
Robustness for Non-Parametric Classification: A Generic Attack and Defense

Yao-Yuan Yang*

Cyrus Rashtchian*

Yizhen Wang

Kamalika Chaudhuri

{yay005, crashtchian, yiw248, kamalika}@eng.ucsd.edu

University of California, San Diego, Computer Science & Engineering

Abstract

Adversarially robust machine learning has received much recent attention. However, prior attacks and defenses for non-parametric classifiers have been developed in an ad-hoc or classifier-specific basis. In this work, we take a holistic look at adversarial examples for non-parametric classifiers, including nearest neighbors, decision trees, and random forests. We provide a general defense method, adversarial pruning, that works by preprocessing the dataset to become well-separated. To test our defense, we provide a novel attack that applies to a wide range of non-parametric classifiers. Theoretically, we derive an optimally robust classifier, which is analogous to the Bayes Optimal. We show that adversarial pruning can be viewed as a finite sample approximation to this optimal classifier. We empirically show that our defense and attack are either better than or competitive with prior work on non-parametric classifiers. Overall, our results provide a strong and broadly-applicable baseline for future work on robust non-parametrics.

1 Introduction

State-of-the-art classifiers have been shown to suffer from substantial drops in accuracy when faced with adversarially modified inputs even if the modifications are imperceptibly slight. Due to the security concerns that this raises, a body of recent research has investigated the construction and prevention of adversarial examples – small perturbations of valid inputs that cause misclassification (Carlini, 2018; Szegedy et al., 2014).

*Equal Contribution.

Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS) 2020, Palermo, Italy. PMLR: Volume 108. Copyright 2020 by the author(s).

Most previous work has looked at parametric methods, i.e., neural networks and linear classifiers (Biggio et al., 2013; Lowd and Meek, 2005; Madry et al., 2018; Papernot et al., 2016b), and there is a mature understanding of what properties can be exploited to design adversarial attacks and defenses for any parametric model. For example, parametric classifiers are based on continuous functions with gradients, which has been used to design gradient-based attacks (Athalye et al., 2018; Carlini and Wagner, 2017). Likewise, parametric models are mostly trained by minimizing a training loss, which has been exploited to build an effective and generic defense – adversarial training, retraining after data augmentation with adversarial examples (Carlini et al., 2019; Madry et al., 2018; Song et al., 2019).

An alternative statistical paradigm is that of non-parametric methods, such as nearest neighbor, decision tree, and random forest classifiers, which typically apply to dense data in lower dimensional spaces. These are local predictors, whose output depends on labeled points close to an input. Surprisingly, these methods behave very differently from parametrics when it comes to adversarial examples. In many cases, they have no gradients, and adversarial examples for parametric models fail to transfer (Papernot et al., 2016a). Generic defenses, such as adversarial training, appear to be ineffective as well (Dubey et al., 2019; Papernot and McDaniel, 2018; Wang et al., 2018).

While prior work has constructed attacks and defenses for some specific classifiers (Chen et al., 2019; Dubey et al., 2019; Kantchelian et al., 2016; Sitawarin and Wagner, 2019; Wang et al., 2018), there appear to be no generic approaches, and no generic principles that can be used to guide the design of attacks and defenses for variety of non-parametric methods.

In this work, we identify two key general principles, and use them to design a generic defense and an attack that apply to a variety of non-parametric methods.

To design defenses, we ask: when do non-parametric methods work well? Figure 1 depicts two variants of

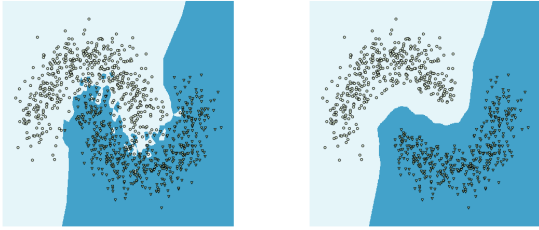


Figure 1: Normal vs. Defended 1-Nearest Neighbor.

random forests. In the left figure, we observe that datasets with nearby oppositely-labeled points may lead to classifiers with convoluted decision boundaries. In the right figure, we see that well-separated data lead to classification regions that are more robust to small perturbations. We will use this low-dimensional intuition as a starting point for generic defense methods.

Figure 1 suggests that since these methods make local predictions, they might work well when data from different classes are well-separated in space. We clearly cannot hope for such separation in most real datasets. Therefore, we propose to preprocess the training data by removing a subset so that different classes are well-separated. To ensure classification accuracy, we propose removing the minimal subset of points that ensure this property. We call our method *Adversarial Pruning*, which can be used as a pre-processing step before training any generic non-parametric classifier.

To evaluate our defense, we propose a new attack that is based on our next key observation: many non-parametric methods divide the instance space into convex polyhedra, and predict in a piecewise constant manner in each. For example, for 1-nearest neighbor, these polyhedra are the Voronoi cells. This suggests the following attack: find the closest polyhedron to an input where the classifier predicts a different label and output the closest point in this region. We implement this strategy by solving a collection of convex programs, and in cases where solution is computationally expensive, we provide a heuristic method for finding an approximate solution. We refer to these attacks as the exact and approximate *region-based attack*.

We next provide some theoretical justification for our methods. For our defense, we show that adversarial pruning can be interpreted as a finite-sample version of a robust analogue to the Bayes Optimal. We formally introduce this robust classifier, that we call the r -optimal, and show that it maximizes *astuteness* (accuracy where it is robust with radius r). For our attack, we show that the exact region-based attack is optimal, in the sense that it yields the closest adversarial example to a test input.

We empirically evaluate the adversarial pruning defense using the region based attack and prior attacks. We provide a general and thorough evaluation, for k -nearest neighbors (k -NN), decision trees, and random forests. We see that adversarial pruning consistently improves robustness, outperforming adversarial training on several datasets and is competitive with classifier-specific defenses. For our attacks, we see that even without any classifier-specific optimization, our new attacks either outperform or are competitive with prior attacks (in terms of perturbation amount). This suggests that both the adversarial pruning defense as well as the region based attack are good generic baselines for evaluating the robustness of non-parametric methods.

2 Preliminaries

We begin with a brief introduction to non-parametric methods that are local classifiers whose output depends on training data close to the test instance. These methods are typically used with dense lower-dimensional data, such as those in Figure 1. Examples are k -nearest neighbor (k -NN) and tree-based classifiers. The k -NN classifier outputs the plurality label among the k training examples closest to \mathbf{x} in an ℓ_p metric. A *tree ensemble* contains T decision trees whose leaves are labeled with vectors in \mathbb{R}^C . Each input \mathbf{x} determines T root-to-leaf paths, corresponding to vectors $\mathbf{u}^1, \dots, \mathbf{u}^T$. The output is the largest coordinate in $\mathbf{u}^1 + \dots + \mathbf{u}^T$. Random forests are a subclass of tree ensembles.

In what follows, $f : \mathbb{R}^d \rightarrow [C]$ denotes a classifier with C classes, where $[C] := \{1, 2, \dots, C\}$. The training data for f is a dataset $\mathcal{S} = \{(\mathbf{x}^j, y^j)\}_{j=1}^n$ of n labeled examples, with $\mathbf{x}^j \in \mathbb{R}^d$ and $y^j \in [C]$.

Robustness. We study robustness in an adversarial model. The adversary’s goal is to modify a true input by a small amount and cause the classifier to output the wrong label. Two main threat models have been proposed. The *black-box* setting restricts the adversary to only querying a classifier f on various inputs. In the *white-box* setting, the adversary has full access to f , including the model structure and parameters.

Fix a classifier f and a norm $\|\cdot\|$ on \mathbb{R}^d . An *adversarial example* for f at \mathbf{x} is any other input $\tilde{\mathbf{x}}$ such that $f(\mathbf{x}) \neq f(\tilde{\mathbf{x}})$. An *optimal adversarial example* for f at \mathbf{x} is an input $\tilde{\mathbf{x}}$ that minimizes $\|\mathbf{x} - \tilde{\mathbf{x}}\|$ subject to $f(\mathbf{x}) \neq f(\tilde{\mathbf{x}})$. In other words, an optimal adversarial example $\tilde{\mathbf{x}}$ is a closest vector to \mathbf{x} that receives a different label. In practice it is not always possible to find the optimal adversarial example, and hence the goal is to find $\tilde{\mathbf{x}}$ that is as close to \mathbf{x} as possible. We also define the robustness radius, the minimum perturbation needed to change the classifier label.

Definition 1. Let $\mathcal{X} \times [C]$ be a labeled space with norm $\|\cdot\|$. The *robustness radius* of f at $\mathbf{x} \in \mathcal{X}$ is

$$\rho(f, \mathbf{x}) := \min_{\tilde{\mathbf{x}} \in \mathcal{X}} \{\|\mathbf{x} - \tilde{\mathbf{x}}\| : f(\mathbf{x}) \neq f(\tilde{\mathbf{x}})\}$$

3 Adversarial Pruning Defense

When are non-parametric methods robust? Since these are local classifiers, Figure 1 suggests that they may be robust when training data from different classes is well-separated, and may fail when they overlap.

The training data may not be separated, so we will preprocess the data. We remove a subset of the training set, so that the remaining data are well-separated. Then, we train a non-parametric classifier on the rest. A remaining question is which subset of points to remove. For high classification accuracy, we remove the minimum subset whose removal ensures this property.

This process of removing examples from training set so that certain properties hold is called *pruning*. In this section, we first introduce the method used to prune the dataset. In Section 5, we justify our method by interpreting it in light of classical results in statistical learning theory (Chaudhuri and Dasgupta, 2014; Cover and Hart, 1967; Devroye et al., 1994).

Formally, given a robustness radius r and training set \mathcal{S} , we propose the following generic way to preprocess the training set and improve the robustness of classifiers:

Adversarial Pruning. Given r and a set \mathcal{S} , compute a maximum subset $\mathcal{S}^{\text{AP}} \subseteq \mathcal{S}$ such that differently-labeled points have distance at least $2r$. Then, train any nonparametric classifier on \mathcal{S}^{AP} .

After computing \mathcal{S}^{AP} once for a dataset, then we may train any classifier on the pruned training set. Our main hypothesis is that this will lead to more robust classifiers when using non-parametric methods. We will demonstrate empirically that this works well, and we will argue that this defense method is a finite-sample approximation to the optimal robust classifier.

Observe that while adversarial pruning is similar to the defense in Wang et al. (2018), they actually retain additional points with confident labels, which ensures that their method converges to being robust where the *Bayes Optimal* is robust. Their work builds on previous results of Gottlieb et al. (2014a) and Kontorovich and Weiss (2015) that sharpen the risk analysis of 1-NN by using pruning. As we explain in Section 5, our method instead can be interpreted as a finite sample version of a different and more appropriate limit.

One drawback of this approach is that the metric must be fine-grained enough to distinguish between close and far pairs. For most datasets and norms (e.g, Euclidean

distance) for which non-parametrics are used, this will be the case. However, for binary features and the ℓ_∞ distance, we have the problem that every pair of different points has distance exactly one, and therefore, the similarity structure is meaningless. To circumvent this, we preprocess the binary feature vectors using standard feature-extraction methods (e.g., PCA), and then operate on the resulting space.

Computing the Robust Dataset. We use known graph algorithms to efficiently compute \mathcal{S}^{AP} . Each training example is a vertex in the graph. Edges connect pairs of differently-labeled examples \mathbf{x} and \mathbf{x}' whenever $\|\mathbf{x} - \mathbf{x}'\| \leq 2r$. We remove as few examples as possible so that no more edges remain. This is equivalent to computing the minimum vertex cover. For binary labels, this graph is bipartite, and a minimum vertex cover can be derived from a maximum matching. The fastest method to solve maximum matching is the Hopcroft-Karp algorithm (Hopcroft and Karp, 1973). For a graph with n vertices and m edges, it takes time $O(m\sqrt{n})$. Fortunately, in practice, the graph of close pairs is quite sparse (for small r and high dimensional feature spaces, with relatively separated classes). For example, if $m = \tilde{O}(n)$ edges, then computing \mathcal{S}^{AP} takes time $\tilde{O}(n^{3/2})$. For large datasets, we note that *linear time* approximation algorithms are known (Duan and Pettie, 2014).

When there are more than two labels, that is $C \geq 3$, it is NP-Hard to compute the optimal pruned subset, but approximation algorithms are known (Gottlieb et al., 2014a; Kontorovich and Weiss, 2015). The greedy algorithm provably generates a 2-approximation. A suboptimal solution still ensures that different classes are separated, and hence, the robustness of the classifier does not require finding the optimal pruned dataset.

4 Region-Based Attack

In this section, we develop a way to evaluate robustness of non-parametric methods. For parametric algorithms, generic gradient-based attacks exist. Our goal is to develop an analogous general attack method, which works well for multiple non-parametrics. Moreover, we aim to develop a white-box attack that will serve as a better baseline than black-box attacks.

The main challenge of finding adversarial examples is that these classifiers have complicated decision regions. The central idea behind our attack is that for many classifiers, such as k -NN or random forests, we can decompose the decision regions into convex sets.

Definition 2. An (s, m) -*decomposition* is a partition of \mathbb{R}^d into convex polyhedra P_1, \dots, P_s such that each P_i can be described by up to m linear constraints, and f is

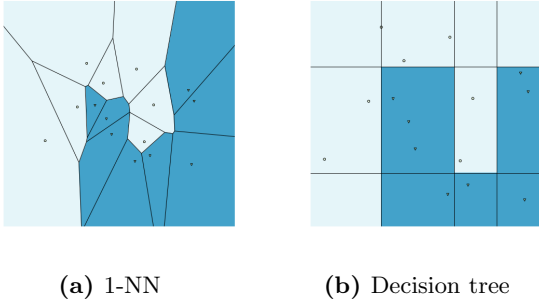


Figure 2: (s, m) -decompositions of two non-parametrics.

(s, m) -decomposable if there is an (s, m) -decomposition such that f is constant on P_i for each $i \in [s]$.

Figure 2 demonstrates the decomposition for two examples. Figure 2(a) shows how 1-NN is decomposed. In particular, a Voronoi diagram for n points is an $(n, n-1)$ -decomposition (P_1, \dots, P_n are Voronoi cells). If $k \geq 1$, then a k -NN classifier is $\binom{n}{k}, k(n-k)$ -decomposable; every k points correspond to polyhedra defined by $k(n-k)$ hyperplanes separating the k points from the other $n-k$ points (Aurenhammer, 1991).

Tree-based classifiers also fit into our framework, and Figure 2(b) shows how a decision tree is decomposed. Any decision tree of depth D with L leaves is (L, D) -decomposable; each root-to-leaf path corresponds to a polyhedron P_i defined by D hyperplanes. Generally, if f is an ensemble of T trees, each with depth D and L leaves, then f is (L^T, DT) -decomposable (proofs in Appendix A). An exponential dependence on T is expected, since the adversarial example problem for tree ensembles is NP-Hard (Kantchelian et al., 2016).

The existence of (s, m) -decompositions suggests the following attack. Given a classifier f and an input \mathbf{x} , suppose we could find the closest polyhedron P_i in the decomposition where f predicts a different label than $f(\mathbf{x})$. Then, the closest point in P_i would be the optimal adversarial example. Our attack implements this strategy by searching over all polyhedra.

Region-Based Attack. Let f be an (s, m) -decomposable classifier with decomposition P_1, \dots, P_s , where $f(\mathbf{z}) = y_i$ when $\mathbf{z} \in P_i$, for labels $y_i \in [C]$. To find an adversarial example for \mathbf{x} , consider all polyhedra P_i such that $f(\mathbf{x}) \neq y_i$. Then, output $\tilde{\mathbf{x}}$ minimizing

$$\min_{i: f(\mathbf{x}) \neq y_i} \min_{\mathbf{z} \in P_i} \|\mathbf{x} - \mathbf{z}\|. \quad (1)$$

Each P_i is described by $\leq m$ linear constraints, and the norm objective is convex (Boyd and Vandenberghe, 2004). Thus, we can solve each inner minimization problem in (1) separately by solving a convex program with $O(m)$ constraints. This results in candi-

dates $\mathbf{z}^i \in P_i$. Taking the outer minimum over i with $f(\mathbf{x}) \neq y_i$ leads to the optimal adversarial example $\tilde{\mathbf{x}} = \operatorname{argmin}_{\mathbf{z}^i} \|\mathbf{x} - \mathbf{z}^i\|$.

Efficiency. The running of the exact attack algorithm depends on two things: (i) the number of regions, which is based on the complexity of the classifier, and (ii) the number of constraints and dimensionality of the polyhedra. Due to advances in linear/quadratic program solvers, finding the adversarial example in a single region is quite efficient, i.e., the inner minimization problem in (1) is easy. We find that the number of regions s dominates the running time, i.e., the outer minimization problem in (1) is hard. For k -NN, the number of convex polyhedra scales with $O(n^k)$. When $k = 1$, this is efficiently solvable, because polyhedra have at most n constraints, and the adversarial examples can be found quickly using a linear program for ℓ_∞ perturbations. Unfortunately, for $k > 1$, this attack does not scale well, and we will develop an approximation algorithm for larger values of k .

For a single decision tree, again the exact attack is very efficient, depending only on the number of nodes in the tree. But for larger tree ensembles (e.g., large random forests), the optimal attack is very slow, as expected.

Speeding Up the Search. The exact attack is computationally intensive when s is large; hence, finding optimal solutions is infeasible for random forests (with many trees) or k -NN (when k is large). We next provide a computationally-efficient algorithm, which searches a constant number of regions.

The region-based attack for an (s, m) -decomposable f requires solving up to s convex programs, one for each polyhedron P_i with a different label. If the number of polyhedra is large, then this may be computationally infeasible. Fortunately, (1) has an obvious subdivision, based on the outer minimum over convex polyhedra. We use a relaxation that considers only a subset of polyhedra. We observe that each training point corresponds to a polyhedron—the one that f uses to predict the label. When finding adversarial examples for \mathbf{x} , the natural choice is to utilize training data close to \mathbf{x} .

Approximate Region-Based Attack. Let \mathcal{S} be the training data. To find an adversarial example under ℓ_p for \mathbf{x} , we first compute the subset $\mathcal{S}' \subseteq \mathcal{S}$ of s' points closest in ℓ_p distance to \mathbf{x} , while having different training labels than $f(\mathbf{x})$. Next, we determine at most s' polyhedra $P_{i_1}, \dots, P_{i_{s'}}$ containing points in \mathcal{S}' (as the polyhedra partition \mathbb{R}^d). We solve the inner optimization problem in (1) for each P_{i_j} to find candidates \mathbf{z}^i for $i \in [s']$. Finally, we output $\tilde{\mathbf{x}} = \operatorname{argmin}_{\mathbf{z}^i} \|\mathbf{x} - \mathbf{z}^i\|$, where the minimum is over these s' candidates.

As we only solve $s' \ll s$ convex programs, the running

time is greatly reduced compared to the optimal region-based attack. Empirically, this approximation finds adversarial examples with low perturbation.

5 Theoretical Justification

We provide some theoretical results to support our methods. To understand the robustness of non-parametric methods, we first derive a theoretically optimal classifier that takes into account robustness as a core objective. Then, we show that adversarial pruning can be interpreted as a finite sample approximation to the optimally robust classifier. Finally, we analyze the exact and approximate region-based attacks.

5.1 Adversarial Pruning vs. Optimal

Under certain conditions, many non-parametric methods converge in the infinite sample limit to the *Bayes Optimal classifier*, the most accurate classifier for a data distribution. In this way, non-parametric classifiers may be viewed as finite-sample approximations to the Bayes Optimal. However, the Bayes Optimal may not be robust to adversarial examples.

We next introduce a novel robust analogue to the Bayes Optimal. For a perturbation amount r , we call it the *r -Optimal classifier*. Surprisingly, to the best of our knowledge, such an analogue seems to be new in the context of adversarial examples.

Let μ denote a distribution on labeled examples $\mathcal{X} \times [C]$ and fix a distance on \mathcal{X} . What is the true objective of a robust classifier? Prior work measures astuteness under μ , which is the probability that the classifier is both r -robust and accurate for a new sample (\mathbf{x}, y) (Madry et al., 2018; Wang et al., 2018).

Definition 3. For distribution μ on $\mathcal{X} \times [C]$, the *astuteness* of a classifier f at radius r is

$$\text{ast}_\mu(f, r) := \Pr_{(\mathbf{x}, y) \sim \mu} [\rho(f, \mathbf{x}) \geq r \text{ and } f(\mathbf{x}) = y].$$

Robust Analogue to Bayes Optimal. We exhibit a classifier, the *r -Optimal classifier*, that achieves optimal astuteness. It is convenient to rewrite astuteness in terms of certain robust subsets of the input space. Then, we define the *r -Optimal classifier* using these subsets. Formally, for a classifier f and label j , let $S_j(f, r) := \{\mathbf{x} \in \mathcal{X} \mid f(\mathbf{x}) = j \text{ and } \rho(f, \mathbf{x}) \geq r\}$. The following lemma expresses astuteness under μ using these subsets (proofs in Appendix B).

Lemma 1. $\text{ast}_\mu(f, r) = \sum_{j=1}^C \int_{\mathbf{x} \in S_j(f, r)} p(y = j \mid \mathbf{x}) d\mu.$

How should we define the classifier that maximizes astuteness? Lemma 1 implies that, to calculate astute-

ness, it suffices to consider the robust regions $S_j(f, r)$ for a classifier. As a consequence, we claim that in order to determine the optimal classifier, it suffices to find the optimal robust regions under μ . We first formalize this intermediate goal using the following maximization problem.

$$\begin{aligned} \max_{S_1, \dots, S_C} \sum_{j=1}^C \int_{\mathbf{x} \in S_j} p(y = j \mid \mathbf{x}) d\mu \quad (2) \\ \text{s.t. } d(S_j, S_{j'}) \geq 2r \text{ for all } j \neq j' \end{aligned}$$

where $d(S_j, S_{j'}) := \min_{u \in S_j, v \in S_{j'}} \|u - v\|$. Notice that for any classifier f , the sets $S_j(f, r)$ for $j \in [C]$ have pairwise distance at least $2r$, implying that they are feasible solutions for (2).

Besides being distance $2r$ apart, an optimal solution S_1^*, \dots, S_C^* to (2) maximizes accuracy in the following sense. The integral measures the probability that $(\mathbf{x}, y) \sim \mu$ has $y = j$ and $\mathbf{x} \in S_j^*$. In other words, S_j^* has the highest frequency of points with label j under μ , subject to the distance constraint. The sets S_j^* form the basis for the optimal classifier’s decision regions. To ensure the separation, we consider the distance r ball around these sets. Formally, we have the following.

Definition 4. Fix r and μ . Let S_1^*, \dots, S_C^* be optimizers of (2). The *r -Optimal classifier* f_{ropt} is any classifier such that $f_{\text{ropt}}(\mathbf{x}) = j$ whenever $d(\mathbf{x}, S_j^*) \leq r$.

We remark that when $r = 0$, the 0-Optimal classifier is the standard Bayes Optimal classifier. Finally, because $S_j(f_{\text{ropt}}, r) = S_j^*$, Lemma 1 then implies that *r -Optimal classifier* maximizes astuteness:

Theorem 1. $f_{\text{ropt}} = \text{argmax}_f \text{ast}_\mu(f, r).$

Finite Sample Approximation. Prior work shows that 1-NN applied to a variant of adversarial pruning leads to provably robust classifiers (Wang et al., 2018). The main difference with our work is their method also selects a subset of confident training examples to keep in the pruned subset - which ensures that the classifier converges to being robust in regions where the Bayes Optimal is robust. In contrast, our aim is to develop generic techniques, for multiple classifiers, and we show that our method can be interpreted as a finite sample approximation to the *r -Optimal classifier* - the optimally astute classifier.

Adversarial pruning works by removing certain training points so that no oppositely labeled pairs of examples remain. We can view this process in the light of the *r -optimal classifier* as follows. To prune the dataset \mathcal{S} , we solve the maximization problem:

$$\begin{aligned} \max_{S_1, \dots, S_C \subseteq \mathcal{S}} \sum_{j=1}^C \sum_{\mathbf{x}^i \in S_j} \mathbf{1}_{\{y^i = j\}} \quad (3) \\ \text{s.t. } d(S_j, S_{j'}) \geq 2r \text{ for all } j \neq j'. \end{aligned}$$

The solution to (3) will be maximum subsets of training data with pairwise distance $2r$. As long as the training set \mathcal{S} is representative of the underlying distribution μ , these subsets will approximate the optimal S_j^* sets. Hence, we posit that a non-parametric method trained on \mathcal{S}^{AP} should approximate the r -Optimal classifier.

5.2 Attack Algorithm Analysis

The run time of the region-based attack depends on the norm. We focus on ℓ_p with $p \in \{1, 2, \infty\}$ as these are the most relevant for adversarial examples. We prove the following theorem in Appendix A.

Theorem 2. *If f is (s, m) -decomposable, then the region-based attack outputs optimal adversarial examples in time $s \cdot \text{poly}(m, d)$, for ℓ_p distance, $p \in \{1, 2, \infty\}$.*

As k -NN and tree ensembles are (s, m) -decomposable, the region-based attack produces an optimal adversarial example for these. Note that an optimal attack *certifies* the robustness radius. Indeed, if on input \mathbf{x} the region-based attack outputs $\tilde{\mathbf{x}}$, then $\|\mathbf{x} - \tilde{\mathbf{x}}\| = \rho(f, \mathbf{x})$.

Approximate Attack Guarantees. We claim that the approximate region-based attack outputs a valid adversarial example when f is (s, m) -decomposable. Each region is defined by m constraints, and f is constant on each region. We search in s' regions, finding the best candidate \mathbf{z}^i from each. Each considered region contains a training example with a different label than $f(\mathbf{x})$. Therefore, the best adversarial example $\tilde{\mathbf{x}}$ in that region receives a different label $f(\tilde{\mathbf{x}}) \neq f(\mathbf{x})$. The analysis of the time complexity for finding candidates is $\text{poly}(m, d)$ for each region P_i . Compared to the exact attack (Theorem 2) we only consider s' regions, so the total time is only $s' \cdot \text{poly}(m, d)$. We find in practice that $s' = 50$ regions suffices for a good attack, and the time only scales with m and d .

6 Experiments

We investigate the effectiveness of our methods by evaluating multiple classifiers on nine datasets. We address the following questions:

1. Does adversarial pruning increase robustness across multiple non-parametric classifiers?
2. How well does the region-based attack perform compared with prior work?

Classifiers and Datasets. We evaluate three non-parametric classifiers: k -nearest neighbor (k -NN), decision tree (DT) and random forest (RF) (Breiman, 2001, 2017; Cover and Hart, 1967). We use nine standard binary classification datasets. All features are scaled to be in $[0, 1]$. We evaluate in ℓ_∞ to be consistent with prior work. We reduce the feature dimension of the

image datasets (f-mnist and mnist) with PCA to 25 dimensions for two reasons: (i) non-parametrics are normally used for low dimensional spaces, (ii) adversarial pruning requires non-binary features for ℓ_∞ . Details are in Appendix C; code in a public repository.¹

Performance Measures. Besides measuring accuracy, we evaluate attacks using empirical robustness, following prior work (Chen et al., 2019; Kantchelian et al., 2016). Intuitively, we want to measure the perturbation distance to the nearest adversarial example (as opposed to fixing r and evaluating error). Formally, the *empirical robustness* for attack A on f at input \mathbf{x} is $\text{ER}(A, f, \mathbf{x}) := \|\mathbf{x} - \tilde{\mathbf{x}}_A\|_\infty$, where A outputs $\tilde{\mathbf{x}}_A$ as the adversarial example for f at \mathbf{x} . Observe that larger empirical robustness means worse attacks, and the minimal empirical robustness of f at \mathbf{x} is the robustness radius $\rho(f, \mathbf{x})$. To fairly compare classifiers having different accuracies, we actually compute $\text{ER}(A, f, S, t)$ over t test inputs. To do so, we draw t random samples S_i from S that are classified correctly by f , and we report the average of $\text{ER}(A, f, \mathbf{x})$ over $\mathbf{x} \in S_i$. We set $t = 100$ to balance efficiency and thoroughness.

Again, for defenses, we use perturbation distance to evaluate robustness. Each defense method D produces a classifier f_D . We evaluate a defense D by assigning it a score, the *defscore*. The *defscore* with respect to an attack A , a test set S and test size t is the ratio

$$\text{defscore}(D, A, f, S, t) = \frac{\text{ER}(A, f_D, S, t)}{\text{ER}(A, f, S, t)},$$

where f is the undefended classifier. A larger *defscore* implies a better defense.

Attack Algorithms. For 1-NN and DT, we apply the exact region-based attack (RBA-Exact). For 3-NN and RF, the RBA-Exact attack is computationally intensive, and we use the approximate region-based attack (RBA-Approx). For 3-NN, it uses $s' = 50$ polyhedra, and for RF, it uses $s' = 100$ polyhedra. We compare RBA-Exact and RBA-Approx against several baselines. A general attack that applies to all methods is the black-box attack (BBox) (Cheng et al., 2019); this attack seems to be the state-of-the-art for non-parametrics. For k -NN, we compare against two white-box attacks, the direct attack (Direct) and kernel substitution attack (Kernel) (Papernot et al., 2016a). The direct attack perturbs the test instance towards the center of the k nearest oppositely-labeled training examples. The kernel substitution attack uses a soft nearest neighbor to build a substitution model and applies the projected gradient descent attack (Kurakin et al., 2016). For DT, the RBA-Exact attack is optimal, and so is the attack

¹<https://github.com/yangarbiter/adversarial-nonparametrics/>

	1-NN					3-NN				DT			RF	
	Direct	BBox	Kernel	RBA Exact	RBA Approx	Direct	BBox	Kernel	RBA Approx	Papernot's	BBox	RBA Exact	BBox	RBA Approx
austr.	.442	.336	.379	.151	.151	.719	.391	.464	.278	.140	.139	.070	.364	.446
cancer	.223	.364	.358	.137	.137	.329	.376	.394	.204	.459	.334	.255	.451	.383
covtype	.130	.199	.246	.066	.067	.200	.259	.280	.108	.254	.083	.051	.233	.214
diabetes	.074	.112	.165	.035	.035	.130	.143	.191	.078	.237	.133	.085	.181	.184
f-mnist06	.080	.140	.187	.029	.030	.129	.169	.202	.051	.189	.134	.079	.206	.188
f-mnist35	.187	.244	.259	.075	.077	.234	.238	.266	.094	.262	.185	.115	.188	.246
fourclass	.109	.124	.137	.090	.090	.101	.113	.134	.096	.288	.197	.137	.159	.133
halfmoon	.070	.129	.102	.058	.058	.105	.132	.115	.096	.098	.148	.085	.182	.149
mnist17	.161	.251	.262	.070	.073	.221	.261	.269	.097	.219	.171	.123	.250	.250

Table 1: The Empirical Robustness for different attacks on four classifiers (lower is better; best is in bold).

	1-NN			3-NN			DT			RF		
	AT	WJC	AP	AT	AP	AT	RS	AP	AT	RS	AP	
aus.	0.64	1.65	1.65	0.68	1.20	2.36	5.86	2.37	1.07	1.12	1.04	
can.	0.82	1.05	1.41	1.06	1.39	0.85	1.09	1.19	0.87	1.54	1.26	
cov.	0.61	4.38	4.38	0.88	3.31	1.47	2.73	4.51	1.02	1.01	2.13	
dia.	0.83	4.69	4.69	0.87	2.97	0.93	1.53	2.22	1.19	1.25	2.22	
f06	0.90	1.93	2.59	0.88	1.75	1.33	2.33	2.57	1.04	1.10	1.77	
f35	0.83	1.05	1.19	0.83	1.15	0.97	3.03	2.06	0.99	1.23	1.41	
fou.	0.93	3.09	3.09	0.89	3.09	1.06	1.23	3.04	1.03	1.92	3.59	
hal.	1.05	2.00	2.78	0.93	1.92	1.54	1.98	2.58	1.04	1.01	1.82	
m17	0.88	1.06	1.39	0.80	1.13	1.11	3.97	1.32	0.88	0.92	1.26	

Table 2: defscore using different defenses (higher is better; best is in bold). The defscore for undefended classifiers is 1.00 (greater than 1.00 is more robust). We use RBA-Exact for 1-NN and DT, and RBA-Approx for 3-NN and RF. We use RBA-Approx for AT on large datasets.

by Kantchelian et al. (2016); we only report RBA-Exact because these achieve the same results. We also evaluate the heuristic DT attack by Papernot et al. (2016a). For RF, both optimal attacks are infeasible, and we only evaluate BBox and RBA-Approx.

Defense Methods. For our defense, we train each classifier on the dataset pre-processed with adversarial pruning (AP); we use ℓ_∞ to determine examples to prune. For the separation r of AP, we found that $r = 0.3$ balances robustness vs. accuracy. We set $r = 0.3$ for all datasets (Appendix C.4 has other r settings). A generic baseline is adversarial training (AT), where the training data is augmented with examples generated by the corresponding attack algorithm. AT has been reported to be ineffective for 1-NN and boosted decision tree (Wang et al., 2018; Chen et al., 2019), but we include it for completeness. For AT, we retrain the classifier after attacking each training point once; we augment the training data with adversarial examples that are distance at most 0.3 from the original input. The parameter 0.3 matches the parameter r for AP. For 1-NN, an available baseline defense is Wang et al. (2018), but for general k -NN, we are not aware of other defenses. For DT and RF, we compare against the best known defense algorithm, Robust Splitting (RS) (Chen et al., 2019). We set the RS parameter to 0.3 as well.

Results. We separately evaluate attacks and defenses, in Tables 1 and 2, respectively. We provide an accuracy vs. perturbation distance experiment in Figure 3.

Effectiveness of Attacks. Table 1 exhibits empirical robustness across four undefended classifiers and nine

datasets. Recall that a smaller empirical robustness implies a more effective attack. For 1-NN, we see that RBA-Exact works as expected, achieving the smallest empirical robustness. For 3-NN, our RBA-Approx attack is more effective than prior attacks, with a much lower empirical robustness. This indicates that RBA-Approx can be a strong attack for $k > 1$, where previously no consistently effective baseline is known. For DT, RBA-Exact again has the best performance. The improvement in many cases shows that the optimal attack for 1-NN and DT can be significantly better than heuristics, which will lead to a more informative defense evaluation. For RF, RBA-Approx wins on five of the nine datasets, and BBox wins on four. Overall, our RBA-Approx attack is competitive with the state-of-the-art attack for RF, and better for 3-NN.

Effectiveness of Defenses. Table 2 shows defscore across four classifiers and several defense methods. For each dataset, the AP defense trains all four classifiers on the same pruned version of the dataset. For all classifiers, we see that AP results in a greater than one defscore, indicating that classifiers trained with AP are more robust. In contrast, AT usually achieves defscore less than one, worse than the undefended classifier; this corroborates previous results (Wang et al., 2018). For 1-NN, observe that AP is slightly better than the defense of Wang et al. (2018). We believe that this is because their method converges to Bayes Optimal, while AP approximates the r -Optimal classifier. For the DT and RF experiments, we see that RS and AP perform competitively, each winning out on some datasets. Overall, AP performs slightly better than RS. We remark that we have evaluated 1-NN and DT against the optimal attack. This provides concrete evidence that AP leads to a more robust classifier.

Discussion. From the results, we see that our generic attack and defense either outperform or perform competitively with prior work on many datasets. We note that there can be a big difference in the perturbation distance depending on the attack algorithms. We also see that our adversarial pruning achieves more robustness compared both to undefended variants and to the classifiers trained using adversarial training. Surprisingly, the pruned subset is computed ahead of time, yet

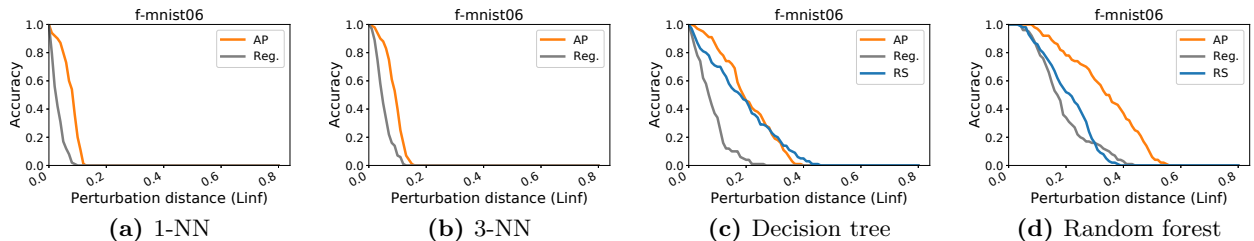


Figure 3: Accuracy (y-axis) vs. perturbation distance (x-axis) for four classifiers on Fashion MNIST classes 0 vs. 6 for the ℓ_∞ distance after applying PCA to 25 dimensions (**larger accuracy is better**). Other datasets appear in Appendix C.4.1. In the legend, Reg. = regular (undefended) classifier, AP = adversarial pruning, and RS = robust splitting.

it improves the robustness of many different classifiers.

The main conclusion from the experiments is that our work provides a new and suitable baseline for many methods. This is analogous to how AT and PGD are generic baselines for parametrics. In particular, if a new non-parametric algorithm is developed, then AP and RBA may be used to evaluate robustness. Our work also opens to the door to combine AP with classifier-specific defenses, e.g. robust boosting (Chen et al., 2019). We note that our methods can sometimes be slow, but we expect that classifier-specific optimizations and techniques will readily improve the running time.

7 Related Work

The bulk of research on robust classifiers has focused on parametric models, with many generic attacks (Carlini and Wagner, 2017; Liu et al., 2017; Papernot et al., 2017b, 2016b; Szegedy et al., 2014), as well as defenses (Hein and Andriushchenko, 2017; Katz et al., 2017; Madry et al., 2018; Papernot et al., 2015; Raghu-nathan et al., 2018; Sinha et al., 2018). In contrast, adversarial examples for non-parametrics have been studied in a more case-by-case basis.

For tree ensembles, Kantchelian et al. (2016) formulate an optimal attack as a Mixed Integer Linear Program (superseding an earlier attack (Papernot et al., 2016a)) and prove NP-Hardness for many trees. Chen et al. (2019) increase the robustness of *boosted* ensembles. Concurrent work also studies the robustness of decision stumps, and we leave it as future work to compare our methods to theirs (Andriushchenko and Hein, 2019).

For k -NN, prior work on adversarial examples only considers suboptimal attacks, such the direct attack and variants thereof (Amsaleg et al., 2017; Sitawarin and Wagner, 2019; Wang et al., 2018). Concurrent work (Khoury and Hadfield-Menell, 2019) on Voronoi-based adversarial training for neural networks also introduces the optimal attack for 1-NN (i.e., Region-Based attack restricted to 1-NN). In terms of defenses, Wang et al. (2018) increase 1-NN robustness by strategically

removing training points. Besides only testing 1-NN against suboptimal attacks, they do not consider other non-parametrics; additionally, their defense is shown to be robust in the large sample limit only where the Bayes Optimal is robust. Our methods are thus more general, and our defense can be interpreted as a finite sample approximation to the r -optimal classifier.

Outside the realm of adversarial examples, pruning has been used to improve the accuracy and generalization (but not robustness) of 1-NN (Gates, 1972; Gottlieb et al., 2014b; Hart, 1968; Kontorovich et al., 2017). Related attacks and defenses have been developed for ReLU networks (Croce et al., 2019; Jordan et al., 2019; Tjeng et al., 2019; Xiao et al., 2019). These results do not directly pertain to non-parametrics, as ReLUs are fundamentally different. The geometric attacks and defenses are similar in spirit to ours. Optimizations based on the dual formulation may improve the efficiency of our methods (Tjeng et al., 2019; Xiao et al., 2019). It would be interesting to explore the relationship between our defense method (adversarial pruning) and the ReLU defense methods and robustness certificates. For example, do robust ReLU networks approximate or converge to the r -Optimal classifier?

8 Conclusion

We consider adversarial examples for non-parametric methods, with a focus on *generic* attacks and defenses. We provide a new attack, the region-based attack, which often outperforms previous attacks. We also provide a new method of defense, adversarial pruning, which should serve as a strong baseline for evaluating the robustness of many classifiers. On the theory side, we prove that the region-based attack outputs the optimal adversarial example. We also introduce and analyze a novel robust analogue to the Bayes Optimal. We prove that the r -Optimal classifier maximizes as-tuteness. On the experimental side, we demonstrate that our methods are better than or competitive with prior work, while being considerably more general.

Acknowledgments. We thank Somesh Jha, Ruslan Salakhutdinov and Michal Moshkovitz for helpful discussions. Part of this research is supported by ONR under N00014-16-1-261, UC Lab Fees under LFR 18-548554 and NSF under 1804829 and 1617157.

References

- Amsaleg, L., Bailey, J., Barbe, D., Erfani, S., Houle, M. E., Nguyen, V., and Radovanović, M. (2017). The vulnerability of learning to adversarial perturbation increases with intrinsic dimensionality. In *WIFS*, pages 1–6.
- Andriushchenko, M. and Hein, M. (2019). Provably robust boosted decision stumps and trees against adversarial attacks. *arXiv preprint arXiv:1906.03526*.
- Athalye, A., Carlini, N., and Wagner, D. (2018). Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples. In *ICML*, pages 274–283.
- Aurenhammer, F. (1991). Voronoi diagrams—a survey of a fundamental geometric data structure. *ACM Computing Surveys*, 23(3):345–405.
- Biggio, B., Corona, I., Maiorca, D., Nelson, B., Šrncić, N., Laskov, P., Giacinto, G., and Roli, F. (2013). Evasion attacks against machine learning at test time. In *ECML-PKDD*, pages 387–402.
- Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge Univ. Press.
- Breiman, L. (2001). Random forests. *Machine learning*, 45(1):5–32.
- Breiman, L. (2017). *Classification and regression trees*. Routledge.
- Carlini, N. (2018). *Evaluation and Design of Robust Neural Network Defenses*. PhD thesis, EECS Department, University of California, Berkeley.
- Carlini, N., Athalye, A., Papernot, N., Brendel, W., Rauber, J., Tsipras, D., Goodfellow, I. J., Madry, A., and Kurakin, A. (2019). On Evaluating Adversarial Robustness. *CoRR*, abs/1902.06705.
- Carlini, N. and Wagner, D. (2017). Towards evaluating the robustness of neural networks. *IEEE Symposium on Security and Privacy*.
- Chaudhuri, K. and Dasgupta, S. (2014). Rates of convergence for nearest neighbor classification. In *NeurIPS*, pages 3437–3445.
- Chen, H., Zhang, H., Boning, D., and Hsieh, C.-J. (2019). Robust Decision Trees Against Adversarial Examples. In *ICML*.
- Cheng, M., Le, T., Chen, P.-Y., Yi, J., Zhang, H., and Hsieh, C.-J. (2019). Query-efficient Hard-label Black-box Attack: An Optimization-based Approach. In *ICLR*.
- Cover, T. M. and Hart, P. E. (1967). Nearest neighbor pattern classification. *IEEE transactions on information theory*, 13(1):21–27.
- Croce, F., Andriushchenko, M., and Hein, M. (2019). Provable robustness of relu networks via maximization of linear regions. In *AISTATS*.
- Devroye, L., Györfi, L., Krzyżak, A., and Lugosi, G. (1994). On the strong universal consistency of nearest neighbor regression function estimates. *The Annals of Statistics*, pages 1371–1385.
- Duan, R. and Pettie, S. (2014). Linear-time approximation for maximum weight matching. *Journal of the ACM (JACM)*, 61(1):1.
- Dubey, A., van der Maaten, L., Yalniz, Z., Li, Y., and Mahajan, D. (2019). Defense against adversarial images using web-scale nearest-neighbor search. *arXiv preprint arXiv:1903.01612*.
- Gates, G. (1972). The Reduced Nearest Neighbor Rule. *IEEE transactions on information theory*, 18(3):431–433.
- Gottlieb, L.-A., Kontorovich, A., and Krauthgamer, R. (2014a). Efficient classification for metric data. *IEEE Transactions on Information Theory*, 60(9):5750–5759.
- Gottlieb, L.-A., Kontorovich, A., and Nisnevitch, P. (2014b). Near-Optimal Sample Compression for Nearest Neighbors. In *NeurIPS*, pages 370–378.
- Gurobi Optimization, L. (2018). Gurobi optimizer reference manual.
- Hart, P. (1968). The Condensed Nearest Neighbor Rule. *IEEE transactions on information theory*, 14(3):515–516.
- Hein, M. and Andriushchenko, M. (2017). Formal guarantees on the robustness of a classifier against adversarial manipulation. In *NeurIPS*, pages 2263–2273.
- Hopcroft, J. E. and Karp, R. M. (1973). An $n^{5/2}$ algorithm for maximum matchings in bipartite graphs. *SIAM Journal on computing*, 2(4):225–231.
- Jordan, M., Lewis, J., and Dimakis, A. G. (2019). Provable Certificates for Adversarial Examples: Fitting a Ball in the Union of Polytopes. *arXiv preprint arXiv:1903.08778*.
- Kantchelian, A., Tygar, J., and Joseph, A. (2016). Evasion and Hardening of Tree Ensemble Classifiers. In *ICML*, pages 2387–2396.
- Katz, G., Barrett, C., Dill, D. L., Julian, K., and Kochenderfer, M. J. (2017). Towards proving the adversarial robustness of deep neural networks. *arXiv preprint arXiv:1709.02802*.

- Khoury, M. and Hadfield-Menell, D. (2019). Adversarial Training with Voronoi Constraints. *Safe Machine Learning workshop at ICLR*.
- Kontorovich, A., Sabato, S., and Weiss, R. (2017). Nearest-neighbor Sample Compression: Efficiency, Consistency, Infinite Dimensions. In *NeurIPS*, pages 1573–1583.
- Kontorovich, A. and Weiss, R. (2015). A Bayes Consistent 1-NN classifier. In *AISTATS*.
- Kurakin, A., Goodfellow, I. J., and Bengio, S. (2016). Adversarial examples in the physical world.
- Liu, Y., Chen, X., Liu, C., and Song, D. (2017). Delving into transferable adversarial examples and black-box attacks. *ICLR*.
- Lowd, D. and Meek, C. (2005). Adversarial learning. In *SIGKDD*, pages 641–647.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. (2018). Towards deep learning models resistant to adversarial attacks. In *ICLR*.
- Manning, C., Raghavan, P., and Schütze, H. (2010). Introduction to Information Retrieval. *Natural Language Engineering*, 16(1):100–103.
- Mulmuley, K. (1991). On levels in arrangements and voronoi diagrams. *Discrete & Computational Geometry*, 6(3):307–338.
- Papernot, N., Carlini, N., Goodfellow, I., Feinman, R., Faghri, F., Matyasko, A., Hambardzumyan, K., Juang, Y.-L., Kurakin, A., Sheatsley, R., Garg, A., and Lin, Y.-C. (2017a). cleverhans v2.0.0: an adversarial machine learning library. *arXiv preprint arXiv:1610.00768*.
- Papernot, N. and McDaniel, P. (2018). Deep k-nearest neighbors: Towards confident, interpretable and robust deep learning. *arXiv preprint arXiv:1803.04765*.
- Papernot, N., McDaniel, P., and Goodfellow, I. (2016a). Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277*.
- Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, B., and Swami, A. (2017b). Practical black-box attacks against deep learning systems using adversarial examples. In *ASIACCS*.
- Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z. B., and Swami, A. (2016b). The limitations of deep learning in adversarial settings. In *EuroS&P*.
- Papernot, N., McDaniel, P., Wu, X., Jha, S., and Swami, A. (2015). Distillation as a defense to adversarial perturbations against deep neural networks. *arXiv preprint arXiv:1511.04508*.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Raghunathan, A., Steinhardt, J., and Liang, P. (2018). Certified defenses against adversarial examples. In *ICLR*.
- Sinha, A., Namkoong, H., and Duchi, J. (2018). Certifiable Distributional Robustness with Principled Adversarial Training. In *ICLR*.
- Sitawarin, C. and Wagner, D. (2019). On the Robustness of Deep K-Nearest Neighbors. *arXiv preprint arXiv:1903.08333*.
- Song, C., He, K., Wang, L., and Hopcroft, J. E. (2019). Improving the Generalization of Adversarial Training with Domain Adaptation. In *ICLR*.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. (2014). Intriguing properties of neural networks. In *ICLR*.
- Tjeng, V., Xiao, K., and Tedrake, R. (2019). Evaluating Robustness of Neural Networks with Mixed Integer Programming. In *ICLR*.
- Wang, Y., Jha, S., and Chaudhuri, K. (2018). Analyzing the Robustness of Nearest Neighbors to Adversarial Examples. In *ICML*, pages 5120–5129.
- Xiao, K. Y., Tjeng, V., Shafiqullah, N. M., and Madry, A. (2019). Training for faster adversarial robustness verification via inducing relu stability. In *ICLR*.