

UNIVERSITY OF CALIFORNIA SAN DIEGO

The Primacy of Applied Privacy

A dissertation submitted in partial satisfaction of the
requirements for the degree Doctor of Philosophy

in

Computer Science

by

Casey Meehan

Committee in charge:

Professor Kamalika Chaudhuri, Chair
Professor Taylor Berg-Kirkpatrick
Professor Sanjoy Dasgupta
Professor Alon Orlitsky

2023

Copyright
Casey Meehan, 2023
All rights reserved.

The Dissertation of Casey Meehan is approved, and it is acceptable in quality and form for publication on microfilm and electronically.

University of California San Diego

2023

DEDICATION

The fact that I have made something I can write a dedication for is owed all to my parents. I cannot imagine following my heart these past few years without their unrelenting support and encouragement.

TABLE OF CONTENTS

| | |
|--|------|
| Dissertation Approval Page | iii |
| Dedication | iv |
| Table of Contents | v |
| List of Figures | viii |
| List of Tables | x |
| Acknowledgements | xi |
| Vita | xii |
| Abstract of the Dissertation | xiii |
| Chapter 1 When are Non-Parametric Methods Robust? | 1 |
| 1.1 Introduction | 1 |
| 1.1.1 Related Work | 3 |
| 1.2 Preliminaries | 4 |
| 1.2.1 Setting | 4 |
| 1.2.2 Notions of Consistency | 5 |
| 1.2.3 Non-parametric Classifiers | 8 |
| 1.3 Warm Up: r -separated distributions | 10 |
| 1.4 General Distributions | 14 |
| 1.4.1 The r -Optimal Classifier and Adversarial Pruning | 14 |
| 1.4.2 Convergence Guarantees | 15 |
| 1.5 Validation | 17 |
| 1.5.1 Experimental Setup | 18 |
| 1.5.2 Results | 19 |
| 1.5.3 Discussion | 19 |
| 1.6 Conclusion | 20 |
| Chapter 2 Robust Empirical Risk Minimization with Tolerance | 21 |
| Bibliography | 22 |

LIST OF FIGURES

| | | |
|-------------|---|----|
| Figure 1.1. | H_S is astute in the green region, but not robust in the red region. | 10 |
| Figure 1.2. | Empirical accuracy/astuteness of different classifiers as a function of training sample size. Accuracy is shown in green, astuteness in purple. Left : Noiseless Setting. Right: Noisy Setting. Top Row: Histogram Classifier, Bottom Row: 1-Nearest Neighbor | 17 |

LIST OF TABLES

ACKNOWLEDGEMENTS

VITA

| | |
|------|---|
| 2015 | Bachelor of Science, Brown University |
| 2018 | Master of Science, Harvard University |
| 2023 | Doctor of Philosophy, University of California, San Diego |

ABSTRACT OF THE DISSERTATION

The Primacy of Applied Privacy

by

Casey Meehan

Doctor of Philosophy in Computer Science

University of California San Diego, 2023

Professor Kamalika Chaudhuri, Chair

As data collection for machine learning (ML) tasks has become more pervasive, it has also become more heterogeneous: we share our writing, images, voices, and location online every day. Naturally, the associated privacy risks are just as complex and variable. My research advances practical data privacy through two avenues: 1) drafting provable privacy definitions and mechanisms for safely sharing data in different ML domains, and 2) empirically quantifying how ML models memorize their sensitive training data and thereby risk disclosing it. This dissertation details the various data domains/tasks considered, and the corresponding privacy methods proposed.

Chapter 1

When are Non-Parametric Methods Robust?

1.1 Introduction

Recent work has shown that many classifiers tend to be highly non-robust and that small strategic modifications to regular test inputs can cause them to misclassify [3, 4, 5]. Motivated by the use of machine learning in safety-critical applications, this phenomenon has recently received considerable interest; however, what exactly causes this phenomenon – known in the literature as *adversarial examples* – still remains a mystery.

Prior work has looked at three plausible reasons why adversarial examples might exist. The first, of course, is the possibility that in real data distributions, different classes are very close together in space – which does not seem plausible in practice. Another possibility is that classification algorithms may require more data to be robust than to be merely accurate; some prior work [6, 7, 8] suggests that this might be true for certain classifiers or algorithms. Finally, others [9, 10, 7] have suggested that better training algorithms may give rise to more robust classifiers – and that in some cases, finding robust classifiers may even be computationally challenging.

In this work, we consider this problem in the context of general non-parametric classifiers. Contrary to parametrics, non-parametric methods are a form of local classifiers, and include a large number of pattern recognition methods such as nearest neighbors, decision trees, random

forests and kernel classifiers. There is a richly developed statistical theory of non-parametric methods [11], which focuses on accuracy, and provides very general conditions under which these methods converge to the Bayes optimal with growing number of samples. We, in contrast, analyze robustness properties of these methods, and ask instead when they converge to the classifier with the highest astuteness at a desired radius r . Recall that the astuteness of a classifier at radius r is the fraction of points from the distribution on which it is accurate and has the same prediction up to a distance r [7, 6].

We begin by looking at the very simple case when data from different classes is well-separated – by at least a distance $2r$. Although achieving astuteness in this case may appear trivial, we show that even in this highly favorable case, not all non-parametric methods provide robust classifiers – and this even holds for methods that converge to the Bayes optimal in the large sample limit.

This raises the natural question – when do non-parametric methods produce astute classifiers? We next provide conditions under which a non-parametric method converges to the most astute classifier in the large sample limit under well-separated data. Our conditions are analogous to the classical conditions for convergence to the Bayes optimal [11, 12], but a little stronger. We show that nearest neighbors and kernel classifiers whose kernel functions decay fast enough, satisfy these conditions, and hence converge to astute classifiers in the large sample limit. In contrast, histogram classifiers, which do converge to the Bayes optimal in the large sample limit, may not converge to the most astute classifier. This indicates that there may be some non-parametric methods, such as nearest neighbors and kernel classifiers, that are more naturally robust when trained on well-separated data, and some that are not.

What happens when different classes in the data are not as well-separated? For this case, [13] proposes a method called Adversarial Pruning that preprocesses the training data by retaining the maximal set of points such that different classes are distance $\geq 2r$ apart, and then trains a non-parametric method on the pruned data. We next prove that if a non-parametric method has certain properties, then the classifier produced by Adversarial Pruning followed by

the method does converge to the most astute classifier in the large sample limit. We show that again nearest neighbors and kernel classifiers whose kernel functions decay faster than inverse polynomials satisfy these properties. Our results thus complement and build upon the empirical results of [13] by providing a performance guarantee.

What can we conclude about the cause for adversarial examples? Our results seem to indicate that at least for non-parametrics, it is mostly the training algorithms that are responsible. With a few exceptions, decades of prior work in machine learning and pattern recognition has largely focussed on designing training methods that provide increasingly accurate classifiers – perhaps to the detriment of other aspects such as robustness. In this context, our results serve to (a) provide a set of guidelines that can be used for designing non-parametric methods that are robust and accurate on well-separated data and (b) demonstrate that when data is not well-separated, preprocessing through adversarial pruning [13] may be used to ensure convergence to optimally astute solutions in the large sample limit.

1.1.1 Related Work

There is a large body of work on adversarial attacks [14, 15, 16, 17, 3] and defenses [18, 19, 2, 20, 21, 22] in the parametric setting, specifically focusing on neural networks. On the other hand, adversarial examples for nonparametric classifiers have mostly been studied in a much more ad-hoc manner, and to our knowledge, there has been no theoretical investigation into general properties of algorithms that promote robustness in non-parametric classifiers.

For nearest neighbors, there has been some prior work on adversarial attacks [23, 24, 7, 13] as well as defenses. Wang et. al. [7] proposes a defense for 1-NN by pruning the input sample. However, their defense learns a classifier whose robustness regions converge towards those of the Bayes optimal classifier, which itself may potentially have poor robustness properties. Yang et. al. [13] accounts for this problem by proposing the notion of the r -optimal classifier, and propose an algorithm called Adversarial Pruning which can be interpreted as a finite sample approximation to the r -optimal. However, they do not provide formal performance guarantees

for Adversarial Pruning, which we do.

For Kernel methods, Hein and Andriushchenko [18] study lower bounds on the norm of the adversarial manipulation that is required for changing a classifiers output. They specifically study bounds for Kernel Classifiers, and propose an empirically based regularization idea that improves robustness. In this work, we improve the robustness properties of kernel classification through adversarial pruning, and show formal guarantees regarding convergence towards the r -optimal classifier.

For decision trees and random forests, attacks and defenses have been provided by [25, 26, 27]. Again, most of the work here is empirical in nature, and convergence guarantees are not provided.

Pruning has a long history of being applied for improving nearest neighbors [28, 29, 30, 31, 32, 33], but this has been entirely done in the context of generalization, without accounting for robustness. In their work, Yang et. al. empirically show that adversarial pruning can improve robustness for nearest neighbor classifiers. However, they do not provide any formal guarantees for their algorithms. In this work, we prove formal guarantees for *adversarial pruning* in the large sample limit, both for nearest neighbors as well as for more general *weight functions*.

There is a long history of literature for understanding the consistency of Kernel classifiers [34, 12], but this has only been done for accuracy and generalization. In this work, we find different conditions are needed to ensure that a Kernel classifier converges in robustness in addition to accuracy.

1.2 Preliminaries

1.2.1 Setting

We consider binary classification where instances are drawn from a totally bounded metric space \mathcal{X} that is equipped with distance metric denoted by d , and the label space is $\{\pm 1\} = \{-1, +1\}$. The classical goal of classification is to build a highly *accurate* classifier,

which we define as follows.

Definition 1. (Accuracy) Let \mathcal{D} be a distribution over $\mathcal{X} \times \{\pm 1\}$, and let $f \in \{\pm 1\}^{\mathcal{X}}$ be a classifier. Then the **accuracy** of f over \mathcal{D} , denoted $A(f, \mathcal{D})$, is the fraction of examples $(x, y) \sim \mathcal{D}$ for which $f(x) = y$. Thus

$$A(f, \mathcal{D}) = P_{(x,y) \sim \mathcal{D}}[f(x) = y].$$

In this work, we consider *robustness* in addition to accuracy. Let $B(x, r)$ denoted the closed ball of radius r centered at x .

Definition 2. (Robustness) A classifier $f \in \{\pm 1\}^{\mathcal{X}}$ is said to be **robust** at x with radius r if $f(x) = f(x')$ for all $x' \in B(x, r)$.

Our goal is to find non-parametric algorithms that output classifiers that are robust, in addition to being accurate. To account for both criteria, we combine them into a notion of *astuteness* [7, 6].

Definition 3. (Astuteness) A classifier $f \in \{\pm 1\}^{\mathcal{X}}$ is said to be **astute** at (x, y) with radius r if f is robust at x with radius r and $f(x) = y$. The **astuteness** of f over \mathcal{D} , denoted $A_r(f, \mathcal{D})$, is the fraction of examples $(x, y) \sim \mathcal{D}$ for which f is astute at (x, y) with radius r . Thus

$$A_r(f, \mathcal{D}) = P_{(x,y) \sim \mathcal{D}}[f(x') = y, \forall x' \in B(x, r)].$$

It is worth noting that $A_0(f, \mathcal{D}) = A(f, \mathcal{D})$, since astuteness with radius 0 is simply the accuracy. For this reason, we will use $A_0(f, \mathcal{D})$ to denote accuracy from this point forwards.

1.2.2 Notions of Consistency

Traditionally, a classification algorithm is said to be consistent if as the sample size grows to infinity, the accuracy of the classifier it learns converges towards the best possible accuracy on the underlying data distribution. We next introduce and formalize an alternative form of consistency, called *r-consistency*, that applies to robust classifiers.

We begin with a formal definition of the Bayes Optimal Classifier – the most accurate classifier on a distribution – and consistency.

Definition 4. (*Bayes Optimal Classifier*) The **Bayes Optimal Classifier** on a distribution \mathcal{D} , denoted by g^* , is defined as follows. Let $\eta(x) = p_{\mathcal{D}}(+1|x)$. Then

$$g^*(x) = \begin{cases} +1 & \eta(x) \geq 0.5 \\ -1 & \eta(x) < 0.5 \end{cases}$$

It can be shown that g^* achieves the highest accuracy over \mathcal{D} over all classifiers.

Definition 5. (*Consistency*) Let M be a classification algorithm over $\mathcal{X} \times \{\pm 1\}$. M is said to be **consistent** if for any \mathcal{D} over $\mathcal{X} \times \{\pm 1\}$, and any ϵ, δ over $(0, 1)$, there exists N such that for $n \geq N$, with probability $1 - \delta$ over $S \sim \mathcal{D}^n$, we have:

$$A(M(S), \mathcal{D}) \geq A(g^*, \mathcal{D}) - \epsilon,$$

where g^* is the Bayes optimal classifier for \mathcal{D} .

How can we incorporate robustness in addition to accuracy in this notion? A plausible way, as used in [7], is that the classifier should converge towards being astute where the Bayes Optimal classifier is astute. However, the Bayes Optimal classifier is not necessarily the most astute classifier and may even have poor astuteness. To see this, consider the following example.

Example 1

Consider \mathcal{D} over $\mathcal{X} = [0, 1]$ such that $\mathcal{D}_{\mathcal{X}}$ is the uniform distribution and

$$p(y = 1|x) = \frac{1}{2} + \sin \frac{4\pi x}{r}.$$

For any point x , there exists $x_1, x_2 \in ([x - r, x + r] \cap [0, 1])$ such that $p(y = 1|x_1) > \frac{1}{2}$ and $p(y = 1|x_2) < \frac{1}{2}$. $A_r(g^*, r) = 0$. However, the classifier that always predicts $f(x) = +1$ does

better. It is robust everywhere, and since $P_{(x,y) \sim \mathcal{D}}[y = +1] = \frac{1}{2}$, it follows that $A_r(f, \mathcal{D}) = \frac{1}{2}$.

This motivates the notion of the r -optimal classifier, introduced by [13], which is the classifier with maximum astuteness.

Definition 6. (*r-optimal classifier*) The ***r-optimal classifier*** of a distribution G denoted by g_r^* is the classifier with maximum astuteness. Thus

$$g_r^* = \arg \max_{f \in \{\pm 1\}^{\mathcal{X}}} A_r(f, \mathcal{D}).$$

We let $A_r^*(\mathcal{D})$ denote $A_r(g_r^*, \mathcal{D})$.

Observe that g_r^* is not necessarily unique. To account for this, we use $A_r^*(\mathcal{D})$ in our definition for r -consistency.

Definition 7. (*r-consistent*) Let M be a classification algorithm over $\mathcal{X} \times \{\pm 1\}$. M is said to be ***r-consistent*** if for any \mathcal{D} , any $\varepsilon, \delta \in (0, 1)$, and $0 < \gamma < r$, there exists N such that for $n \geq N$, with probability $1 - \delta$ over $S \sim \mathcal{D}^n$,

$$A_{r-\gamma}(M(S), \mathcal{D}) \geq A_r^*(\mathcal{D}) - \varepsilon.$$

if the above conditions hold for a specific distribution \mathcal{D} , we say that M is *r-consistent with respect to \mathcal{D}* .

Observe that in addition to the usual ε and δ , there is an extra parameter γ which measures the gap in the robustness radius. We may need this parameter as when classes are exactly $2r$ apart, we may not be able to find the exact robust boundary with only finite samples.

Our analysis will be centered around understanding what kinds of algorithms M provide highly astute classifiers for a given radius r . We begin by first considering the special case of

r -separated distributions.

Definition 8. (r -separated distributions) A distribution \mathcal{D} is said to be **r -separated** if there exist subsets $T^+, T^- \subset \mathcal{X}$ such that

1. $\mathbb{P}_{(x,y) \sim \mathcal{D}}[x \in T^+] = 1$.
2. $\forall x_1 \in T^+, \forall x_2 \in T^-, d(x_1, x_2) > 2r$.

Observe that if \mathcal{D} is r -separated, $A_r(g_r^*, \mathcal{D}) = 1$.

1.2.3 Non-parametric Classifiers

Many non-parametric algorithms classify points by averaging labels over a local neighborhood from their training data. A very general form of this idea is encapsulated in *weight functions* – which is the general form we will use.

Definition 9. [11] A **weight function** W is a non-parametric classifier with the following properties.

1. Given input $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \sim \mathcal{D}^n$, W constructs functions $w_1^S, \dots, w_n^S : \mathcal{X} \rightarrow [0, 1]$ such that for all $x \in \mathcal{X}$, $\sum_1^n w_i^S(x) = 1$. The functions w_i^S are allowed to depend on x_1, x_2, \dots, x_n but must be independent of y_1, y_2, \dots, y_n .
2. W has output W_S defined as

$$W_S(x) = \begin{cases} +1 & \sum_1^n w_i^S(x) y_i > 0 \\ -1 & \sum_1^n w_i^S(x) y_i \leq 0 \end{cases}$$

As a result, $w_i^S(x)$ can be thought of as the weight that (x_i, y_i) has in classifying x .

Weight functions encompass a fairly extensive set of common non-parametric classifiers, which is the motivation for considering them. We now define several common non-parametric algorithms that can be construed as weight functions.

Definition 10. A *histogram classifier*, H , is a non-parametric classification algorithm over $\mathbb{R}^d \times \{\pm 1\}$ that works as follows. For a distribution \mathcal{D} over $\mathbb{R} \times \{\pm 1\}$, H takes $S = \{(x_i, y_i) : 1 \leq i \leq n\} \sim \mathcal{D}^n$ as input. Let k_i be a sequence with $\lim_{i \rightarrow \infty} k_i = \infty$ and $\lim_{i \rightarrow \infty} \frac{k_i}{i} = 0$. H constructs a set of hypercubes $C = \{c_1, c_2, \dots, c_m\}$ as follows:

1. Initially $C = \{c\}$, where $S \subset c$.
2. For $c \in C$, if c contains more than k_n points of S , then partition c into 2^d equally sized hypercubes, and insert them into C .
3. Repeat step 2 until all cubes in C have at most k_n points.

For $x \in \mathbb{R}$ let $c(x)$ denote the unique cell in C containing x . If $c(x)$ doesn't exist, then $H_S(x) = -1$ by default. Otherwise,

$$H_S(x) = \begin{cases} +1 & \sum_{x_i \in c(x)} y_i > 0 \\ -1 & \sum_{x_i \in c(x)} y_i \leq 0 \end{cases}.$$

Histogram classifiers are weight functions in which all x_i contained within the same cell as x are given the same weight $w_i^S(x)$ in predicting x , while all other x_i are given weight 0.

Definition 11. A *kernel classifier* is a weight function W over $\mathcal{X} \times \{\pm 1\}$ constructed from function $K : \mathbb{R}^+ \cup \{0\} \rightarrow \mathbb{R}^+$ and some sequence $\{h_n\} \subset \mathbb{R}^+$ in the following manner. Given $S = \{(x_i, y_i)\} \sim \mathcal{D}^n$, we have

$$w_i^S(x) = \frac{K\left(\frac{d(x, x_i)}{h_n}\right)}{\sum_{j=1}^n K\left(\frac{d(x, x_j)}{h_n}\right)}.$$

Then, as above, W has output

$$W_S(x) = \begin{cases} +1 & \sum_1^n w_i^S(x) y_i > 0 \\ -1 & \sum_1^n w_i^S(x) y_i \leq 0 \end{cases}$$

Finally, we note that k_n -nearest neighbors is also a weight function; $w_i^S(x) = \frac{1}{k_n}$ if x_i is one of the k_n closest neighbors of x and 0 otherwise.

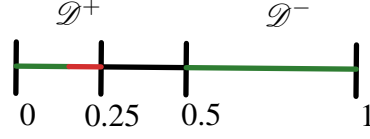


Figure 1.1. H_S is astute in the green region, but not robust in the red region.

1.3 Warm Up: r -separated distributions

We begin by considering the case when the data distribution is r -separated; the more general case is considered in Section 1.4. While classifying r -separated distributions robustly may appear almost trivial, learning an arbitrary classifier does not necessarily produce an astute result. To see this, consider the following example of a histogram classifier – which is known to be consistent.

We let H denote the histogram classifier over \mathbb{R} .

Example 2

Consider the data distribution $\mathcal{D} = \mathcal{D}^+ \cup \mathcal{D}^-$ where \mathcal{D}^+ is the uniform distribution over $[0, \frac{1}{4})$ and \mathcal{D}^- is the uniform distribution over $(\frac{1}{2}, 1]$, $p(+1|x) = 1$ for $x \in \mathcal{D}^+$, and $p(-1|x) = 1$ for $x \in \mathcal{D}^-$.

We make the following observations (refer to Figure 1.1).

1. \mathcal{D} is 0.1-separated, since the supports of \mathcal{D}^+ and \mathcal{D}^- have distance $0.25 > 0.2$.
2. If n is sufficiently large, H will construct the cell $[0.25, 0.5)$, which will not be split because it will never contain any points.
3. $H_S(x) = -1$ for $x \in [0.25, 0.5)$.
4. H_S is not astute at $(x, 1)$ for $x \in (0.15, 0.25)$. Thus $A_{0.1}(H_S, \mathcal{D}) = 0.8$.

Example 2 shows that histogram classifiers do not always learn astute classifiers even

when run on r -separated distributions. This motivates the question: which non-parametric classifiers do?

We answer this question in the following theorem, which gives sufficient conditions for a weight function (definition 9) to be r -consistent over an r -separated distribution.

Theorem 12. *Let \mathcal{D} be a distribution over $\mathcal{X} \times \{\pm 1\}$, and let W be a weight function. Let X be a random variable with distribution $\mathcal{D}_{\mathcal{X}}$, and $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \sim \mathcal{D}^n$. Suppose that for any $0 < a < b$,*

$$\lim_{n \rightarrow \infty} \mathbb{E}_{X, S} \left[\sup_{x' \in B(X, a)} \sum_{i=1}^n w_i^S(x') I_{\|x_i - x'\| > b} \right] = 0.$$

Then if \mathcal{D} is r -separated, W is r -consistent with respect to \mathcal{D} .

First, we compare Theorem 12 to Stone's theorem [12], which gives sufficient conditions for a weight function to be consistent (i.e. converge in accuracy towards the Bayes optimal). For convenience, we include a statement of Stone's theorem.

Theorem 13. [12] *Let W be weight function over $\mathcal{X} \times \{\pm 1\}$. Suppose the following conditions hold for any distribution \mathcal{D} over $\mathcal{X} \times \{\pm 1\}$. Let X be a random variable with distribution $\mathcal{D}_{\mathcal{X}}$, and $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \sim \mathcal{D}^n$. All expectations are taken over X and S .*

1. *There is a constant c such that, for every nonnegative measurable function f satisfying*

$$\mathbb{E}[f(X)] < \infty,$$

$$\mathbb{E} \left[\sum_{i=1}^n w_i^S(X) f(x_i) \right] \leq c \mathbb{E}[f(x)].$$

2. *For all $a > 0$,*

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\sum_{i=1}^n w_i^S(x) I_{\|x_i - X\| > a} \right] = 0,$$

where $I_{\|x_i - X\| > a}$ is an indicator variable.

- 3.

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\max_{1 \leq i \leq n} w_i^S(X) \right] = 0.$$

Then W is consistent.

There are two main differences between Theorem 12 and Stone's theorem.

1. Conditions 1. and 3. of Stone's theorem are no longer necessary. This is because r -separated distributions are well-separated and thus have simpler conditions for consistency. In fact, a slight modification of the arguments of [12] shows that for r -separated distributions, condition 2. alone is sufficient for consistency.
2. Condition 2. is strengthened. Instead of requiring the weight of x_i 's outside of a given radius to go to 0 for $X \sim \mathcal{D}$, we require the same to *uniformly* hold over a ball centered at X .

Theorem 12 provides a general condition that allows us to verify the r -consistency of non-parametric methods. We now show below that two common non-parametric algorithms – k_n -nearest neighbors and kernel classifiers with rapidly decaying kernel functions – satisfy the conditions of Theorem 12.

Corollary 14. *Let \mathcal{D} be any r -separated distribution. Let k_n be any sequence such that $\lim_{n \rightarrow \infty} \frac{k_n}{n} = 0$, and let M be the k_n -nearest neighbors classifier on a sample $S \sim \mathcal{D}^n$. Then M is r -consistent with respect to \mathcal{D} .*

Remarks:

1. Because the data distribution is r -separated, $k_n = 1$ will be r -consistent. Also observe that for r -separated distributions, $k_n = 1$ will converge towards the Bayes Optimal classifier.
2. In general, M converges towards the Bayes Optimal classifier provided that $k_n \rightarrow \infty$ in addition to $k_n/n \rightarrow 0$. This condition is not necessary for r -consistency– because the distribution is r -separated.

We next show that kernel classifiers are also r -consistent on r -separated data distributions, provided the kernel function decreases rapidly enough.

Corollary 15. *Let W be a kernel classifier over $\mathcal{X} \times \{\pm 1\}$ constructed from K and h_n . Suppose the following properties hold for K and h_n .*

1. *For any $c > 1$, $\lim_{x \rightarrow \infty} \frac{K(cx)}{K(x)} = 0$.*
2. *$\lim_{n \rightarrow \infty} h_n = 0$.*

If \mathcal{D} is an r -separated distribution over $\mathcal{X} \times \{\pm 1\}$, then W is r -consistent with respect to \mathcal{D} .

Observe that Condition 1. is satisfied for any $K(x)$ that decreases more rapidly than an inverse polynomial – and is hence satisfied by most popular kernels like the Gaussian kernel. Is the condition on K in Corollary 15 necessary? The following example illustrates that a kernel classifier with any arbitrary K is not necessarily r -consistent. This indicates that some sort of condition needs to be imposed on K to ensure r -consistency; finding a tight necessary condition however is left for future work.

Example 3

Let $\mathcal{X} = [-1, 1]$ and let \mathcal{D} be a distribution with $p_{\mathcal{D}}(-1, -1) = 0.1$ and $p_{\mathcal{D}}(1, 1) = 0.9$. Clearly, \mathcal{D} is 0.3-separated. Let $K(x) = e^{-\min(|x|, 0.2)^2}$. Let h_n be any sequence with $\lim_{n \rightarrow \infty} h_n = 0$ and $\lim_{n \rightarrow \infty} nh_n = \infty$. Let W be the weight classifier with input $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ such that

$$w_i^S(x) = \frac{K\left(\frac{|x - x_i|}{h_n}\right)}{\sum_{j=1}^n K\left(\frac{|x - x_j|}{h_n}\right)}.$$

W can be shown to satisfy all the conditions of Theorem 13 (the proof is analogous to the case for a Gaussian Classifier), and is therefore consistent. However, W does not learn a robust classifier on \mathcal{D} for $r = 0.3$.

Consider $x = -0.7$. For any $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \sim \mathcal{D}^n$, all x_i will either be -1 or 1 . Therefore, since $K(|x - (-1)|) = K(|x - 1|)$, it follows that $w_i^S(x) = \frac{1}{n}$ for all $1 \leq i \leq n$. Since $x_i = 1$ with probability 0.9, it follows that with high probability x will be classified as 1 which means that f , the output of W , is not robust at $x = -1$. Thus f has astuteness at most 0.9 which means that W is *not* r -consistent for $r = 0.3$.

1.4 General Distributions

We next consider more general data distributions, where data from different classes may be close together in space, and may even overlap. Observe that unlike the r -separated case, here there may be no classifier with astuteness one. Thus, a natural question is: what does the optimally astute classifier look like, and how can we build non-parametric classifiers to this limit?

1.4.1 The r -Optimal Classifier and Adversarial Pruning

[13] propose a large-sample limit – called the r -optimal – and show that it is analogous to the Bayes Optimal classifier for robustness. More specifically, given a data distribution D , to find the r -optimal classifier, we solve the following optimization problem.

$$\begin{aligned} \max_{S_{+1}, S_{-1}} & \int_{x \in S_{+1}} p(y = +1|x) d\mu_{\mathcal{D}}(x) + \\ & \int_{x \in S_{-1}} p(y = -1|x) d\mu_{\mathcal{D}}(x) \\ \text{subject to} & d(S_{+1}, S_{-1}) > 2r \end{aligned} \tag{1.1}$$

Then, the r -optimal classifier is defined as follows.

Definition 16. [13] Fix r, \mathcal{D} . Let S_{+1}^* and S_{-1}^* be any optimizers of (1.1). Then the r -optimal classifier, g_r^* is any classifier such that $g_r^*(x) = j$ whenever $d(S_j^*, x) \leq r$.

[13] show that the r -optimal classifier achieves the optimal astuteness – out of all classifiers on the data distribution \mathcal{D} ; hence, it is a robustness analogue to the Bayes Optimal Classifier. Therefore, for general distributions, the goal in robust classification is to find non-parametric algorithms that output classifiers that converge towards g_r^* .

To find robust classifiers, [13] propose Adversarial Pruning – a defense method that preprocesses the training data by making it better separated. More specifically, Adversarial

Pruning takes as input a training dataset S and a radius r , and finds the largest subset of the training set where differently labeled points are at least distance $2r$ apart.

Definition 17. A set $S_r \subset \mathcal{X} \times \{\pm 1\}$ is said to be *r -separated* if for all $(x_1, y_1), (x_2, y_2) \in S_r$, if $y_1 \neq y_2$, then $d(x_1, x_2) > 2r$. To *adversarially prune* a set S is to return its largest r -separated subset. We let $\text{AdvPrun}(S, r)$ denote the result of adversarially pruning S .

Once an r -separated subset S_r of the training set is found, a standard non-parametric method is trained on S_r . While [13] show good empirical performance of such algorithms, no formal guarantees are provided. We next formally characterize when adversarial pruning followed by a non-parametric method results in a classifier that is provably r -consistent.

Specifically, we consider analyzing the general algorithm provided in Algorithm 1.

Algorithm 1: RobustNonPar

- 1 **Input:** $S \sim \mathcal{D}^n$, weight function W , robustness radius r ;
 - 2 $S_r \leftarrow \text{AdvPrun}(S, r)$;
 - 3 **Output:** W_{S_r} ;
-

1.4.2 Convergence Guarantees

We begin with some notation. For any weight function W and radius $r > 0$, we let $\text{RobustNonPar}(W, r)$ represent the weight function that outputs weights for $S \sim \mathcal{D}^n$ according to $\text{RobustNonPar}(S, W, r)$. In particular, this can be used to convert any weight function algorithm into a new weight function which takes robustness into account. A natural question is, for which weight functions W is $\text{RobustNonPar}(W, r)$ r -consistent? Our next theorem provides sufficient conditions for this.

Theorem 18. Let W be a weight function over $\mathcal{X} \times \{\pm 1\}$, and let \mathcal{D} be a distribution over $\mathcal{X} \times \{\pm 1\}$. Fix $r > 0$. Let $S_r = \text{AdvPrun}(S, r)$. For convenience, relabel x_i, y_i so that $S_r =$

$\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$. Suppose that for any $0 < a < b$,

$$\lim_{n \rightarrow \infty} \mathbb{E}_{S \sim \mathcal{D}^n} \left[\frac{1}{m} \sum_{i=1}^m \sup_{x \in B(x_i, a)} \sum_{j=1}^m w_j^{S_r}(x) I_{||x_j - x|| > b} \right] = 0.$$

Then $\text{RobustNonPar}(W, r)$ is r -consistent with respect to \mathcal{D} .

Remark:

There are two important differences between the conditions in Theorem 18 and Theorem 12.

1. We replace S with S_r .
2. The expectation over $X \sim \mathcal{D}_{\mathcal{X}}$ is replaced with an average over $\{x_1, x_2, \dots, x_m\}$. The intuition here is that we are replacing \mathcal{D} with a uniform distribution over S_r . While \mathcal{D} may not be r -separated, the uniform distribution over S_r is, and represents the region of points where our classifier is astute.

A natural question is what satisfies the conditions in Theorem 18. We next show that k_n -nearest neighbors and kernel classifiers with rapidly decaying kernel functions continue to satisfy the conditions in Theorem 18; this means that these classifiers, when combined with Adversarial Pruning, will converge to r -optimal classifiers in the large sample limit.

Corollary 19. *Let k_n be a sequence with $\lim_{n \rightarrow \infty} \frac{k_n}{n} = 0$, and let M denote the k_n -nearest neighbor algorithm. Then for any $r > 0$, $\text{RobustNonPar}(M, r)$ is r -consistent.*

Remark:

Corollary 19 gives a formal guarantee in the large sample limit for the modified nearest-neighbor algorithm proposed by [13].

Corollary 20. *Let W be a kernel classifier over $\mathcal{X} \times \{\pm 1\}$ constructed from K and h_n . Suppose the following properties hold for K and h_n .*

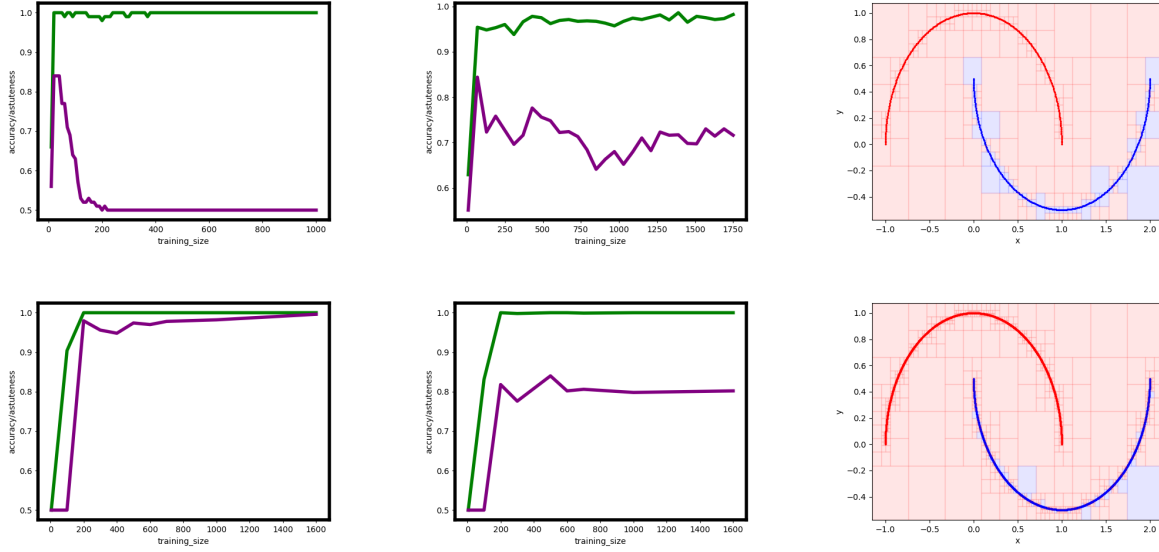


Figure 1.2. Empirical accuracy/astuteness of different classifiers as a function of training sample size. Accuracy is shown in green, astuteness in purple. Left : Noiseless Setting. Right: Noisy Setting. Top Row: Histogram Classifier, Bottom Row: 1-Nearest Neighbor

1. For any $c > 1$, $\lim_{x \rightarrow \infty} \frac{K(cx)}{K(x)} = 0$.

2. $\lim_{n \rightarrow \infty} h_n = 0$.

Then for any $r > 0$, $\text{RobustNonPar}(W, r)$ is r -consistent.

Observe again that Condition 1. is satisfied by any K that decreases more rapidly than an inverse polynomial kernel; it is thus satisfied by most popular kernels, such as the Gaussian kernel.

1.5 Validation

Our theoretical results are, by nature, large sample; we next validate how well they apply to the finite sample case by trying them out on a simple example. In particular, we ask the following question:

How does the robustness of non-parametric classifiers change with increasing sample size?

This question is considered in the context of two simple non-parametric classifiers – one nearest neighbor (which is guaranteed to be r -consistent) and histograms (which is not). To be able to measure performance with increasing data size, we look at a simple synthetic dataset – the Half Moons.

1.5.1 Experimental Setup

Classifiers and Dataset.

We consider two different classification algorithms – one nearest neighbor (NN) and a Histogram Classifier (HC). We use the Halfmoon dataset with two settings of the gaussian noise parameter σ , $\sigma = 0$ (Noiseless) and $\sigma = 0.08$ (Noisy). For the Noiseless setting, observe that the data is already 0.1-separated; for the Noisy setting, we use Adversarial Pruning (Algorithm 1) with parameter $r = 0.1$ for both classification methods.

Performance Measure.

We evaluate robustness with respect to the ℓ_∞ metric, that is commonly used in the adversarial examples literature. Specifically, for each classifier, we calculate the *empirical astuteness*, which is the fraction of test examples on which it is astute.

Observe that computing the empirical astuteness of a classifier around an input x amounts to finding the adversarial example that is *closest to x* according to the ℓ_∞ norm. For the 1-nearest neighbor, we do this using the optimal attack algorithm proposed by Yang et. al. [13]. For the histogram classifier, we use the optimal attack framework proposed by [13], and show that the structure of the classifier can be exploited to solve the convex program efficiently. Details are in Appendix C.

We use an attack radius of $r = 0.1$ for the Noiseless setting, and $r = 0.09$ for the Noisy setting. For all classification algorithms, we plot the empirical astuteness as a function of the training set size. As a baseline, we also plot their standard accuracy on the test set.

1.5.2 Results

The results are presented in Figure 1.2; the left two panels are for the Noiseless setting while the two center ones are for the Noisy setting.

The results show that as predicted by our theory, for the Noiseless setting, the empirical astuteness of nearest neighbors converges to 1 as the training set grows. For Histogram Classifiers, the astuteness converges to 0.5 – indicating that the classifier may grow less and less astute with higher sample size even for well-separated data. This is plausibly because the cell size induced by the histogram grows smaller with growing training data; thus, the classifier that outputs the default label -1 in empty cells is incorrect on adversarial examples that are close to a point with $+1$ label, but belongs to a different, empty cell. The rightmost panels in Figure 1.2 provide a visual illustration of this process.

For the Noisy setting, the empirical astuteness of adversarial pruning followed by nearest neighbors converges to 0.8. For histograms with adversarial pruning, the astuteness converges to 0.7, which is higher than the noiseless case but still clearly sub-optimal.

1.5.3 Discussion

Our results show that even though our theory is asymptotic, our predictions continue to be relevant in finite sample regimes. In particular, on well-separated data, nearest neighbors that we theoretically predict to be intrinsically robust is robust; histogram classifiers, which do not satisfy the conditions in Theorem 12 are not. Our predictions continue to hold for data that is not well-separated. Nearest neighbors coupled with Adversarial Pruning continues to be robust with growing sample size, while histograms continue to be non-robust. Thus our theory is confirmed by practice.

1.6 Conclusion

In conclusion, we rigorously analyze when non-parametric methods provide classifiers that are robust in the large sample limit. We provide a general condition that characterizes when non-parametric methods are robust on well-separated data, and show that Adversarial Pruning of [13] works on data that is not well-separated.

Our results serve to provide a set of guidelines that can be used for designing non-parametric methods that are robust and accurate on well-separated data; additionally, we demonstrate that when data is not well-separated, preprocessing by adversarial pruning [13] does lead to optimally astute solutions in the large sample limit.

Chapter 2

Robust Empirical Risk Minimization with Tolerance

Bibliography

- [1] Chen Dan, Yuting Wei, and Pradeep Ravikumar. Sharp statistical guarantees for adversarially robust gaussian classification. *CoRR*, abs/2006.16384, 2020. URL <https://arxiv.org/abs/2006.16384>.
- [2] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
- [3] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- [4] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples, March 20 2014. URL <http://arxiv.org/abs/1412.6572>.
- [5] Daniel Lowd and Christopher Meek. Adversarial learning. In Robert Grossman, Roberto J. Bayardo, and Kristin P. Bennett, editors, *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Chicago, Illinois, USA, August 21-24, 2005*, pages 641–647. ACM, 2005. ISBN 1-59593-135-X.
- [6] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially robust generalization requires more data. In *Advances in Neural Information Processing Systems 31*, pages 5014–5026. Curran Associates, Inc., 2018.
- [7] Yizhen Wang, Somesh Jha, and Kamalika Chaudhuri. Analyzing the robustness of nearest neighbors to adversarial examples. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, pages 5120–5129, 2018.
- [8] Omar Montasser, Steve Hanneke, and Nathan Srebro. VC classes are adversarially robustly learnable, but only improperly. In *Conference on Learning Theory, COLT 2019, 25-28 June 2019, Phoenix, AZ, USA*, pages 2512–2530, 2019.

- [9] Hadi Salman, Greg Yang, Jerry Li, Pengchuan Zhang, Huan Zhang, Ilya P. Razenshteyn, and Sébastien Bubeck. Provably robust deep learning via adversarially trained smoothed classifiers. *CoRR*, abs/1906.04584, 2019. URL <http://arxiv.org/abs/1906.04584>.
- [10] Dmitrii Avdiukhin, Slobodan Mitrovic, Grigory Yaroslavtsev, and Samson Zhou. Adversarially robust submodular maximization under knapsack constraints. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019.*, pages 148–156, 2019. doi: 10.1145/3292500.3330911.
- [11] Luc Devroye, László Györfi, and Gábor Lugosi. *A Probabilistic Theory of Pattern Recognition*, volume 31 of *Stochastic Modelling and Applied Probability*. Springer, 1996.
- [12] Charles Stone. Consistent nonparametric regression. *The Annals of Statistics*, 5(4):595–645, 1977.
- [13] Yao-Yuan Yang, Cyrus Rashtchian, Yizhen Wang, and Kamalika Chaudhuri. Adversarial examples for non-parametric methods: Attacks, defenses and large sample limits. *CoRR*, abs/1906.03310, 2019. URL <http://arxiv.org/abs/1906.03310>.
- [14] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*, pages 39–57, 2017.
- [15] Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- [16] Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. Practical black-box attacks against deep learning systems using adversarial examples. *ASIACCS*, 2017.
- [17] Nicolas Papernot, Patrick D. McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In *IEEE European Symposium on Security and Privacy, EuroS&P 2016, Saarbrücken, Germany, March 21-24, 2016*, pages 372–387, 2016.
- [18] Matthias Hein and Maksym Andriushchenko. Formal guarantees on the robustness of a classifier against adversarial manipulation. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 2266–2276. Curran Associates, Inc., 2017.
- [19] Guy Katz, Clark W. Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer.

- Towards proving the adversarial robustness of deep neural networks. In *Proceedings First Workshop on Formal Verification of Autonomous Vehicles, FVAV@iFM 2017, Turin, Italy, 19th September 2017.*, pages 19–26, 2017.
- [20] Nicolas Papernot, Patrick D. McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *IEEE Symposium on Security and Privacy, SP 2016, San Jose, CA, USA, May 22-26, 2016*, pages 582–597, 2016.
 - [21] Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified defenses against adversarial examples. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
 - [22] Aman Sinha, Hongseok Namkoong, and John C. Duchi. Certifying some distributional robustness with principled adversarial training. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
 - [23] Laurent Amsaleg, James Bailey, Dominique Barbe, Sarah M. Erfani, Michael E. Houle, Vinh Nguyen, and Milos Radovanovic. The vulnerability of learning to adversarial perturbation increases with intrinsic dimensionality. In *2017 IEEE Workshop on Information Forensics and Security, WIFS 2017, Rennes, France, December 4-7, 2017*, pages 1–6, 2017.
 - [24] Chawin Sitawarin and David A. Wagner. On the robustness of deep k-nearest neighbors. In *2019 IEEE Security and Privacy Workshops, SP Workshops 2019, San Francisco, CA, USA, May 19-23, 2019*, pages 1–7, 2019.
 - [25] Maksym Andriushchenko and Matthias Hein. Provably robust boosted decision stumps and trees against adversarial attacks. In H. Wallach, H. Larochelle, A. Beygelzimer, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 12997–13008. Curran Associates, Inc., 2019.
 - [26] Alex Kantchelian, J. D. Tygar, and Anthony D. Joseph. Evasion and hardening of tree ensemble classifiers. *CoRR*, abs/1509.07892, 2015. URL <http://arxiv.org/abs/1509.07892>.
 - [27] Hongge Chen, Huan Zhang, Duane S. Boning, and Cho-Jui Hsieh. Robust decision trees against adversarial examples. *CoRR*, abs/1902.10660, 2019.
 - [28] Geoffrey W. Gates. The reduced nearest neighbor rule (corresp.). *IEEE Trans. Information Theory*, 18(3):431–433, 1972.
 - [29] Lee-Ad Gottlieb, Aryeh Kontorovich, and Pinhas Nisnevitch. Near-optimal sample compression for nearest neighbors. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014*,

Montreal, Quebec, Canada, pages 370–378, 2014.

- [30] Peter E. Hart. The condensed nearest neighbor rule (corresp.). *IEEE Trans. Information Theory*, 14(3):515–516, 1968.
- [31] Aryeh Kontorovich, Sivan Sabato, and Roi Weiss. Nearest-neighbor sample compression: Efficiency, consistency, infinite dimensions. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 1573–1583, 2017.
- [32] Aryeh Kontorovich and Roi Weiss. A bayes consistent 1-nn classifier. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2015, San Diego, California, USA, May 9-12, 2015*, 2015.
- [33] Steve Hanneke, Aryeh Kontorovich, Sivan Sabato, and Roi Weiss. Universal bayes consistency in metric spaces. *CoRR*, abs/1906.09855, 2019.
- [34] Ingo Steinwart. Consistency of support vector machines and other regularized kernel classifiers. *IEEE Trans. Information Theory*, 51(1):128–142, 2005.