

1 We thank all reviewers for the constructive comments. Due to the space limitation, we only address major review
2 concerns and will incorporate other suggestions in a revised version of the current work.

3 **R1-Q1: Motivation of using the SoRR loss.** Traditional approaches to handle outliers focus on the design of robust
4 *individual losses* that apply to individual training samples, notable examples include the Huber loss and the capped
5 hinge loss as pointed out by the reviewer, or give pre-determined weights to individual training samples. We take a
6 different approach to this work. Instead of changing the definition of individual losses, our method builds in robustness
7 at the aggregate loss level, using the AoRR loss. The resulting learning algorithm is more flexible and allows the user to
8 choose an individual loss form that is relevant to the learning problem. On the other hand, the weights on the individual
9 training samples under AoRR loss is determined automatically.

10 **R1-Q2: Related work of (Chang, Yu, & Yang, KDD 2017).** This work addresses the outliers for multiclass SVM by
11 using a hyperparameter to cap the values of the individual losses. This approach is *different* from ours since we directly
12 address the original top- k multiclass SVM problem using our TKML loss *without* introducing this convex surrogate.
13 Furthermore, it has been shown that the loss used in this work is not multi-class top- k consistent as shown in [34], while
14 our TKML loss is consistent (see lines 265-270).

15 **R1-Q3: Comparison with maximum loss.** Although the sensitivity to outliers of the maximum loss can be alleviated
16 with the tricks used in [27], our work on the AoRR aggregate loss provides a general mechanism to exclude the influence
17 of potential outliers in the training data. We use the straightforward maximum loss, as it is a special case of the AoRR
18 loss ($m = 0, k = 1$) and can better reflect the continuum of performance change while we adjust k and m . As such, we
19 do not want to give any special treatment to this special case of AoRR loss.

20 **R1-Q4: Performance comparison with the AT_k loss.** The advantage of the AoRR aggregate loss is its improved
21 robustness to outliers. On datasets that do not have a significant fraction of outliers, its performance is expected to be
22 similar to the AT_k loss. On the other hand, when the dataset contains outliers, AoRR can significantly outperform
23 existing aggregate losses, as demonstrated on datasets Monk and Phoneme (Table 1, with an average difference/standard
24 deviation in the accuracy of 4.07%/0.05% and 2.22%/0.18%, respectively).

25 **R1-Q5: Value of k when comparing with the AT_k loss.** Experiments on the synthetic
26 data aim to show that the AoRR loss can exclude the influence of outliers. To exemplify
27 this case, we chose to have one outlier in the synthetic data, and correspondingly, use
28 $k = 2$. Furthermore, we perform an additional experiment to the request of R1 showing
29 misclassification rate w.r.t different values of k in Figure 1. This result shows that for
30 values other than 2, the AoRR loss still exhibits an advantage over the AT_k loss.

31 **R1-Q6: Computation complexity of optimizing AoRR.** For large datasets, the original
32 DC algorithm may take a longer time to run. However, we can use a stochastic
33 version of DCA to optimize our problem according to Thi, Hoai An Le, et al. "Stochastic
34 DCA for minimizing a large sum of DC functions with application to Multi-class Logistic
35 Regression." arXiv preprint arXiv:1911.03992 (2019). The authors have proved
36 it is much more efficient and the convergence of it to a critical point is guaranteed with
37 probability 1.

38 **R1-Q7: Grid search for m and k .** Our experiments are based on a grid search for
39 selecting the value of k and m which may not be ideal in practice for large-scale
40 datasets. For this reason, we provided a well-known adaptive setting method there, and its feasibility had been discussed
41 in [13]. Since k and m are hyper-parameters which are chosen by cross-validation (CV), we did not include the time of
42 searching for them into the overall time complexity, as commonly practiced in machine learning.

43 **R1-Q8: Strategy to choose k .** We suggested a strategy to choose k , but for the experiments we performed, we found
44 that a grid search is simpler and often yields comparable performance.

45 **R2-Q1: Novelty of this work.** Although Theorem 1 is related with Lemma 1 in [10], there are some fundamental
46 differences. The average top- k function in [10] is convex with regards to its inputs. In our work, we show that the SoRR
47 can be expressed as the difference of two sum of top- k functions, which shows that it is not convex but can be solved
48 with the DC algorithm. This connection has not been studied previously. In addition, based on the SoRR approach, we
49 describe in details of new types of aggregate loss that are more robust to the presence of outliers in training data, and
50 new types of individual losses for multi-class multi-label learning.

51 **R4-Q1: Related work of (Nie, et al., AAI 2017).** The capped hinge loss in this work is designed as an individual
52 loss. Since the AoRR loss is at the aggregate level, we did not compare it as related work.

53 **R4-Q2: Related work of (Z. Zhou, NSR 2018).** The outliers in Figure 3 are different from the case of inaccurate
54 supervision as described in (Z.-H Zhou, NSR, 2018), which assumes errors in the labels of training data. In our study of
55 aggregate losses, we assume the training labels are *accurate* for all training samples, including the outliers.

56 **R4-Q3: Value of k when comparing with the AT_k loss.** Please refer to our answers to R1-Q5.

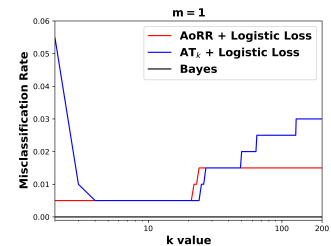


Figure 1: Misclassification rates of the AT_k loss (blue) and AoRR loss (red) for different k values on a balanced but multi-modal synthetic dataset.