

# Ordinal Losses for Classification of Cervical Cancer Risk

**Tomé Albuquerque**<sup>Corresp., 1, 2</sup>, **Ricardo Cruz**<sup>1, 2</sup>, **Jaime S Cardoso**<sup>1, 2</sup>

<sup>1</sup> Institute for Systems and Computer Engineering, Technology and Science, Porto, Portugal

<sup>2</sup> Faculty of Engineering of the University of Porto, Porto, Portugal

Corresponding Author: Tomé Albuquerque  
Email address: tome.m.albuquerque@inesctec.pt

Cervical cancer is the fourth leading cause of cancer-related deaths in women, especially in low to middle-income countries. Despite the outburst of recent scientific advances, there is no totally effective treatment, especially when diagnosed in an advanced stage. Screening tests, such as cytology or colposcopy, have been responsible for a strong decrease in cervical cancer deaths. Cervical cancer automatic screening via Pap smear is a highly valuable cell imaging-based detection tool, where cells must be classified as being within one a multitude of ordinal classes, ranging from abnormal to normal. Current approaches at ordinal inference for neural networks are found to not take advantage of the ordinal problem or to be too uncompromising. A non-parametric ordinal loss for neuronal networks is proposed that promotes the output probabilities to follow a unimodal distribution. This is done by imposing a set of different constraints over all pairs of consecutive labels which allows for a more flexible decision boundary relative to approaches from the literature. Our proposed loss is contrasted against other methods from the literature by using a plethora of deep architectures. A first conclusion is the benefit of using non-parametric ordinal losses against parametric losses in cervical cancer risk prediction. Additionally, the proposed loss is found to be the top-performer in several cases. The best performing model scores an accuracy of 75.6% for 7 classes and 81.3% for 4 classes.

# Ordinal Losses for Classification of Cervical Cancer Risk

Tomé Albuquerque<sup>1,2</sup>, Ricardo Cruz<sup>1,2</sup>, and Jaime S. Cardoso<sup>1,2</sup>

<sup>1</sup>INESC TEC, Porto, Portugal

<sup>2</sup>Faculty of Engineering, University of Porto, Porto, Portugal

Corresponding author:

Tomé Albuquerque<sup>1,2</sup>

Email address: tome.m.albuquerque@inesctec.pt

## ABSTRACT

Cervical cancer is the fourth leading cause of cancer-related deaths in women, especially in low to middle-income countries. Despite the outburst of recent scientific advances, there is no totally effective treatment, especially when diagnosed in an advanced stage. Screening tests, such as cytology or colposcopy, have been responsible for a strong decrease in cervical cancer deaths. Cervical cancer automatic screening via Pap smear is a highly valuable cell imaging-based detection tool, where cells must be classified as being within one a multitude of ordinal classes, ranging from abnormal to normal. Current approaches at ordinal inference for neural networks are found to not take advantage of the ordinal problem or to be too uncompromising. A non-parametric ordinal loss for neuronal networks is proposed that promotes the output probabilities to follow a unimodal distribution. This is done by imposing a set of different constraints over all pairs of consecutive labels which allows for a more flexible decision boundary relative to approaches from the literature. Our proposed loss is contrasted against other methods from the literature by using a plethora of deep architectures. A first conclusion is the benefit of using non-parametric ordinal losses against parametric losses in cervical cancer risk prediction. Additionally, the proposed loss is found to be the top-performer in several cases. The best performing model scores an accuracy of 75.6% for 7 classes and 81.3% for 4 classes.

## INTRODUCTION

The survival rate for women with cervical cancer is disturbing – in the USA, the 5-year survival rate for all women with cervical cancer is just 66% and is responsible for around 10 deaths per week in women aged 20 to 39 years (Siegel et al., 2020). The main factor for the high mortality rate is the asymptomatic characteristic of cervical cancer in its initial stages which justifies the need for early diagnosis. Screening tests have been responsible for a strong decrease in cervical cancer deaths. The screening programs are implemented in most developed countries and the process includes Human papillomavirus (HPV) test, cytology test (or Pap smear), colposcopy, and biopsy (WHO, 2019). HPV is a group of viruses known to influence the risk of cervical cancer – some types of HPV viruses produce dysplastic changes in cells that can progressively lead to the development of cancer (WHO, 2019).

A cervical cytology test is used to detect potentially abnormal cells from the uterine cervix. These premalignant dysplastic changes of cells are classified in progressive stages: 7 stages by the World Health Organization classification (WHO) system or 4 stages by The Bethesda classification system (TBS) (DeMay, 2007).

The risk of developing cancer is especially pronounced for the later stages. Therefore, distinguishing between the stages can be crucial for diagnosis. Yet, most of the literature focuses on binary classification (normal or abnormal), ignoring the fine-grained classification of cervical cells into different stages.

The classification of observations into naturally ordered classes, as the stages of the premalignant dysplastic changes, are traditionally handled by conventional methods intended for the classification of nominal classes where the order relation is ignored. This paper introduces a new machine learning paradigm intended for multi-class classification problems where the classes are ordered. A non-parametric loss for ordinal data classification is proposed whose goal is to promote unimodality in the prediction

**Table 1.** The 7 Different Pap Smear Classes in the Herlev dataset.

	WHO	TBS	Type of cell	Quantity
Normal	1	1	Superficial squamous epithelial	74 cells
	2	1	Intermediate squamous epithelial	70 cells
	3	1	Columnar epithelial	98 cells
Abnormal	4	2	Mild squamous non-keratinizing dysplasia	182 cells
	5	3	Moderate squamous non-keratinizing dysplasia	146 cells
	6	3	Severe squamous non-keratinizing dysplasia	197 cells
	7	4	Squamous cell carcinoma in situ intermediate	150 cells

distributions produced by the neural network; e.g., it would be inconsistent to predict that stage 1 and stage 3 are both more likely than stage 2. Yet, this loss is more flexible than other losses from the literature which force a binomial distribution in the output (Costa and Cardoso, 2005). This loss is also contrasted with the standard cross-entropy loss and networks that predict classes in the form of an ordinal encoding (Cheng et al., 2008). The Herlev dataset, which comprises 917 images of individual cervical cells in different stages of the disease, is used in the experiments (Jantzen and Dounias, 2006) together with a plethora of CNN architectures.

In the next section, the problem and dataset at hand are presented. Other work for Pap smear cell classification is then reviewed in the “Related Work” section. The proposed loss is elaborated on the “Proposal” section, and the experimental details are described in “Experiments” with results and discussion presented in “Results”. The study finished with a “Conclusion” section.

## BACKGROUND

According to the WHO classification system, there are seven different types of Pap smear cells in cervical cancer progression. This system rules the existence of three different types of normal cells and four different types of abnormal cells. From suspicious cells to carcinoma in situ (CIS), the premalignant dysplastic changes of cells can also include four stages, which are mild, moderate, severe dysplasia, and carcinoma in situ (Suhrlund, 2000). However, nowadays the most used classification system is the TBS classification system, which is widely accepted among the medical society. According to the TBS system, the Pap smear cells can be divided into four classes: normal, Low-grade Squamous Intraepithelial Lesion (LSIL), High-grade Squamous Intraepithelial Lesion (HSIL) and Carcinoma in situ (Nayar and Wilbur, 2015).

The different stages of cervical cytology abnormalities are associated with different morphological changes in the cells including the cytoplasm and nucleus. However, the small visual differences between some stages of cervical cells make the construction of a multi-class autonomous classification system a true challenge.

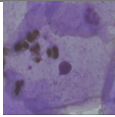
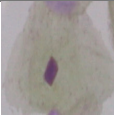
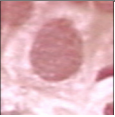
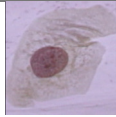
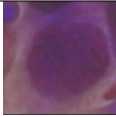
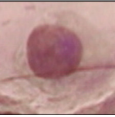
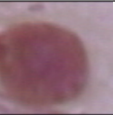
The dataset used in this work was the Herlev Dataset which is a publicly available dataset<sup>1</sup> collected at the Herlev University Hospital (Denmark) using a digital camera and microscope with an image resolution of 0.201  $\mu\text{m}$  per pixel (Jantzen and Dounias, 2006). The preparation of the specimens followed the traditional Pap smear and Pap staining. To amplify the certainty of diagnosis, two cyto-technicians and a doctor characterized the cervical images in Herlev dataset into seven classes. The Herlev dataset is composed of a total of 917 images of individual cervical cells. Each image contains ground truth segmentation and classification label. Table 1 shows the nomenclature of the 7 different classes from the dataset, wherein classes 1–3 correspond to types of normal cells and classes 4–7 to different levels of abnormal cells. Illustrations of these classes are then displayed in Table 2.

In most cases, the abnormal cells present a nucleus size bigger than healthy cells. However, the difference between normal columnar nucleus and severe and/or carcinoma nucleus is not easy to differentiate, which makes the classification between these different types of cells a challenge.

There is some imbalance in the class distribution of the dataset: 8%, 7%, 11%, 19%, 16%, 22%, and 17%, whereas 14% would be expected if the distribution was uniform.

<sup>1</sup><http://mde-lab.aegean.gr/index.php/downloads>

**Table 2.** Image examples of the 7 Different Pap Smear Classes in the Herlev dataset.

Normal			Abnormal			
						
$k = 1$	$k = 2$	$k = 3$	$k = 4$	$k = 5$	$k = 6$	$k = 7$
TBSWHO			$k = 2$	$k = 3$		$k = 4$

## RELATED WORK

In most literature, the classification of Pap smear images consists in a binary separation between normal and abnormal cell (two classes), using different methodologies such as: Support Vector Machines (SVM) (Chen et al., 2014; Chankong et al., 2014; Kashyap et al., 2016; Bora et al., 2017),  $k$ -Nearest Neighbours (kNN) (Chankong et al., 2014; Bora et al., 2017; Marinakis et al., 2009; Fekri Ershad, 2019), Fuzzy  $c$ -Means Algorithm (FCM) (Chankong et al., 2014; William et al., 2019),  $k$ -Means clustering (Paul et al., 2015), Artificial Neural Networks (ANN) (Chankong et al., 2014), and, more recently, Convolutional Neural Networks (CNN) (Zhang et al., 2017; Lin et al., 2019; Kurnianingsih et al., 2019).

However, all this work consists of binary classification, which is useful for screening, but not enough for a confident diagnosis. Fewer works explore the multi-class classification of cervical cells on Herlev dataset:

Chankong et al. (2014) proposed a multi-class automatic cervical cancer cell classification system using different classifiers, such as FCM, ANN, and kNN. However, this system is based only on 9 cell-based features. The Chankong's approach applies feature extraction from the nucleus and cytoplasm in each image and requires manual selection of the best threshold to minimize the error when applying the FCM to construct the cell mask. More recently, Kurnianingsih et al. (2019) perform feature extraction in a more autonomous way using a CNN. The use of a CNN simplifies the pre-processing steps that were necessary for the previous Chankong's approach. Ghoneim et al. (2019) proposed a new approach for multi-class cervical cancer cell detection and classification, using in the first step CNNs to extract deep-learned features and in second step, extreme learning machine (ELM)-based classifiers to classify the input cell images. Lin et al. (2019) proposed a new CNN-based method that combines cell image appearance with cell morphology for multi-class classification of cervical cells in the Herlev dataset. In all these cases, cross-entropy is adopted for ordinal data classification.

Assume that examples in a classification problem come from one of  $K$  classes, labeled from  $\mathcal{C}^{(1)}$  to  $\mathcal{C}^{(K)}$ , corresponding to their natural order in ordinal classes, and arbitrarily for nominal classes.

**Cross-Entropy (CE):** Traditionally, a CNN would perform multi-class classification by minimizing cross-entropy, averaged over the training set,

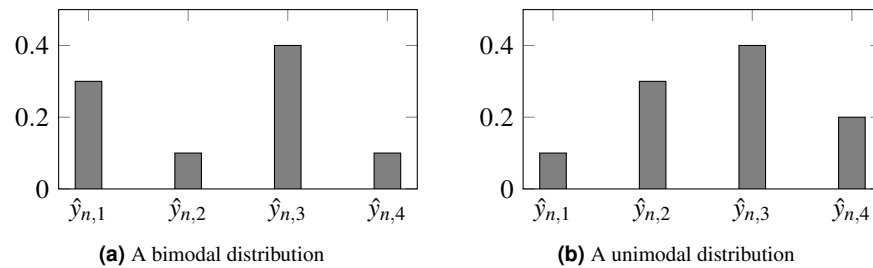
$$CE(\mathbf{y}_n, \hat{\mathbf{y}}_n) = - \sum_{k=1}^K y_{nk} \log(\hat{y}_{nk}),$$

where  $\mathbf{y}_n = [y_{n1} \cdots y_{nk} \cdots y_{nK}] \in R^K$  represents the one-hot encoding of the class of the  $n$ -th observation and  $\hat{\mathbf{y}}_n = [\hat{y}_{n1} \cdots \hat{y}_{nk} \cdots \hat{y}_{nK}] \in R^K$  is the output probability vector given by the neural network for observation  $n$ . Note that  $y_{nk} \in \{0, 1\}$ ,  $\hat{y}_{nk} \in [0, 1]$  and  $\sum_{k=1}^K y_{nk} = \sum_{k=1}^K \hat{y}_{nk} = 1$ .

However, CE has limitations when applied to ordinal data. Defining  $k_n^* \in \{1, \dots, K\}$  as the index of the true class of observation  $\mathbf{x}_n$  (the position where  $y_{nk} = 1$ ), it is then clear that

$$CE(\mathbf{y}_n, \hat{\mathbf{y}}_n) = - \log(\hat{y}_{nk_n^*}).$$

Intuitively, CE is just trying to maximize the probability in the output corresponding to the true class, ignoring all the other probabilities. For this loss, an error between classes  $\mathcal{C}^{(1)}$  and  $\mathcal{C}^{(2)}$  is treated as the same as an error between  $\mathcal{C}^{(1)}$  and  $\mathcal{C}^{(K)}$ , which is undesirable for ordinal problems.



**Figure 1.** Probabilities produced by two different models for observation  $n$ . CE is unable to distinguish both scenarios, setting the same loss for both. For ordinal problems, a unimodal distribution, peaking in the true class, is, arguably, preferable. In this example,  $k_n^* = 3$  is the assumed true class.

Furthermore, the loss does not constrain the model to produce unimodal probabilities, so inconsistencies can be produced such as  $\hat{y}_{nj} > \hat{y}_{n\ell} < \hat{y}_{ni}$ , even when  $1 \leq j < \ell < i \leq K$ . It would be preferable for output probabilities to follow a unimodal distribution, as depicted by Figure 1.

Cross-entropy is a fair approach for nominal data, where no additional information is available. However, for ordinal data, the order can be explored to further regularize learning.

**Ordinal Encoding (OE):** A model agnostic way to introduce ordinality is by training binary classifiers, in the form of an ensemble, where each classifier tries to distinguish between each pair of adjacent classes,  $\mathcal{C}^{(i)}$  and  $\mathcal{C}^{(i+1)}$  (Frank and Hall, 2001). An adaptation for neural networks consists of training a single neural network to produce  $K - 1$  outputs, where each output makes a binary decision between each pair of adjacent classes. The information on the ordinal distribution can, therefore, be encoded in the  $\mathbf{y}$  labels themselves (Cheng et al., 2008).

In traditional one-hot encoding, classes are encoded using the indicator function  $\mathbb{1}(k = k^*)$ , so that  $y_{nm}$  is represented by 1 if  $k = k_n^*$  and 0 otherwise. In ordinal encoding, classes are encoded using a cumulative distribution – the indicator function used is  $\mathbb{1}(k < k^*)$  so that  $y_{nm}$  is represented by 1 if  $k < k_n^*$  and 0 otherwise. Each output represents the incremental neighbor probability, and the inverse operation (during inference) is performed by summing up these outputs,  $p_{nk} = \sum_{m=1}^{K-1} y_{nm}$ .

**Binomial Unimodal (BU):** A different approach is to constrain the output of the network directly, approaching the problem under a regression setting. Instead of several outputs, the output predicts a single output representing the probability along the classes, with  $y_n = 0$  representing  $k_n^* = 1$  and  $y_n = 1$  representing  $k_n^* = K$  (Costa and Cardoso, 2005). Thus, this model has only one output unit as the final layer. The model's sigmoid output is converted into class probabilities using Binomial's probability mass function.

The goal of this approach is to maintain the ordinality of the classes by applying a parametric model for the output probabilities.

**Poisson Unimodal (PU):** The Poisson probability mass function (PMF) is used to enforce a discrete unimodal probability distribution (Beckham and Pal, 2017). As a final layer, the log Poisson PMF transform is applied together with a softmax to normalize the output as a probability distribution.

These parametric approaches sometimes sacrifice accuracy to ensure the ordinality assumption. This sacrifice might sometimes prove too much, especially given the fact that modern deep learning datasets are very big and have a significant number of mislabels. A loss is now proposed to stimulate a unimodal output without modifying the network architecture.

## PROPOSAL

As already explored, CE presents drawbacks when applied to ordinal data. By focusing only on the mode of the distribution and ignoring all the other values in the output probability vector, one is not leveraging the ordinal information intrinsic to the data.

### Fixing CE with an Ordinal Loss Term

A possible fix for CE is to a regularization term that penalizes the deviations from the unimodal setting. Defining  $\mathbb{1}(x)$  as the indicator function of  $x$  and  $\text{ReLU}(x) = x\mathbb{1}(x > 0) = \max(0, x)$ , a tentative solution

for an order-aware loss could be

$$\begin{aligned} \text{CO}(\mathbf{y}_n, \hat{\mathbf{y}}_n) &= \text{CE}(\mathbf{y}_n, \hat{\mathbf{y}}_n) \\ &+ \lambda \sum_{k=1}^{K-1} \mathbb{1}(k \geq k_n^*) \text{ReLU}(\hat{y}_{n(k+1)} - \hat{y}_{n(k)}) + \lambda \sum_{k=1}^{K-1} \mathbb{1}(k \leq k_n^*) \text{ReLU}(\hat{y}_{n(k)} - \hat{y}_{n(k+1)}), \quad (1) \end{aligned}$$

where  $\lambda \geq 0$  is controlling the relative importance of the extra terms favoring unimodal distributions. Predicted probability values are expected to decrease monotonously as we depart left and right from the true class. The added terms penalize any deviation from this expected unimodal distribution, with a penalty proportional to the difference of the consecutive probabilities. The additional terms, although promoting uni-modality, still allow flat distributions. A generalization of the previous idea is to add a margin of  $\delta > 0$  to the ReLU, imposing that the difference between consecutive probabilities is at least  $\delta$ . This leads us to a second CE loss, CO2, suitable for ordinal classes:

$$\begin{aligned} \text{CO2}(\mathbf{y}_n, \hat{\mathbf{y}}_n) &= \text{CE}(\mathbf{y}_n, \hat{\mathbf{y}}_n) \\ &+ \lambda \sum_{k=1}^{K-1} \mathbb{1}(k \geq k_n^*) \text{ReLU}(\delta + \hat{y}_{n(k+1)} - \hat{y}_{n(k)}) + \lambda \sum_{k=1}^{K-1} \mathbb{1}(k \leq k_n^*) \text{ReLU}(\delta + \hat{y}_{n(k)} - \hat{y}_{n(k+1)}). \quad (2) \end{aligned}$$

A value of  $\delta = 0.05$  has been empirically found to provide a sensible margin. This loss is aligned with the proposal present in Belharbi et al. (2019).

### Beyond CO2: Ordinal Entropy Loss Function

In CO2, the CE term by itself is only trying to maximize the probability estimated in the true output class (while ignoring the remaining probabilities); the ordinal terms are promoting unimodality but not penalizing (almost) flat distributions. This also explains why the ordinal terms by themselves (especially the version without margin) are not enough to promote strong learning: the model could converge to solutions where the predicted probability in the true class is only slightly above the neighboring probabilities, which will not, most likely, provide a strong generalization for new observations.

However, the extreme nature of CE, ignoring almost everything in the predicted distribution  $\hat{\mathbf{y}}_n$  is equivalent to assume that the perfect probability distribution is one on the true class and zero everywhere else. This assumes a strong belief and dependence on the chosen one-hot encoding, which is often a crude approximation to the true probability class distribution. Seldom, for a fixed observation  $\mathbf{x}_n$ , the class is deterministically known; rather, we expect a class distribution with a few non-zero values. This is particularly true for observations close to the boundaries between classes. A softer assumption is that the distribution should have a low entropy, only.

This leads us to propose the ordinal entropy loss, HO2, for ordinal data as

$$\begin{aligned} \text{HO2}(\mathbf{y}_n, \hat{\mathbf{y}}_n) &= H(\hat{\mathbf{y}}_n) \\ &+ \lambda \sum_{k=1}^{K-1} \mathbb{1}(k \geq k_n^*) \text{ReLU}(\delta + \hat{y}_{n(k+1)} - \hat{y}_{n(k)}) + \lambda \sum_{k=1}^{K-1} \mathbb{1}(k \leq k_n^*) \text{ReLU}(\delta + \hat{y}_{n(k)} - \hat{y}_{n(k+1)}), \quad (3) \end{aligned}$$

where  $H(\mathbf{p})$  denotes the entropy of the distribution  $\mathbf{p}$ .

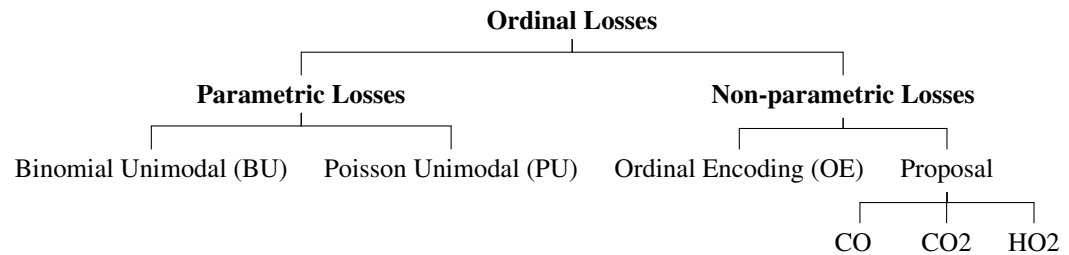
## EXPERIMENTS

Several neural network architectures are now trained using the aforementioned losses for the dataset at hand. In this work it was also evaluated the performance differences between parametric and non-parametric losses for ordinal classification (Figure 2). All the experiments were implemented in PyTorch and are available online<sup>2</sup>.

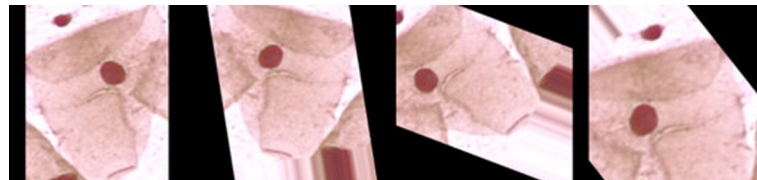
### Data Pre-processing

Given that all images from the Herlev dataset were of different sizes, all images were resized to  $224 \times 224$  pixels; however, before the resize of cytological images a zero-padding must be done to avoid the

<sup>2</sup><https://github.com/tomealbuquerque/ordinal-losses>



**Figure 2.** Schematic representation of the used and proposed ordinal losses.



**Figure 3.** Examples of data augmentation on the Herlev database. The original zero-padding image (left) and random transformations.

loss of essential information regarding cells shape. The last preprocessing step was to apply the same normalization as used by ImageNet (Simonyan and Zisserman, 2014).

Since the Herlev database has a relatively small number of observations (917), the dataset was augmented by a series of random transformations: 10% of width and height shift, 10% of zoom, image rotation, horizontal and vertical flips, and color saturation. These transformations are illustrated in Figure 3.

### Convolutional Neural Networks

A convolutional neural network (CNN) is a neural network that successively applies convolutions of filters to the image. These filters are learned and consist of quadrilateral patches that are convolved across the whole input image – unlike previous fully-connected networks, only local inputs are connected at each layer. Typically, each convolution is intertwined with downsampling operations, such as max-pooling, that successively reduces the size of the original image. In the end, a series of outputs produce the desired classification.

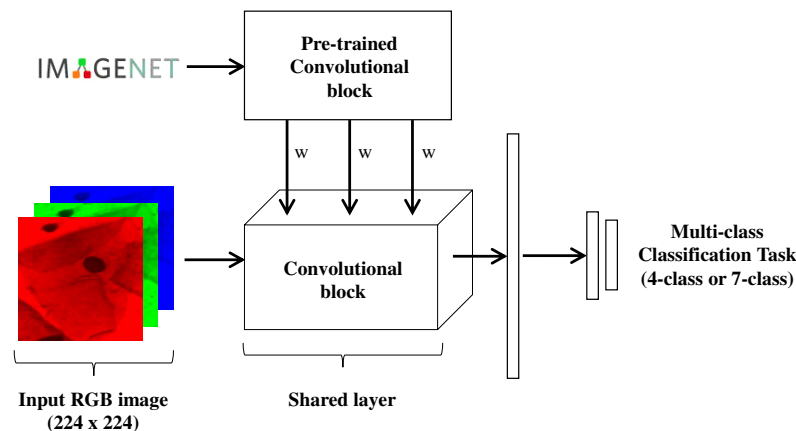
The final layers are fully-connected and then the final output is processed by a soft-max for multiclass problems or a logistic function for binary classification. Dropout was used to reduce overfitting by constraining these fully-connected layers.

### Network Architectures

Two different models were trained and tested in this work for multi-class (4-class and 7-class) classification of Pap smear cells images (Figure 4). Both models were trained and tested with eight different convolutional network architectures: AlexNet (Krizhevsky et al., 2012), GoogLeNet (Szegedy et al., 2015), MobileNet\_V2 (Howard et al., 2017), ResNet18 (He et al., 2016), ResNeXt50\_32X4D (Xie et al., 2017), ShuffleNet\_V2\_X1\_0 (Zhang et al., 2018), SqueezeNet1\_0 (Iandola et al., 2016), VGG-16 (Simonyan and Zisserman, 2014), and Wide\_ResNet50.2 (Zagoruyko and Komodakis, 2016). The goal in testing these different architectures is to evaluate how well the proposed loss behaves in a wide range of architectures. These 9 different architectures were chosen as they are often used in the literature and came pre-trained with PyTorch on ImageNet<sup>3</sup>. The last block of each architecture was replaced by these layers: dropout with  $p=20\%$ , 512 dense layer with ReLU, dropout with  $p=20\%$ , 256 dense layer with ReLU, followed by  $K$  neuron outputs.

A brief introduction of each architecture is now presented. AlexNet, based on LeNet, formalized the Convolutional Neural Network (CNN) as is known today: a series of convolutions intertwined by downsampling blocks. Max-pooling was used for downsampling and ReLU was used as the activation

<sup>3</sup><https://pytorch.org/docs/stable/torchvision/models.html>



**Figure 4.** Schematic representation of the model used for multi-class classification of Pap smear cells.

function. It became famous for winning ImageNet, the first CNN to do so (Krizhevsky et al., 2012). The following ImageNet competitions were also won by other CNNs – VGG and GoogLeNet – which were evolutions on top of AlexNet which consisted mostly in an explodingly higher number of parameters (Simonyan and Zisserman, 2014; Szegedy et al., 2015). Then, MobileNet (Howard et al., 2017) introduced hyperparameters to help the user choose between latency and accuracy trade-offs. An attempt was then made at curbing the number of parameters with ShuffleNet (Zhang et al., 2018) by approximating convolution operators using fewer parameters.

Finally, an attempt was made at curbing the number of parameters, which had been exploding, while keeping the accuracy of these early CNNs with SqueezeNet (Iandola et al., 2016).

In another line of research, ResNet (He et al., 2016) introduced residual blocks whose goal was to make the optimization process easier for gradient descent. Each residual block learns  $a = f(x) + x$  instead of  $a = f(x)$ . Given that weights are initialized randomly around zero and most activation functions are also centered in zero (an exception would be the logistic activation function), then, in expectation, all neurons output zero before any training. Therefore, when using residual blocks, at time=0,  $a = x$ , i.e. activations produce the identity function. This greatly helps gradient descent focus on finding improvements (residuals) on top of the identity function. While this model allowed for deeper neural networks, each percent of improved accuracy cost nearly doubling the number of layers, which motivated WideResNet (Zagoruyko and Komodakis, 2016) and ResNeXt (Xie et al., 2017) to improve the residual architecture to improve learning time.

## Training

The weights of the architectures mentioned in the previous section are already initialized by pre-training on ImageNet. Adam was used as the optimizer and starts with a learning rate of  $10^{-4}$ . The learning rate is reduced by 10% whenever the validation loss is stagnant for 10 epochs. The training process is completed after 100 epochs, and the best model is saved as judged by the lowest validation loss. The dataset was divided into 10 different k-folds using stratified cross-validation, maintaining the ratio among the different classes.

The loss hyperparameter  $\lambda$  was chosen using nested k-fold cross-validating ( $K=5$ ) for each architecture, in order to reduce the bias in combined hyperparameter tuning and model selection.

## Evaluation Metrics

The most popular classification metric is accuracy (Acc). For  $N$  observations, taking  $k_i$  and  $\hat{k}_i$  to be the label and prediction of the  $n$ -th observation, respectively, then  $\text{Acc} = \frac{1}{N} \sum_{n=1}^N \mathbb{1}(\hat{k}_n^* = k_n^*)$ , where  $\mathbb{1}$  is the indicator function.

However, this metric treats all class errors as the same, whether the error is between adjacent classes or between classes in the extreme. If we have  $K$  classes represented by a set  $\mathcal{C} = \{\mathcal{C}^{(1)}, \mathcal{C}^{(2)}, \dots, \mathcal{C}^{(K)}\}$ , then accuracy will treat an error between  $\mathcal{C}^{(1)}$  and  $\mathcal{C}^{(2)}$  with the same magnitude as an error between  $\mathcal{C}^{(1)}$  and  $\mathcal{C}^{(K)}$  which is clearly worse. As an illustration, in a medical setting, a misdiagnosis between



Stage II and Stage III of a disease, while bad, is not as bad as a misdiagnosis between Healthy and Stage III. For that reason, a popular metric for ordinal classification is the Mean Absolute Error (MAE),  $MAE = \frac{1}{N} \sum_i |k_i^* - \hat{k}_i^*|$ . This metric is not perfect since it treats an ordinal variable as a cardinal variable. An error between classes  $\mathcal{C}^{(1)}$  and  $\mathcal{C}^{(3)}$  will be treated as two times worse than an error between classes  $\mathcal{C}^{(1)}$  and  $\mathcal{C}^{(2)}$ . Naturally, the assumption of ordinality is not always warranted.

To evaluate the models' performance we also used a specific metric for ordinal classification, Uniform Ordinal Classification Index (UOC) which takes into account accuracy and ranking in the performance assessment and it is also robust against imbalanced classes (Silva et al., 2018). The better the performance, the lower the UOC.

By combining a quality assessment (accuracy) with a quantity assessment (MAE) and also with a specific metric for ordinality (UOC) we hope to provide a balanced view of the performance of the methods.

The two other metrics used are the AUC of ROC or AUROC (Area Under the Receiver Operating Characteristic) and Kendall's  $\tau$  rank correlation coefficient. AUROC measures how well-calibrated are the probabilities produced by the model. This first metric is used in the binary classification context (two classes) and are extended for multi-class by comparing each class against the rest (one vs rest strategy) and performing an overall average, known as macro averaging. On the other hand, Kendall's Tau is a non-parametric evaluation of relationships between columns of ranked data, so it is a measure of ordinal association between data. The  $\tau$  correlation coefficient returns a value that ranges from -1 to 1, with 0 being no correlation and 1 perfect correlation.

## RESULTS

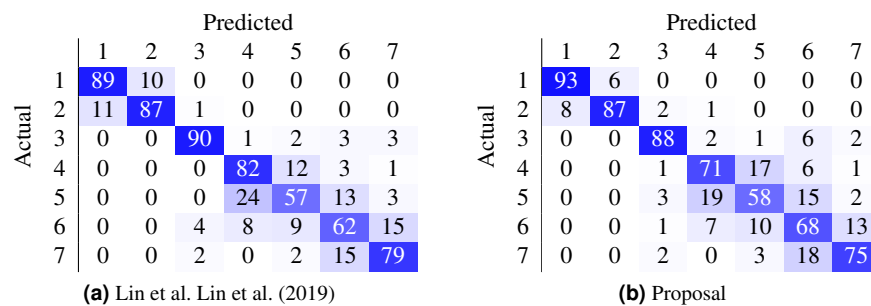
The average performance for the 10-folds of 9 different architectures are presented in Tables 3–8, for both the 7-class and 4-class classification problems, with the seven different learning losses – conventional Cross-Entropy (CE), Binomial Unimodal (BU) (Costa and Cardoso, 2005), Poisson Unimodal (PU) (Beckham and Pal, 2017), Ordinal Encoding (OE) (Cheng et al., 2008) and our proposed losses (CO, CO2 and HO2), as measured by MAE, accuracy, UOC index and Kendall's coefficient (Appendix - table A1–A2) detailed in the previous section. In all the Tables of results, the best results for each loss are in bold. Furthermore, the other results are compared to the best result and a statistical test is used with the hypothesis of them being the same or not. A  $p$ -value of 0.1 is used with a one-sided paired  $t$ -test due to the small sample size (10 folds).

For the 7-class classification problem, Table 3 shows the results for MAE, which confirm the influence of ordinal losses in promoting ordinality when comparing to nominal loss (CE). OE loss achieved the best performance across the different architectures but it is also notable the good performance of our loss, which in 67% of cases, the models trained with our proposed loss provide better MAE results. The MAE results present in Table 4 for 4-class classification are consistent with the 7-class Table 3, with ordinal losses winning over nominal CE.

In Table 5 and Table 6 are presented the accuracy results for 7-class and 4-class classification problems respectively. Regarding this metric, the results between nominal and ordinal losses are more balanced. CE loss performance is above ordinal losses in 11% for 7-class problem and is tied for 4-class problem. This can be explained by the lower role of ordinality in the CE loss, as also confirmed by the MAE results. This means that when misclassification occurs, ordinal losses tend to classify Pap smear images as being closer to the real class. Results for UOC index (Table 7–8) are also consistent with MAE metric, with 78% of the models presenting a lowest UOC index when using the ordinal losses. Table A1 and Table A2 in appendix represents the results for Kendall's  $\tau$  coefficient test in 4-class and 7-class classification problems. This results are also aligned with the results of MAE and UOC metrics, being the ordinal losses in advantage when comparing with nominal CE.

Adding the margin (CO  $\rightarrow$  CO2) influences positively most of metrics for 7 and 4 classes. Using entropy (CO2 or HO2), instead of cross-entropy, promotes better results on the metrics intrinsically connected with ordinality (MAE, UOC and Kendall's  $\tau$  coefficient).

The average results for all losses across the 9 different architectures for MAE, accuracy, UOC, AUROC, Kendall's  $\tau$  coefficient and Gini index metrics are present in appendix Table A3 and Table A4 for 4 and 7 class classification respectively. In both Tables are represented the results using the classical mode (softmax) to aggregate the probabilities and also using mean (expectation trick) Beckham and Pal (2017). Concerning the sparsity of the prediction probabilities, as measured by the Gini index, it is notable



**Figure 5.** Comparison of state-of-the-art confusion matrix (7 classes) against WideResNet50 trained using the HO2 loss.

that as the loss is made more ordinal-aware then the predicted probabilities tend to be more spread across the classes. This could already be seen in Figure 6. Interestingly, the OE distribution is almost identical to the CE distribution and has been omitted from the figure for eligibility.

In average, in most metrics, non-parametric losses outperformed parametric losses. This difference can be justified with the greater flexibility in boundary decisions provided by non-parametric losses. OE, CO2 and HO2 provided better results across the different metrics when comparing to BU and PU.

Most work from the literature concerns the binary case using Herlev dataset (normal vs abnormal); only a couple concern themselves with the 7-class and 4-class ordinal classification problem. Table 9 contrasts the best performing models from two recent works against the proposed method. In our case, the non-parametric loss (CO2) was able to beat the state-of-the-art by 11.1% (7 classes) and by 10% (4 classes) in accuracy metric. Furthermore, the confusion matrix from Figure 5 contrasts the proposal against Lin et al. (2019).

There are classes of cells easier to be classified than others, as shown by the confusion matrix in Figure 5 (b). Columnar cells are sometimes inappropriately classified as severe dysplasia cells since severe dysplasia cells have similar characteristics in appearance and morphology with columnar cells (e.g., small cytoplasm, dark nuclei).

The main challenge occurs in the classification of abnormal cells (i.e., mild dysplasia, moderate dysplasia, severe dysplasia, and carcinoma) where the characteristics of these kinds of cells are very similar. The fact is that the abnormal classes correspond to different levels of evolution of structures, with a progressive change in its characteristics which leads them to present characteristics common to two levels, being a hard task even for cytopathologists to classify them correctly. Thus, the right multiclass classification of abnormal cells is highly desirable, and with substantial clinical value.

Finally, the influence of the losses on the output probabilities is illustrated in Figure 6 when predicting two classes for the 7-class case. Contrasting this to Figure 1, it is clear that the proposed loss tends to promote a unimodal distribution of probabilities relative to the CE loss, which tends to maximize the probability in the output corresponding to the true class and ignoring all the other probabilities distribution,

**Table 3.** Results in terms of Mean Absolute Error (MAE) for 7 class problem, averaged for 10 folds. Lower is better

	CE	BU	PU	OE	CO	CO2	HO2
<b>alexnet</b>	0.46 ± 0.08	0.52 ± 0.09	0.50 ± 0.09	<b>0.44 ± 0.08</b>	0.90 ± 0.19	<b>0.41 ± 0.08</b>	0.45 ± 0.10
<b>googlenet</b>	<b>0.39 ± 0.05</b>	0.41 ± 0.07	0.42 ± 0.08	<b>0.38 ± 0.09</b>	0.53 ± 0.10	<b>0.37 ± 0.07</b>	<b>0.36 ± 0.06</b>
<b>mobilenet_v2</b>	0.34 ± 0.05	0.36 ± 0.04	<b>0.31 ± 0.04</b>	<b>0.33 ± 0.05</b>	0.52 ± 0.26	0.34 ± 0.06	<b>0.34 ± 0.05</b>
<b>resnet18</b>	<b>0.34 ± 0.09</b>	<b>0.36 ± 0.06</b>	<b>0.35 ± 0.06</b>	<b>0.35 ± 0.10</b>	0.49 ± 0.11	<b>0.34 ± 0.07</b>	<b>0.35 ± 0.10</b>
<b>resnext50_32x4d</b>	0.34 ± 0.07	<b>0.33 ± 0.05</b>	<b>0.33 ± 0.03</b>	0.34 ± 0.06	0.41 ± 0.08	<b>0.33 ± 0.06</b>	<b>0.31 ± 0.07</b>
<b>shufflenet_v2_x1.0</b>	0.41 ± 0.07	0.49 ± 0.07	<b>0.41 ± 0.05</b>	<b>0.38 ± 0.07</b>	0.47 ± 0.08	<b>0.40 ± 0.05</b>	<b>0.38 ± 0.06</b>
<b>squeezenet1.0</b>	<b>0.38 ± 0.07</b>	0.45 ± 0.05	0.46 ± 0.07	<b>0.40 ± 0.09</b>	0.97 ± 0.31	0.41 ± 0.08	0.45 ± 0.09
<b>vgg16</b>	<b>0.37 ± 0.09</b>	0.44 ± 0.05	0.44 ± 0.10	<b>0.37 ± 0.06</b>	0.67 ± 0.15	<b>0.36 ± 0.06</b>	<b>0.36 ± 0.07</b>
<b>wide_resnet50_2</b>	0.33 ± 0.06	0.37 ± 0.05	0.32 ± 0.06	<b>0.30 ± 0.04</b>	0.45 ± 0.13	0.33 ± 0.06	0.35 ± 0.09
<b>Avg</b>	0.37	0.41	0.39	0.36	0.60	0.37	0.37
<b>Winners</b>	4	2	4	8	0	6	6

**Table 4.** Results in terms of **Mean Absolute Error (MAE)** for **4 class** problem, averaged for 10 folds. Lower is better

	CE	BU	PU	OE	CO	CO2	HO2
<b>alexnet</b>	0.31 ± 0.06	0.32 ± 0.04	<b>0.28 ± 0.04</b>	<b>0.29 ± 0.06</b>	0.47 ± 0.19	<b>0.29 ± 0.05</b>	<b>0.31 ± 0.06</b>
<b>googlenet</b>	<b>0.24 ± 0.04</b>	0.25 ± 0.03	0.25 ± 0.05	<b>0.24 ± 0.05</b>	0.38 ± 0.17	<b>0.22 ± 0.05</b>	0.25 ± 0.06
<b>mobilenet_v2</b>	<b>0.22 ± 0.06</b>	<b>0.21 ± 0.03</b>	0.24 ± 0.05	<b>0.22 ± 0.06</b>	<b>0.23 ± 0.04</b>	0.24 ± 0.05	<b>0.22 ± 0.05</b>
<b>resnet18</b>	0.24 ± 0.03	0.26 ± 0.05	<b>0.24 ± 0.05</b>	<b>0.22 ± 0.04</b>	0.29 ± 0.11	<b>0.22 ± 0.04</b>	0.26 ± 0.06
<b>resnext50_32x4d</b>	<b>0.21 ± 0.03</b>	0.22 ± 0.04	0.23 ± 0.03	<b>0.20 ± 0.04</b>	0.28 ± 0.07	<b>0.21 ± 0.03</b>	0.22 ± 0.05
<b>shufflenet_v2_x1.0</b>	<b>0.28 ± 0.05</b>	0.33 ± 0.05	<b>0.27 ± 0.05</b>	0.31 ± 0.06	0.36 ± 0.09	<b>0.28 ± 0.06</b>	<b>0.28 ± 0.04</b>
<b>squeezenet1.0</b>	<b>0.28 ± 0.06</b>	0.30 ± 0.05	0.30 ± 0.06	<b>0.27 ± 0.07</b>	0.66 ± 0.17	<b>0.29 ± 0.04</b>	0.31 ± 0.05
<b>vgg16</b>	0.27 ± 0.06	0.28 ± 0.06	<b>0.26 ± 0.05</b>	<b>0.24 ± 0.03</b>	0.53 ± 0.18	0.26 ± 0.05	0.27 ± 0.05
<b>wide_resnet50_2</b>	0.23 ± 0.05	0.22 ± 0.04	<b>0.20 ± 0.06</b>	<b>0.22 ± 0.05</b>	0.43 ± 0.22	<b>0.21 ± 0.05</b>	0.22 ± 0.03
<b>Avg</b>	0.25	0.27	0.25	0.24	0.40	0.25	0.26
<b>Winners</b>	5	1	5	8	1	7	3

**Table 5.** Results in terms of **Accuracy** for **7 class** problem, averaged for 10 folds. Higher is better

	CE	BU	PU	OE	CO	CO2	HO2
<b>alexnet</b>	<b>71.1 ± 5.1</b>	60.6 ± 3.7	64.8 ± 5.4	<b>70.1 ± 5.1</b>	44.2 ± 7.6	<b>70.8 ± 5.1</b>	67.9 ± 5.4
<b>googlenet</b>	<b>72.5 ± 3.7</b>	66.1 ± 4.3	68.5 ± 4.5	<b>71.5 ± 5.3</b>	59.7 ± 8.2	<b>72.4 ± 4.9</b>	<b>72.4 ± 3.7</b>
<b>mobilenet_v2</b>	<b>75.0 ± 4.4</b>	69.0 ± 3.5	<b>74.2 ± 2.8</b>	<b>74.4 ± 3.8</b>	64.4 ± 16.5	73.1 ± 3.7	<b>74.1 ± 3.9</b>
<b>resnet18</b>	<b>74.4 ± 6.1</b>	69.5 ± 3.7	<b>73.3 ± 4.3</b>	<b>73.6 ± 6.4</b>	64.6 ± 6.5	<b>73.3 ± 4.5</b>	<b>73.3 ± 6.4</b>
<b>resnext50_32x4d</b>	<b>74.4 ± 3.7</b>	72.4 ± 4.3	72.8 ± 2.8	<b>74.0 ± 4.2</b>	68.0 ± 5.9	<b>75.5 ± 3.5</b>	<b>75.7 ± 5.3</b>
<b>shufflenet_v2_x1.0</b>	<b>71.9 ± 5.5</b>	61.0 ± 4.5	67.7 ± 4.6	70.7 ± 4.9	65.5 ± 4.5	<b>70.7 ± 3.1</b>	<b>71.3 ± 3.7</b>
<b>squeezenet1.0</b>	<b>73.0 ± 4.3</b>	63.3 ± 2.4	67.3 ± 3.6	<b>71.8 ± 5.3</b>	40.5 ± 13.3	70.8 ± 4.5	67.1 ± 5.0
<b>vgg16</b>	<b>73.1 ± 4.7</b>	63.9 ± 4.6	67.6 ± 6.2	<b>72.6 ± 3.8</b>	54.4 ± 8.5	<b>71.8 ± 3.3</b>	<b>72.0 ± 3.7</b>
<b>wide_resnet50_2</b>	<b>75.7 ± 3.2</b>	69.7 ± 3.1	74.5 ± 4.3	<b>76.8 ± 1.9</b>	66.1 ± 7.8	<b>75.6 ± 4.0</b>	74.3 ± 5.7
<b>Avg</b>	73.4	66.2	70.1	72.8	58.6	72.6	72.0
<b>Winners</b>	9	0	2	8	0	7	6

**Table 6.** Results in terms of **Accuracy** for **4 class** problem, averaged for 10 folds. Higher is better

	CE	BU	PU	OE	CO	CO2	HO2
<b>alexnet</b>	<b>76.1 ± 3.8</b>	72.8 ± 2.7	<b>75.7 ± 4.0</b>	<b>76.8 ± 3.6</b>	63.9 ± 12.5	<b>75.9 ± 3.5</b>	74.9 ± 3.9
<b>googlenet</b>	<b>79.9 ± 1.8</b>	<b>78.3 ± 2.6</b>	77.3 ± 3.1	<b>79.2 ± 4.0</b>	69.4 ± 12.0	<b>80.0 ± 3.8</b>	<b>78.4 ± 4.0</b>
<b>mobilenet_v2</b>	<b>81.8 ± 4.3</b>	<b>80.7 ± 2.5</b>	78.8 ± 3.4	<b>81.2 ± 4.9</b>	<b>79.8 ± 3.7</b>	79.2 ± 3.2	<b>80.8 ± 3.7</b>
<b>resnet18</b>	<b>79.8 ± 2.6</b>	77.2 ± 2.3	<b>78.5 ± 4.1</b>	<b>80.7 ± 4.1</b>	75.2 ± 8.4	<b>80.4 ± 3.8</b>	78.0 ± 4.3
<b>resnext50_32x4d</b>	<b>82.0 ± 3.1</b>	80.0 ± 3.5	79.5 ± 3.2	<b>82.3 ± 4.3</b>	76.2 ± 5.1	80.8 ± 2.8	79.9 ± 3.9
<b>shufflenet_v2_x1.0</b>	<b>77.1 ± 3.7</b>	72.1 ± 3.5	<b>76.1 ± 3.5</b>	75.0 ± 4.4	70.4 ± 6.6	<b>76.9 ± 3.9</b>	<b>76.2 ± 2.3</b>
<b>squeezenet1.0</b>	<b>77.2 ± 4.2</b>	73.5 ± 3.1	74.9 ± 5.1	<b>77.3 ± 5.3</b>	49.9 ± 12.2	75.5 ± 3.3	74.3 ± 4.5
<b>vgg16</b>	<b>77.9 ± 4.8</b>	74.4 ± 4.7	<b>77.5 ± 3.8</b>	<b>79.4 ± 2.5</b>	58.1 ± 11.8	77.0 ± 3.9	<b>77.4 ± 3.7</b>
<b>wide_resnet50_2</b>	80.8 ± 3.2	79.3 ± 3.3	<b>82.2 ± 4.2</b>	<b>81.0 ± 3.9</b>	64.0 ± 15.3	<b>81.3 ± 4.2</b>	80.6 ± 2.6
<b>Avg</b>	79.2	76.5	77.8	79.2	67.4	78.5	77.8
<b>Winners</b>	8	2	5	8	1	5	4

**Table 7.** Results in terms of **Uniform Ordinal Classification Index (UOC)** for **7 class** problem, averaged for 10 folds. Lower is better.

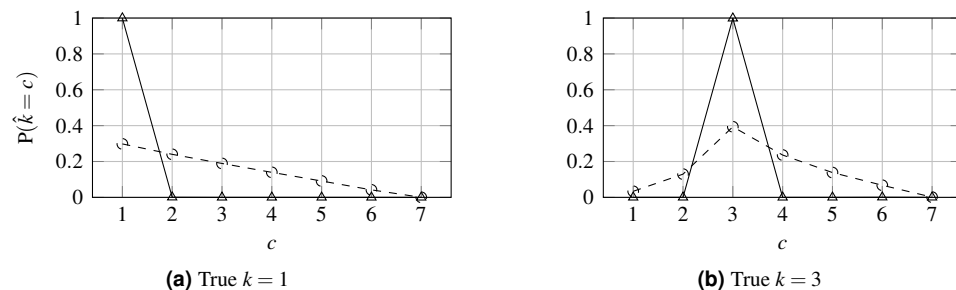
	CE	BU	PU	OE	CO	CO2	HO2
<b>alexnet</b>	<b>45.1 ± 6.5</b>	51.7 ± 5.7	49.8 ± 6.6	<b>44.0 ± 6.9</b>	70.3 ± 7.8	<b>42.8 ± 7.3</b>	46.4 ± 7.8
<b>googlenet</b>	<b>38.9 ± 6.0</b>	44.2 ± 5.7	44.6 ± 7.3	<b>39.0 ± 7.2</b>	51.3 ± 9.1	<b>38.8 ± 6.9</b>	<b>38.1 ± 4.7</b>
<b>mobilenet_v2</b>	36.0 ± 5.7	39.7 ± 4.9	<b>33.6 ± 4.5</b>	35.4 ± 5.6	46.7 ± 15.0	36.2 ± 6.4	36.2 ± 6.1
<b>resnet18</b>	<b>36.2 ± 9.3</b>	40.1 ± 5.7	<b>37.2 ± 6.3</b>	<b>37.3 ± 9.1</b>	46.9 ± 6.8	<b>37.1 ± 7.6</b>	<b>37.8 ± 8.7</b>
<b>resnext50_32x4d</b>	36.9 ± 6.8	<b>37.0 ± 5.2</b>	37.6 ± 4.6	36.8 ± 6.1	42.2 ± 6.7	<b>35.3 ± 6.7</b>	<b>34.0 ± 7.2</b>
<b>shufflenet_v2_x1.0</b>	41.8 ± 7.1	49.6 ± 6.4	43.6 ± 4.9	<b>40.3 ± 6.3</b>	46.3 ± 6.0	42.4 ± 4.1	<b>40.3 ± 4.9</b>
<b>squeezenet1.0</b>	<b>40.4 ± 6.0</b>	47.9 ± 3.8	47.5 ± 4.8	<b>42.4 ± 8.1</b>	73.6 ± 13.6	42.7 ± 7.4	46.8 ± 7.0
<b>vgg16</b>	<b>38.5 ± 8.2</b>	47.2 ± 4.9	45.5 ± 8.6	<b>39.0 ± 6.4</b>	60.3 ± 10.0	<b>40.2 ± 6.1</b>	<b>39.6 ± 6.8</b>
<b>wide_resnet50_2</b>	35.7 ± 5.2	40.8 ± 5.4	35.6 ± 6.3	<b>33.5 ± 4.5</b>	44.2 ± 9.1	<b>34.8 ± 6.5</b>	36.6 ± 8.4
<b>Avg</b>	38.8	44.2	41.7	38.6	53.5	39.0	39.5
<b>Winners</b>	5	1	2	7	0	6	5

**Table 8.** Results in terms of **Uniform Ordinal Classification Index (UOC)** for **4 class** problem, averaged for 10 folds. Lower is better.

	CE	BU	PU	OE	CO	CO2	HO2
<b>alexnet</b>	<b>38.2 ± 5.1</b>	39.5 ± 3.4	<b>37.1 ± 4.3</b>	<b>37.0 ± 4.9</b>	52.7 ± 14.2	<b>37.4 ± 5.8</b>	<b>38.9 ± 5.8</b>
<b>googlenet</b>	<b>31.6 ± 3.1</b>	<b>31.7 ± 3.6</b>	34.4 ± 5.6	<b>32.5 ± 5.7</b>	44.7 ± 14.6	<b>30.8 ± 5.5</b>	<b>32.9 ± 6.3</b>
<b>mobilenet_v2</b>	<b>30.1 ± 6.9</b>	<b>29.2 ± 3.7</b>	32.8 ± 5.2	<b>30.6 ± 7.5</b>	<b>31.0 ± 4.8</b>	32.5 ± 5.5	<b>30.5 ± 5.4</b>
<b>resnet18</b>	31.4 ± 4.6	33.1 ± 3.7	<b>32.3 ± 5.5</b>	<b>29.4 ± 6.0</b>	36.7 ± 11.0	<b>30.3 ± 4.1</b>	33.2 ± 6.7
<b>resnext50_32x4d</b>	<b>28.7 ± 4.7</b>	<b>29.8 ± 4.9</b>	32.0 ± 3.9	<b>27.5 ± 5.3</b>	35.9 ± 4.8	<b>28.8 ± 4.6</b>	31.0 ± 5.2
<b>shufflenet_v2_x1.0</b>	<b>35.8 ± 5.3</b>	38.6 ± 4.7	<b>36.7 ± 4.4</b>	39.0 ± 6.5	43.5 ± 9.0	<b>36.4 ± 6.9</b>	<b>35.9 ± 4.7</b>
<b>squeezenet1.0</b>	<b>36.6 ± 5.8</b>	<b>37.3 ± 4.3</b>	38.2 ± 6.8	<b>35.3 ± 6.9</b>	65.1 ± 9.4	37.6 ± 4.1	39.6 ± 4.6
<b>vgg16</b>	35.3 ± 6.4	36.2 ± 6.4	<b>34.6 ± 4.7</b>	<b>32.3 ± 3.8</b>	55.1 ± 10.5	34.7 ± 5.5	35.1 ± 6.0
<b>wide_resnet50_2</b>	30.2 ± 5.7	<b>29.9 ± 4.9</b>	<b>28.2 ± 5.0</b>	<b>30.5 ± 6.2</b>	47.7 ± 14.4	<b>29.1 ± 5.6</b>	30.7 ± 4.3
<b>Avg</b>	33.1	33.9	34.0	32.7	45.8	33.1	34.2
<b>Winners</b>	6	5	5	8	1	6	4

**Table 9.** Accuracy comparison of different models with literature for 7 and 4 classes.

	<b>7 classes</b> <b>Accuracy (%)</b>	<b>4 classes</b> <b>Accuracy (%)</b>
<b>Jantzen et al.</b>	61.1	—
<b>Lin et al.</b>	64.5	71.3
<b>Proposal</b>	<b>75.6</b>	<b>81.3</b>



**Figure 6.** Probability distribution for WideResNet50 contrasting losses CE (solid line) and HO2 (dashed line).

and even in contrast to OE.

## CONCLUSION

Non-parametric losses achieved better results when comparing with parametric losses. This type of losses does not limit the learned representation to a specific parametric model, which allows during the train to explore different and larger spaces of solutions avoiding the need for ad hoc choices.

A new non-parametric loss is proposed for multi-class Pap smear cell-classification based on convolutional neural networks. This new loss demonstrated to be competitive with state-of-the-art results and more flexible than the existing in deep ordinal classification techniques when imposing uni-modality in probabilities distribution. The use of the proposed loss in training popular architectures from the literature outperforms the state-of-the-art by over 10%.

Furthermore, the proposed loss is a convenient way of introducing ordinality to the optimization problem without the major changes in architecture or data format required by other techniques from the literature. On the other hand, the proposed loss is pestered by two new hyperparameters, albeit the suggested values have been found to be robust. While motivated by this dataset, the proposed loss could potentially be used by other applications of ordinal classification or even time-series problems.

In any case, there is a lot to improve in multi-class classification of cervical cells to achieve better accuracy performances since results are still short of 75.6% accuracy. The Herlev data set is mainly composed of expert-selected “typical” cells, however, in real-life circumstances, it is more complex because a cytology image contains lots of cells and not only a single cropped cell, so further work are needed before moving the results of this work to practice. Another important detail is the effect of overlapping nuclei and cell clumps, which has not been taken into account in this work. The presence of artifacts on the images also interferes with classification accuracy.

## ACKNOWLEDGMENTS

The authors would like to acknowledge for the access to Herlev Pap smear dataset collected by Herlev University Hospital (Denmark) and the Technical University of Denmark.

This work was financed by the ERDF – European Regional Development Fund through the Operational Programme for Competitiveness and Internationalisation – COMPETE 2020 Programme and by National Funds through the Portuguese funding agency, FCT – Fundação para a Ciência e a Tecnologia within project “POCI-01-0145-FEDER-028857”. Ricardo Cruz was supported by Ph.D. grant SFRH/BD/122248/2016, also provided by FCT. There was no additional external funding received for this study.

## REFERENCES

- Beckham, C. and Pal, C. (2017). Unimodal probability distributions for deep ordinal classification. In Precup, D. and Teh, Y. W., editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 411–419, International Convention Centre, Sydney, Australia. PMLR.

- 361 Belharbi, S., Ayed, I. B., McCaffrey, L., and Granger, E. (2019). Non-parametric uni-modality constraints  
362 for deep ordinal classification.
- 363 Bora, K., Chowdhury, M., Mahanta, L., Kundu, M., and Das, A. (2017). Automated classification of pap  
364 smear image to detect cervical dysplasia. *Computer Methods and Programs in Biomedicine*, 138:31–47.
- 365 Chankong, T., Theera-Umpon, N., and Auephanwiriyakul, S. (2014). Automatic cervical cell segmentation  
366 and classification in PAP smears. *Computer methods and programs in biomedicine*, 113.
- 367 Chen, Y.-F., Huang, P.-C., Lin, K.-C., Lin, H.-H., Wang, L.-E., Cheng, C.-C., Chen, T.-P., Chan, Y.-K.,  
368 and Y. Chiang, J. (2014). Semi-automatic segmentation and classification of pap smear cells. *IEEE*  
369 *journal of biomedical and health informatics*, 18:94–108.
- 370 Cheng, J., Wang, Z., and Pollastri, G. (2008). A neural network approach to ordinal regression. In *2008*  
371 *IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational*  
372 *Intelligence)*, pages 1279–1284. IEEE.
- 373 Costa, J. and Cardoso, J. (2005). Classification of ordinal data using neural networks. In *European*  
374 *Conference on Machine Learning*, pages 690–697.
- 375 DeMay, R. M. (2007). *Practical Principles of Cytopathology*. American Society for Clinical Pathology  
376 Press.
- 377 Fekri Ershad, S. (2019). Pap smear classification using combination of global significant value, texture  
378 statistical features and time series features. *Multimedia Tools and Applications*.
- 379 Frank, E. and Hall, M. (2001). A simple approach to ordinal classification. In *European Conference on*  
380 *Machine Learning*, pages 145–156.
- 381 Ghoneim, A., Muhammad, G., and Hossain, M. S. (2019). Cervical cancer classification using con-  
382 volutional neural networks and extreme learning machines. *Future Generation Computer Systems*,  
383 102.
- 384 He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In  
385 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- 386 Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., and Adam,  
387 H. (2017). MobileNets: Efficient convolutional neural networks for mobile vision applications. *arXiv*  
388 *preprint arXiv:1704.04861*.
- 389 Iandola, F. N., Han, S., Moskewicz, M. W., Ashraf, K., Dally, W. J., and Keutzer, K. (2016).  
390 SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and 0.5 MB model size. *arXiv*  
391 *preprint arXiv:1602.07360*.
- 392 Jantzen, J. and Dounias, G. (2006). Analysis of Pap-smear image data. In *Proceedings of the Nature-*  
393 *Inspired Smart Information Systems 2nd Annual Symposium*, volume 10.
- 394 Kashyap, D., Somani, A., Shekhar, J., Bhan, A., Dutta, M. K., Burget, R., and Riha, K. (2016). Cervical  
395 cancer detection and classification using independent level sets and multi SVMs. In *2016 39th*  
396 *international conference on telecommunications and signal processing (TSP)*, pages 523–528. IEEE.
- 397 Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). ImageNet classification with deep convolutional  
398 neural networks. In *Advances in neural information processing systems*, pages 1097–1105.
- 399 Kurnianingsih, K., Allehaibi, K., Nugroho, L., Widyawan, W., Lazuardi, L., Prabuwo, A. S., and  
400 Mantoro, T. (2019). Segmentation and classification of cervical cells using deep learning. *IEEE Access*,  
401 PP:1–1.
- 402 Lin, H., Hu, Y., Chen, S., Yao, J., and Zhang, L. (2019). Fine-grained classification of cervical cells using  
403 morphological and appearance based convolutional neural networks. *IEEE Access*, PP:1–1.
- 404 Marinakis, Y., Dounias, G., and Jantzen, J. (2009). Pap smear diagnosis using a hybrid intelligent scheme  
405 focusing on genetic algorithm based feature selection and nearest neighbor classification. *Computers in*  
406 *biology and medicine*, 39:69–78.
- 407 Nayar, R. and Wilbur, D. (2015). The pap test and Bethesda 2014. *Journal of the American Society of*  
408 *Cytopathology*, 123.
- 409 Paul, P. R., Bhowmik, M. K., and Bhattacharjee, D. (2015). *Automated cervical cancer detection using*  
410 *Pap smear images*, pages 267–278. Springer.
- 411 Siegel, R. L., Miller, K. D., and Jemal, A. (2020). Cancer statistics, 2020. *CA: A Cancer Journal for*  
412 *Clinicians*.
- 413 Silva, W., Pinto, J. R., and Cardoso, J. S. (2018). A uniform performance index for ordinal classification  
414 with imbalanced classes. In *2018 International Joint Conference on Neural Networks (IJCNN)*, pages  
415 1–8. IEEE.

- 416 Simonyan, K. and Zisserman, A. (2014). Very deep convolutional networks for large-scale image  
417 recognition. *arXiv preprint arXiv:1409.1556*.
- 418 Suhrland, M. (2000). Practical principles of cytopathology: Author: Mac demay ascp press, chicago,  
419 1999. *Diagnostic Cytopathology - DIAGN CYTOPATHOL*, 23:213–213.
- 420 Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and  
421 Rabinovich, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE conference on*  
422 *computer vision and pattern recognition*, pages 1–9.
- 423 WHO (2019). Human papillomavirus (hvp) and cervical cancer. [https://www.](https://www.who.int/news-room/fact-sheets/detail/human-papillomavirus-(hvp)-and-cervical-cancer)  
424 [who.int/news-room/fact-sheets/detail/human-papillomavirus-\(hvp\)](https://www.who.int/news-room/fact-sheets/detail/human-papillomavirus-(hvp)-and-cervical-cancer)  
425 [-and-cervical-cancer](https://www.who.int/news-room/fact-sheets/detail/human-papillomavirus-(hvp)-and-cervical-cancer). [Online; accessed 2020-02-19].
- 426 William, W., Ware, A., Basaza-Ejiri, A. H., and Obungoloch, J. (2019). Cervical cancer classification  
427 from pap-smears using an enhanced fuzzy c-means algorithm. *Informatics in Medicine Unlocked*,  
428 14:23 – 33.
- 429 Xie, S., Girshick, R., Dollár, P., Tu, Z., and He, K. (2017). Aggregated residual transformations for deep  
430 neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
431 pages 1492–1500.
- 432 Zagoruyko, S. and Komodakis, N. (2016). Wide residual networks. *arXiv preprint arXiv:1605.07146*.
- 433 Zhang, L., Lu, L., Nogues, I., Summers, R., Liu, S., and Yao, J. (2017). Deeppap: Deep convolutional  
434 networks for cervical cell classification. *IEEE Journal of Biomedical and Health Informatics*, PP:1–1.
- 435 Zhang, X., Zhou, X., Lin, M., and Sun, J. (2018). ShuffleNet: An extremely efficient convolutional neural  
436 network for mobile devices. In *Proceedings of the IEEE conference on computer vision and pattern*  
437 *recognition*, pages 6848–6856.

## 438 APPENDIX

439 Some extra results are made available in this appendix.

**Table A1.** Results in terms of **Kendall's  $\tau$**  for **7 class** problem, averaged for 10 folds. Higher is better.

	CE	BU	PU	OE	CO	CO2	HO2
<b>alexnet</b>	75.9±4.1	76.9±4.8	74.5±5.8	<b>77.9±4.8</b>	57.3±9.3	<b>79.9±4.7</b>	77.5±5.5
<b>googlenet</b>	<b>80.9±2.7</b>	<b>81.9±4.1</b>	79.3±4.4	<b>81.9±4.3</b>	76.9±4.0	<b>81.7±4.2</b>	<b>82.8±3.7</b>
<b>mobilenet_v2</b>	83.2±2.9	84.5±1.9	<b>85.6±2.5</b>	<b>84.4±2.1</b>	75.6±11.5	83.7±3.2	83.8±3.0
<b>resnet18</b>	<b>83.1±5.1</b>	<b>84.7±3.1</b>	<b>83.4±3.0</b>	<b>83.5±4.9</b>	76.9±5.1	<b>83.8±4.0</b>	<b>83.4±5.1</b>
<b>resnext50_32x4d</b>	83.2±4.5	<b>85.2±2.8</b>	<b>84.7±1.4</b>	83.2±3.6	81.5±3.4	<b>84.0±3.6</b>	<b>85.4±3.1</b>
<b>shufflenet_v2_x1.0</b>	78.8±4.0	78.9±3.6	<b>80.9±2.4</b>	<b>81.7±3.1</b>	77.8±4.5	80.1±2.6	<b>81.6±3.4</b>
<b>squeezenet1.0</b>	<b>81.0±3.9</b>	<b>80.4±2.6</b>	77.0±4.0	<b>79.6±5.1</b>	54.1±14.6	<b>79.4±4.1</b>	78.1±4.8
<b>vgg16</b>	<b>81.6±5.1</b>	<b>81.3±2.2</b>	78.3±6.4	<b>81.9±3.5</b>	68.8±6.9	<b>82.6±3.2</b>	<b>82.8±4.6</b>
<b>wide_resnet50_2</b>	83.4±3.1	83.8±2.8	84.9±3.4	<b>85.9±2.7</b>	79.8±6.1	84.1±3.6	82.7±4.3
<b>Avg</b>	81.2	82.0	80.9	82.2	72.1	82.2	82.0
<b>Winners</b>	4	5	4	8	0	6	5

**Table A2.** Results in terms of **Kendall's  $\tau$**  for **4 class** problem, averaged for 10 folds. Higher is better.

	CE	BU	PU	OE	CO	CO2	HO2
<b>alexnet</b>	73.9±6.3	75.4±4.0	<b>77.4±3.8</b>	<b>76.0±5.7</b>	58.7±19.3	<b>76.5±5.7</b>	75.0±5.9
<b>googlenet</b>	80.3±4.2	<b>81.3±3.0</b>	<b>80.4±4.9</b>	<b>81.2±4.4</b>	70.1±14.2	<b>82.5±4.2</b>	80.2±5.1
<b>mobilenet_v2</b>	81.3±5.1	<b>83.9±2.0</b>	81.3±3.7	<b>81.8±6.2</b>	<b>82.3±3.7</b>	80.4±4.3	82.3±4.4
<b>resnet18</b>	81.0±2.7	80.3±4.4	<b>81.2±4.4</b>	<b>82.8±3.5</b>	77.5±7.7	<b>82.5±3.2</b>	79.5±5.1
<b>resnext50_32x4d</b>	<b>83.2±3.1</b>	<b>83.2±2.9</b>	81.8±2.8	<b>84.3±3.4</b>	78.1±5.1	<b>83.9±2.4</b>	82.3±3.9
<b>shufflenet_v2_x1.0</b>	77.1±4.7	76.1±3.7	<b>78.8±4.1</b>	74.8±6.2	71.1±9.1	<b>77.3±5.6</b>	<b>78.0±4.5</b>
<b>squeezenet1.0</b>	<b>76.2±5.2</b>	<b>77.5±4.5</b>	<b>75.8±5.6</b>	<b>78.2±5.9</b>	48.5±11.0	<b>76.6±3.8</b>	74.8±3.8
<b>vgg16</b>	77.9±5.4	<b>79.2±4.5</b>	<b>80.0±4.4</b>	<b>81.0±2.4</b>	63.2±10.3	79.2±4.2	78.3±5.4
<b>wide_resnet50_2</b>	81.7±4.6	<b>83.4±2.8</b>	<b>84.4±4.6</b>	<b>82.2±5.1</b>	67.8±18.6	<b>83.2±4.3</b>	<b>82.6±2.8</b>
<b>Avg</b>	79.2	80.0	80.1	80.2	68.6	80.3	79.2
<b>Winners</b>	2	6	7	8	1	7	2

**Table A3.** Aggregate results for 7 class problem, averaged for 10 folds.

	CE	BU	PU	OE	CO	CO2	HO2
	<b>Mode</b>						
<b>uoc</b>	<b>38.8 ± 7.5</b>	44.2 ± 7.2	41.1 ± 7.7	<b>38.6 ± 7.5</b>	53.5 ± 14.7	<b>39.0 ± 7.3</b>	39.5 ± 8.1
<b>mae</b>	0.37 ± 0.08	0.41 ± 0.09	0.40 ± 0.10	<b>0.36 ± 0.08</b>	0.60 ± 0.26	<b>0.37 ± 0.07</b>	<b>0.37 ± 0.09</b>
<b>acc</b>	<b>73.4 ± 4.8</b>	66.2 ± 5.5	71.5 ± 4.9	72.8 ± 5.1	58.6 ± 13.4	72.6 ± 4.5	72.0 ± 5.6
<b>tau</b>	81.2 ± 4.7	<b>82.0 ± 4.2</b>	80.0 ± 5.8	<b>82.2 ± 4.5</b>	72.1 ± 12.4	<b>82.2 ± 4.2</b>	<b>82.0 ± 4.9</b>
<b>auroc</b>	<b>95.9 ± 1.4</b>	93.0 ± 2.2	95.5 ± 1.4	95.5 ± 1.5	82.3 ± 9.9	94.5 ± 1.8	93.9 ± 1.9
<b>gini</b>	85.1 ± 0.2	64.0 ± 1.6	84.8 ± 0.5	84.8 ± 0.4	28.2 ± 33.3	50.1 ± 6.3	45.0 ± 4.2
	<b>Mean</b>						
<b>uoc</b>	<b>39.2 ± 7.4</b>	42.1 ± 7.5	41.7 ± 8.2	<b>39.7 ± 7.4</b>	75.9 ± 20.6	79.8 ± 4.9	83.7 ± 1.7
<b>mae</b>	<b>0.37 ± 0.08</b>	0.39 ± 0.08	0.39 ± 0.09	<b>0.37 ± 0.08</b>	1.17 ± 0.50	0.91 ± 0.13	1.03 ± 0.04
<b>acc</b>	<b>72.1 ± 5.1</b>	67.7 ± 5.8	70.1 ± 5.6	<b>71.6 ± 5.1</b>	33.6 ± 19.9	28.3 ± 5.4	25.0 ± 3.1
<b>tau</b>	<b>82.3 ± 4.3</b>	<b>82.5 ± 4.1</b>	80.9 ± 5.5	<b>82.6 ± 4.3</b>	<i>nan ± nan</i>	78.3 ± 3.6	76.0 ± 3.2
<b>auroc</b>	<b>95.9 ± 1.4</b>	93.0 ± 2.2	95.5 ± 1.4	95.5 ± 1.5	82.3 ± 9.9	94.5 ± 1.8	93.9 ± 1.9
<b>gini</b>	85.1 ± 0.2	64.0 ± 1.6	84.8 ± 0.5	84.8 ± 0.4	28.2 ± 33.3	50.1 ± 6.3	45.0 ± 4.2

**Table A4.** Aggregate results for 4 class problem, averaged for 10 folds.

	CE	BU	PU	OE	CO	CO2	HO2
	<b>Mode</b>						
<b>uoc</b>	<b>33.1 ± 6.3</b>	33.9 ± 5.9	33.6 ± 5.7	<b>32.7 ± 7.0</b>	45.8 ± 14.9	<b>33.1 ± 6.3</b>	34.2 ± 6.4
<b>mae</b>	0.25 ± 0.06	0.27 ± 0.06	<b>0.25 ± 0.06</b>	<b>0.24 ± 0.06</b>	0.40 ± 0.20	<b>0.25 ± 0.06</b>	0.26 ± 0.06
<b>acc</b>	<b>79.2 ± 4.1</b>	76.5 ± 4.5	<b>78.7 ± 4.0</b>	<b>79.2 ± 4.8</b>	67.4 ± 13.7	78.5 ± 4.2	77.8 ± 4.3
<b>tau</b>	79.2 ± 5.5	<b>80.0 ± 4.7</b>	<b>79.7 ± 5.4</b>	<b>80.2 ± 5.8</b>	68.6 ± 15.7	<b>80.3 ± 5.1</b>	79.2 ± 5.4
<b>auroc</b>	<b>94.5 ± 1.7</b>	92.7 ± 1.9	<b>94.6 ± 1.5</b>	<b>94.5 ± 1.7</b>	83.5 ± 10.5	93.2 ± 2.1	92.9 ± 2.0
<b>gini</b>	74.3 ± 0.3	58.4 ± 2.0	74.2 ± 0.4	74.2 ± 0.3	31.5 ± 31.4	44.5 ± 16.3	36.6 ± 11.2
	<b>Mean</b>						
<b>uoc</b>	<b>33.4 ± 6.3</b>	34.2 ± 5.8	<b>34.0 ± 5.9</b>	<b>33.2 ± 7.1</b>	57.0 ± 17.3	52.2 ± 13.5	58.4 ± 8.5
<b>mae</b>	<b>0.25 ± 0.06</b>	0.26 ± 0.05	<b>0.25 ± 0.06</b>	<b>0.25 ± 0.07</b>	0.50 ± 0.20	0.42 ± 0.14	0.47 ± 0.09
<b>acc</b>	<b>78.5 ± 4.4</b>	77.0 ± 4.2	77.8 ± 4.4	<b>78.4 ± 5.1</b>	54.7 ± 17.8	60.1 ± 13.4	54.6 ± 9.3
<b>tau</b>	79.7 ± 5.1	<b>80.4 ± 4.6</b>	<b>80.1 ± 4.9</b>	<b>80.5 ± 5.8</b>	63.8 ± 16.6	73.4 ± 7.3	70.5 ± 5.6
<b>auroc</b>	<b>94.5 ± 1.7</b>	92.7 ± 1.9	<b>94.6 ± 1.5</b>	<b>94.5 ± 1.7</b>	83.5 ± 10.5	93.2 ± 2.1	92.9 ± 2.0
<b>gini</b>	74.3 ± 0.3	58.4 ± 2.0	74.2 ± 0.4	74.2 ± 0.3	31.5 ± 31.4	44.5 ± 16.3	36.6 ± 11.2