

---

# Rank4Class: A Ranking Formulation for Multiclass Classification

---

Nan Wang<sup>1,2</sup> Zhen Qin<sup>3</sup> Le Yan<sup>3</sup> Honglei Zhuang<sup>3</sup> Xuanhui Wang<sup>3</sup> Michael Bendersky<sup>3</sup> Marc Najork<sup>3</sup>

## Abstract

Multiclass classification (MCC) is a fundamental machine learning problem of classifying each instance into one of a predefined set of classes. In the deep learning era, extensive efforts have been spent on developing more powerful neural embedding models to better represent the instance for improving MCC performance. In this paper, we do not aim to propose new neural models for instance representation learning, but to show that it is promising to boost MCC performance with a novel formulation through the lens of ranking. In particular, by viewing MCC as to rank classes for an instance, we first argue that ranking metrics, such as Normalized Discounted Cumulative Gain, can be more informative than the commonly used Top- $K$  metrics. We further demonstrate that the dominant neural MCC recipe can be transformed to a neural ranking framework. Based on such generalization, we show that it is intuitive to leverage advanced techniques from the learning to rank literature to improve the MCC performance out of the box. Extensive empirical results on both text and image classification tasks with diverse datasets and backbone neural models show the value of our proposed framework.

## 1. Introduction

Multiclass classification (MCC) is the problem of classifying each instance into one of a predefined set of classes (Hastie et al., 2001). It is one of the most fundamental machine learning problems that has broad applications in many fields such as natural language processing (Sun et al., 2019) and computer vision (He et al., 2016). For example, deciding the category of a news article or the subject of an image can be formulated as an MCC problem.

---

<sup>1</sup>Department of Computer Science, University of Virginia, VA, USA <sup>2</sup>Work was done when the author was a research intern at Google Research <sup>3</sup>Google Research, Mountain View, CA, USA. Correspondence to: Nan Wang <nw6a@virginia.edu>, Zhen Qin <zhenqin@google.com>.

Numerous MCC models have been proposed in the past, ranging from linear models to nonlinear decision trees and neural models (Aly, 2005). In the modern deep learning era, while there are significant advances in neural architectures, dominating MCC methods share the same recipe: an input instance, being it a feature vector, an image, or a text sentence, is fed into a neural embedding model to produce its vector representation, which is followed by a classification layer to score against the candidate classes. The model is trained by using a loss function, typically the softmax cross entropy loss, between the labels and scores over all candidate classes (Goodfellow et al., 2016). During inference, the classes are sorted after an instance is scored against them. Metrics such as Top- $K$  Accuracy (the percentage of test instances whose correct class label is in the top  $K$  predicted classes, also known as simply “classification accuracy” when  $K=1$ ) are usually used for evaluating MCC performance. To better display the headroom of the models, Top- $K$  Error ( $1 - \text{Top-}K \text{ Accuracy}$ ) is widely adopted to compare different classification models (Krizhevsky et al., 2012). Following this recipe, most efforts in the MCC literature focus on designing more powerful neural network architectures for representation learning of the input instances (He et al., 2016; Huang et al., 2017; Dosovitskiy et al., 2021; Minaee et al., 2021). Few work has studied a formulation different from the recipe above.

Paralleling with developing more powerful embedding models, we seek for a novel formulation for MCC by examining it through the lens of ranking, or more specifically, learning to rank (LTR), a rich research field stemmed from information retrieval (Liu, 2009). As shown in the rest of the paper, firstly, such a formulation allows us to better evaluate model performances by borrowing more informative ranking metrics. Secondly, it improves MCC performance out of the box, agnostic to the embedding architecture used, by training models with advanced ranking losses and use more flexible matching schemes between the input instance and candidate classes.

In particular, we first formalize MCC as a ranking problem. As hinted above, dominating neural MCC methods score a given instance on a predefined set of classes and sort the classes during inference. This is equivalent to a LTR setting where a set of items (i.e., classes) are ranked given a query (i.e., the input instance). In fact, a Top- $K$  Accuracy

measure itself can be viewed as a ranking metric, similar to Precision@K used in information retrieval. However, other ranking metrics such as Normalized Discounted Cumulative Gain (NDCG) are more commonly used in LTR because they are more informative than Precision@K in real-world user-facing applications. They can be borrowed to enrich the MCC evaluation metrics but have not been widely adopted.

For modeling, we show a general “equivalent view” for MCC from ranking perspectives, where the dominant MCC recipe is equivalent to a neural ranking pipeline with *a specific set of design choices*. This LTR view gives us insights into existing MCC models’ groundings as well as limitations, which then inspires more design options. We further propose a general framework with several intuitive approaches under the LTR formulation to improve MCC performance. We mainly focus on two aspects: loss functions and model architectures. For loss functions, we leverage the rich literature in LTR of advanced ranking losses specifically designed for certain ranking metrics. For model architecture, we realize that the vast MCC literature focuses on only one component in the neural ranking architecture, i.e. the input instance embedding. With the LTR view, we can enhance the modeling capacity of other components as well. In this paper, we specifically study the effect of enhancing the interactions between instances and classes, a popular setup in LTR to match queries and items (Li & Xu, 2014).

We report experimental results on a variety of MCC tasks, including different datasets and backbone models for different modalities (image and text). Results show that the proposed framework outperform or perform competitively with the baselines in all settings. We expect that the promising results can encourage the community to further examine MCC from LTR perspectives.

## 2. Ranking Metrics for MCC

### 2.1. Metrics for Classification

The basic binary classification problem classifies instances into positive and negative classes. MCC extends binary classification to more than two classes. Going from binary to multiclass is not trivial for evaluation. Binary classification metrics are usually *class-oriented*. For example, metrics such as AUC and Accuracy are based on measures like true positives (TP) and false negatives (FN), which are computed with respect to positive and negative classes (Sokolova & Lapalme, 2009). These metrics are not directly used in MCC for more than two classes. In contrast, MCC metrics are usually *instance-oriented*. The commonly used metrics are the Top- $K$  Accuracy/Error metrics, popularized by the ImageNet competition (Russakovsky et al., 2015). Some earlier works like (Crammer & Singer, 2002) defined the “empirical error” when working on multiclass SVM algorithms,

which is equivalent to Top-1 Error.

### 2.2. Relation to Ranking Metrics

For easier illustration, we use Top- $K$  Accuracy in the rest of this section, which is simply  $1 - \text{Top-}K \text{ Error}$ . As discussed above, MCC metrics are instance-oriented, and the metric value is simply averaged over all instances in a given evaluation set. Thus we study the metric calculated on a single instance to be more concise. In particular, given an instance and  $n$  candidate classes, let  $\mathbf{y}$  be the true labels of the classes and  $\mathbf{y}[c] \in \{0, 1\}$  be the label for class  $c$ . We consider the standard MCC setting where there is exactly one class with label 1 per instance (Aly, 2005). Let  $\pi_A$  be a ranking of classes produced by model A, with the  $i$ -th most likely class being  $\pi_A(i)$ . Then the Top- $K$  Accuracy ( $1 \leq K \leq n$ , abbreviated as  $\text{acc}@K$  in Eq 1) is

$$\text{acc}@K(\pi_A, \mathbf{y}) = \sum_{i=1}^K \mathbf{y}[\pi_A(i)]. \quad (1)$$

The Top- $K$  Accuracy can be thought of a type of a ranking metric. It is very close to the Precision@ $K$ : Top- $K$  Accuracy is the same as  $\min(1, K \cdot \text{Precision}@K)$  and Top-1 Accuracy is the same as Precision@1. Besides Precision@ $K$ , there are other commonly used ranking metrics such as Normalized Discounted Cumulative Gains (NDCG) and Mean Reciprocal Rank (MRR) (Järvelin & Kekäläinen, 2002). To the best of our knowledge, they are not commonly used for MCC evaluation, but can be more informative.

We use NDCG@ $K$  (simply called NDCG when  $K = n$ ) as an example. Consider the instance as a query, and all the classes as candidate items with relevance labels  $\mathbf{y}$ , we have

$$\text{NDCG}@K(\pi_A, \mathbf{y}) = \frac{\text{DCG}@K(\pi_A, \mathbf{y})}{\text{DCG}@K(\pi_*, \mathbf{y})}, \quad (2)$$

$$\text{and } \text{DCG}@K(\pi, \mathbf{y}) = \sum_{i=1}^K \frac{2^{\mathbf{y}[\pi(i)]} - 1}{\log_2(1 + i)}.$$

where  $\pi_*$  is the ideal ranking sorted by  $\mathbf{y}$ . It’s easy to prove that NDCG@1 is equivalent to Top-1 Accuracy.

**Theorem 1.** *Given an instance with class labels  $\mathbf{y} \in \{0, 1\}^n$  and  $\sum_{c=1}^n \mathbf{y}[c] = 1$ , for any  $1 < K \leq n$ , NDCG@ $K$  preserves more information than Top- $K$  Accuracy in evaluating the class rankings.*

The proof of Theorem 1 is in Appendix A. The basic idea is to measure the information of a metric by the entropy of its value evaluated on any class ranking as a random variable. This information measure indicates that, given the metric’s value, how much can we say about its input class ranking.

Theorem 1 shows that for  $K > 1$ , NDCG@ $K$  is strictly more informative than Top- $K$  Accuracy. To better illustrate this, let us consider two models A and B, which rank

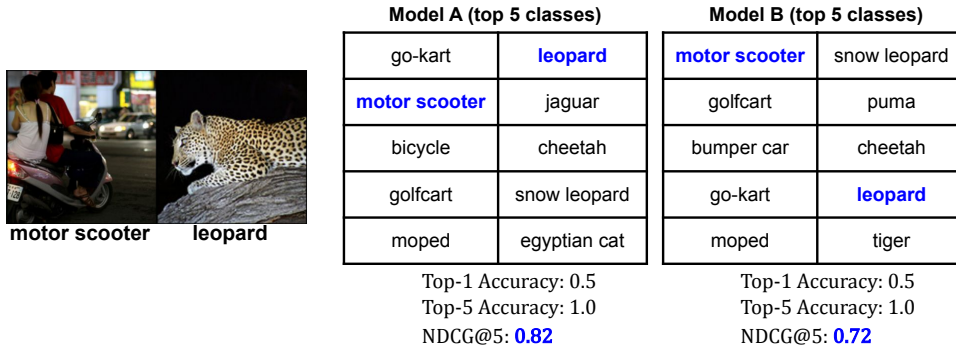


Figure 1. We show an simulated example of evaluating image classification performance of two models with Top-1/5 Accuracy and NDCG@5. The Top-1/5 Accuracy of two models are exactly the same, and can not differentiate the two models. While NDCG@5 can successfully tell that model A has a better performance than model B.

the correct class at position  $p_A$  and  $p_B$ , respectively, and  $p_A < p_B \leq K, 2 \leq K \leq n$ . Thus, both models can rank the correct class in the top  $K$  positions while model A ranks it higher than model B. With Top- $K$  Accuracy, we have  $\text{acc}@K(\pi_A, \mathbf{y}) = \text{acc}@K(\pi_B, \mathbf{y}) = 1$ . But with  $\text{NDCG}@K$ , we have  $\text{NDCG}@K(\pi_A, \mathbf{y}) = 1/\log_2(1 + p_A)$ , and  $\text{NDCG}@K(\pi_B, \mathbf{y}) = 1/\log_2(1 + p_B)$ , which gives us  $\text{NDCG}@K(\pi_A, \mathbf{y}) > \text{NDCG}@K(\pi_B, \mathbf{y})$ . Therefore, when model A consistently ranks the correct class higher than model B in top  $K$  positions, Top- $K$  Accuracy may not be able to reflect the better performance of A, while  $\text{NDCG}@K$  will always detect it. More importantly, although both metrics ignore the correct classes ranked below the  $K$ 'th position, we can simply set  $K = n$  to evaluate the whole ranked list of classes. In this case,  $\text{NDCG}$  is still a meaningful metric that reflects the position of the correct class in the full list, while Top- $K$  Accuracy becomes meaningless (always 1).

In Figure 1, we simulate an image MCC application in a *user-facing interface* and compare the results from two models with different metrics.  $\text{NDCG}@5$  is clearly more capable of distinguishing the two models, while both Top-1 Accuracy and Top-5 fails to. The reason is that there is a position discounting function in  $\text{NDCG}@K$ , while Top- $K$  metrics impose a simple hard cut at position  $K$ , which loses information about the exact ranked position of the correct class. Other ranking metrics such as MRR share similar characteristics with  $\text{NDCG}$ . In the rest of this paper, we choose  $\text{NDCG}@K$  as our representatives for ranking metrics and use it in addition to Top- $K$  metrics in experiments.

### 3. MCC from a Ranking Perspective

#### 3.1. Classical MCC Model Architectures

The general model architecture for MCC based on deep neural networks (DNN) is composed of three parts: an input instance to be classified, an encoder to extract latent representation (embedding) of the input, and a classification

layer for generating scores on candidate classes, as shown in Figure 2 Left. We use  $\mathbf{x}$  to represent the input instance such as a textual sentence or an image. The encoder can have different structure designs based on the modality of the input, such as transformers (Vaswani et al., 2017) to encode textual sentences; and convolutional neural networks (Krizhevsky et al., 2012) to encode images. It can be represented by a function  $\mathcal{H}(\cdot)$  that maps  $\mathbf{x}$  to a  $d$ -dimensional embedding vector  $\mathbf{h} = \mathcal{H}(\mathbf{x})$ . The classification layer is in most cases a dense layer with weight matrix  $W \in \mathbb{R}^{n \times (d+1)}$ . The classification scores are calculated by  $\mathbf{s} = W\mathbf{h}'$ , where  $\mathbf{h}' := [\mathbf{h}, 1]$  with an added bias dimension of value 1 for the bias.  $\mathbf{s}$  is an  $n$ -dimensional score vector for  $n$  classes with  $s_i = \mathbf{e}_i^\top W\mathbf{h}'$  for the  $i$ -th class, where  $\mathbf{e}_i$  is a  $n$ -dimensional one-hot vector with the  $i$ -th dimension being 1.

For training neural MCC models, the softmax cross entropy loss (SoftmaxCE) is used by default in almost all prior work (Goodfellow et al., 2016; Zhang et al., 2021). However, whether it is the most suitable loss for optimizing the evaluation metrics of interest is not carefully studied. Very recently, empirical results show that the mean squared error (MSE) can sometimes outperform SoftmaxCE in MCC tasks (Hui & Belkin, 2021), but it is sensitive to an extra rescaling parameter for tasks with many classes.

#### 3.2. An Equivalent View from Neural Ranking Models

LTR learns a ranking model to rank a set of items (e.g., documents, news, etc.) based on their relevance to a query. Neural ranking models (Guo et al., 2020) adopt DNNs to match the query and items with their latent representations.

**Claim 1.** *The classical DNN model for MCC as described in Section 3.1 is equivalent to a specific neural ranking model trained with softmax cross entropy loss.*

To prove Claim 1, we show the equivalent view of a classical MCC model as a neural ranking model in Figure 2 Middle. In particular, we treat the input instance  $\mathbf{x}$  as the query, and the  $n$  candidate classes as input items to be ranked

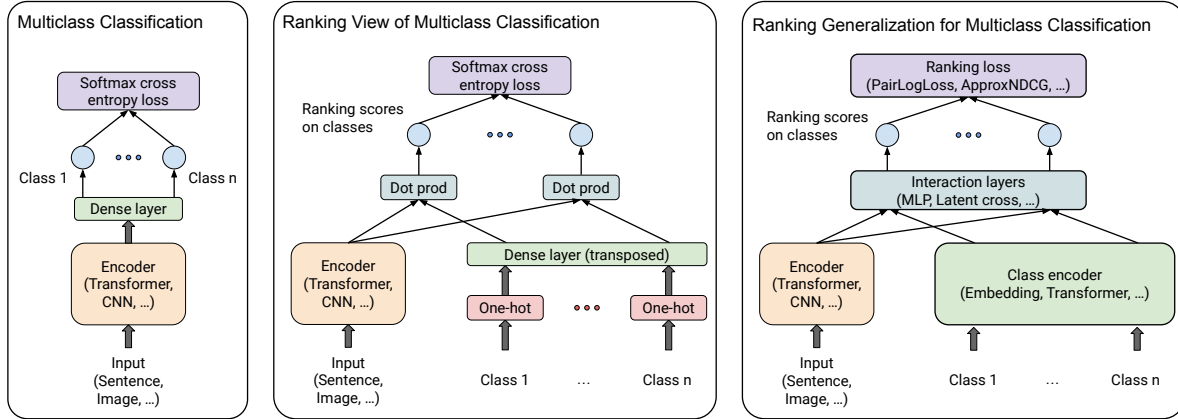


Figure 2. **Left:** the classical DNN architecture for MCC; **Middle:** the equivalent neural ranking architecture; **Right:** The ranking generalization for MCC which allows the exploration of different class encoders, interaction layers, and loss functions.

by the neural ranking model. The same as in the classical MCC model, we apply the encoder  $\mathcal{H}(\cdot)$  to  $\mathbf{x}$  to get a  $d$ -dimensional embedding  $\mathbf{h} = \mathcal{H}(\mathbf{x})$ , and  $\mathbf{h}' = [\mathbf{h}, 1]$ . For representing the classes, we use one-hot vector  $\mathbf{e}_i$  for the  $i$ -th class and obtain its embedding  $\mathbf{c}_i = W^\top \mathbf{e}_i$ , where  $W$  is the weight matrix from the classification layer of the classical MCC model. Finally, the ranking score of the  $i$ -th class can be calculated by simple *dot product*  $s'_i \equiv \mathbf{c}_i^\top \mathbf{h}' = \mathbf{e}_i^\top W \mathbf{h}' = s_i$ . In this way, we show the classical MCC model can be transformed to an equivalent neural ranking model when both are trained by the SoftmaxCE.

## 4. Ranking Architectures for MCC

In this section, we discuss the possible new designs of different components in the equivalent ranking view of MCC models. We call the resulting new framework ranking for multiclass classification, or *Rank4Class*.

As discussed in Section 1, existing work mostly focus on representation learning in the instance encoder. From the ranking perspective, however, we can see that there are several other promising directions to improve MCC performance. First, existing work mainly use SoftmaxCE, while there exist many advanced ranking losses that can be leveraged. Second, the interaction between the instance and class embeddings is simply a dot product, while richer interaction patterns can be explored. Third, while existing work focus on instance encoder, the class encoder is a naive linear projection that uses one-hot encoding of classes as input, where more powerful class encoders with richer inputs (e.g., class metadata) can be deployed. We illustrate the potential improvements in Figure 2 **Right**. In this paper, we focus on ranking losses and interaction patterns for improving MCC performance, and leave the rest for future work.

### 4.1. Ranking Losses for MCC

Although SoftmaxCE is widely adopted for training MCC models, whether it is the best option to optimize existing evaluation metrics is not clear (Hui & Belkin, 2021). On the other hand, we have shown that ranking metrics are better for evaluating MCC tasks. Thus, ranking losses are natural considerations since many of them are theoretically grounded to directly optimize certain ranking metrics. Next, we first discuss the soundness of using SoftmaxCE in MCC tasks with respect to ranking metrics. Then we present two example ranking losses, i.e., the pairwise logistic loss (PairLogLoss) and the approximate NDCG loss (ApproxNDCG).

We re-denote the label of the  $i$ -th class for an instance by  $y_i \in \{0, 1\}$ ,  $i \in \{1, \dots, n\}$ , where only the correct class has label 1. To compute the SoftmaxCE, the classification scores  $s_i$ ,  $i \in \{1, \dots, n\}$  produced by the model are first projected to the probability simplex  $\mathbf{p}$  by softmax activation as  $p_i = \frac{e^{s_i}}{\sum_{j=1}^n e^{s_j}}$ . Then SoftmaxCE is defined as

$$\ell_{ce}(\mathbf{y}, \mathbf{p}) = - \sum_{i=1}^n y_i \log p_i.$$

Intuitively, the SoftmaxCE is promoting the correct class against all other classes since only the term with  $y_i = 1$  is counted in it. It can be viewed as a listwise ranking loss, which aims to rank the correct class above all other classes.

**Theorem 2 (Bruch et al. (2019)).** *Softmax cross entropy loss is a bound on mean reciprocal rank and mean normalized discounted cumulative gain in log-scale.*

The proof of Theorem 2 can be found in Bruch et al. (2019) (Theorems 1-3). It explains the promising performance of SoftmaxCE in MCC tasks. To the best of our knowledge, there is no such bounds for MSE, so the findings in (Hui & Belkin, 2021) may need further theoretical investigation.



On the other hand, there is rich literature of developing ranking losses for optimizing ranking metrics (Burges et al., 2005; 2006; Burges, 2010). Such ranking losses are usually directly derived from the target ranking metric, and bound the metric by approximation techniques (Qin et al., 2010). One of the most historical and popularly used ranking losses is the PairLogLoss.

$$\ell_{pl}(\mathbf{y}, \mathbf{s}) = \sum_{i=1}^n \sum_{j=1}^n \mathbb{I}_{y_i > y_j} \log(1 + e^{-\sigma(s_i - s_j)}), \quad (3)$$

where  $\mathbb{I}$  is the indicator and  $\sigma$  is a hyper-parameter. The PairLogLoss is proved to be able to minimize the rank of the relevant item (Wang et al., 2018).

Another popular ranking loss is the ApproxNDCG (Qin et al., 2010), which directly optimizes the NDCG metric in Eq 2. The rank of an item  $i$  can be computed as  $\pi_s(i) = 1 + \sum_{j \neq i} \mathbb{I}_{s_i < s_j}$ , where the indicator  $\mathbb{I}_{s < t}$  is discrete but can be approximated by a sigmoid function to be smooth:

$$\mathbb{I}_{s < t} = \mathbb{I}_{t-s > 0} \approx \frac{1}{1 + e^{-\alpha(t-s)}}, \quad (4)$$

where  $\alpha > 0$  is a parameter to control how tightly the indicator is approximated.

To summarize, SoftmaxCE can be viewed as a ranking loss that bounds specific ranking metrics. Other ranking losses can also be valuable. We further investigate the empirical use of different ranking losses for MCC in Section 5.

#### 4.2. Enhancing Instance-Class Interactions

After obtaining the embeddings  $\mathbf{h}$  for the input instance and  $\mathbf{c}_i$  for class  $i$ , it is also important to design the matching mechanism for producing the score  $s_i$  between the embeddings of the instance and the class. For this purpose, we can add different interactions between  $\mathbf{h}$  and  $\mathbf{c}_i$  rather than simple dot-product for producing the score,  $s_i = \mathcal{I}(\mathbf{h}, \mathbf{c}_i)$ .  $\mathcal{I}(\cdot, \cdot)$  is a function to represent the interaction between instance and class embeddings which produces a ranking score. In this paper, we consider the following two patterns as examples to enhance the interaction, but much richer interactions can be deployed based on the LTR formulation:

- **LC+MLP:** We first apply element-wise multiplication, also known as the latent cross (LC) operation (Beutel et al., 2018) on  $\mathbf{h}$ ,  $\mathbf{c}_i$ , and then follow with a multilayer perceptron (MLP) to get the score  $s_i$ . LC has shown to be a simple and efficient way to generate higher-order matching interactions in DNNs.
- **Concat+MLP:** We first concatenate the two embeddings  $\mathbf{h}$ ,  $\mathbf{c}_i$  and follow with an MLP to compute  $s_i$ .

Table 1. Statistics of datasets used in experiments.

Dataset	#classes	Train	Validation	Test
GoEmotions	28	36.3K	4.5K	4.6K
MIND	18	78.8K	25.7K	25.9K
ImageNet	1000	1.28M	50K	-
CIFAR-10	10	50K	-	10K

## 5. Experiments

We study Rank4Class on both text classification and image classification tasks in comparison to classical MCC models. We consider several datasets and instance encoders for evaluation. The datasets are summarized in Table 1. For text classification, we include two recent large-scale datasets, GoEmotions (Demszky et al., 2020) and MIND (Wu et al., 2020). The GoEmotions dataset contains instances with multiple labels, which are filtered out, since we focus on single-label MCC tasks. We adopt ELECTRA (Clark et al., 2020) and BERT (Devlin et al., 2019) as text encoders in the experiments. For image classification, we use the popular ImageNet (Russakovsky et al., 2015) and CIFAR-10 (Krizhevsky & Hinton, 2009). We adopt ResNet50 (He et al., 2016) as image encoder for ImageNet and VGG16 (Simonyan & Zisserman, 2014) for CIFAR-10. More details of the datasets and instance encoders are given in Appendix B.

For text datasets with both validation and test sets provided, we tune hyper-parameters on the validation set and report results on the test set. For ImageNet/CIFAR-10, we tune hyper-parameters and report the performance on the validation/test set respectively, which is the norm in the literature. More details on experimental settings such as data processing and hyper-parameter tuning are included in Appendix C. We use Top-1 Error (equivalent to  $1 - \text{NDCG}@1$ , lower the better), Top-5 Error, and NDCG@5 (higher the better) for evaluation. All result points are multiplied by 100 for better illustration as commonly done in the literature.

In Section 5.1, we summarize the overall MCC performance of Rank4Class with respect to different configurations of loss functions and interaction patterns. In Section 5.2, we study different ranking losses in optimizing MCC tasks with respect to different evaluation metrics. In Section 5.3, we examine the effect of different interaction patterns between the instance and class embeddings. We discuss the effectiveness of ranking metrics in MCC evaluation in both sections.

### 5.1. Overall Performance of Rank4Class

Besides SoftmaxCE, PairLogLoss, and ApproxNDCG discussed in Section 4.1, we also include Gumbel-ApproxNDCG (Bruch et al., 2020) loss and MSE (Hui & Belkin, 2021) in experiments. We study different combinations between these five losses and the two interaction

Table 2. Results from classical MCC models and Rank4Class evaluated by Top-1/5 Error and NDCG@5. Bold font indicates the best value in each row. Relative Improvement means relative reduction in Top-1/5 Error and relative increase in NDCG@5.

Dataset	Encoder	Metrics	Classical MCC	Rank4Class	Relative Improvement
GoEmotions	BERT	Top-1 Error	41.00	<b>40.65</b>	0.85%
		Top-5 Error	13.35	<b>11.85</b>	11.24%
		NDCG@5	73.90	<b>74.62</b>	0.97%
	ELECTRA	Top-1 Error	38.45	<b>37.36</b>	2.83%
		Top-5 Error	10.45	<b>8.65</b>	17.22%
		NDCG@5	76.96	<b>78.30</b>	1.74%
MIND	BERT	Top-1 Error	30.90	<b>30.20</b>	2.27%
		Top-5 Error	6.15	<b>5.05</b>	17.89%
		NDCG@5	82.78	<b>83.56</b>	0.94%
	ELECTRA	Top-1 Error	27.09	<b>26.40</b>	2.55%
		Top-5 Error	4.25	<b>3.70</b>	12.94%
		NDCG@5	85.73	<b>86.13</b>	0.47%
ImageNet	ResNet50	Top-1 Error	23.74	<b>23.58</b>	0.67%
		Top-5 Error	6.90	<b>6.75</b>	2.17%
		NDCG@5	85.74	<b>85.85</b>	0.13%
CIFAR-10	VGG16	Top-1 Error	6.56	<b>6.40</b>	2.44%
		Top-5 Error	0.20	<b>0.15</b>	25.00%
		NDCG@5	97.19	<b>97.23</b>	0.04%

patterns introduced in Section 4.2 in the Rank4Class framework. In Table 2, we report the best performance of different configurations for Rank4Class under each metric in comparison to the baseline MCC models. The complete results of all combinations are included in Appendix D, where we use \* to mark the combinations that achieve the best performance in each task under each metric. As shown in the table, Rank4Class can improve the MCC performance in virtually all tasks through specific combinations of losses and interaction patterns. This shows the increased capacity from Rank4Class in MCC tasks evaluated on different metrics, which is achieved by adding more flexible design options in different components from LTR perspectives.

## 5.2. Effect of Ranking Losses

In this section, we study the use of ranking losses for optimizing MCC performance. In particular, we use PairLogLoss and ApproxNDCG as two most representative ranking losses in comparison to the SoftmaxCE. The results on other losses can be found in Appendix D. We only vary the loss function in the base Rank4Class structure in Figure 2 (Middle), so the ‘‘SoftmaxCE’’ method is the baseline that is equivalent to classical MCC models.

The results are shown in Table 3 and 4 for text and image classification tasks respectively. Overall, PairLogLoss or ApproxNDCG can outperform SoftmaxCE in nearly all tasks and metrics except for the Top-1 Error on GoEmo-

tions and MIND datasets with BERT model. Besides, PairLogLoss is generally good at reducing Top-5 Error than SoftmaxCE and ApproxNDCG, achieving the lowest Top-5 Error in all tasks. ApproxNDCG performs well on both Top-1 Error and NDCG@5 (achieves the top in four out of six tasks), which shows its effectiveness in directly optimizing NDCG metrics. Moreover, we see that the improvements from PairLogLoss and ApproxNDCG are more significant on text classification tasks than that on image classification tasks. Our observation is similar to that in (Hui & Belkin, 2021), which hypothesized the reason being that structures and parameters of popular image encoders are all heavily tuned with the SoftmaxCE.

Finally, we observe that different metrics are not always consistent in evaluating MCC tasks. For example, on GoEmotions with ELECTRA, PairLogLoss performs better than ApproxNDCG on Top-5 Error, while ApproxNDCG outperforms PairLogLoss on NDCG@5. This means that the PairLogLoss tends to rank the correct class in top 5 positions in more instances than ApproxNDCG, but ApproxNDCG can put the correct class relatively higher in the rankings. Furthermore, on MIND with BERT, Top-5 Error can not tell if PairLogLoss or ApproxNDCG is better. But NDCG@5 can successfully differentiate them since it also takes the absolute rank of the correct class in top 5 positions into account and thus is more informative.

Table 3. Results on text classification tasks trained with different losses.

Dataset	Encoder	Metrics	SoftmaxCE	PairLogLoss	ApproxNDCG
GoEmotions	BERT	Top-1 Error	<b>41.00</b>	41.79	41.42
		Top-5 Error	13.35	<b>12.00</b>	13.00
		NDCG@5	73.90	<b>74.27</b>	74.02
	ELECTRA	Top-1 Error	38.45	39.78	<b>37.67</b>
		Top-5 Error	10.45	<b>8.80</b>	9.25
		NDCG@5	76.96	77.17	<b>78.03</b>
MIND	BERT	Top-1 Error	<b>30.90</b>	31.21	30.97
		Top-5 Error	6.15	<b>5.25</b>	<b>5.25</b>
		NDCG@5	82.78	83.25	<b>83.36</b>
	ELECTRA	Top-1 Error	27.09	27.72	<b>26.65</b>
		Top-5 Error	4.25	<b>3.75</b>	4.30
		NDCG@5	85.73	85.74	<b>85.89</b>

Table 4. Results on image classification tasks trained with different losses.

Dataset	Encoder	Metrics	SoftmaxCE	PairLogLoss	ApproxNDCG
ImageNet	ResNet50	Top-1 Error	23.74	23.64	<b>23.62</b>
		Top-5 Error	6.90	<b>6.80</b>	6.85
		NDCG@5	85.74	<b>85.82</b>	85.80
CIFAR-10	VGG16	Top-1 Error	6.56	6.43	<b>6.41</b>
		Top-5 Error	0.20	<b>0.15</b>	<b>0.15</b>
		NDCG@5	97.19	97.21	<b>97.23</b>

### 5.3. Effect of Instance and Class Interactions

In this section, we study the effect of different interaction patterns between instance and class embeddings for producing ranking scores. Specifically, we use the two interaction patterns in section 4.2 with simple 2-layer MLPs.

We compare these two types of interactions with the dot-product baseline between the instance and class embeddings by fixing the loss function. In particular, we use ApproxNDCG in text classification tasks and SoftmaxCE in image classification tasks. The study of different interactions on other losses are included in Appendix D. The results are shown in Table 5 and 6 for text and image classification tasks. The results demonstrate that the two added interactions can outperform or achieve competitive performance than simple dot-product in all tasks and metrics. This shows the effectiveness of adding enhanced interactions based on the Rank4Class framework. In particular, latent-cross embedding tends to perform better than the concatenation of instance and class embeddings, achieving top performance in four out of six on Top-1 Error and five out of six tasks on NDCG@5. Again, we see that different evaluation metrics are not always consistent with each other, and NDCG@5 can be more informative than Top-5 Error, as observed and discussed in Section 5.2.

## 6. Related Work

Modern deep neural networks for MCC converge to the same recipe: given an input instance, a neural encoder is used to output its vector representation for generating scores on a set of classes. The vast research literature focus on developing more effective encoders in diverse domains, such as computer vision (He et al., 2016; Krizhevsky et al., 2012; Huang et al., 2017; Dosovitskiy et al., 2021; Tay et al., 2021a), natural language processing (Sun et al., 2019; Tay et al., 2021b; Minaee et al., 2021), and automatic speech recognition (Moritz et al., 2019), among others. The softmax cross entropy loss is the dominant loss function discussed in these papers. Only very recently, Hui & Belkin (2021) study the mean squared error as another loss for MCC. Our work is orthogonal to the extensive research on neural encoders in that we provide a new formulation from the LTR perspective. Such a perspective inspires more diverse loss functions and model architectures that can model interactions between inputs and classes more effectively.

Learning to rank (LTR) is a long-established interdisciplinary research area at the intersection of machine learning and information retrieval (Liu, 2009). Neural rankers are dominating in ranking virtually all modalities recently, including text ranking (Lin et al., 2020), image

Table 5. Results on text classification tasks of different interactions trained with ApproxNDCG.

Dataset	Encoder	Metrics	dot product	LC+MLP	Concat+MLP
GoEmotions	BERT	Top-1 Error	41.42	<b>40.65</b>	40.92
		Top-5 Error	13.00	12.40	<b>12.20</b>
		NDCG@5	74.02	<b>74.62</b>	74.61
	ELECTRA	Top-1 Error	<b>37.67</b>	<b>37.67</b>	37.80
		Top-5 Error	9.25	<b>8.65</b>	9.15
		NDCG@5	78.03	<b>78.30</b>	77.97
MIND	BERT	Top-1 Error	30.97	30.81	<b>30.77</b>
		Top-5 Error	5.25	5.25	<b>5.05</b>
		NDCG@5	83.36	<b>83.43</b>	83.35
	ELECTRA	Top-1 Error	26.65	26.74	<b>26.40</b>
		Top-5 Error	4.30	<b>3.70</b>	4.30
		NDCG@5	85.89	<b>86.13</b>	86.05

Table 6. Results on image classification tasks of different interactions trained with SoftmaxCE.

Dataset	Encoder	Metrics	dot product	LC+MLP	Concat+MLP
ImageNet	ResNet50	Top-1 Error	23.74	<b>23.62</b>	23.66
		Top-5 Error	6.90	6.80	<b>6.75</b>
		NDCG@5	85.74	<b>85.84</b>	85.83
CIFAR-10	VGG16	Top-1 Error	6.56	<b>6.40</b>	6.43
		Top-5 Error	<b>0.20</b>	<b>0.20</b>	<b>0.20</b>
		NDCG@5	97.19	97.22	<b>97.23</b>

retrieval (Gordo et al., 2016), and tabular data ranking (Qin et al., 2021). Many LTR papers focus on more effective loss functions (Qin et al., 2010; Bruch et al., 2020) to rank items with respect to a query. One of the focuses of this paper is to introduce new techniques stemmed from LTR to solve MCC problems.

LTR has been explored in multi-label classification (MLC) (Agrawal et al., 2013; Zhang et al., 2018; Jain et al., 2019), where each instance can belong to multiple classes. MLC is generally treated as a different problem from classical MCC: the number of labels assigned to an instance could be arbitrary and one research focus is to decide the threshold to cutoff the prediction list. For example, (Azarbondy et al., 2021) focuses on feature engineering in learning a linear scoring function for ranking classes, and studies how the top-k threshold affect the MLC performance. Instead, we considers the recent neural ranking models for MCC, and do not consider multi-label scenario. An important related direction in MLC is extreme classification (EC) (Zhang et al., 2018), which formulates document retrieval or recommendation as MLC problems, where each document/item is treated as a candidate class. EC considers a huge number of classes and the main focus is on the scalability and efficiency (Weston et al., 2011; 2013). LTR techniques have

been proposed for MLC in the past (Yang & Gopal, 2012). In contrast, standard MCC has a smaller number of classes and each instance has a single correct class label. LTR have not been well studied for MCC, and we are the first to propose a unified ranking formulation for it.

## 7. Conclusion

In this paper we examine the classical MCC problem through the lens of ranking. Such a perspective brings benefits to MCC from three aspects: ranking metrics, ranking losses, and ranking architectures. We first show that ranking metrics can be more informative for MCC evaluations. Then, in the deep learning setting, we show an equivalent view of MCC in LTR setting. Such a connection provides new perspectives for MCC with respect to loss functions and model architectures. We studies these new formulations of MCC on various datasets and observe promising results.

Our work opens up several research directions. First, the new ranking architectures allow to take more class information such as class metadata into account and it is interesting to study how this additional information can improve MCC. Second, it is also possible to apply the proposed framework Rank4Class to binary classification. Third, classes are usu-



ally not independent and our framework can incorporate the relationship between classes into the MCC through attention mechanisms, which is worth studying.

## References

- Agrawal, R., Gupta, A., Prabhu, Y., and Varma, M. Multi-label learning with millions of labels: Recommending advertiser bid phrases for web pages. In *Proceedings of the 22nd International Conference on World Wide Web, WWW '13*, New York, NY, USA, 2013. Association for Computing Machinery. ISBN 9781450320351.
- Aly, M. Survey on multiclass classification methods. *Neural Netw*, 19:1–9, 2005.
- Azarbonyad, H., Dehghani, M., Marx, M., and Kamps, J. Learning to rank for multi-label text classification: Combining different sources of information. *Natural Language Engineering*, 2021.
- Beutel, A., Covington, P., Jain, S., Xu, C., Li, J., Gatto, V., and Chi, E. H. Latent cross: Making use of context in recurrent recommender systems. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, 2018.
- Bruch, S., Wang, X., Bendersky, M., and Najork, M. An analysis of the softmax cross entropy loss for learning-to-rank with binary relevance. In *Proceedings of the 2019 ACM SIGIR International Conference on Theory of Information Retrieval*, pp. 75–78, 2019.
- Bruch, S., Han, S., Bendersky, M., and Najork, M. A stochastic treatment of learning to rank scoring functions. In *Proceedings of the 13th International Conference on Web Search and Data Mining*, pp. 61–69, 2020.
- Burges, C., Shaked, T., Renshaw, E., Lazier, A., Deeds, M., Hamilton, N., and Hullender, G. Learning to rank using gradient descent. In *Proceedings of the 22nd International Conference on Machine Learning*, pp. 89–96, 2005.
- Burges, C., Ragno, R., and Le, Q. Learning to rank with nonsmooth cost functions. In *Advances in Neural Information Processing Systems*, 2006.
- Burges, C. J. C. From RankNet to LambdaRank to LambdaMART: An overview. Technical report, Microsoft Research, 2010.
- Clark, K., Luong, M.-T., Le, Q. V., and Manning, C. D. ELECTRA: Pre-training text encoders as discriminators rather than generators. In *International Conference on Learning Representations*, 2020.
- Crammer, K. and Singer, Y. On the algorithmic implementation of multiclass kernel-based vector machines. *J. Mach. Learn. Res.*, 2:265–292, March 2002.
- Demszky, D., Movshovitz-Attias, D., Ko, J., Cowen, A., Nemade, G., and Ravi, S. GoEmotions: A Dataset of Fine-Grained Emotions. In *58th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, 2019.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations*, 2021.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(61): 2121–2159, 2011.
- Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, 2016.
- Gordo, A., Almazán, J., Revaud, J., and Larlus, D. Deep image retrieval: Learning global representations for image search. In *European Conference on Computer Vision*, pp. 241–257. Springer, 2016.
- Guo, J., Fan, Y., Pang, L., Yang, L., Ai, Q., Zamani, H., Wu, C., Croft, W. B., and Cheng, X. A deep look into neural ranking models for information retrieval. *Information Processing & Management*, 57(6):102067, 2020.
- Hastie, T., Tibshirani, R., and Friedman, J. *The Elements of Statistical Learning*. 2001.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2261–2269, 2017.
- Hui, L. and Belkin, M. Evaluation of neural architectures trained with square loss vs cross entropy in classification tasks. In *International Conference on Learning Representations*, 2021.

- Jain, H., Balasubramanian, V., Chunduri, B., and Varma, M. Slice: Scalable linear extreme classifiers trained on 100 million labels for related searches. In *WSDM '19, February 11–15, 2019, Melbourne, VIC, Australia*. ACM, February 2019.
- Järvelin, K. and Kekäläinen, J. Cumulated gain-based evaluation of ir techniques. *ACM Trans. Inf. Syst.*, pp. 422–446, October 2002.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical report, University of Toronto, Toronto, Ontario, 2009.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems*, pp. 1097–1105, 2012.
- Li, H. and Xu, J. Semantic matching in search. *Foundations and Trends in Information Retrieval*, 7(5):343–469, 2014.
- Lin, J., Nogueira, R., and Yates, A. Pretrained transformers for text ranking: BERT and beyond. *arXiv preprint arXiv:2010.06467*, 2020.
- Liu, T.-Y. Learning to rank for information retrieval. *Found. Trends Inf. Retr.*, 3(3):225–331, March 2009.
- Minaee, S., Kalchbrenner, N., Cambria, E., Nikzad, N., Chenaghlu, M., and Gao, J. Deep learning-based text classification: A comprehensive review. *ACM Computing Surveys (CSUR)*, 54(3):1–40, 2021.
- Moritz, N., Hori, T., and Le Roux, J. Triggered attention for end-to-end speech recognition. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5666–5670, 2019.
- Pasumathi, R. K., Bruch, S., Wang, X., Li, C., Bendersky, M., Najork, M., Pfeifer, J., Golbandi, N., Anil, R., and Wolf, S. Tf-ranking: Scalable tensorflow library for learning-to-rank. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.
- Qin, T., Liu, T.-Y., and Li, H. A general approximation framework for direct optimization of information retrieval measures. *Information Retrieval*, 13(4):375–397, 2010.
- Qin, Z., Yan, L., Zhuang, H., Tay, Y., Pasumathi, R. K., Wang, X., Bendersky, M., and Najork, M. Are neural rankers still outperformed by gradient boosted decision trees? In *Proceedings of the 9th International Conference on Learning Representations*, 2021.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Sokolova, M. and Lapalme, G. A systematic analysis of performance measures for classification tasks. *Information Processing and Management*, 45(4):427–437, 2009.
- Sun, C., Qiu, X., Xu, Y., and Huang, X. How to fine-tune bert for text classification? In *China National Conference on Chinese Computational Linguistics*, pp. 194–206. Springer, 2019.
- Tay, Y., Dehghani, M., Aribandi, V. K., Gupta, J. P., Pham, P., Qin, Z., Bahri, D., Juan, D.-C., and Metzler, D. Omninet: Omnidirectional representations from transformers. In *International Conference on Machine Learning*, 2021a.
- Tay, Y., Dehghani, M., Gupta, J., Bahri, D., Aribandi, V., Qin, Z., and Metzler, D. Are pre-trained convolutions better than pre-trained transformers? *arXiv preprint arXiv:2105.03322*, 2021b.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.
- Wang, X., Li, C., Golbandi, N., Bendersky, M., and Najork, M. The lambdaloss framework for ranking metric optimization. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pp. 1313–1322, 2018.
- Weston, J., Bengio, S., and Usunier, N. Wsabie: Scaling up to large vocabulary image annotation. In *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence - Volume Volume Three, IJCAI'11*. AAAI Press, 2011.
- Weston, J., Makadia, A., and Yee, H. Label partitioning for sublinear ranking. In *Proceedings of the 30th International Conference on Machine Learning*, Proceedings of Machine Learning Research, Atlanta, Georgia, USA, 2013. PMLR.
- Wu, F., Qiao, Y., Chen, J.-H., Wu, C., Qi, T., Lian, J., Liu, D., Xie, X., Gao, J., Wu, W., and Zhou, M. MIND: A large-scale dataset for news recommendation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020.

Yang, Y. and Gopal, S. Multilabel classification with meta-level features in a learning-to-rank framework. *Machine Learning*, 88(1-2):47–68, 2012.

Zhang, A., Lipton, Z. C., Li, M., and Smola, A. J. Dive into deep learning. *arXiv preprint arXiv:2106.11342*, 2021.

Zhang, W., Yan, J., Wang, X., and Zha, H. Deep extreme multi-label learning. In *Proceedings of the 2018 ACM on International Conference on Multimedia Retrieval*, pp. 100–107, 2018.

## A. Proof of Theorem 1

**Theorem 1.** *Given an instance with class labels  $\mathbf{y} \in \{0, 1\}^n$  and  $\sum_{c=1}^n \mathbf{y}[c] = 1$ , for any  $1 < K \leq n$ ,  $NDCG@K$  preserves more information than Top- $K$  Accuracy in evaluating the class rankings.*

*Proof.* Consider  $NDCG@K$  and Top- $K$  Accuracy as functions of any input class ranking  $\pi \in \Pi$ , where  $\Pi$  is the space of all possible rankings (of size  $n!$ ). Let  $P$  be a distribution over  $\Pi$  from which the evaluation set is sampled. Let the probability of the correct class being at position  $i$  in the evaluation set be  $p_i, p_i \geq 0, i \in [n]$ , and  $\sum_{i=1}^n p_i = 1$ . Thus  $\mathbb{E}[p_i] = \sum_{\mathcal{Y}[\pi(i)]=1} P(\pi)$ , where  $P(\pi)$  is the probability that a class ranking  $\pi$  being sampled from  $\Pi$ .

Note that  $NDCG@K$  has  $K$  distinct values corresponding to the correct class being ranked at the top  $K$  positions plus 0 for not in top  $K$  positions. Therefore, the entropy of the values of  $NDCG@K$  can be obtained as:

$$H(NDCG@K) = - \sum_{i=1}^K p_i \log p_i - p_- \log p_-, \quad (5)$$

where  $p_- = \sum_{i=K+1}^n p_i$  is the probability that the correct class is not ranked in top  $K$  positions.  $p_- \log p_- = 0$  when  $p_- = 0$  as a convention in information theory.

However, Top- $K$  Accuracy only has two possible values: 1 when the correct class is in top  $K$  positions, and 0 otherwise. Hence the entropy of Top- $K$  Accuracy values is

$$H(\text{acc}@K) = -p_+ \log p_+ - p_- \log p_-, \quad (6)$$

where  $p_+ = \sum_{i=1}^K p_i$  is the probability that the correct class is ranked in top  $K$  positions.

Finally, we have

$$\begin{aligned} & H(NDCG@K) - H(\text{acc}@K) \\ &= - \sum_{i=1}^K p_i \log p_i + p_+ \log p_+ \\ &= \sum_{i=1}^K p_i (\log p_+ - \log p_i) \geq 0. \end{aligned} \quad (7)$$

Therefore,  $NDCG@K$  has more information than Top- $K$  Accuracy in terms of the entropy of the evaluation results.  $\square$

## B. Datasets and Encoders

### B.1. Datasets

- **GoEmotions** (Demszky et al., 2020) is the largest manually annotated dataset for fine-grained emotion classification. It contains 58k English Reddit comments, labeled with 27 emotion categories and a default ‘‘Neutral’’ category. On top of the raw data, the authors also provided a version that only includes comments with two or more raters agreeing on at least one label, split into train/test/validation sets. In our experiments, we filter out comments with more than 1 emotion labels in all three sets.
- **MIND** (Wu et al., 2020) is a large-scale dataset for news recommendation. MIND contains about 130k English news articles. Every news article contains rich textual content including title, abstract, body, category and entities. We use the concatenation of title and abstract as the content of the news and use the category of the news as the classification label. We split all news instances into train/test/validation sets by roughly 60/20/20 for experiments.
- **ImageNet** (Krizhevsky et al., 2012) is an image dataset with around 1.28 million images in the training set and 50k images in the validation set. Each image is labeled by one of 1,000 classes. The images are cropped and resized to  $224 \times 224 \times 3$  pixels as the input following the preprocessing in [https://github.com/tensorflow/models/tree/master/official/vision/image\\_classification/resnet](https://github.com/tensorflow/models/tree/master/official/vision/image_classification/resnet). We follow the common practices of using ImageNet and report the results on the validation set.
- **CIFAR-10** (Krizhevsky & Hinton, 2009) contains 50k training images and 10k test images. There are 10 classes and each class has 6k images. All images have the same size of  $32 \times 32 \times 3$ .



## B.2. Instance Encoders

- **BERT** (Vaswani et al., 2017) is a model based on transformers pretrained on a large corpus of English data. We use the BERT-Base-uncased model from <https://github.com/google-research/bert> for finetuning in our experiments. The maximum sequence length is set to 32 for GoEmotions and 128 for MIND.
- **ELECTRA** (Clark et al., 2020) is another pretrained transformer model. We use the ELECTRA-Base-uncased model from <https://github.com/google-research/electra> in our experiments. The maximum sequence length is set in the same way as BERT.
- **ResNet50** (He et al., 2016) is a 50 layers deep convolutional neural network with residual connections. We use the implementation of ResNet50 from the tensorflow official models at [https://github.com/tensorflow/models/tree/master/official/vision/image\\_classification/resnet](https://github.com/tensorflow/models/tree/master/official/vision/image_classification/resnet). We adopt the pretrained network on ImageNet and finetune it in our Rank4Class framework.
- **VGG16** (Simonyan & Zisserman, 2014) is a convolutional neural network with 16 layers. We use the implementation of VGG16 from <https://github.com/geifmany/cifar-vgg> with input size  $32 \times 32$  in our experiments on CIFAR-10.

## C. Experimental Settings

We provide the implementation and hyper-parameter tuning details for reproducibility. For text classification tasks, we tokenize the raw sentences into word ids based on BERT vocabulary, and create the input masks and segment ids following standard BERT input formats. For MIND, we concatenate the title and abstract of each news by adding a “[SEP]” token between them. ELECTRA uses exactly the same input formats as BERT. The data processing of image datasets follow the standard methods as given in the original papers. We implement the Rank4Class pipeline based on TF-ranking (Pasumarthi et al., 2019). We adopt TF-ranking’s implementation of SoftmaxCE, MSE and all the ranking losses.

For instance encoders, we use the default hyper-parameters suggested in the original implementations listed in Appendix B without further tuning. We use the same number of hidden units in the two layers of MLP for interactions, and tune the number of hidden units in  $\{64, 128, 256, 512\}$ . For all experiments, we use Adam (Kingma & Ba, 2014) and Adagrad (Duchi et al., 2011) as optimizers for training models. We tune the initial learning rate for all experiments in the range of  $1e^{-7}$  to 0.1 with a multiplicative step size of 3. Adam has a slight edge over Adagrad in most experiments and the best learning rate for Adam are  $3e^{-6}$  and  $1e^{-5}$  for most experiments. We use a batch size of 32 for text classification tasks, and a batch size of 64 for image classification tasks. For all configurations of the models and datasets, we train the model for 100,000 steps in text classification tasks, and 50 epochs in image classification tasks. We pick the best checkpoint on the validation set (if provided) for evaluation.

## D. Results on different combinations of ranking losses and interaction patterns

Here we provide the experiment results on different combinations of loss functions and interaction patterns. Besides PairLogLoss and ApproxNDCG included in Section 5, we also include results from Gumbel-ApproxNDCG and MSE here. Note that we use the rescaled MSE as suggested in (Hui & Belkin, 2021) on ImageNet dataset, since there is only 1 correct class in 1,000 classes, and regular MSE performs poorly on such imbalanced data. In particular, as the number of combinations is large, we display the results with respect to each dataset and instance encoder in each table for better visualization and comparison. Then in each table, we group the results from different combinations of loss functions and interaction patterns according to the three evaluation metrics. We underline the top 3 combinations (ties are included) under each metric, and use \* to mark the best. The results of all datasets and instance encoders are shown in Tables 7 - 12.

Firstly, the conclusions on loss functions, interactions and evaluation metrics from the tables are the same as those in Sections 5.2 and 5.3. The Gumbel-ApproxNDCG performs well when more complicated interactions such as latent cross and concatenating embeddings are applied. MSE achieves the best top-1 Error in GoEmotions with BERT and CIFAR-10 with VGG16. It indicates that MSE can be effective in optimizing top-1 metrics in certain tasks, which aligns with the observation in Hui & Belkin (2021). However, MSE is not suitable for direct application in problems with highly imbalanced correct and incorrect classes as on the ImageNet classification task. A few rescaling techniques with hyper-parameters need to be applied for MSE to perform properly, which also increase the burden of hyper-parameters tuning. Besides, adding different combinations of losses and interaction patterns can always improve the performance, as indicated by \* under each

**Rank4Class: A Ranking Formulation for Multiclass Classification**

Table 7. Results on GoEmotions with BERT as text encoder.

GoEmotions + BERT						
Metric	Interaction	SoftmaxCE	PairLogLoss	ApproxNDCG	Gumbel ApproxNDCG	MSE
Top-1 Error	dot product	41.00	41.79	41.42	41.11	41.20
	LC+MLP	41.35	41.98	<u>40.65*</u>	<u>40.70</u>	41.31
	Concat+MLP	41.02	41.55	<u>40.92</u>	<u>41.59</u>	41.63
Top-5 Error	dot product	13.35	<u>12.00</u>	13.00	13.80	15.00
	LC+MLP	12.85	<u>12.25</u>	12.40	12.85	14.30
	Concat+MLP	12.40	<u>11.85*</u>	<u>12.20</u>	12.90	13.85
NDCG@5	dot product	73.90	74.27	74.02	73.71	73.00
	LC+MLP	74.11	74.22	<u>74.62*</u>	<u>74.53</u>	73.50
	Concat+MLP	74.36	74.19	<u>74.61</u>	74.05	73.55

Table 8. Results on GoEmotions with ELECTRA as text encoder.

GoEmotions + ELECTRA						
Metric	Interaction	SoftmaxCE	PairLogLoss	ApproxNDCG	Gumbel ApproxNDCG	MSE
Top-1 Error	dot product	38.45	39.78	<u>37.67</u>	38.63	37.97
	LC+MLP	38.69	39.13	<u>37.67</u>	38.45	<u>37.36*</u>
	Concat+MLP	39.00	39.48	<u>37.80</u>	38.06	38.15
Top-5 Error	dot product	10.45	<u>8.80</u>	9.25	10.65	10.30
	LC+MLP	9.30	<u>8.90</u>	<u>8.65*</u>	10.00	10.15
	Concat+MLP	9.75	<u>8.70</u>	<u>9.15</u>	9.90	10.40
NDCG@5	dot product	76.96	77.17	<u>78.03</u>	76.72	77.37
	LC+MLP	77.45	77.49	<u>78.30*</u>	77.41	77.69
	Concat+MLP	77.06	77.47	<u>77.97</u>	77.35	77.24

metric. Note that we did not further tune the hyper-parameters of instance encoders for different architectures of Rank4Class. It is possible that the hyper-parameters of encoders are more suitable for the classical MCC models, and fine-tuning may further boost the performance of Rank4Class. Last but not least, we see that NDCG is more informative in evaluating MCC performance than Top-5 Error which creates many ties. For example, In Table 9 for MIND with BERT, Top-5 Error fails to find the best performing model, while NDCG@5 can successfully differentiate them.

**Rank4Class: A Ranking Formulation for Multiclass Classification**

*Table 9. Results on MIND with BERT as text encoder.*

MIND + BERT						
Metric	Interaction	SoftmaxCE	PairLogLoss	ApproxNDCG	Gumbel ApproxNDCG	MSE
Top-1 Error	dot product	30.90	31.21	30.97	30.88	30.71
	LC+MLP	<u>30.20*</u>	<u>30.51</u>	30.81	30.79	30.56
	Concat+MLP	<u>30.51</u>	31.29	30.77	30.92	30.71
Top-5 Error	dot product	6.15	<u>5.25</u>	<u>5.25</u>	5.75	6.95
	LC+MLP	5.35	5.30	<u>5.25</u>	5.45	7.15
	Concat+MLP	5.55	<u>5.25</u>	<u>5.05*</u>	5.65	7.30
NDCG@5	dot product	82.78	83.25	83.36	83.11	82.59
	LC+MLP	<u>83.56*</u>	<u>83.54</u>	<u>83.43</u>	83.28	82.48
	Concat+MLP	83.31	83.27	83.35	83.16	82.37

*Table 10. Results on MIND with ELECTRA as text encoder.*

MIND + ELECTRA						
Metric	Interaction	SoftmaxCE	PairLogLoss	ApproxNDCG	Gumbel ApproxNDCG	MSE
Top-1 Error	dot product	27.09	27.72	26.65	26.77	26.99
	LC+MLP	26.89	27.79	26.74	<u>26.55</u>	26.92
	Concat+MLP	26.78	27.53	<u>26.40*</u>	<u>26.55</u>	26.65
Top-5 Error	dot product	4.25	<u>3.75</u>	4.30	4.25	5.35
	LC+MLP	4.04	<u>3.75</u>	<u>3.70*</u>	3.90	5.20
	Concat+MLP	4.25	<u>3.70*</u>	4.30	3.95	5.55
NDCG@5	dot product	85.73	85.74	85.89	85.89	85.26
	LC+MLP	85.67	85.71	<u>86.13*</u>	<u>86.09</u>	85.31
	Concat+MLP	85.81	85.85	<u>86.05</u>	85.99	85.25

*Table 11. Results on ImageNet with ResNet50 as image encoder.*

ImageNet + ResNet50						
Metric	Interaction	SoftmaxCE	PairLogLoss	ApproxNDCG	Gumbel ApproxNDCG	MSE (rescaled)
Top-1 Error	dot product	23.74	23.64	23.62	23.65	26.46
	LC+MLP	23.62	23.70	23.65	<u>23.60</u>	24.18
	Concat+MLP	23.66	23.79	<u>23.61</u>	<u>23.58*</u>	23.83
Top-5 Error	dot product	6.90	6.80	6.85	6.80	8.95
	LC+MLP	6.80	<u>6.75*</u>	6.80	6.80	6.90
	Concat+MLP	<u>6.75*</u>	<u>6.75*</u>	6.80	6.80	6.80
NDCG@5	dot product	85.74	85.82	85.80	85.82	83.24
	LC+MLP	<u>85.84</u>	85.80	85.82	<u>85.84</u>	85.44
	Concat+MLP	85.83	85.77	85.82	<u>85.85*</u>	85.72

Table 12. Results on CIFAR-10 with VGG16 as image encoder.

CIFAR-10 + VGG16						
Metric	Interaction	SoftmaxCE	PairLogLoss	ApproxNDCG	Gumbel ApproxNDCG	MSE
Top-1 Error	dot product	6.56	6.43	6.41	6.45	6.48
	LC+MLP	<u>6.40*</u>	6.42	6.46	<u>6.40*</u>	6.42
	Concat+MLP	6.43	6.44	6.44	6.47	<u>6.40*</u>
Top-5 Error	dot product	0.20	<u>0.15*</u>	<u>0.15*</u>	0.20	<u>0.15*</u>
	LC+MLP	0.20	<u>0.15*</u>	0.25	0.20	0.25
	Concat+MLP	0.20	0.20	0.25	<u>0.15*</u>	0.35
NDCG@5	dot product	97.19	97.21	<u>97.23*</u>	97.21	97.16
	LC+MLP	<u>97.22</u>	97.21	97.17	97.17	97.14
	Concat+MLP	<u>97.23*</u>	97.21	97.19	<u>97.22</u>	97.16