

A Optimal value of triplet problem

The triplet problem is formalized as below:

$$\min_{\delta} \|\delta\| \quad \text{s.t.} \quad d_M(\mathbf{x} + \delta, \mathbf{x}^-) \leq d_M(\mathbf{x} + \delta, \mathbf{x}^+). \quad (15)$$

It is equivalent to the optimization

$$\min_{\delta} \delta^\top \delta \quad \text{s.t.} \quad \mathbf{a}^\top \delta \leq b, \quad (16)$$

where we have

$$\mathbf{a} = \mathbf{M}(\mathbf{x}^+ - \mathbf{x}^-), \quad (17)$$

$$b = \frac{1}{2} (d_M(\mathbf{x}, \mathbf{x}^+) - d_M(\mathbf{x}, \mathbf{x}^-)). \quad (18)$$

The dual function is

$$g(\lambda) = \inf_{\delta} \delta^\top \delta + \lambda(\mathbf{a}^\top \delta - b) \quad (19)$$

$$= -\frac{1}{4} \mathbf{a}^\top \mathbf{a} \lambda^2 - b\lambda, \quad (20)$$

where inf holds for $\delta = -\lambda \mathbf{a}/2$. Then the dual problem is

$$\max_{\lambda \geq 0} -\frac{1}{4} \mathbf{a}^\top \mathbf{a} \lambda^2 - b\lambda. \quad (21)$$

The optimal point is

$$\left[-\frac{2b}{\mathbf{a}^\top \mathbf{a}} \right]_+, \quad (22)$$

and the optimal value is

$$\begin{cases} 0 & \text{if } b \geq 0 \\ \frac{b^2}{\mathbf{a}^\top \mathbf{a}} & \text{otherwise.} \end{cases} \quad (23)$$

By the Slater's condition, if $\mathbf{x}^+ \neq \mathbf{x}^-$ holds, we have the strong duality. Therefore, the optimal value of (15) is

$$\left[\frac{-b}{\sqrt{\mathbf{a}^\top \mathbf{a}}} \right]_+ = \frac{[d_M(\mathbf{x}, \mathbf{x}^-) - d_M(\mathbf{x}, \mathbf{x}^+)]_+}{2\sqrt{(\mathbf{x}^+ - \mathbf{x}^-)^\top \mathbf{M}^{-1} \mathbf{M}(\mathbf{x}^+ - \mathbf{x}^-)}}. \quad (24)$$

In fact, it is easy to verify that even if $\mathbf{x}^+ = \mathbf{x}^-$ obtains, the optimal value also holds.

B Proof of Theorem 1

Proof. Let $\epsilon^{(\mathbb{I}, \mathbb{J})}$ denote the optimal value of the inner minimization problem of (6). By relaxing the constraint via replacing the universal quantifier, we have

$$\epsilon^{(\mathbb{I}, \mathbb{J})} \geq \max_{i \in \{i: y_i = y_{\text{test}}\} - \mathbb{I}, j \in \mathbb{J}} \tilde{\epsilon}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_{\text{test}}; \mathbf{M}). \quad (25)$$

Substitute it in (6) and then we have

$$\epsilon^* \geq \min_{\mathbb{I}, \mathbb{J}} \epsilon^{(\mathbb{I}, \mathbb{J})} \quad (26)$$

$$\geq \min_{\mathbb{I}, \mathbb{J}} \max_{i \in \{i: y_i = y_{\text{test}}\} - \mathbb{I}, j \in \mathbb{J}} \tilde{\epsilon}(\mathbf{x}_i^+, \mathbf{x}_j^-, \mathbf{x}_{\text{test}}) \quad (27)$$

$$\geq \min_{\mathbb{I}, \mathbb{J}} \max_{j \in \mathbb{J}} \max_{i \in \{i: y_i = y_{\text{test}}\} - \mathbb{I}} \tilde{\epsilon}(\mathbf{x}_i^+, \mathbf{x}_j^-, \mathbf{x}_{\text{test}}) \quad (28)$$

$$\geq \min_{\mathbb{I}, \mathbb{J}} \max_{j \in \mathbb{J}} k\text{th max}_{i \in \{i: y_i = y_{\text{test}}\}} \tilde{\epsilon}(\mathbf{x}_i^+, \mathbf{x}_j^-, \mathbf{x}_{\text{test}}) \quad (29)$$

$$\geq \min_{\mathbb{I}, \mathbb{J}} k\text{th min}_{j \in \{j: y_j \neq y_{\text{test}}\}} k\text{th max}_{i \in \{i: y_i = y_{\text{test}}\}} \tilde{\epsilon}(\mathbf{x}_i^+, \mathbf{x}_j^-, \mathbf{x}_{\text{test}}) \quad (30)$$

$$= k\text{th min}_{j \in \{j: y_j \neq y_{\text{test}}\}} k\text{th max}_{i \in \{i: y_i = y_{\text{test}}\}} \tilde{\epsilon}(\mathbf{x}_i^+, \mathbf{x}_j^-, \mathbf{x}_{\text{test}}) \quad (31)$$

□

C Details of computing exact minimal adversarial perturbation of Mahalanobis 1-NN

The overall algorithm is displayed in Algorithm 2. We denote $\epsilon^{(j)}$ as the optimal value of the inner minimization problem with respect to j , and denote $\underline{\epsilon}^{(j)}$ as its lower bound. We first sort the subproblems according to the ascending order of $\|\mathbf{x}_j - \mathbf{x}_{\text{test}}\|$ for $\{j : y_j \neq y_{\text{test}}\}$. For every subproblem, we compute the lower bound of its optimal value. If the optimal value is too large, we just screen the subproblem safely without solving it exactly.

Algorithm 2: Computing the minimal adversarial perturbation for Mahalanobis 1-NN

Input: Test instance $(\mathbf{x}_{\text{test}}, y_{\text{test}})$, dataset $\mathbb{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$.

Output: Perturbation norm ϵ .

```

1 Initialize  $\epsilon = \infty$ ;
2 Sort  $\{j : y_j \neq y_{\text{test}}\}$  by the ascending order of  $d_M(\mathbf{x}_{\text{test}}, \mathbf{x}_j)$ ;
3 for  $j : y_j \neq y_{\text{test}}$  according to the ascending order do
4   Compute a lower bound  $\underline{\epsilon}^{(j)}$  of the inner minimization corresponding to  $j$ ;
5   if  $\underline{\epsilon} < \epsilon$  then
6     Solve the inner minimization problem exactly via the greedy coordinate ascent method
       and derive the optimal value  $\epsilon^{(j)}$ ;
7     if  $\epsilon^{(j)} < \epsilon$  then
8        $\epsilon = \epsilon^{(j)}$ 
9     end
10  end
11 end

```

C.1 Greedy coordinate ascent (descent)

For the subproblem we have to solve exactly, we employ the greedy coordinate ascent method. Note that the inner minimization problem of (13) is a convex quadratic programming problem. We solve the problem by dealing with its dual formulation. The greedy coordinate ascent method is used because the optimal dual variables are very sparse. The algorithm is shown in Algorithm 3. At every iteration, only one dual variable is updated.

Algorithm 3: Greedy coordinate descent for QP: $\min_{\mathbf{x} \geq 0} \frac{1}{2} \mathbf{x}^\top \mathbf{P} \mathbf{x} + \mathbf{q}^\top \mathbf{x}$

Input: $\mathbf{P}, \mathbf{q}, \epsilon, T$.

```

1  $\mathbf{x} \leftarrow \mathbf{0}, \mathbf{g} \leftarrow \mathbf{P} \mathbf{x} + \mathbf{q}$ ;
2 for  $t = 0$  to  $T - 1$  do
3    $\forall i, y_i \leftarrow \max\left(x_i - \frac{g_i}{p_{i,i}}, 0\right) - x_i$ ;
4    $i^* \leftarrow \arg \max_i |y_i|$ ; // choose a coordinate
5   if then
6     | break;
7   end
8    $x_{i^*} \leftarrow x_{i^*} + y_{i^*}$ ; // update the solution
9    $\mathbf{g} \leftarrow \mathbf{g} + y_{i^*} \mathbf{p}_{i^*}$ ; // update the gradient
10 end
Output:  $\mathbf{x}$ .

```

C.2 Lower bound of inner minimization problem

The following theorem is dependent on the solution of the triplet problem.

Theorem 2. The optimal value $\epsilon^{(j)}$ of the inner minimization of (13) with respect to j is lower bounded as

$$\epsilon^{(j)} \geq \max_{i: y_i = y_{\text{test}}} [\tilde{\epsilon}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_{\text{test}}; \mathbf{M})]_+. \quad (32)$$

Proof. Relaxing the constraint of (13) by means of replacing the universal quantifier, we know $\epsilon^{(j)}$ is lower bounded by the optimal value of the following optimization problem

$$\max_{i: y_i = y_{\text{test}}} \min_{\delta_{i,j}} \|\delta_{i,j}\| \quad (33)$$

$$\text{s.t. } d_{\mathbf{M}}(\mathbf{x}_{\text{test}} + \delta_{i,j}, \mathbf{x}_j) \leq d_{\mathbf{M}}(\mathbf{x}_{\text{test}} + \delta_{i,j}, \mathbf{x}_i). \quad (34)$$

Obviously, the optimal value of the inner problem is $[\tilde{\epsilon}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_{\text{test}}; \mathbf{M})]_+$. \square

In this way, we could derive a lower bound of the optimal value in closed form.

D Experimental details

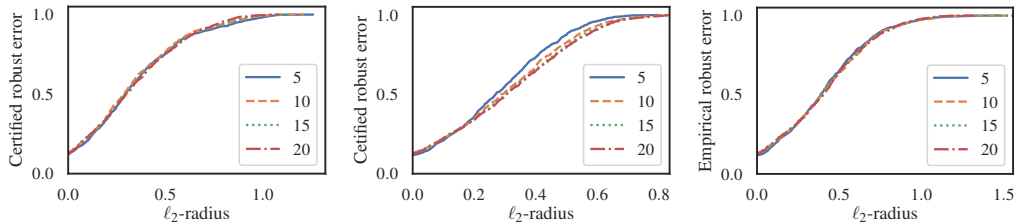
Datasets Dataset statistics and test errors (on all test instances) of Euclidean K -NN with $K = 11$ are shown in Table 3. All training data are used to learn metrics, and 1,000 instances are randomly sampled to compute certified robust errors.

Table 3: Dataset statistics.

	# features	# classes	# train	# test	1-NN test error	K -NN test error
MNIST	784	10	60,000	10,000	0.031	0.033
Fashion-MNIST	784	10	60,000	10,000	0.150	0.150
Splice	60	2	1,000	2,175	0.295	0.291
Pendigits	16	10	7,494	3,498	0.023	0.027
Satimage	36	6	4,435	2,000	0.112	0.106
USPS	256	10	7,291	2,007	0.049	0.060

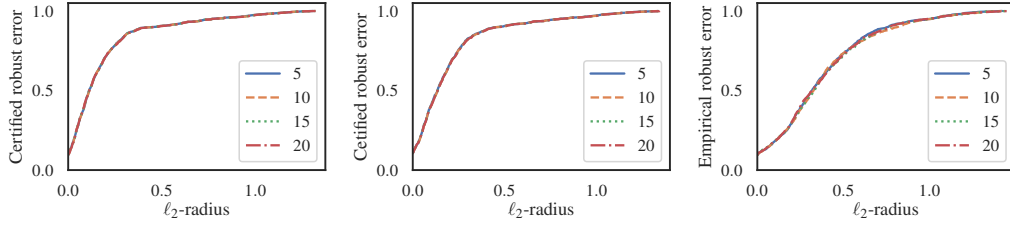
Hyperparameters Hyperparameters of our ARML algorithm are fixed across all datasets. Specifically, the size of the neighborhood where randnear^+ and randnear^- sample random instances is 10. In other words, at every iteration, we sample one instance from the nearest 10 instances in the same class with the test instance, and sample one instance from the nearest 10 instances in the different classes from the test instance. We employ the Adam algorithm [24] to update parameters with gradients and the parameters is in the default setting (learning rate: 0.001, betas: (0.9, 0.999)). The number of epochs is 1,000. The loss function is the negative loss.

Hyperparameter sensitivity We investigate the sensitivity of the size of neighborhood used for randnear^+ and randnear^- . We plot the robust error curves against the ℓ_2 radius for different neighborhood sizes in Figure 3 and Figure 4. It suggests that ARML is not very sensitive to this hyperparameter in terms of certified and empirical robust errors.



(a) 1-NN certified robust error. (b) K -NN certified robust error. (c) K -NN empirical robust error.

Figure 3: Sensitivity to neighborhood size on Splice.



(a) 1-NN certified robust error. (b) K -NN certified robust error. (c) K -NN empirical robust error.

Figure 4: Sensitivity to neighborhood size on Satimage.

Implementations of NCA, LMNN, ITML and LFDA We use the implementations of the metric-learn library [14] for NCA, LMNN, ITML and LFDA. Similar to ARML, hyperparameters are fixed across all datasets and are in the default setting. In particular, the maximum numbers of iterations for NCA, LMNN and ITML are 1,00, 1000 and 1,000 respectively.