

1 We appreciate all the reviewers for the constructive comments, and our responses are as below:

2 **R1.1: concerns on compared methods and datasets.** We fully understand the concern about our baselines since
3 we are the first to improve certified robustness of metric learning. Therefore, as Reviewer 4 suggested, we add
4 experiments comparing with neural networks certification methods, including ordinary neural networks certified by
5 CROWN [48] and randomized-smoothing neural networks [11]. The results are shown in Figure i. Our method
6 outperforms randomized-smoothing networks for a wide range of ℓ_2 perturbations, and we will add these baselines to
7 the paper. As mentioned in the review, these two datasets are often used for evaluating robustness of neural networks, so
8 surpassing a very recently proposed defensive network (RandSmooth) suggests the potential high impact of this work.
9 Regarding the datasets, such UCI datasets are actually commonly used for performance evaluation in previous metric
10 learning papers in various settings (Perrot and Habrard, NeurIPS’15; Zadeh *et al.*, ICML’16; Chen *et al.*, IJCAI’18).

11 **R2.1: computational cost.** In general, computational
12 cost is not an issue for ARML. The average runtime
13 (of 5 trials) of LMNN, ITML, NCA and ARML, are
14 66.4s, 95.2s, 480.9s and 146.1s respectively, on USPS
15 for 100 iterations. To make the comparison fair, all of
16 the methods are run on CPU (Xeon(R) E5-2620 v4 @
17 2.10GHz). In fact, ARML is highly parallelizable and
18 our implementation also supports GPU: when running
19 on GPU (one Nvidia TITAN Xp), the average runtime of
20 ARML is only 10.6s. We will add detailed discussions.

21 **R2.2: classification error.** We did report the classifi-
22 cation error. The classification error is equivalent to the
23 empirical robust error at radius 0 in both Table 1 and Table 2. (Note that the certified robust error in Table 1 is equivalent
24 to the empirical robust error.) One of the major advantages of ARML over existing metric learning methods is that it
25 could improve the classification error and robust error simultaneously. We will highlight them in these two tables.

26 **R2.3 & R3.1: missing papers and existing metric learning methods against adversarial examples.** Thanks for
27 mentioning these related work. Both Chen *et al.* (2018) and Duan *et al.* (2018) introduced the adversarial framework
28 as a mining strategy aiming to improve classification accuracy of metric learning, and Zheng *et al.* (2019) made
29 an improvement upon them; Li *et al.* (2019) proposed an attack method for K -NN via a differentiable substitute;
30 Mao *et al.* (2019) used a metric learning loss as a regularization to improve *empirical* robustness of deep learning
31 models; A concurrent paper (Yang *et al.*, 2020), made public after the NeurIPS deadline, proposed a similar training
32 method to ours, but still lacks formal robustness verification like our Theorem 1 and Algorithm 2, and therefore did
33 not *certify* robustness formally. In contrast, as we underscore in our paper, we propose the first *certified* defense based
34 on metric learning. We will add a paragraph in the related work section to talk further about adversarial robustness of
35 metric learning with additional references.

36 **R3.2: loss function in Eq. (12).** The loss functions in Eq. (11) and Eq. (12) are the same to the one in Eq. (5): a
37 monotonically non-increasing function. In our implementation, we simply use the “negative” loss, i.e., $\ell(\epsilon) = -\epsilon$. We
38 will make it clear in the paper.

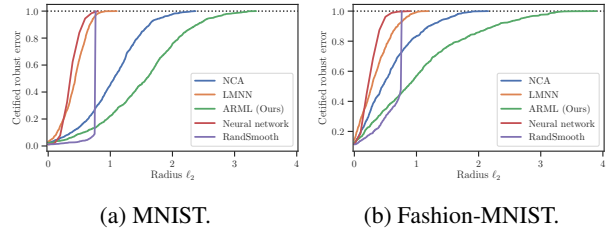
39 **R3.3: efficiency issue concerning sampling instances.** Sampling from neighborhood instead of computing k th max
40 and k th min is indeed a crucial technique of ARML to help improve training efficiency. We will clarify it.

41 **R4.1: compare with robust deep learning.** We have compared with robust neural networks (see R1.1).

42 **R4.2: Eq. (6) and Eq. (13).** It is noteworthy that Eq. (6) involves two layers of optimization, and it is correct that
43 \mathbb{I} and \mathbb{J} are fixed for the inner optimization. The constraints assure that *at most* $k - 1 = (K + 1)/2 - 1$ “positive”
44 training instances, i.e., instances in \mathbb{I} , are in the K -size neighborhood of $\mathbf{x}_{\text{test}} + \delta_{(\mathbb{I}, \mathbb{J})}$, by means of *not* constraining
45 $(\mathbf{x}_{\text{test}} + \delta_{(\mathbb{I}, \mathbb{J})})$ ’s distances from these “positive” training instances; it is a *necessary* condition for a successful attack.
46 Interestingly and not surprisingly, for Mahalanobis 1-NN — the case where we have $K = 1$ and $k = (K + 1)/2 = 1$,
47 and hence \mathbb{I} is an empty set, and \mathbb{J} has only one element — Eq. (6) is exactly reduced to Eq. (13).

48 **R4.3: selection of \mathbf{x}^+ in Eq. (7) (or \mathbf{x}_i in Eq. (10)).** It is an insightful question. In theory, enumerating every
49 $\{i : y_i = y_{\text{test}}\}$ in Eq. (10) will derive the *tightest* lower bound of the minimal adversarial perturbation. Nevertheless, in
50 practice, only selecting a subset of $\{i : y_i = y_{\text{test}}\}$ of which \mathbf{x}_i is “close” to \mathbf{x}_{test} usually suffices to derive a satisfying
51 lower bound. In the robustness evaluation phase, to derive certification bounds as tight as possible, we did not employ
52 this strategy. We will add more discussions for Theorem 1.

53 **R4.4: how the radius is determined in experiments.** The radius is only used to show the experimental results, and is
54 *not* a hyperparameter. In fact, we could also plot the certified robust error curve across “all” radii as in Figure i.



(a) MNIST. (b) Fashion-MNIST.
Figure i: Certified robust errors of Mahalanobis 1-NN compared with ordinary neural networks and randomized-smoothing neural networks (RandSmooth).