

# MEASURING ADVERSARIAL ROBUSTNESS USING A VORONOI-EPSILON ADVERSARY

Hyeongji Kim<sup>1,2</sup>, Pekka Parviainen<sup>1</sup>, and Ketil Malde<sup>1,2</sup>

<sup>1</sup>Department of Informatics, University of Bergen, Norway

<sup>2</sup>Institute of Marine Research, Bergen, Norway

kim.hyeongji@hi.no

## ABSTRACT

Previous studies on robustness have argued that there is a tradeoff between accuracy and adversarial accuracy. The tradeoff can be inevitable even when we neglect generalization. We argue that the tradeoff is inherent to the commonly used definition of adversarial accuracy, which uses an adversary that can construct adversarial points constrained by  $\epsilon$ -balls around data points. As  $\epsilon$  gets large, the adversary may use real data points from other classes as adversarial examples. We propose a Voronoi-epsilon adversary which is constrained both by Voronoi cells and by  $\epsilon$ -balls. This adversary balances between two notions of perturbation. As a result, adversarial accuracy based on this adversary avoids a tradeoff between accuracy and adversarial accuracy on training data even when  $\epsilon$  is large. Finally, we show that a nearest neighbor classifier is the maximally robust classifier against the proposed adversary on the training data.

## 1 INTRODUCTION

By applying a carefully crafted, but imperceptible perturbation to input images, so-called adversarial examples can be constructed that cause classifiers to misclassify the perturbed inputs (Szegedy et al., 2013). Defense methods like adversarial training (Madry et al., 2017) and certified defenses (Wong & Kolter, 2018) against adversarial examples have often resulted in decreased accuracies on clean samples (Tsipras et al., 2018). Previous studies have argued that the tradeoff between accuracy and adversarial accuracy may be inevitable in classifiers (Tsipras et al., 2018; Dohmatob, 2018; Zhang et al., 2019).

### 1.1 PROBLEM SETTINGS

**Problem setting.** Let  $\mathcal{X} \subset \mathbb{R}^{\text{dim}}$  be a nonempty input space and  $\mathcal{Y}$  be a set of possible classes. Data points  $x \in \mathcal{X}$  and corresponding classes  $c_x \in \mathcal{Y}$  are sampled from a joint distribution  $\mathcal{D}$ . The distribution  $\mathcal{D}$  should satisfy the condition that  $c_x$  is unique for all  $x$ . The set of the data points is denoted as  $X$ .  $X$  is a nonempty finite set. A classifier  $f$  assigns a class label from  $\mathcal{Y}$  for each point  $x \in \mathcal{X}$ .  $L(x, y)$  is a classification loss of the classifier  $f$  provided an input  $x \in \mathcal{X}$  and a label  $y \in \mathcal{Y}$ .

More notations are summarized in A.1. Abbreviations are summarized in A.2. We focus on situations that we neglect generalization to simplify the analysis.

### 1.2 ADVERSARIAL ACCURACY (AA)

Adversarial accuracy is a commonly used measure of adversarial robustness of classifiers (Madry et al., 2017; Tsipras et al., 2018). It is defined by an adversary region  $R(x) \subset \mathcal{X}$ , which is an allowed region of the perturbations for a data point  $x$ .

**Definition 1 (Adversarial accuracy).** Given an adversary that is constrained to an adversary region  $R(x)$ , adversarial accuracy  $a$  is defined as follows.

$$a = \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f(x^*) = c_x)] \text{ where } x^* = \arg \max_{x' \in R(x)} L(x', c_x)$$

The choice of  $R(x)$  will determine the adversarial accuracy that we are measuring. Commonly considered adversary region is  $\mathbb{B}(x, \epsilon)$ , which is a  $\epsilon$ -ball around a data point  $x$  based on a distance metric  $d$  (Biggio et al., 2013; Madry et al., 2017; Tsipras et al., 2018; Zhang et al., 2019).

**Definition 2 (Standard adversarial accuracy).** When the adversary region is  $\mathbb{B}(x, \epsilon)$ , we refer to the adversarial accuracy  $a$  as standard adversarial accuracy (SAA)  $a_{std}(\epsilon)$ . For SAA, we denote  $R(x)$  as  $R_{std}(\epsilon; x)$ .

$$a_{std}(\epsilon) = \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f(x^*) = c_x)] \text{ where } x^* = \arg \max_{x' \in R_{std}(\epsilon; x)} L(x', c_x)$$

This adversary region  $\mathbb{B}(x, \epsilon)$  is based on an implicit assumption that there might be an adequate single epsilon  $\epsilon$  that perturbed samples do not change their classes. However, this assumption has some limitations. We explain that in the next section.

### 1.3 THE TRADEOFF BETWEEN ACCURACY AND STANDARD ADVERSARIAL ACCURACY

The usage of  $\epsilon$ -ball-based adversary can cause the tradeoff between accuracy and adversarial accuracy. When the two clean samples  $x_1$  and  $x_2$  with  $d(x_1, x_2) \leq \epsilon$  have different classes, the increase of standard adversarial accuracy requires misclassification. We illustrate this with a toy example.

#### 1.3.1 TOY EXAMPLE

Let us consider an example visualized in Figure 1a. The input space is  $\mathbb{R}^2$ . There are only two classes  $A$  and  $B$ , i.e.,  $\mathcal{Y} = \{A, B\}$ . We use the  $l_2$  norm as a distance metric in this example.

Let us consider a situation when  $\epsilon = 1.0$  (see Figure 1c). In this case, clean samples can also be considered as adversarial examples. For example, the point  $(2, 1)$  can be considered as an adversarial example originating from the point  $(1, 1)$ . If we choose a robust model based on SAA, we might choose a model with excessive invariance. For example, we might choose a model that predicts points belong to  $\mathbb{B}((1, 1), 1)$  (including the point  $(2, 1)$ ) have class A. Or, we can choose a model that predicts points belong to  $\mathbb{B}((2, 1), 1)$  (including the point  $(1, 1)$ ) have class B. In either case, the accuracy of the chosen model is smaller than 1. This situation explains the tradeoff between accuracy and standard adversarial accuracy when large  $\epsilon$  is used. It originates from the overlapping adversary regions from the samples with different classes.

To avoid the tradeoff between accuracy and adversarial accuracy, one can use small  $\epsilon$  values. Actually, a previous study has argued that commonly used  $\epsilon$  values are small enough to avoid the tradeoff (Yang et al., 2020b). However, when small  $\epsilon$  values are used, we can only analyze local robustness, and we need to ignore robustness beyond the chosen  $\epsilon$ . For instance, let us consider our example when  $\epsilon = 0.5$  (see Figure 1b). In this case, we ignore robustness on  $\mathbb{B}((-2, 1), 1.0) - \mathbb{B}((-2, 1), 0.5)$ . Models with local but without global robustness enable attackers to use large  $\epsilon$  values to fool the models. Ghiasi et al. (2019) have experimentally shown that even models with certified local robustness can be attacked by attacks with large  $\epsilon$  values. Note that their attack applies little semantic perturbations even though the perturbation norms measured by  $l_p$  norms are large.

These limitations motivate us to find an alternative way to measure robustness. **The contributions of this paper are as follows.**

- We propose Voronoi-epsilon adversarial accuracy (VAA) that avoids the tradeoff between accuracy and adversarial accuracy. This allows the adversary regions to scale to cover most of the input space without incurring a tradeoff. To our best knowledge, this is the first work to achieve this without an external classifier. (In Appendix A.3, we introduce formulas for adversary regions that can be used to estimate VAA.)
- We explain the connection between SAA and VAA. We define global Voronoi-epsilon robustness as a limit of the Voronoi-epsilon adversarial accuracy. We show that a nearest neighbor (1-NN) classifier maximizes global Voronoi-epsilon robustness.

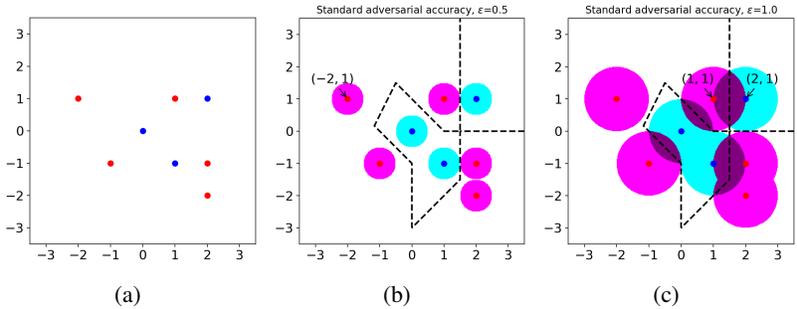


Figure 1: (a): Plot of the two-dimensional toy example. Data points are colored based on their classes (class A: red and class B: blue). (b): Visualization of the adversary regions for SAA when  $\epsilon = 0.5$ . The regions are colored differently depending on their classes (class A: magenta and class B: cyan). The decision boundary of a single nearest neighbor classifier is shown as a dashed black curve. (c): Visualization of the adversary regions for SAA when  $\epsilon = 1.0$ . The overlapping adversary regions from the samples with different classes are colored in purple.

## 2 VORONOI-EPSILON ADVERSARIAL ACCURACY (VAA)

Our approach restricts the allowed region of the perturbations to avoid the tradeoff originating from the definition of standard adversarial accuracy. This is achieved without limiting the magnitude of  $\epsilon$  and without using an external model. We want to have the following property to avoid the tradeoff.

$$\forall x_i, x_j \in X, x_i \neq x_j \implies R(x_i) \cap R(x_j) = \emptyset \tag{1}$$

When Property (1) holds for the adversary region, we no longer have the tradeoff as  $x_i \notin R(x_j)$  for  $x_i \neq x_j$ . In other words, a clean sample cannot be an adversarial example originating from another clean sample. We propose a new adversary called a Voronoi-epsilon adversary that combines the Voronoi-adversary introduced by Khoury & Hadfield-Menell (2019) with an  $\epsilon$ -ball-based adversary. This adversary is constrained to an adversary region  $Vor(x) \cap \mathbb{B}(x, \epsilon)$  where  $Vor(x)$  is the (open) Voronoi cell around a data point  $x \in X$ .  $Vor(x)$  consists of every point in  $\mathcal{X}$  that is closer than any  $x_{clean} \in X - \{x\}$ . Then, Property (1) holds as  $Vor(x_i) \cap Vor(x_j) = \emptyset$  for  $x_i \neq x_j$ .

Based on a Voronoi-epsilon adversary, we define Voronoi-epsilon adversarial accuracy (VAA).

**Definition 3 (Voronoi-epsilon adversarial accuracy).** When a Voronoi-epsilon adversary is used for the adversary, we refer to the adversarial accuracy as Voronoi-epsilon adversarial accuracy (VAA)  $a_{Vor}(\epsilon)$ . For VAA, we denote  $R(x)$  as  $R_{Vor}(\epsilon; x)$ .

$$a_{Vor}(\epsilon) = \mathbb{E}_{x \in X} [\mathbb{1}(f(x^*) = c_x)] \text{ where } x^* = \arg \max_{x' \in R_{Vor}(\epsilon; x)} L(x', c_x)$$

Note that VAA is only defined on a fixed set of data points  $X$ . As we do not know the distribution  $\mathcal{D}$ , in practice, the fact that VAA is not defined on the whole input space does not matter.

Figure 2 shows the adversary regions for VAA with varying  $\epsilon$  values. When  $\epsilon = 0.5$ , the regions are same with SAA except for the points  $(1.5, 1)$ ,  $(1.5, -1)$  and  $(2, -1.5)$ . Even when  $\epsilon$  is large ( $\epsilon > 0.5$ ), there is no overlapping adversary region, which was a source of the tradeoff in SAA. Therefore, when we choose a robust model based on VAA, we can get a model that is both accurate and robust. Figure 2c shows the single nearest neighbor (1-NN) classifier would maximize VAA. The adversary regions cover most of the points in  $\mathbb{R}^2$  for large  $\epsilon$ .

**Observation 1.** Let  $d_{min}$  be the nearest distance of the data point pairs, i.e.,  $d_{min} = \min_{x_i, x_j \in X, x_i \neq x_j} d(x_i, x_j)$ . Then, the following equivalence holds.

$$a_{Vor}(\epsilon) = a_{std}(\epsilon) \text{ when } \epsilon < \frac{1}{2}d_{min} \tag{2}$$

Observation 1 shows that VAA is equivalent to SAA for sufficiently small  $\epsilon$  values. This indicates that VAA is an extension of SAA that avoids the tradeoff when  $\epsilon$  is large. The proof of the obser-

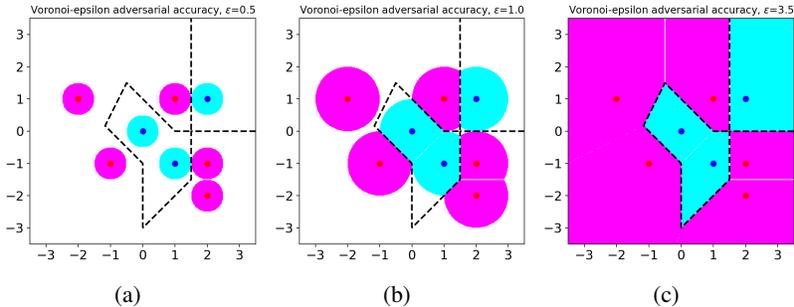


Figure 2: Visualization of the adversary regions for VAA with varying  $\epsilon$  values. The data points and regions are colored as in Figure 1. (a): When  $\epsilon = 0.5$ . (b): When  $\epsilon = 1.0$ . (c): When  $\epsilon = 3.5$ .

vation is in Appendix A.5. We point out that equivalent findings were also mentioned in Yang et al. (2020a;b); Khoury & Hadfield-Menell (2019).

As explained in Section 1.3.1, studying the local robustness of classifiers has a limitation. Attackers can attack models with only local robustness by using large  $\epsilon$  values. The absence of a tradeoff between accuracy and VAA enables us to increase  $\epsilon$  values and to study global robustness. We define a measure for global robustness using VAA.

**Definition 4 (Global Voronoi-epsilon robustness).** *Global Voronoi-epsilon robustness  $a_{global}$  is defined as*

$$a_{global} = \lim_{\epsilon \rightarrow \infty} a_{Vor}(\epsilon).$$

Global Voronoi-epsilon robustness considers the robustness of classifiers for most points in  $\mathcal{X}$  (all points except for Voronoi boundary  $VB(X)$ , which is the complement set of the unions of Voronoi cells.). We derive the following theorem from global Voronoi-epsilon robustness.

**Theorem 1.** *A single nearest neighbor (1-NN) classifier maximizes global Voronoi-epsilon robustness  $a_{global}$  on training data. 1-NN classifier is a unique classifier that satisfies this except for Voronoi boundary  $VB(X)$ .*

Note that Theorem 1 only holds for exactly the same data under the exclusive class condition as mentioned in the problem settings 1.1. It does not take into account generalization. The proof of the theorem is in A.6.

### 3 DISCUSSION

In this work, we address the tradeoff between accuracy and adversarial robustness by introducing the Voronoi-epsilon adversary. Another way to address this tradeoff is to use a Bayes optimal classifier (Suggala et al., 2019; Kim & Wang, 2020). Since this is not available in practice, a reference model must be used as an approximation. In that case, the meaning of adversarial robustness is dependent on the choice of the reference model. VAA removes the need for a reference model by using the data point set  $X$  and the distance metric  $d$  to construct adversary. This is in contrast to Khoury & Hadfield-Menell (2019) who used Voronoi cell-based constraints (without  $\epsilon$ -balls) for an adversarial training purpose, but not for measuring adversarial robustness.

By avoiding the tradeoff with VAA, we can extend the study of local robustness to global robustness. Also, Theorem 1 implies that VAA is a measure of agreement with the 1-NN classifier. For sufficiently small  $\epsilon$  values, SAA is also a measure of agreement with the 1-NN classifier because SAA is equivalent to VAA as in Observation 1. This implies that many defenses (Goodfellow et al., 2014; Madry et al., 2017; Zhang et al., 2019; Wong & Kolter, 2018; Cohen et al., 2019) with small  $\epsilon$  values unknowingly try to make locally the same predictions with a 1-NN classifier.

In our analysis, we do not consider generalization, and robust models are known to often generalize poorly (Raghunathan et al., 2020). The close relationship between adversarially robust models and the 1-NN classifier revealed by Theorem 1 highlights a possible avenue to explore this phenomenon.

## ACKNOWLEDGMENTS

We thank Dr. Nils Olav Handegard, Dr. Yi Liu, and Jungeum Kim for the helpful feedback. We also thank Dr. Wieland Brendel for the helpful discussions.

## REFERENCES

- Battista Biggio, Igino Corona, Davide Maiorca, Blaine Nelson, Nedim Šrđić, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. Evasion attacks against machine learning at test time. In *Joint European conference on machine learning and knowledge discovery in databases*, pp. 387–402. Springer, 2013.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning*, pp. 1310–1320. PMLR, 2019.
- Elvis Dohmatob. Limitations of adversarial robustness: strong No Free Lunch Theorem. *arXiv:1810.04065 [cs, stat]*, October 2018. URL <http://arxiv.org/abs/1810.04065>. arXiv: 1810.04065.
- Amin Ghiasi, Ali Shafahi, and Tom Goldstein. Breaking certified defenses: Semantic adversarial examples with spoofed robustness certificates. In *International Conference on Learning Representations*, 2019.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Marc Khoury and Dylan Hadfield-Menell. Adversarial training with Voronoi constraints. *arXiv preprint arXiv:1905.01019*, 2019.
- Jungeum Kim and Xiao Wang. Sensible adversarial learning, 2020. URL [https://openreview.net/forum?id=rJlf\\_RVKwr](https://openreview.net/forum?id=rJlf_RVKwr).
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Aditi Raghunathan, Sang Michael Xie, Fanny Yang, John Duchi, and Percy Liang. Understanding and mitigating the tradeoff between robustness and accuracy. *arXiv preprint arXiv:2002.10716*, 2020.
- Arun Sai Suggala, Adarsh Prasad, Vaishnavh Nagarajan, and Pradeep Ravikumar. Revisiting adversarial risk. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2331–2339. PMLR, 2019.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- Eric Wong and Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. In *International Conference on Machine Learning*, pp. 5286–5295. PMLR, 2018.
- Yao-Yuan Yang, Cyrus Rashtchian, Yizhen Wang, and Kamalika Chaudhuri. Robustness for non-parametric classification: A generic attack and defense. In *International Conference on Artificial Intelligence and Statistics*, pp. 941–951. PMLR, 2020a.
- Yao-Yuan Yang, Cyrus Rashtchian, Hongyang Zhang, Russ R Salakhutdinov, and Kamalika Chaudhuri. A closer look at accuracy vs. robustness. *Advances in Neural Information Processing Systems*, 33, 2020b.
- Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. Theoretically principled trade-off between robustness and accuracy. *arXiv preprint arXiv:1901.08573*, 2019.

## A APPENDIX

## A.1 LIST OF NOTATION

$\epsilon$	A perturbation budget.
$\dim$	The dimension of the input space.
$\mathcal{X}$	The nonempty input space. $\mathcal{X} \subset \mathbb{R}^{\dim}$ .
$\mathcal{Y}$	The set of possible classes.
$c_x$	A corresponding class of a clean data point $x \in \mathcal{X}$ .
$\mathcal{D}$	The joint distribution. $\mathcal{D} \subset \mathcal{X} \times \mathcal{Y}$ .
$X$	The set of data points. We assume it is a nonempty finite set.
$f$	The classifier that we want to analyze. $f : \mathcal{X} \rightarrow \mathcal{Y}$ .
$L(x, y)$	A classification loss of the classifier $f$ provided an input $x \in \mathcal{X}$ and a label $y \in \mathcal{Y}$ .
$R(x)$	An adversary region which is an allowed region of the perturbations for a data point $x$ . It can be depend on a perturbation budget $\epsilon$ .
$\mathbb{1}()$	The indicator function. $\mathbb{1}(True) = 1$ and $\mathbb{1}(False) = 0$ .
$a$	Adversarial accuracy.
$d$	The distance metric that is used for measuring adversarial robustness. It is not limited to $l_p$ norms. It can be a learned metric or more complex distance.
$\mathbb{B}(x, \epsilon)$	An $\epsilon$ -ball around a sample $x$ . Mathematically, $\mathbb{B}(x, \epsilon) = \{x' \in \mathcal{X}   d(x, x') \leq \epsilon\}$ .
$R_{std}(\epsilon; x)$	The allowed regions of the perturbations for standard adversarial accuracy around a data point $x$ . $R_{std}(\epsilon; x) = \mathbb{B}(x, \epsilon)$ .
$a_{std}(\epsilon)$	Standard adversarial accuracy using a perturbation budget $\epsilon$ . In other words, the adversarial accuracy when the adversary region is $R_{std}(\epsilon; x) = \mathbb{B}(x, \epsilon)$ .
$HS(x, x_{clean})$	The (open) half-space closer to $x \in X$ than $x_{clean} \in X - \{x\}$ . Mathematically, $HS(x, x_{clean}) = \{x' \in \mathcal{X}   d(x, x') < d(x_{clean}, x')\}$ .
$Vor(x)$	The (open) Voronoi cell of a sample $x \in X$ . Mathematically, $Vor(x) = \{x' \in \mathcal{X}   d(x, x') < d(x_{clean}, x'), \forall x_{clean} \in X - \{x\}\} = \bigcap_{x_{clean} \in X - \{x\}} HS(x, x_{clean})$ .
$R_{Vor}(\epsilon; x)$	The allowed regions of the perturbations for Voronoi-epsilon adversarial accuracy around a data point $x$ . $R_{Vor}(\epsilon; x) = Vor(x) \cap \mathbb{B}(x, \epsilon)$ .
$a_{Vor}(\epsilon)$	The Voronoi-epsilon adversarial accuracy using perturbation budget $\epsilon$ . In other words, the adversarial accuracy when the adversary region is $R_{Vor}(\epsilon; x) = Vor(x) \cap \mathbb{B}(x, \epsilon)$ .
$S^c$	The complement set of a set $S$ . For $S \subset \mathcal{X}$ , $S^c = \mathcal{X} - S$ .
$VB(X)$	Voronoi boundary based on $X$ . It is the complement set of the unions of Voronoi cells. $VB(X) = \left( \bigcup_{x \in X} Vor(x) \right)^c = \bigcap_{x \in X} Vor(x)^c$ .
$a_{global}$	Global Voronoi-epsilon robustness.
$N$	The number of data points.
$R_{Vor;LB}(\epsilon; x)$	The allowed regions of the perturbations for the lower bound of Voronoi-epsilon adversarial accuracy around a data point $x$ . When $\epsilon < \frac{1}{2}d(x, x_{m+2})$ , $R_{Vor;LB}(\epsilon; x) = R_{Vor}(\epsilon; x)$ . When $\epsilon \geq \frac{1}{2}d(x, x_{m+2})$ , $R_{Vor;LB}(\epsilon; x) = \mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^{m+1} HS(x, x_i) \right)$ .

$R_{Vor;UB}(\epsilon; x)$	The allowed regions of the perturbations for the upper bound of Voronoi-epsilon adversarial accuracy around a sample $x$ . When $\epsilon < \frac{1}{2}d(x, x_{m+2})$ , $R_{Vor;UB}(\epsilon; x) = R_{Vor}(\epsilon; x)$ . When $\epsilon \geq \frac{1}{2}d(x, x_{m+2})$ , $R_{Vor;UB}(\epsilon; x) = \mathbb{B}(x, \frac{1}{2}d(x, x_{m+2}) - \tau) \cap \left( \bigcap_{i=2}^{m+1} HS(x, x_i) \right)$ .
$a_{Vor;LB}(\epsilon)$	The lower bound of Voronoi-epsilon adversarial accuracy using perturbation budget $\epsilon$ . It is defined as the adversarial accuracy when the adversary region for a data point $x$ is $R_{Vor;LB}(\epsilon; x)$ .
$a_{Vor;UB}(\epsilon)$	The upper bound of Voronoi-epsilon adversarial accuracy using perturbation budget $\epsilon$ . It is defined as the adversarial accuracy when the adversary region for a data point $x$ is $R_{Vor;UB}(\epsilon; x)$ .

## A.2 LIST OF ABBREVIATION

AA	Adversarial accuracy.
SAA	Standard adversarial accuracy.
VAA	Voronoi-epsilon adversarial accuracy.
1-NN	Single nearest neighbor.
LB	Lower bound.
UB	Upper bound.

## A.3 ADVERSARY REGION $R_{Vor}(\epsilon; x)$

Voronoi-epsilon adversarial accuracy (VAA) uses  $R_{Vor}(\epsilon; x) = Vor(x) \cap \mathbb{B}(x, \epsilon)$ . We introduce upper and lower bounds of  $R_{Vor}(\epsilon; x)$  using  $m + 1$  nearest neighbors of a data point  $x$ . These bounds enable to calculate approximate upper and lower bounds of VAA.

**Lemma 1.** *When  $N$  is the number of data points, let  $x_2, \dots, x_N \in X - \{x\}$  be the sorted neighbors of a data point  $x \in X$ . Mathematically,  $d(x, x_2) \leq d(x, x_3) \leq \dots \leq d(x, x_N)$ . Then, the following relations hold for a fixed number  $m \leq N - 2$ .*

$$R_{Vor}(\epsilon; x) = \mathbb{B}(x, \epsilon) \text{ when } \epsilon < \frac{1}{2}d(x, x_2) \quad (3)$$

$$R_{Vor}(\epsilon; x) = \mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^j HS(x, x_i) \right) \text{ when } \frac{1}{2}d(x, x_j) \leq \epsilon < \frac{1}{2}d(x, x_{j+1}) \quad (4)$$

$(j = 2, \dots, m + 1)$

$$\mathbb{B}(x, \frac{1}{2}d(x, x_{m+2}) - \tau) \cap \left( \bigcap_{i=2}^{m+1} HS(x, x_i) \right) \subset R_{Vor}(\epsilon; x) \subset \mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^{m+1} HS(x, x_i) \right) \quad (5)$$

*when  $\epsilon \geq \frac{1}{2}d(x, x_{m+2})$  and  $\tau > 0$*

When  $\epsilon < \frac{1}{2}d(x, x_{m+2})$ , we can calculate VAA using relations (3) and (4). The relation (5) of Lemma 1 enables to calculate the lower and upper bound of VAA when  $\epsilon \geq \frac{1}{2}d(x, x_{m+2})$ . When  $\epsilon \geq \frac{1}{2}d(x, x_{m+2})$ , we denote the leftmost set in the relation (5) as  $R_{Vor;UB}(\epsilon; x)$  and the rightmost set as  $R_{Vor;LB}(\epsilon; x)$ . (When  $\epsilon < \frac{1}{2}d(x, x_{m+2})$ , we set  $R_{Vor;LB}(\epsilon; x) = R_{Vor;UB}(\epsilon; x) = R_{Vor}(\epsilon; x)$ .) Figure 3 visualizes the relationship  $R_{Vor;UB}(\epsilon; x) \subset R_{Vor}(\epsilon; x) \subset R_{Vor;LB}(\epsilon; x) \subset R_{std}(\epsilon; x)$ . The proof of the lemma is in Appendix A.4.

**Proposition 1.**  *$a_{Vor;LB}(\epsilon)$  is defined as the adversarial accuracy when the allowed regions of perturbation is  $R_{Vor;LB}(\epsilon; x)$ .  $a_{Vor;UB}(\epsilon)$  is defined as the adversarial accuracy when the allowed regions of perturbation is  $R_{Vor;UB}(\epsilon; x)$ . Then, the following relation holds.*

$$a_{std}(\epsilon) \leq a_{Vor;LB}(\epsilon) \leq a_{Vor}(\epsilon) \leq a_{Vor;UB}(\epsilon) \quad (6)$$

The proof of Proposition 1 is in Appendix A.5.

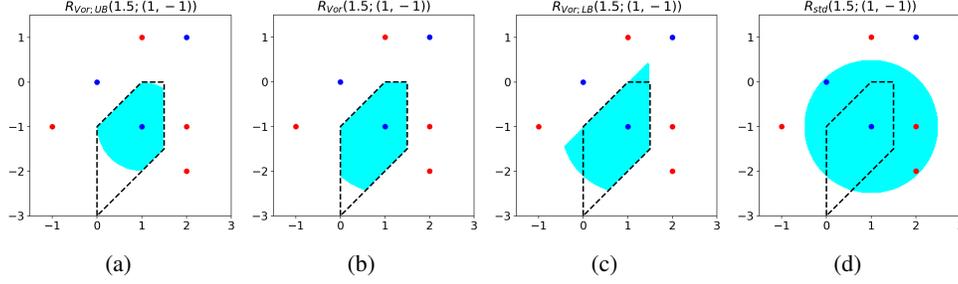


Figure 3: Visualization of the adversary region for the point  $(1, -1)$  when  $m = 3$  and  $\epsilon = 1.5$  on our example 1.3.1. (a):  $R_{Vor;UB}(1.5; (1, -1))$ . (b):  $R_{Vor}(1.5; (1, -1))$ . (c):  $R_{Vor;LB}(1.5; (1, -1))$ . (d):  $R_{Std}(1.5; (1, -1))$ .

#### A.4 PROOF OF LEMMA 1

##### *Proof.* **Relation (3)**

First, we consider when  $\epsilon < \frac{1}{2}d(x, x_2)$ .

Let  $x' \in \mathbb{B}(x, \epsilon)$ . Then,  $d(x, x') \leq \epsilon$ .

$\frac{1}{2}d(x, x_2) \leq \frac{1}{2}d(x, x_{clean}), \forall x_{clean} \in X - \{x\}$ .

Due to the triangle inequality,  $\frac{1}{2}d(x, x_{clean}) \leq \frac{1}{2}d(x, x') + \frac{1}{2}d(x', x_{clean})$ .

When we combine the above inequalities,  $d(x, x') \leq \epsilon < \frac{1}{2}d(x, x_2) \leq \frac{1}{2}d(x, x_{clean}) \leq \frac{1}{2}d(x, x') + \frac{1}{2}d(x', x_{clean}), \forall x_{clean} \in X - \{x\}$ .

Then,  $\frac{1}{2}d(x, x') < \frac{1}{2}d(x', x_{clean}) = \frac{1}{2}d(x_{clean}, x'), \forall x_{clean} \in X - \{x\}$ . Thus,  $x' \in Vor(x)$ .

Hence,  $\mathbb{B}(x, \epsilon) \subset Vor(x)$  and  $R_{Vor}(\epsilon; x) = \mathbb{B}(x, \epsilon) \cap Vor(x) = \mathbb{B}(x, \epsilon)$ .

##### **Relation (4)**

Now, we consider when  $\frac{1}{2}d(x, x_j) \leq \epsilon < \frac{1}{2}d(x, x_{j+1})$  ( $j = 2, \dots, m+1$ ).

$R_{Vor}(\epsilon; x) = \mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^{N-1} HS(x, x_i) \right) \subset \mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^j HS(x, x_i) \right)$  is obvious as  $j \leq N-1$ .

We only need to proof  $\mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^j HS(x, x_i) \right) \subset R_{Vor}(\epsilon; x)$ .

Let  $x' \in \mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^j HS(x, x_i) \right)$ . Then,  $d(x, x') \leq \epsilon, d(x, x') < d(x_2, x'), \dots, d(x, x') < d(x_j, x')$ .

$\frac{1}{2}d(x, x_{j+1}) \leq \frac{1}{2}d(x, x_k)$  for  $k = j+1, \dots, N-1$ .

Due to the triangle inequality,  $\frac{1}{2}d(x, x_k) \leq \frac{1}{2}d(x, x') + \frac{1}{2}d(x', x_k)$ .

When we combine the above inequalities,  $d(x, x') \leq \epsilon < \frac{1}{2}d(x, x_{j+1}) \leq \frac{1}{2}d(x, x_k) \leq \frac{1}{2}d(x, x') + \frac{1}{2}d(x', x_k)$  for  $k = j+1, \dots, N-1$ .

Then,  $\frac{1}{2}d(x, x') < \frac{1}{2}d(x', x_k) = \frac{1}{2}d(x_k, x')$  for  $k = j+1, \dots, N-1$ .

We got  $d(x, x') \leq \epsilon, d(x, x') < d(x_2, x'), \dots, d(x, x') < d(x_{N-1}, x')$  and we proved  $\mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^j HS(x, x_i) \right) \subset R_{Vor}(\epsilon; x)$ .

##### **Relation (5)**

Finally, we consider when  $\epsilon \geq \frac{1}{2}d(x, x_{m+2})$ .

(i)  $\mathbb{B}(x, \frac{1}{2}d(x, x_{m+2}) - \tau) \cap \left( \bigcap_{i=2}^{m+1} HS(x, x_i) \right) \subset R_{Vor}(\epsilon; x)$  for  $\tau > 0$ :

Let  $x' \in \mathbb{B}(x, \frac{1}{2}d(x, x_{m+2}) - \tau) \cap \left( \bigcap_{i=2}^{m+1} HS(x, x_i) \right)$ . Then,  $d(x, x') \leq \frac{1}{2}d(x, x_{m+2}) - \tau <$

$\frac{1}{2}d(x, x_{m+2}) \leq \epsilon, d(x, x') < d(x_2, x'), \dots, d(x, x') < d(x_{m+1}, x')$ .

Through similar process used in the proof of **Relation (3)** and **Relation (4)**, we have  $d(x, x') < \frac{1}{2}d(x, x_{m+2}) \leq \frac{1}{2}d(x, x_k) \leq \frac{1}{2}d(x, x') + \frac{1}{2}d(x', x_k)$  for  $k = m+2, \dots, N-1$ .

Then,  $\frac{1}{2}d(x, x') < \frac{1}{2}d(x', x_k) = \frac{1}{2}d(x_k, x')$  for  $k = m+2, \dots, N-1$ .

We got  $d(x, x') < \epsilon, d(x, x') < d(x_2, x'), \dots, d(x, x') < d(x_{N-1}, x')$  and we proved (i).

(ii)  $R_{Vor}(\epsilon; x) \subset \mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^{m+1} HS(x, x_i) \right)$ :

It is obvious as  $R_{Vor}(\epsilon; x) = \mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^{N-1} HS(x, x_i) \right)$  and  $m+1 \leq N-1$ .  $\square$

#### A.5 PROOF OF OBSERVATION 1 AND PROPOSITION 1

##### *Proof.* **Observation 1**

$d_{min} \leq d(x, x_i), \forall x, x_i \in X, x \neq x_i$ .

When  $\epsilon < \frac{1}{2}d_{min}, \epsilon < \frac{1}{2}d_{min} \leq \frac{1}{2}d(x, x_i), \forall x, x_i \in X, x \neq x_i$ . Thus,  $R_{Vor}(\epsilon; x) = \mathbb{B}(x, \epsilon), \forall x \in X$  due to the relation (3) in Lemma 1.

Then,  $a_{Vor}(\epsilon)$  is same with  $a_{std}(\epsilon)$  as  $R_{Vor}(\epsilon; x) = \mathbb{B}(x, \epsilon) = R_{std}(\epsilon; x), \forall x \in X$ .

##### **Proposition 1**

First, we consider a data point  $x \in X$  and let  $x_2, \dots, x_N \in X - \{x\}$  be the sorted neighbors of  $x$ .

Let  $x^{*1} = \arg \max_{x' \in R_{std}(\epsilon; x)} L(x', c_x), x^{*2} = \arg \max_{x' \in R_{Vor; LB}(\epsilon; x)} L(x', c_x), x^{*3} = \arg \max_{x' \in R_{Vor}(\epsilon; x)} L(x', c_x)$ , and

$x^{*4} = \arg \max_{x' \in R_{Vor; UB}(\epsilon; x)} L(x', c_x)$ .

(i) When  $\epsilon < \frac{1}{2}d(x, x_{m+2})$ :

$R_{Vor; UB}(\epsilon; x) = R_{Vor}(\epsilon; x) = R_{Vor; LB}(\epsilon; x)$  from the definition.

$R_{Vor; LB}(\epsilon; x) = R_{Vor}(\epsilon; x) \subset \mathbb{B}(x, \epsilon) = R_{std}(\epsilon; x)$  from the relations (3) and (4).

Then,  $\mathbb{1}(f(x^{*1}) = c_x) \leq \mathbb{1}(f(x^{*2}) = c_x) = \mathbb{1}(f(x^{*3}) = c_x) = \mathbb{1}(f(x^{*4}) = c_x)$  as  $R_{Vor; UB}(\epsilon; x) = R_{Vor}(\epsilon; x) = R_{Vor; LB}(\epsilon; x) \subset R_{std}(\epsilon; x)$ .

(ii) When  $\epsilon \geq \frac{1}{2}d(x, x_{m+2})$ :

$R_{Vor; UB}(\epsilon; x) \subset R_{Vor}(\epsilon; x) \subset R_{Vor; LB}(\epsilon; x)$  from the relation (5).

$R_{Vor; LB}(\epsilon; x) = \mathbb{B}(x, \epsilon) \cap \left( \bigcap_{i=2}^{m+1} HS(x, x_i) \right) \subset \mathbb{B}(x, \epsilon) = R_{std}(\epsilon; x)$  from the definition.

Then,  $\mathbb{1}(f(x^{*1}) = c_x) \leq \mathbb{1}(f(x^{*2}) = c_x) \leq \mathbb{1}(f(x^{*3}) = c_x) \leq \mathbb{1}(f(x^{*4}) = c_x)$  as  $R_{Vor; UB}(\epsilon; x) \subset R_{Vor}(\epsilon; x) \subset R_{Vor; LB}(\epsilon; x) \subset R_{std}(\epsilon; x)$ .

From (i) and (ii),  $\mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f(x^{*1}) = c_x)] \leq \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f(x^{*2}) = c_x)] \leq \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f(x^{*3}) = c_x)] \leq \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f(x^{*4}) = c_x)]$ .

We finished the proof of the relation (6) as  $a_{std}(\epsilon) = \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f(x^{*1}) = c_x)]$ ,  $a_{Vor; LB}(\epsilon) = \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f(x^{*2}) = c_x)]$ ,  $a_{Vor}(\epsilon) = \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f(x^{*3}) = c_x)]$ , and  $a_{Vor; UB}(\epsilon) = \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f(x^{*4}) = c_x)]$ .  $\square$

#### A.6 PROOF OF THEOREM 1

To proof Theorem 1, we introduce the following lemma.

**Lemma 2.** By changing  $\epsilon$  and  $x \in X, x'$  that satisfies  $x' \in R_{Vor}(\epsilon; x)$  can fill up  $\mathcal{X}$  except for  $VB(X)$ . In other words,  $VB(X)^c = \mathcal{X} - VB(X) \subset \bigcup_{\epsilon \geq 0} \left( \bigcup_{x \in X} R_{Vor}(\epsilon; x) \right)$ .

##### *Proof.* **Lemma 2**

Let  $x' \in VB(X)^c$ .

$VB(X)^c = \mathcal{X} - VB(X) = \mathcal{X} - \left( \bigcup_{x \in X} Vor(x) \right)^c = \mathcal{X} \cap \left( \bigcup_{x \in X} Vor(x) \right) = \bigcup_{x \in X} Vor(x)$ .

$\exists x \in X$  such that  $x' \in Vor(x)$ .

Let  $\epsilon^* = d(x, x')$ . Then,  $d(x, x') \leq \epsilon^*$  and  $x' \in Vor(x)$ .

$x' \in \mathbb{B}(x, \epsilon^*) \cap Vor(x) = R_{Vor}(\epsilon^*; x) \subset \bigcup_{\epsilon \geq 0} \left( \bigcup_{x \in X} R_{Vor}(\epsilon; x) \right)$ .

We proved  $VB(X)^c \subset \bigcup_{\epsilon \geq 0} \left( \bigcup_{x \in X} R_{Vor}(\epsilon; x) \right)$ .  $\square$

Now, we proof Theorem 1.

**Proof. Part 1**

First, we prove that a 1-NN classifier maximizes global Voronoi-epsilon robustness. We denote the 1-NN classifier as  $f_{1-NN}$  and calculate its global Voronoi-epsilon robustness.

For a data point  $x \in X$ , let  $x' \in R_{Vor}(\epsilon; x) = \mathbb{B}(x, \epsilon) \cap Vor(x)$ .

$x' \in Vor(x) \iff d(x, x') < d(x_{clean}, x'), \forall x \in X - \{x\}$ .

As  $x' \in R_{Vor}(\epsilon; x) \subset Vor(x)$ ,  $x$  is unique nearest data point in  $X$  and thus  $f_{1-NN}(x') = c_x$ .

When  $x^* = \arg \max_{x' \in R_{Vor}(\epsilon; x)} L(x', c_x)$ ,  $a_{Vor}(\epsilon) = \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f_{1-NN}(x^*) = c_x)] = \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [1] =$

1.

$a_{global} = \lim_{\epsilon \rightarrow \infty} a_{Vor}(\epsilon) = \lim_{\epsilon \rightarrow \infty} 1 = 1$ . Thus,  $f_{1-NN}$  takes the maximum global Voronoi-epsilon robustness 1.

**Part 2**

Now, we prove that if  $f^*$  maximizes global Voronoi-epsilon robustness, then  $f^*$  becomes the 1-NN classifier except for Voronoi boundary  $VB(X)$ .

Let  $f^{*1}$  be a function that maximizes global Voronoi-epsilon robustness.

From the last part of the part 1, when we calculate global Voronoi-epsilon robustness of  $f^{*1}$ , it should satisfy the equation  $a_{global} = 1$ .

For a data point  $x \in X$  and  $\epsilon_1 < \epsilon_2$ ,  $R_{Vor}(\epsilon_1; x) = \mathbb{B}(x, \epsilon_1) \cap Vor(x) \subset \mathbb{B}(x, \epsilon_2) \cap Vor(x) = R_{Vor}(\epsilon_2; x)$ .

Thus, for a data point  $x \in X$  and  $\epsilon_1 < \epsilon_2$ ,  $\mathbb{1}(f^{*1}(x^{*1}) = c_x) \geq \mathbb{1}(f^{*1}(x^{*2}) = c_x)$  where  $x^{*1} = \arg \max_{x' \in R_{Vor}(\epsilon_1; x)} L(x', c_x)$  and  $x^{*2} = \arg \max_{x' \in R_{Vor}(\epsilon_2; x)} L(x', c_x)$ .

$a_{Vor}(\epsilon_1) = \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f^{*1}(x^{*1}) = c_x)] \geq \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f^{*1}(x^{*2}) = c_x)] = a_{Vor}(\epsilon_2)$  for  $\epsilon_1 < \epsilon_2$ . In other words,  $a_{Vor}(\epsilon)$  is a decreasing function.

$a_{Vor}(\epsilon) = 1, \forall \epsilon \geq 0$  ( $\because a_{Vor}(\epsilon^*) < 1$  for a  $\epsilon^* > 0$ , then it is a contradictory to  $a_{global} = 1$  as  $a_{Vor}(\epsilon)$  is a decreasing function.).

$1 = a_{Vor}(\epsilon) = \mathbb{E}_{(x, c_x) \sim \mathcal{D}} [\mathbb{1}(f^{*1}(x^*) = c_x)]$  where  $x^* = \arg \max_{x' \in R_{Vor}(\epsilon; x)} L(x', c_x)$ .

As the calculation is based on the finite set  $X$ ,  $f^{*1}(x^*) = c_x$  ( $\because \mathbb{1}(f^{*1}(x^*) = c_x) = 1$ ) where  $x^* = \arg \max_{x' \in R_{Vor}(\epsilon; x)} L(x', c_x)$ .

As  $x^*$  are the worst case adversarially perturbed samples, i.e., samples that output mostly different from  $c_x$ ,  $f^{*1}(x') = c_x = f_{1-NN}(x')$  where  $x' \in R_{Vor}(\epsilon; x)$ .

By changing  $\epsilon$  and  $x \in X$ ,  $x'$  that satisfies  $x' \in R_{Vor}(\epsilon; x)$  can fill up  $\mathcal{X}$  except for  $VB(X)$  ( $\because$  Lemma 2). Hence,  $f^{*1}$  is equivalent to  $f_{1-NN}$  except for Voronoi boundary  $VB(X)$ .

□