors significantly affects the results obtained by the algorithm.

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 groups, an HSD Tukey's test has been applied. The results of this test over the test set are shown in Table

TABLE VIII
TUKEY'S TEST RESULTS FOR THE ADIANCE DATASET.

LF	LF	Mean diff.	Sig.
logit	probit	0.046	0.000
	clog-log	0.084	0.000
	logit	-0.046	0.000
	clog-log	0.038	0.000
clog-log	logit	-0.084	0.000
	probit	-0.038	0.000
LR	LR	Mean diff.	Sig.
10^{-2}	10^{-3}	-0.046	0.000
	10^{-4}	0.148	0.000
10^{-3}	10^{-2}	0.046	0.000
	10^{-4}	0.194	0.000
10^{-4}	10^{-2}	-0.148	0.000
	10^{-3}	-0.194	0.000
BS	BS	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

VIII.

The results over the test set show that the best link function is the **logit** one, the best learning rate is 10^{-3} and the best batch size is 128. However, the interactions between these factors made the configuration that uses a **logit** link, $\eta = 10^{-3}$ and batch size 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 optimal value for the batch size and the optimal link function are not the same for Retinopathy and Adience datasets. These results highlight the importance of adjusting the hyper-parameters for each problem instead of trying to find an optimal configuration for all the datasets. However, the best learning rate for both datasets were 10^{-3} . It is recommended to use this value for future datasets. The best batch size 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 the Diabetic Retinopathy one. This highlights the importance of making experimental designs associated with each dataset to determine the best value for each factor.

D. Comparison with nominal method and previous works

The same experiments described in Section IV were repeated using the cross-entropy instead of the QWK as loss function for the optimizer and the softmax function instead of the Cumulative Link Model for the output of the network. The evaluation metric remains the same in order to be able to compare. There are some parameter configurations where the training process gets stagnated and a very low QWK is obtained. As we saw in Sections

TABLE IX
COMPARISON BETWEEN THE BEST RESULTS OF NOMINAL, ORDINAL AND PREVIOUS WORKS FOR THE DIABETIC RETINOPATHY DATASET.

Method	$\overline{QWK}_{(SD)}$	$\overline{CCR}_{(SD)}$	$\overline{1-off}_{(SD)}$
Ordinal network	0.582 _(0.016)	0.723 _(0.004)	0.868 _(0.002)
Nominal network	0.498 _(0.011)	0.692 _(0.012)	0.854 _(0.006)
J. Torre et al. [19]	0.537 ₍₋₎	-	-
Å. Nebot et al. [45]	0.555 ₍₋₎	-	-

TABLE X
COMPARISON BETWEEN THE BEST RESULTS OF NOMINAL, ORDINAL AND PREVIOUS WORKS FOR THE ADIANCE DATASET.

Method	$\overline{QWK}_{(SD)}$	$\overline{CCR}_{(SD)}$	$\overline{1-off}_{(SD)}$
Ordinal network	0.881 _(0.005)	0.519 _(0.013)	0.894 _(0.005)
Nominal network	0.787 _(0.004)	0.458 _(0.008)	0.800 _(0.007)
E. Eidinger et al. [27]	-	0.451 _(0.026)	0.807 _(0.011)
J.-C. Chen et al. [28]	-	0.529 _(0.060)	0.885 _(0.022)
G. Levi et al. [29]	-	0.507 _(0.051)	0.847 _(0.022)

V-A and V-B, this problem is not found when using the ordinal method.

For the Diabetic Retinopathy dataset, the best mean value of QWK was 0.497 and was obtained when using a batch size of 10 and a learning rate of 10^{-4} .

In the case of Adience dataset, the highest QWK was 0.787 and was achieved with a batch size of 64 and a learning rate of 10^{-3} .

Finally, a comparison of the best results for each dataset for ordinal and nominal cases and previous works is shown in tables IX and X. All the results are given for the test set, except those from [19] (DR dataset), because the authors only provided validation results for 128×128 images (however, validation results are usually better than test results). The proposed ordinal model outperforms all the other alternatives in terms of QWK. The performance gain of CLM over the softmax reaches 16.8% for Diabetic Retinopathy and 11.9% for Adience dataset. The improvement of the ordinal method for Retinopathy dataset is higher than for Adience dataset. It seems that the method proposed in this work offers a more significant improvement as the given problem complexity increases.

VI. CONCLUSIONS

The first conclusion obtained from our results is that the optimal values for the different parameters considered are problem-dependant.

The complementary log-log function offers the best results in Diabetic Retinopathy dataset while the **logit** link is the best option for the Adience dataset. These results provide an opportunity for exploring new generalised link functions that could be dynamically adapted to any problem.

The best value for the learning rate parameter for both datasets is $\eta = 10^{-3}$. It can be considered a good value for this parameter when training the model with new datasets.

Both datasets have obtained the best performance with an intermediate batch size: 10 for Diabetic Retinopathy and 128 for Adience.

Also, the statistical tests reported that there are relevant interactions between the three factors that we have taken into account. The results highlight the importance of making an experimental design where all of these parameters are adjusted for each problem.

The proposed CLM has improved the performance of the deep network compared to the model that uses the softmax function and the models proposed in previous works. Also, it reduces the chance that the model gets stuck when training with some parameter configurations. So, the most significant improvements of these link functions are the performance increase, the reduction of the number of parameters configurations that should be tried to find the best one and the prevention of the over-fitting and the stagnation.

REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] D. Cireřan, U. Meier, and J. Schmidhuber, "Multi-column deep neural networks for image classification," *arXiv preprint arXiv:1202.2745*, 2012.
- [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, "iprivacy: 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] 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.
- [11] 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.
- [12] 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.
- [13] 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.
- [14] 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.
- [15] 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.
- [16] 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.
- [17] N. Mehdiyev, D. Enke, P. Fetteke, and P. Loos, "Evaluating forecasting methods by considering different accuracy measures," *Procedia Computer Science*, vol. 95, pp. 264–271, 2016.
- [18] J. Cohen, "A coefficient of agreement for nominal scales," *Educational and psychological measurement*, vol. 20, no. 1, pp. 37–46, 1960.
- [19] 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.
- [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] R. G. Miller Jr, *Beyond ANOVA: basics of applied statistics*. Chapman and Hall/CRC, 1997.
- [23] J. W. Tukey, "Comparing individual means in the analysis of variance," *Biometrics*, pp. 99–114, 1949.
- [24] C. Beckham and C. Pal, "Unimodal probability distributions for deep ordinal classification," *arXiv preprint arXiv:1705.05278*, 2017.
- [25] 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.
- [26] C. Beckham and C. Pal, "A simple squared-error reformulation for ordinal classification," *arXiv preprint arXiv:1612.00775*, 2016.
- [27] 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.
- [28] 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.
- [29] 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.
- [30] 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.
- [31] 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.
- [32] 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.
- [33] 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.
- [34] 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.
- [35] 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.
- [36] 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.
- [37] 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.
 - [38] 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.
 - [39] 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.
 - [40] 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.
 - [41] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, “Fast and accurate deep network learning by exponential linear units (elus),” *arXiv preprint arXiv:1511.07289*, 2015.
 - [42] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” *arXiv preprint arXiv:1502.03167*, 2015.
 - [43] P. McCullagh, “Regression models for ordinal data,” *Journal of the royal statistical society. Series B (Methodological)*, pp. 109–142, 1980.
 - [44] R. Herbrich, “Large margin rank boundaries for ordinal regression,” *Advances in large margin classifiers*, pp. 115–132, 2000.
 - [45] 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.
 - [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] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
 - [48] 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.
 - [49] A. Ben-David, “Comparison of classification accuracy using cohen’s weighted kappa,” *Expert Systems with Applications*, vol. 34, no. 2, pp. 825–832, 2008.
 - [50] Z. Wang and A. C. Bovik, “Mean squared error: Love it or leave it? a new look at signal fidelity measures,” *IEEE signal processing magazine*, vol. 26, no. 1, pp. 98–117, 2009.
 - [51] L. N. Smith, “Cyclical learning rates for training neural networks,” in *Applications of Computer Vision (WACV)*, 2017 IEEE Winter Conference on. IEEE, 2017, pp. 464–472.
 - [52] A. Senior, G. Heigold, K. Yang et al., “An empirical study of learning rates in deep neural networks for speech recognition,” in *Acoustics, Speech and Signal Processing (ICASSP)*, 2013 IEEE International Conference on. IEEE, 2013, pp. 6724–6728.