

# Adversarially Robust Learning with Tolerance

**Hassan Ashtiani\***  
McMaster University

ZOKAEIAM@MCMASTER.CA

**Vinayak Pathak**  
Layer 6 AI

VINAYAK@LAYER6.AI

**Ruth Urner†**  
York University

RUTH@ECS.YORKU.CA

**Editors:** Shipra Agrawal and Francesco Orabona

## Abstract

We initiate the study of tolerant adversarial PAC-learning with respect to metric perturbation sets. In adversarial PAC-learning, an adversary is allowed to replace a test point  $x$  with an arbitrary point in a closed ball of radius  $r$  centered at  $x$ . In the tolerant version, the error of the learner is compared with the best achievable error with respect to a slightly larger perturbation radius  $(1 + \gamma)r$ . This simple tweak helps us bridge the gap between theory and practice and obtain the first PAC-type guarantees for algorithmic techniques that are popular in practice.

Our first result concerns the widely-used “perturb-and-smooth” approach for adversarial learning. For perturbation sets with doubling dimension  $d$ , we show that a variant of these approaches PAC-learns any hypothesis class  $\mathcal{H}$  with VC-dimension  $v$  in the  $\gamma$ -tolerant adversarial setting with  $O\left(\frac{v(1+\gamma)^{O(d)}}{\epsilon}\right)$  samples. This is in contrast to the traditional (non-tolerant) setting in which, as we show, the perturb-and-smooth approach can provably fail.

Our second result shows that one can PAC-learn the same class using  $\tilde{O}\left(\frac{O(d)\text{VC}(\mathcal{H})\log(1+\gamma)}{\epsilon^2}\right)$  samples even in the agnostic setting. This result is based on a novel compression-based algorithm, and achieves a linear dependence on the doubling dimension as well as the VC-dimension. This is in contrast to the non-tolerant setting where there is no known sample complexity upper bound that depends polynomially on the VC-dimension.

**Keywords:** Adversarial Perturbations, PAC-learning, Tolerant Adversarial Robustness, Perturb-and-S Sample Compression

## 1. Introduction

Several empirical studies (Szegeedy et al., 2014; Goodfellow et al., 2018) have demonstrated that models trained to have a low accuracy on a data set often have the undesirable property that a small perturbation to an input instance can change the label outputted by the model. For most domains this does not align with human perception and thus indicates that the learned models are not representing the ground truth despite obtaining good accuracy on test sets.

The theory of PAC-learning characterizes the conditions under which learning is possible. For binary classification, the following conditions are sufficient: a) unseen data should arrive from the same distribution as training data, and b) the class of models should have a low capacity (as measured,

\* Hassan Ashtiani is also a faculty affiliate at Toronto’s Vector Institute and was supported by an NSERC Discovery Grant.

† Ruth Urner is also a faculty affiliate at Toronto’s Vector Institute and was supported by an NSERC Discovery Grant.

for example, by its VC-dimension). If these conditions are met, an *Empirical Risk Minimizer* (ERM) that simply optimizes model parameters to maximize accuracy on the training set learns successfully.

Recent work has studied test-time adversarial perturbations under the PAC-learning framework. If an adversary is allowed to perturb data during test time then the conditions above do not hold, and we cannot hope for the model to learn to be robust just by running ERM. Thus, the goal here is to bias the learning process towards finding models where label-changing perturbations are rare. This is achieved by defining a loss function that combines both classification error and the probability of seeing label-changing perturbations, and learning models that minimize this loss on unseen data. It has been shown that even though (robust) ERM can fail in this setting, PAC-learning is still possible as long as we know during training the kinds of perturbations we want to guard against at test time (Montasser et al., 2019). This result holds for all perturbation sets. However, the learning algorithm is significantly more complex than robust ERM and requires a large number of samples (with the best known sample complexity bounds potentially being exponential in the VC-dimension).

We study a *tolerant* version of the adversarially robust learning framework and restrict the perturbations to balls in a general metric space with finite doubling dimension. We show this slight shift in the learning objective yields significantly improved sample complexity bounds through a simpler learning paradigm than what was previously known. In fact, we show that a version of the common “perturb-and-smooth” paradigm successfully PAC-learns any class of bounded VC-dimension in this setting.

**Learning in general metric spaces.** What kinds of perturbations should a learning algorithm guard against? Any transformation of the input that we believe should not change its label could be a viable perturbation for the adversary to use. The early works in this area considered perturbations contained within a small  $\ell_p$ -ball of the input. More recent work has considered other transformations such as a small rotation, or translation of an input image (Engstrom et al., 2019; Fawzi and Frossard, 2015; Kanbak et al., 2018; Xiao et al., 2018), or even adding small amounts of fog or snow (Kang et al., 2019). It has also been argued that small perturbations in some *feature space* should be allowed as opposed to the input space (Inkawhich et al., 2019; Sabour et al., 2016; Xu et al., 2020; Song et al., 2018; Hosseini and Poovendran, 2018). This motivates the study of more general perturbations.

We consider a setting where the input comes from a domain that is equipped with a distance metric and allows perturbations to be within a small metric ball around the input. Earlier work on general perturbation sets (for example, (Montasser et al., 2019)) considered arbitrary perturbations. In this setting one does not quantify the magnitude of a perturbation and thus cannot talk about small versus large perturbations. Modeling perturbations using a metric space enables us to do that while also keeping the setup general enough to be able to encode a large variety of perturbation sets by choosing appropriate distance functions.

**Learning with tolerance.** In practice, we often believe that small perturbations of the input should not change its label but we do not know *precisely* what small means. However, in the PAC-learning framework for adversarially robust classification, we are required to define a precise perturbation set and learn a model that has error arbitrarily close to the smallest error that can be achieved with respect to that perturbation set. In other words, we aim to be arbitrarily close to a target that was picked somewhat arbitrarily to begin with. Due to the uncertainty about the correct perturbation size, it is more meaningful to allow for a wider range of error values. To achieve this, we introduce the concept of tolerance. In the tolerant setting, in addition to specifying a perturbation size  $r$ , we introduce a tolerance parameter  $\gamma$  that encodes our uncertainty about the size of allowed perturbations. Then, for any given  $\epsilon > 0$ , we aim to learn a model whose error with

respect to perturbations of size  $r$  is at most  $\epsilon$  more than the smallest error achievable with respect to perturbations of size  $r(1 + \gamma)$ .

## 2. Our results

In this paper we formalize and initiate the study of the problem of adversarially robust learning in the tolerant setting for general metric spaces and provide two algorithms for the task. Both of our algorithms rely on: 1) modifying the training data by randomly sampling points from the perturbation sets around each data point, and 2) smoothing the output of the model by taking a majority over the labels returned by the model for nearby points.

Our first algorithm starts by modifying the training set by randomly perturbing each training point using a certain distribution (see Section 5 for details). It then trains a (non-robust) PAC-learner (such as ERM) on the perturbed training set to find a hypothesis  $h$ . Finally, it outputs a smooth version of  $h$ . The smoothing step replaces  $h(x)$  at each point  $x$  with the a majority label outputted by  $h$  on the points around  $x$ . We show that for metric spaces of a fixed doubling dimension, this algorithm successfully learns in the tolerant setting assuming tolerant realizability.

**Theorem 1 (Informal version of Theorem 10)** *Let  $(X, \text{dist})$  be a metric space with doubling dimension  $d$  and  $\mathcal{H}$  a hypothesis class. Assuming tolerant realizability,  $\mathcal{H}$  can be learned tolerantly in the adversarially robust setting using  $O\left(\frac{(1+1/\gamma)^{O(d)} \text{VC}(\mathcal{H})}{\epsilon}\right)$  samples, where  $\gamma$  encodes the amount of allowed tolerance, and  $\epsilon$  is the desired accuracy.*

An interesting feature of the above result is the linear dependence of the sample complexity with respect to  $\text{VC}(\mathcal{H})$ . This is in contrast to the best known upper bound for non-tolerant adversarial setting (Montasser et al., 2019) which depends on the *dual VC-dimension* of the hypothesis class and in general is exponential in  $\text{VC}(\mathcal{H})$ . Moreover, this is the first PAC-type guarantee for the general perturb-and-smooth paradigm, indicating that the tolerant adversarial learning is the “right” learning model for studying these approaches. While the above method enjoys simplicity and can be computationally efficient, one downside is that its sample complexity grows exponentially with the doubling dimension. For instance, such algorithm cannot be used on high-dimensional data in the Euclidean space. Another limitation is that the guarantee holds only in the (robustly) realizable setting.

In the second main part of our submission (Section 6) we show that, surprisingly, these limitations can be overcome by incorporating ideas from the tolerant framework and perturb-and-smooth algorithms into a novel compression scheme for robust learning. The resulting algorithm improves the dependence on the doubling dimension, and works in the general agnostic setting.

**Theorem 2 (Informal version of Corollary 16)** *Let  $(X, \text{dist})$  be a metric space with doubling dimension  $d$  and  $\mathcal{H}$  a hypothesis class. Then  $\mathcal{H}$  can be learned tolerantly in the adversarially robust setting using  $\tilde{O}\left(\frac{O(d)\text{VC}(\mathcal{H})\log(1+1/\gamma)}{\epsilon^2}\right)$  samples, where  $\tilde{O}$  hides logarithmic factors,  $\gamma$  encodes the amount of allowed tolerance, and  $\epsilon$  is the desired accuracy.*

This algorithm exploits the connection between sample compression and adversarially robust learning Montasser et al. (2019). However, unlike Montasser et al. (2019), our new compression scheme sidesteps the dependence on the dual VC-dimension (refer to the discussion at the end of Section 6 for more details). As a result, we get an exponential improvement over the best known (nontolerant) sample complexity in terms of dependence on VC-dimension.

### 3. Related work

PAC-learning for adversarially robust classification has been studied extensively in recent years (Cullina et al., 2018; Awasthi et al., 2019; Montasser et al., 2019; Feige et al., 2015; Attias et al., 2019; Montasser et al., 2020a; Ashtiani et al., 2020). These works provide learning algorithms that guarantee low generalization error in the presence of adversarial perturbations in various settings. The most general result is due to Montasser et al. (2019), and is proved for general hypothesis classes and perturbation sets. All of the above results assume that the learner knows the kinds of perturbations allowed for the adversary. Some more recent papers have considered scenarios where the learner does not even need to know that. Goldwasser et al. (2020) allow the adversary to perturb test data in unrestricted ways and are still able to provide learning guarantees. The catch is that it only works in the transductive setting and only if the learner is allowed to abstain from making a prediction on some test points. Montasser et al. (2021a) consider the case where the learner needs to infer the set of allowed perturbations by observing the actions of the adversary.

Tolerance was introduced by Ashtiani et al. (2020) in the context of certification. They provide examples where certification is not possible unless we allow some tolerance. Montasser et al. (2021b) study transductive adversarial learning and provide a “tolerant” guarantee. Note that unlike our work, the main focus of that paper is on the transductive setting. Moreover, they do not specifically study tolerance with respect to metric perturbation sets. Without a metric, it is not meaningful to expand perturbation sets by a factor  $(1 + \gamma)$  (as we do in our definition of tolerance). Instead, they expand their perturbation sets by applying two perturbations in succession, which is akin to setting  $\gamma = 1$ . In contrast, our results hold in the more common inductive setting, and capture a more realistic setting where  $\gamma$  can be any (small) real number larger than zero.

Subsequent to our work, Bhattacharjee et al. (2022) study adversarially robust learning with tolerance for “regular” VC-classes and show that a simple modification of robust ERM achieves a sample complexity polynomial in both VC-dimension and doubling dimension. In a similar vein, Raman et al. (2022) identify a more general property of hypothesis classes for which robust ERM is sufficient for adversarially robust learning with tolerance.

Like many recent adversarially robust learning algorithms (Feige et al., 2015; Attias et al., 2019), our first algorithm relies on calls to a non-robust PAC-learner. Montasser et al. (2020b) formalize the question of reducing adversarially robust learning to non-robust learning and study finite perturbation sets of size  $k$ . They show a reduction that makes  $O(\log^2 k)$  calls to the non-robust learner and also prove a lower bound of  $\Omega(\log k)$ . It will be interesting to see if our algorithms can be used to obtain better bounds for the tolerant setting. Our first algorithm makes one call to the non-robust PAC-learner at training time, but needs to perform potentially expensive smoothing for making actual predictions (see Theorem 10).

A related line of work studies smallest achievable robust loss for various distributions and hypothesis classes. For example, Bubeck and Sellke (2021) show that hypothesis classes with low robust loss must be overparametrized. Yang et al. (2020b) explore real-world datasets and provide evidence that they are separable and therefore there must exist locally Lipschitz hypotheses with low robust loss. Note that the existence of such hypotheses does not immediately imply that PAC-learning is possible.

The techniques of randomly perturbing the training data and smoothing the output classifier has been extensively used in practice and has shown good empirical success. Augmenting the training data with some randomly perturbed samples was used for handwriting recognition as early as by

Yaeger et al. (1996). More recently, “stability training” was introduced by Zheng et al. (2016) for state of the art image classifiers where training data is augmented with Gaussian perturbations. Empirical evidence was provided that the technique improved the accuracy against naturally occurring perturbations. Augmentations with non-Gaussian perturbations of a large variety were considered by Hendrycks et al. (2019).

Smoothing the output classifier using random samples around the test point is a popular technique for producing *certifiably* robust classifiers. A certification, in this context, is a guarantee that given a test point  $x$ , all points within a certain radius of  $x$  receive the same label as  $x$ . Several papers have provided theoretical analyses to show that smoothing produces certifiably robust classifiers (Cao and Gong, 2017; Cohen et al., 2019; Lecuyer et al., 2019; Li et al., 2019; Liu et al., 2018; Salman et al., 2019; Levine and Feizi, 2020), whereas others have identified cases where smoothing does not work Yang et al. (2020a); Blum et al. (2020).

However, to the best of our knowledge, a PAC-type guarantee has not been shown for any algorithm that employs training data perturbations or output classifier smoothing, and our paper provides the first such analysis.

#### 4. Notations and setup

We denote by  $X$  the input domain and by  $Y = \{0, 1\}$  the binary label space. We assume that  $X$  is equipped with a metric  $\text{dist}$ . A hypothesis  $h : X \rightarrow Y$  is a function that assigns a label to each point in the domain. A hypothesis class  $\mathcal{H}$  is a set of such hypotheses. For a sample  $S = ((x_1, y_1), \dots, (x_n, y_n)) \in (X \times Y)^n$ , we use the notation  $S^X = \{x_1, x_2, \dots, x_n\}$  to denote the collection of domain points  $x_i$  occurring in  $S$ . The binary (also called 0-1) loss of  $h$  on data point  $(x, y) \in X \times Y$  is defined by

$$\ell^{0/1}(h, x, y) = \mathbb{1}[h(x) \neq y],$$

where  $\mathbb{1}[\cdot]$  is the indicator function. Let  $P$  be a probability distribution over  $X \times Y$ . Then the *expected binary loss* of  $h$  with respect to  $P$  is defined by

$$\mathcal{L}_P^{0/1}(h) = \mathbb{E}_{(x,y) \sim P}[\ell^{0/1}(h, x, y)]$$

Similarly, the *empirical binary loss* of  $h$  on sample  $S = ((x_1, y_1), \dots, (x_n, y_n)) \in (X \times Y)^n$  is defined as  $\mathcal{L}_S^{0/1}(h) = \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(h, x_i, y_i)$ . We also define the *approximation error* of  $\mathcal{H}$  with respect to  $P$  as  $\mathcal{L}_P^{0/1}(\mathcal{H}) = \inf_{h \in \mathcal{H}} \mathcal{L}_P^{0/1}(h)$ .

A *learner*  $\mathcal{A}$  is a function that takes in a finite sequence of labeled instances  $S = ((x_1, y_1), \dots, (x_n, y_n))$  and outputs a hypothesis  $h = \mathcal{A}(S)$ . The following definition abstracts the notion of PAC-learning Vapnik and Chervonenkis (1971); Valiant (1984).

**Definition 3 (PAC-learner)** *Let  $\mathcal{P}$  be a set of distributions over  $X \times Y$  and  $\mathcal{H}$  a hypothesis class. We say  $\mathcal{A}$  PAC-learns  $(\mathcal{H}, \mathcal{P})$  with  $m_{\mathcal{A}} : (0, 1)^2 \rightarrow \mathbb{N}$  samples if the following holds: for every distribution  $P \in \mathcal{P}$  over  $X \times Y$ , and every  $\epsilon, \delta \in (0, 1)$ , if  $S$  is an i.i.d. sample of size at least  $m_{\mathcal{A}}(\epsilon, \delta)$  from  $P$ , then with probability at least  $1 - \delta$  (over the randomness of  $S$ ) we have*

$$\mathcal{L}_P(\mathcal{A}(S)) \leq \mathcal{L}_P(\mathcal{H}) + \epsilon.$$

$\mathcal{A}$  is called an *agnostic learner* if  $\mathcal{P}$  is the set of all distributions on  $X \times Y$ , and a *realizable learner* if  $\mathcal{P} = \{P : \mathcal{L}_P(\mathcal{H}) = 0\}$ .

The smallest function  $m : (0, 1)^2 \rightarrow \mathbb{N}$  for which there exists a learner  $\mathcal{A}$  that satisfies the above definition with  $m_{\mathcal{A}} = m$  is referred to as the (realizable or agnostic) *sample complexity* of learning  $\mathcal{H}$ .

The existence of sample-efficient PAC-learners for VC classes is a standard result [Vapnik and Chervonenkis \(1971\)](#). We state the results formally in [Appendix A](#).

#### 4.1. Tolerant adversarial PAC-learning

Let  $\mathcal{U} : X \rightarrow 2^X$  be a function that maps each point in the domain to the set of its ‘‘admissible’’ perturbations. We call this function the *perturbation type*. The adversarial loss of  $h$  with respect to  $\mathcal{U}$  on  $(x, y) \in X \times Y$  is defined by

$$\ell^{\mathcal{U}}(h, x, y) = \max_{z \in \mathcal{U}(x)} \{\ell^{0/1}(h, z, y)\}$$

The *expected adversarial loss* with respect to  $P$  is defined by  $\mathcal{L}_P^{\mathcal{U}}(h) = \mathbb{E}_{(x,y) \sim P} \ell^{\mathcal{U}}(h, x, y)$ . The *empirical adversarial loss* of  $h$  on sample  $S = ((x_1, y_1), \dots, (x_n, y_n)) \in (X \times Y)^n$  is defined by  $\mathcal{L}_S^{\mathcal{U}}(h) = \frac{1}{n} \sum_{i=1}^n \ell^{\mathcal{U}}(h, x_i, y_i)$ . Finally, the *adversarial approximation error* of  $\mathcal{H}$  with respect to  $\mathcal{U}$  and  $P$  is defined by  $\mathcal{L}_P^{\mathcal{U}}(\mathcal{H}) = \inf_{h \in \mathcal{H}} \mathcal{L}_P^{\mathcal{U}}(h)$ .

The following definition generalizes the setting of PAC adversarial learning to what we call the *tolerant* setting, where we consider two perturbation types  $\mathcal{U}$  and  $\mathcal{V}$ . We say  $\mathcal{U}$  is *contained in*  $\mathcal{V}$  and and write it as  $\mathcal{U} \prec \mathcal{V}$  if  $\mathcal{U}(x) \subset \mathcal{V}(x)$  for all  $x \in X$ .

**Definition 4 (Tolerant Adversarial PAC-learner)** *Let  $\mathcal{P}$  be a set of distributions over  $X \times Y$ ,  $\mathcal{H}$  a hypothesis class, and  $\mathcal{U} \prec \mathcal{V}$  two perturbation types. We say  $\mathcal{A}$  tolerantly PAC-learns  $(\mathcal{H}, \mathcal{P}, \mathcal{U}, \mathcal{V})$  with  $m_{\mathcal{A}} : (0, 1)^2 \rightarrow \mathbb{N}$  samples if the following holds: for every distribution  $P \in \mathcal{P}$  and every  $\epsilon, \delta \in (0, 1)$ , if  $S$  is an i.i.d. sample of size at least  $m_{\mathcal{A}}(\epsilon, \delta)$  from  $P$ , then with probability at least  $1 - \delta$  (over the randomness of  $S$ ) we have*

$$\mathcal{L}_P^{\mathcal{U}}(\mathcal{A}(S)) \leq \mathcal{L}_P^{\mathcal{V}}(\mathcal{H}) + \epsilon.$$

*We say  $\mathcal{A}$  is a tolerant PAC-learner in the agnostic setting if  $\mathcal{P}$  is the set of all distributions over  $X \times Y$ , and in the tolerantly realizable setting if  $\mathcal{P} = \{P : \mathcal{L}_P^{\mathcal{V}}(\mathcal{H}) = 0\}$ .*

In the above context, we refer to  $\mathcal{U}$  as the *actual perturbation type* and to  $\mathcal{V}$  as the *reference perturbation type*. The case where  $\mathcal{U}(x) = \mathcal{V}(x)$  for all  $x \in X$  corresponds to the usual adversarial learning scenario (with no tolerance).

#### 4.2. Tolerant adversarial PAC-learning in metric spaces

If  $X$  is equipped with a metric  $\text{dist}(\cdot, \cdot)$ , then  $\mathcal{U}(x)$  can be naturally defined by a ball of radius  $r$  around  $x$ , i.e.,  $\mathcal{U}(x) = \mathcal{B}_r(x) = \{z \in X \mid \text{dist}(x, z) \leq r\}$ . To simplify the notation, we sometimes use  $\ell^r(h, x, y)$  instead of  $\ell^{\mathcal{B}_r}(h, x, y)$  to denote the adversarial loss with respect to  $\mathcal{B}_r$ .

In the tolerant setting, we consider the perturbation sets  $\mathcal{U}(x) = \mathcal{B}_r(x)$  and  $\mathcal{V}(x) = \mathcal{B}_{(1+\gamma)r}(x)$ , where  $\gamma > 0$  is called the *tolerance parameter*. Note that  $\mathcal{U} \prec \mathcal{V}$ . We now define PAC-learning with respect to the metric space.

**Definition 5 (Tolerant Adversarial Learning in metric spaces)** *Let  $(X, \text{dist})$  be a metric space,  $\mathcal{H}$  a hypothesis class, and  $\mathcal{P}$  a set of distributions of  $X \times Y$ . We say  $(\mathcal{H}, \mathcal{P}, \text{dist})$  is tolerantly PAC-learnable with  $m : (0, 1)^3 \rightarrow \mathbb{N}$  samples when for every  $r, \gamma > 0$  there exist a PAC-learner  $A_{r, \gamma}$  for  $(\mathcal{H}, \mathcal{P}, B_r, B_{r(1+\gamma)})$  that uses  $m(\epsilon, \delta, \gamma)$  samples.*

**Remark 6** *In this definition the learner receives  $\gamma$  and  $r$  as input but its sample complexity does not depend on  $r$  (but can depend on  $\gamma$ ). Also, as in Definition 4, the tolerantly realizable setting corresponds to  $\mathcal{P} = \{P : \mathcal{L}_P^{r(1+\gamma)}(\mathcal{H}) = 0\}$  while in the agnostic setting  $\mathcal{P}$  is the set of all distributions over  $X \times Y$ .*

The doubling dimension and the doubling measure of the metric space will play important roles in our analysis. We refer the reader to Appendix B for their definitions.

We will use the following lemma in our analysis, whose proof can be found in Appendix B:

**Lemma 7** *For any family  $\mathcal{M}$  of complete, doubling metric spaces, there exist constants  $c_1, c_2 > 0$  such that for any metric space  $(X, \text{dist}) \in \mathcal{M}$  with doubling dimension  $d$ , there exists a measure  $\mu$  such that if a ball  $\mathcal{B}_r$  of radius  $r > 0$  is completely contained inside a ball  $\mathcal{B}_{\alpha r}$  of radius  $\alpha r$  (with potentially a different center) for any  $\alpha > 1$ , then  $0 < \mu(\mathcal{B}_{\alpha r}) \leq (c_1 \alpha)^{c_2 d} \mu(\mathcal{B}_r)$ . Furthermore, if we have a constant  $\alpha_0 > 1$  such that we know that  $\alpha \geq \alpha_0$  then the bound can be simplified to  $0 < \mu(\mathcal{B}_{\alpha r}) \leq \alpha^{\zeta d} \mu(\mathcal{B}_r)$ , where  $\zeta$  depends on  $\mathcal{M}$  and  $\alpha_0$ .*

Later, we will set  $\alpha = 1 + 1/\gamma$  where  $\gamma$  is the tolerance parameter. Since we are mostly interested in small values of  $\gamma$ , suppose we decide on some loose upper bound  $\Gamma \gg \gamma$ . This corresponds to saying that there exists some  $\alpha_0 > 1$  such that  $\alpha \geq \alpha_0$ .

It is worth noting that in the special case of Euclidean metric spaces, we can set both  $c_1$  and  $c_2$  to be 1. In the rest of the paper, we will assume we have a loose upper bound  $\Gamma \gg \gamma$  and use the simpler bound from Lemma 24 extensively.

Given a metric space  $(X, d)$  and a measure  $\mu$  defined over it, for any subset  $Z \subseteq X$  for which  $\mu(Z)$  is non-zero and finite,  $\mu$  induces a probability measure  $P_Z^\mu$  over  $Z$  as follows. For any set  $Z' \subseteq Z$  in the  $\sigma$ -algebra over  $Z$ , we define  $P_Z^\mu(Z') = \mu(Z')/\mu(Z)$ . With a slight abuse of notation, we write  $z \sim Z$  to mean  $z \sim P_Z^\mu$  whenever we know  $\mu$  from the context.

Our learners rely on being able to sample from  $P_Z^\mu$ . Thus we define the following oracle, which can be implemented efficiently for  $\ell_p$  spaces.

**Definition 8 (Sampling Oracle)** *Given a metric space  $(X, \text{dist})$  equipped with a doubling measure  $\mu$ , a sampling oracle is an algorithm that when queried with a  $Z \subseteq X$  such that  $\mu(Z)$  is finite, returns a sample drawn from  $P_Z^\mu$ . We will use the notation  $z \sim Z$  for queries to this oracle.*

## 5. The perturb-and-smooth approach for tolerant adversarial learning

In this section we focus on tolerant adversarial PAC-learning in metric spaces (Definition 5), and show that VC classes are tolerantly PAC-learnable in the tolerantly realizable setting. Interestingly, we prove this result using an approach that resembles the ‘‘perturb-and-smooth’’ paradigm which is used in practice (for example by Cohen et al. (2019)). The overall idea is to ‘‘perturb’’ each training point  $x$ , train a classifier on the ‘‘perturbed’’ points, and ‘‘smooth out’’ the final hypothesis using a certain majority rule.

We employ three perturbation types:  $\mathcal{U}$  and  $\mathcal{V}$  play the role of the *actual* and the *reference* perturbation type respectively. Additionally, we consider a perturbation type  $\mathcal{W} : X \rightarrow 2^X$ , which is used for smoothing. We assume  $\mathcal{U} \prec \mathcal{V}$  and  $\mathcal{W} \prec \mathcal{V}$ . For this section, we will use metric balls for the three types. Specifically, if  $\mathcal{U}$  consists of balls of radius  $r$  for some  $r > 0$ , then  $\mathcal{W}$  will consist of balls of radius  $\gamma r$  and  $\mathcal{V}$  will consist of balls of radius  $(1 + \gamma)r$ .

**Definition 9 (Smoothed classifier)** For a hypothesis  $h : X \rightarrow \{0, 1\}$ , and perturbation type  $\mathcal{W} : X \rightarrow 2^X$ , we let  $\bar{h}_{\mathcal{W}}$  denote the classifier resulting from replacing the label  $h(x)$  with the average label over  $\mathcal{W}(x)$ , that is

$$\bar{h}_{\mathcal{W}}(x) = \mathbb{1} \left[ \mathbb{E}_{x' \sim \mathcal{W}(x)} h(x') \geq 1/2 \right]$$

For metric perturbation types, where  $\mathcal{W}$  is a ball of some radius  $r$ , we also use the notation  $\bar{h}_r$  and when the type  $\mathcal{W}$  is clear from context, we may omit the subscript altogether and simply write  $\bar{h}$  for the smoothed classifier.

**The tolerant perturb-and-smooth algorithm** We propose the following learning algorithm, TPaS, for tolerant learning in metric spaces. Let the perturbation radius be  $r > 0$  for the actual type  $\mathcal{U} = \mathcal{B}_r$ , and let  $S = ((x_1, y_1), \dots, (x_m, y_m))$  be the training sample. For each  $x_i \in S^X$ , the learner samples a point  $x'_i \sim \mathcal{B}_{r \cdot (1+\gamma)}(x_i)$  (using the sampling oracle) from the expanded reference perturbation set  $\mathcal{V}(x_i) = \mathcal{B}_{(1+\gamma)r}(x_i)$ . Let  $S' = ((x'_1, y_1), \dots, (x'_m, y_m))$ . TPaS then invokes a (standard, non-robust) PAC-learner  $\mathcal{A}_{\mathcal{H}}$  for the hypothesis class  $\mathcal{H}$  on the perturbed data  $S'$ . We let  $\hat{h} = \mathcal{A}_{\mathcal{H}}(S')$  denote the output of this PAC-learner. Finally, TPaS outputs the  $\mathcal{W}$ -smoothed version of  $\bar{h}_{\gamma r}$  for  $\mathcal{W} = \mathcal{B}_{\gamma r}$ . That is,  $\bar{h}_{\gamma r}(x)$  is simply the majority label in a ball of radius  $\gamma r$  around  $x$  with respect to the distribution defined by  $\mu$ , see also Definition 9. We will prove below that this  $\bar{h}_{\gamma r}$  has a small  $\mathcal{U}$ -adversarial loss. Algorithm 1 below summarizes our learning procedure.

---

**Algorithm 1** Tolerant Perturb and Smooth (TPaS)

---

**Input:** Radius  $r$ , tolerance parameter  $\gamma$ , data  $S = ((x_1, y_1), \dots, (x_m, y_m))$ , access to sampling oracle  $\mathcal{O}$  for  $\mu$  and PAC-learner  $\mathcal{A}_{\mathcal{H}}$ .

Initialize  $S' = \emptyset$

**for**  $i = 1$  to  $m$  **do**

Sample  $x'_i \sim \mathcal{B}_{(1+\gamma)r}(x_i)$

Add  $(x'_i, y_i)$  to  $S'$

**end for**

Set  $\hat{h} = \mathcal{A}_{\mathcal{H}}(S')$

**Output:**  $\bar{h}_{\gamma r}$  defined by

$$\bar{h}_{\gamma r}(x) = \mathbb{1} \left[ \mathbb{E}_{x' \sim \mathcal{B}_{\gamma r}(x)} \hat{h}(x') \geq 1/2 \right]$$


---

The following is the main result of this section.

**Theorem 10** Let  $(X, \text{dist})$  be an any metric space with doubling dimension  $d$  and doubling measure  $\mu$ . Let  $\mathcal{O}$  be a sampling oracle for  $\mu$ . Let  $\mathcal{H}$  be a hypothesis class and  $\mathcal{P}$  a set of distributions over  $X \times Y$ . Assume  $\mathcal{A}_{\mathcal{H}}$  PAC-learns  $\mathcal{H}$  with  $m_{\mathcal{H}}(\epsilon, \delta)$  samples in the realizable setting. Then there exists a learner  $\mathcal{A}$ , namely TPaS, that



- *Tolerantly PAC-learns*  $(\mathcal{H}, \mathcal{P}, \text{dist})$  in the tolerantly realizable setting with sample complexity bounded by  $m(\epsilon, \delta, \gamma) = O(m_{\mathcal{H}}(\epsilon, \delta) \cdot (1 + 1/\gamma)^{\zeta d}) = O\left(\frac{\text{VC}(\mathcal{H}) + \log 1/\delta}{\epsilon} \cdot (1 + 1/\gamma)^{\zeta d}\right)$ , where  $\gamma$  is the tolerance parameter and  $d$  is the doubling dimension.
- *Makes only one query to*  $\mathcal{A}_{\mathcal{H}}$
- *Makes*  $m(\epsilon, \delta, \gamma)$  *queries to sampling oracle*  $\mathcal{O}$

The proof of this theorem uses the following key technical lemma (its proof can be found in Appendix C):

**Lemma 11** *Let*  $r > 0$  *be a perturbation radius,*  $\gamma > 0$  *a tolerance parameter, and*  $g : X \rightarrow Y$  *a classifier. For*  $x \in X$  *and*  $y \in Y = \{0, 1\}$ , *we define*

$$\Sigma_{g,y}(x) = \mathbb{E}_{z \sim \mathcal{B}_{r(1+\gamma)}(x)} \mathbb{1}[g(z) \neq y] \quad \text{and} \quad \sigma_{g,y}(x) = \mathbb{E}_{z \sim \mathcal{B}_{r\gamma}(x)} \mathbb{1}[g(z) \neq y].$$

Then  $\Sigma_{g,y}(x) \leq \frac{1}{3} \cdot \left(\frac{1+\gamma}{\gamma}\right)^{-\zeta d}$  implies that  $\sigma_{g,y}(z) \leq 1/3$  for all  $z \in \mathcal{B}_r(x)$ .

**Proof** [Proof of Theorem 10] Consider some  $\epsilon_0 > 0$  and  $0 < \delta < 1$  to be given (we will pick a suitable value of  $\epsilon_0$  later), and assume the PAC-learner  $\mathcal{A}_{\mathcal{H}}$  was invoked on the perturbed sample  $S'$  of size at least  $m_A(\epsilon_0, \delta)$ . According to definition 3, this implies that with probability  $1 - \delta$ , the output  $\hat{h} = \mathcal{A}_{\mathcal{H}}(S)$  has (binary) loss at most  $\epsilon_0$  with respect to the data-generating distribution. Note that the relevant distribution here is the two-stage process of the original data generating distribution  $P$  and the perturbation sampling according to  $\mathcal{V} = \mathcal{B}_{(1+\gamma)r}$ . Since  $P$  is  $\mathcal{V}$ -robustly realizable, the two-stage process yields a realizable distribution with respect to the standard 0/1-loss. Thus, we have

$$\mathbb{E}_{(x,y) \sim P} \mathbb{E}_{z \sim \mathcal{B}_{r(1+\gamma)}(x)} \mathbb{1}[\hat{h}(z) \neq y] \leq \epsilon_0.$$

With Lemma 11, this becomes  $\mathbb{E}_{(x,y) \sim P} \Sigma_{\hat{h},y}(x) \leq \epsilon_0$ . For  $\lambda > 0$ , Markov's inequality then yields :

$$\mathbb{E}_{(x,y) \sim P} \mathbb{1}[\Sigma_{\hat{h},y}(x) \leq \lambda] > 1 - \epsilon_0/\lambda \tag{1}$$

Thus setting  $\lambda = \frac{1}{3} \cdot \left(\frac{1+\gamma}{\gamma}\right)^{-\zeta d}$  and plugging in the result of the Lemma 11 to equation (1), we get

$$\mathbb{E}_{(x,y) \sim P} \mathbb{1}[\forall z \in \mathcal{B}_r(x), \sigma_{\hat{h},y}(z) \leq 1/3] > 1 - \epsilon_0/\lambda.$$

Since  $\sigma_{\hat{h},y}(z) \leq 1/3$  implies that  $\mathbb{1}[\mathbb{E}_{z' \sim \mathcal{B}_{\gamma r}(z)} \hat{h}(z') \geq 1/2] = y$ , using the definition of the smoothed classifier  $\bar{h}_{\gamma r}$  we get

$$\mathbb{E}_{(x,y) \sim P} \mathbb{1}[\exists z \in \mathcal{B}_r(x), \bar{h}_{\gamma r}(z) \neq y] \leq \epsilon_0/\lambda, \tag{2}$$

which implies  $\mathcal{L}_P^r(\bar{h}_{\gamma r}) \leq \epsilon_0/\lambda$ . Thus, for the robust learning problem, if we are given a desired accuracy  $\epsilon$  and we want  $\mathcal{L}_P^r(\bar{h}_{\gamma r}) \leq \epsilon$ , we can pick  $\epsilon_0 = \lambda\epsilon$ . Putting it all together, we get

sample complexity  $m \leq O\left(\frac{\text{VC}(\mathcal{H}) + \log 1/\delta}{\epsilon_0}\right)$  where  $\epsilon_0 = \lambda\epsilon$ , and  $\lambda = \frac{1}{3} \cdot \left(\frac{1+\gamma}{\gamma}\right)^{-\zeta d}$ . Therefore,  $m \leq O\left(\frac{\text{VC}(\mathcal{H}) + \log 1/\delta}{\epsilon} \cdot (1 + 1/\gamma)^{\zeta d}\right)$ . ■

**Computational complexity of the learner.** Assuming we have access to  $\mathcal{O}$  and an efficient algorithm for non-robust PAC-learning in the realizable setting, we can compute  $\hat{h}$  efficiently. Therefore, the learning can be done efficiently in this case. However, at the prediction time, we need to compute  $\bar{h}(x)$  on new test points which requires us to compute an expectation. We can instead *estimate* the expectations using random samples from the sampling oracle. For a single test point  $x$ , if the number of samples we draw is  $\Omega(\log 1/\delta)$  then with probability at least  $1 - \delta$  we get the same result as that of the optimal  $\bar{h}(x)$ . Using more samples we can boost this probability to guarantee a similar output to that of  $\bar{h}$  on a larger set of test points.

**The traditional non-tolerant framework does not justify the use of perturb-and-smooth-type approaches.** The introduction of the tolerance in the adversarial learning framework is crucial for being able to prove guarantees for perturb-and-smooth-type algorithms. To see why, consider a simple case where the domain is the real line, the perturbation set is an open Euclidean ball of radius 1, and the hypothesis class is the set of all thresholds. Assume that the underlying distribution is supported only on two points:  $\mathbb{P}(x = -1, y = 1) = \mathbb{P}(x = 1, y = 0) = 0.5$ . This distribution is robustly realizable, but the threshold should be set exactly to  $x = 0$  to get a small error. However, the perturb-and-smooth method will fail because the only way the PAC-learner  $\mathcal{A}_{\mathcal{H}}$  sets the threshold to  $x = 0$  is if it receives a (perturbed) sample exactly at  $x = 0$ , whose probability is 0.

## 6. Improved tolerant learning guarantees through sample compression

The perturb-and-smooth approach discussed in the previous section offers a general method for tolerant robust learning. However, one shortcoming of this approach is the exponential dependence of its sample complexity with respect to the doubling dimension of the metric space. Furthermore, the tolerant robust guarantee relied on the data generating distribution being tolerantly realizable. In this section, we propose another approach that addresses both of these issues. The idea is to adopt the perturb-and-smooth approach within a sample compression argument. We introduce the notion of a  $(\mathcal{U}, \mathcal{V})$ -tolerant sample compression scheme and present a learning bound based on such a compression scheme, starting with the realizable case. We then show that this implies learnability in the agnostic case as well. Remarkably, this tolerant compression based analysis will yield bounds on the sample complexity that avoid the exponential dependence on the doubling dimension.

For a compact representation, we will use the general notation  $\mathcal{U}, \mathcal{V}$ , and  $\mathcal{W}$  for the three perturbation types (actual, reference and smoothing type) in this section and will assume that they satisfy the Property 1 below for some parameter  $\beta > 0$ . Lemma 11 implies that, in the metric setting, for any radius  $r$  and tolerance parameter  $\gamma$  the perturbation types  $\mathcal{U} = \mathcal{B}_r$ ,  $\mathcal{V} = \mathcal{B}_{(1+\gamma)r}$ , and  $\mathcal{W} = \mathcal{B}_{\gamma r}$  have this property for  $\beta = \frac{1}{3} \left(\frac{1+\gamma}{\gamma}\right)^{-\zeta d}$ .

**Property 1** For a fixed  $0 < \beta < 1/2$ , we assume that the perturbation types  $\mathcal{V}, \mathcal{U}$  and  $\mathcal{W}$  are so that for any classifier  $h$  and any  $x \in X$ , any  $y \in \{0, 1\}$  if

$$\mathbb{E}_{z \sim \mathcal{V}(x)}[h(z) = y] \geq 1 - \beta$$

then  $\mathcal{W}$ -smoothed class classifier  $\bar{h}_{\mathcal{W}}$  satisfies  $\bar{h}_{\mathcal{W}}(z) = y$  for all  $z \in \mathcal{U}(x)$ .

A compression scheme of size  $k$  is a pair of functions  $(\kappa, \rho)$ , where the compression function  $\kappa : \bigcup_{i=1}^{\infty} (X \times Y)^i \rightarrow \bigcup_{i=1}^k (X \times Y)^i$  maps samples  $S = ((x_1, y_1), (x_2, y_2), \dots, (x_m, y_m))$  of arbitrary size to sub-samples of  $S$  of size at most  $k$ , and  $\rho : \bigcup_{i=1}^k (X \times Y)^i \rightarrow Y^X$  is a decompression function that maps samples to classifiers. The pair  $(\kappa, \rho)$  is a sample compression scheme for loss  $\ell$  and class  $\mathcal{H}$ , if for any samples  $S$  realizable by  $\mathcal{H}$ , we recover the correct labels for all  $(x, y) \in S$ , that is,  $\mathcal{L}_S(H) = 0$  implies that  $\mathcal{L}_S(\kappa \circ \rho(S)) = 0$ .

For tolerant learning, we introduce the following generalization of compression schemes:

**Definition 12 (Tolerant sample compression scheme)** *A sample compression scheme  $(\kappa, \rho)$  is a  $\mathcal{U}, \mathcal{V}$ -tolerant sample compression scheme for class  $\mathcal{H}$ , if for any samples  $S$  that are  $\ell^{\mathcal{V}}$  realizable by  $\mathcal{H}$ , that is  $\mathcal{L}_S^{\mathcal{V}}(\mathcal{H}) = 0$ , we have  $\mathcal{L}_S^{\mathcal{U}}(\kappa \circ \rho(S)) = 0$ .*

The next lemma establishes that the existence of a sufficiently small tolerant compression scheme for the class  $\mathcal{H}$  yields bounds on the sample complexity of tolerantly learning  $\mathcal{H}$ . The proof of the lemma is based on a modification of a standard compression based generalization bound. Appendix Section D provides more details.

**Lemma 13** *Let  $\mathcal{H}$  be a hypothesis class and  $\mathcal{U}$  and  $\mathcal{V}$  be perturbation types with  $\mathcal{U}$  included in  $\mathcal{V}$ . If the class  $\mathcal{H}$  admits a  $(\mathcal{U}, \mathcal{V})$ -tolerant compression scheme of size bounded by  $k \ln(m)$  for sample of size  $m$ , then the class is  $(\mathcal{U}, \mathcal{V})$ -tolerantly learnable in the realizable case with sample complexity bounded by  $m(\epsilon, \delta) = \tilde{O}\left(\frac{k + \ln(1/\delta)}{\epsilon}\right)$ .*

We next establish a bound on the tolerant compression size for general VC-classes, which will then immediately yield the improved sample complexity bounds for tolerant learning in the realizable case. The proof is sketched here; its full version has been moved to the Appendix.

**Lemma 14** *Let  $\mathcal{H} \subseteq Y^X$  be some hypothesis class with finite VC-dimension  $\text{VC}(\mathcal{H}) < \infty$ , and let  $\mathcal{U}, \mathcal{V}, \mathcal{W}$  satisfy the conditions in Property 1 for some  $\beta > 0$ . Then there exists a  $(\mathcal{U}, \mathcal{V})$ -tolerant sample compression scheme for  $\mathcal{H}$  of size  $\tilde{O}\left(\text{VC}(\mathcal{H}) \ln\left(\frac{m}{\beta}\right)\right)$ .*

**Proof [Proof Sketch]** We will employ a boosting-based approach to establish the claimed compression sizes. Let  $S = ((x_1, y_1), (x_2, y_2), \dots, (x_m, y_m))$  be a data-set that is  $\ell^{\mathcal{V}}$ -realizable with respect to  $\mathcal{H}$ . We let  $S_{\mathcal{V}}$  denote an ‘‘inflated data-set’’ that contains all domain points in the  $\mathcal{V}$ -perturbation sets of the  $x_i \in S^X$ , that is  $S_{\mathcal{V}}^X := \bigcup_{i=1}^m \mathcal{V}(x_i)$ . Every point  $z \in S_{\mathcal{V}}^X$  is assigned the label  $y = y_i$  of the minimally-indexed  $(x_i, y_i) \in S$  with  $z \in \mathcal{V}(x_i)$ , and we set  $S_{\mathcal{V}}$  to be the resulting collection of labeled data-points.

We then use the boost-by-majority method to encode a classifier  $g$  that (roughly speaking) has error bounded by  $\beta/m$  over (a suitable measure over)  $S_{\mathcal{V}}$ . This boosting method outputs a  $T$ -majority vote  $g(x) = \mathbb{1}\left[\sum_{i=1}^T h_i(x)\right] \geq 1/2$  over weak learners  $h_i$ , which in our case will be hypotheses from  $\mathcal{H}$ . We prove that this error can be achieved with  $T = 18 \ln\left(\frac{2m}{\beta}\right)$  rounds of boosting. We prove that each weak learner that is used in the boosting procedure can be encoded with  $n = \tilde{O}(\text{VC}(\mathcal{H}))$  many sample points from  $S$ . The resulting compression size is thus  $n \cdot T = \tilde{O}\left(\text{VC}(\mathcal{H}) \ln\left(\frac{m}{\beta}\right)\right)$ .

Finally, the error bound  $\beta/m$  of  $g$  over  $S_{\mathcal{V}}$  implies that the error in each perturbation set  $\mathcal{V}(x_i)$  of a sample point  $(x_i, y_i) \in S$  is at most  $\beta$ . Property 1 then implies  $\mathcal{L}_S^{\mathcal{U}}(\bar{g}_{\mathcal{W}}) = 0$  for the  $\mathcal{W}$ -smoothed classifier  $\bar{g}_{\mathcal{W}}$ , establishing the  $(\mathcal{U}, \mathcal{V})$ -tolerant correctness of the compression scheme. ■

This yields the following result

**Theorem 15** *Let  $\mathcal{H}$  be a hypothesis class of finite VC-dimension and  $\mathcal{V}, \mathcal{U}, \mathcal{W}$  be three perturbation types (actual, reference and smoothing) satisfying Property 1 for some  $\beta > 0$ . Then the sample complexity (omitting log-factors) of  $(\mathcal{U}, \mathcal{V})$ -tolerantly learning  $\mathcal{H}$  is bounded by*

$$m(\epsilon, \delta) = \tilde{O} \left( \frac{\text{VC}(\mathcal{H}) \ln(1/\beta) + \ln(1/\delta)}{\epsilon} \right)$$

*in the realizable case, and in the agnostic case by*

$$m(\epsilon, \delta) = \tilde{O} \left( \frac{\text{VC}(\mathcal{H}) \ln(1/\beta) + \ln(1/\delta)}{\epsilon^2} \right)$$

**Proof** The bound for the realizable case follows immediately from Lemma 14 and the subsequent discussion (in the Appendix). For the agnostic case, we employ a reduction from agnostic robust learnability to realizable robust learnability (Montasser et al., 2019; Moran and Yehudayoff, 2016). The reduction is analogous to the one presented in Appendix C of Montasser et al. (2019) for usual (non-tolerant) robust learnability with some minor modifications. Namely, for a sample  $S$ , we choose the largest subsample  $S'$  that is  $\ell^{\mathcal{V}}$ -realizable (this will result in competitiveness with a  $\ell^{\mathcal{V}}$ -optimal classifier), and we will use the boosting procedure described there for the  $\ell^{\mathcal{U}}$  loss. For the sample sizes employed for the weak learners in that procedure, we can use the sample complexity for  $\epsilon = \delta = 1/3$  of an optimal  $(\mathcal{U}, \mathcal{V})$ -tolerant learner in the realizable case (note that each learning problem during the boosting procedure is a realizable  $(\mathcal{U}, \mathcal{V})$ -tolerant learning task). These modifications result in the stated sample complexity for agnostic tolerant learnability. ■

In particular, for the doubling measure scenario (as considered in the previous section), we obtain

**Corollary 16** *For metric tolerant learning with tolerance parameter  $\gamma$  in doubling dimension  $d$  the sample complexity of adversarially robust learning with tolerance in the realizable case is bounded by  $m(\epsilon, \delta) = \tilde{O} \left( \frac{\text{VC}(\mathcal{H}) \zeta d \ln(1+1/\gamma) + \ln(1/\delta)}{\epsilon} \right)$  and in the agnostic case by  $m(\epsilon, \delta) = \tilde{O} \left( \frac{\text{VC}(\mathcal{H}) \zeta d \ln(1+1/\gamma) + \ln(1/\delta)}{\epsilon^2} \right)$ .*

**Discussion of linear dependence on  $\text{VC}(\mathcal{H})$**  Earlier, general compression based sample complexity bounds for robust learning with arbitrary perturbation sets exhibit a dependence on the dual VC-dimension of the hypothesis class and therefore potentially an exponential dependence on  $\text{VC}(\mathcal{H})$  (Montasser et al., 2019). In our setting, we show that it is possible to avoid the dependence on dual-VC by exploiting both the metric structure of the domain set and the tolerant framework. In the full proof of Lemma 14, we show that *if we can encode a classifier with small error* (exponentially small with respect to the doubling dimension for the metric case) on the perturbed distribution w.r.t. larger perturbation sets, then we can *use smoothing to get a classifier that correctly classifies every point in the inner inflated sets*. And, as for TPaS, the tolerant perspective is crucial to exploit a smoothing step in the compression approach (through the guarantee from Property 1 or Lemma 11).

More specifically, we define a tolerant compression scheme (Definition 12) that naturally extends the classic definition of compression to the tolerant framework. The compression scheme we establish in the proof of Lemma 14 then borrows ideas from our perturb-and-smooth algorithm. Within the compression argument, we define the perturbed distribution over the sample that we want to compress with respect to the larger perturbation sets. We then use boosting to build a classifier with very small error with respect to this distribution. The nice property of boosting is that its error decreases

exponentially with the number of iterations. As a result, we also get linear dependence on the doubling dimension. This classifier can be encoded using  $\tilde{O}(TVC(\mathcal{H}))$  samples ( $T$  rounds of boosting, and each weak classifier can be encoded using  $O(VC(\mathcal{H}))$  samples, since we can here use simple  $\epsilon$ -approximations rather than invoking VC-theory in the dual space). Our decoder receives the description of these weak classifiers, combines them, and performs a final smoothing step. The smoothing step translates the exponentially small error with respect to the perturbed distribution to zero error with respect to the (inner) inflated set, thereby satisfying the requirement of a tolerant compression scheme.

## References

- Hassan Ashtiani, Vinayak Pathak, and Ruth Uerner. Black-box certification and learning under adversarial perturbations. In *International Conference on Machine Learning*, pages 388–398. PMLR, 2020.
- Idan Attias, Aryeh Kontorovich, and Yishay Mansour. Improved generalization bounds for robust learning. In *Algorithmic Learning Theory, ALT*, pages 162–183, 2019.
- Pranjal Awasthi, Abhratanu Dutta, and Aravindan Vijayaraghavan. On robustness to adversarial examples and polynomial optimization. In *Advances in Neural Information Processing Systems, NeurIPS*, pages 13760–13770, 2019.
- Robi Bhattacharjee, Max Hopkins, Akash Kumar, Hantao Yu, and Kamalika Chaudhuri. Robust empirical risk minimization with tolerance. *arXiv preprint arXiv:2210.00635*, 2022.
- Avrim Blum, Travis Dick, Naren Manoj, and Hongyang Zhang. Random smoothing might be unable to certify  $l_\infty$  robustness for high-dimensional images. *Journal of machine learning research*, 21(211), 2020.
- Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K Warmuth. Learnability and the vapnik-chervonenkis dimension. *Journal of the ACM (JACM)*, 36(4):929–965, 1989.
- Sébastien Bubeck and Mark Sellke. A universal law of robustness via isoperimetry. *Advances in Neural Information Processing Systems*, 34:28811–28822, 2021.
- Xiaoyu Cao and Neil Zhenqiang Gong. Mitigating evasion attacks to deep neural networks via region-based classification. In *Proceedings of the 33rd Annual Computer Security Applications Conference*, pages 278–287, 2017.
- Jeremy M. Cohen, Elan Rosenfeld, and J. Zico Kolter. Certified adversarial robustness via randomized smoothing. In *Proceedings of the 36th International Conference on Machine Learning, ICML*, pages 1310–1320, 2019.
- Daniel Cullina, Arjun Nitin Bhagoji, and Prateek Mittal. Pac-learning in the presence of adversaries. In *Advances in Neural Information Processing Systems, NeurIPS*, pages 230–241, 2018.
- Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. Exploring the landscape of spatial robustness. In *International Conference on Machine Learning*, pages 1802–1811. PMLR, 2019.

- Alhussein Fawzi and Pascal Frossard. Manitest: Are classifiers really invariant? In *British Machine Vision Conference (BMVC)*, number CONF, 2015.
- Uriel Feige, Yishay Mansour, and Robert Schapire. Learning and inference in the presence of corrupted inputs. In *Conference on Learning Theory, COLT*, pages 637–657, 2015.
- Shafi Goldwasser, Adam Tauman Kalai, Yael Tauman Kalai, and Omar Montasser. Beyond perturbations: Learning guarantees with arbitrary adversarial test examples. *arXiv preprint arXiv:2007.05145*, 2020.
- Ian J. Goodfellow, Patrick D. McDaniel, and Nicolas Papernot. Making machine learning robust against adversarial inputs. *Commun. ACM*, 61(7):56–66, 2018.
- Steve Hanneke. The optimal sample complexity of pac learning. *The Journal of Machine Learning Research*, 17(1):1319–1333, 2016.
- David Haussler. Decision theoretic generalizations of the pac model for neural net and other learning applications. *Information and computation*, 100(1):78–150, 1992.
- David Haussler and Emo Welzl. epsilon-nets and simplex range queries. *Discret. Comput. Geom.*, 2: 127–151, 1987.
- Dan Hendrycks, Norman Mu, Ekin Dogus Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. In *International Conference on Learning Representations*, 2019.
- Hossein Hosseini and Radha Poovendran. Semantic adversarial examples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1614–1619, 2018.
- Nathan Inkawhich, Wei Wen, Hai Helen Li, and Yiran Chen. Feature space perturbations yield more transferable adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7066–7074, 2019.
- Can Kanbak, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. Geometric robustness of deep networks: Analysis and improvement. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4441–4449. IEEE, 2018.
- Daniel Kang, Yi Sun, Dan Hendrycks, Tom Brown, and Jacob Steinhardt. Testing robustness against unforeseen adversaries. *arXiv preprint arXiv:1908.08016*, 2019.
- Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified robustness to adversarial examples with differential privacy. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 656–672. IEEE, 2019.
- Alexander Levine and Soheil Feizi. Robustness certificates for sparse adversarial attacks by randomized ablation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 4585–4593, 2020.
- Bai Li, Changyou Chen, Wenlin Wang, and Lawrence Carin. Certified adversarial robustness with additive noise. *Advances in Neural Information Processing Systems*, 32:9464–9474, 2019.

- Xuanqing Liu, Minhao Cheng, Huan Zhang, and Cho-Jui Hsieh. Towards robust neural networks via random self-ensemble. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 369–385, 2018.
- Jouni Luukkainen and Eero Saksman. Every complete doubling metric space carries a doubling measure. *Proceedings of the American Mathematical Society*, 126(2):531–534, 1998.
- Omar Montasser, Steve Hanneke, and Nathan Srebro. VC classes are adversarially robustly learnable, but only improperly. In *Conference on Learning Theory, COLT*, pages 2512–2530, 2019.
- Omar Montasser, Surbhi Goel, Ilias Diakonikolas, and Nathan Srebro. Efficiently learning adversarially robust halfspaces with noise. *arXiv preprint arXiv:2005.07652*, 2020a.
- Omar Montasser, Steve Hanneke, and Nati Srebro. Reducing adversarially robust learning to non-robust pac learning. In *NeurIPS*, 2020b.
- Omar Montasser, Steve Hanneke, and Nathan Srebro. Adversarially robust learning with unknown perturbation sets. *arXiv preprint arXiv:2102.02145*, 2021a.
- Omar Montasser, Steve Hanneke, and Nathan Srebro. Transductive robust learning guarantees. *arXiv preprint arXiv:2110.10602*, 2021b.
- Shay Moran and Amir Yehudayoff. Sample compression schemes for vc classes. *Journal of the ACM (JACM)*, 63(3):1–10, 2016.
- Vinod Raman, Unique Subedi, and Ambuj Tewari. Probabilistically robust pac learning. *arXiv preprint arXiv:2211.05656*, 2022.
- Sara Sabour, Yanshuai Cao, Fartash Faghri, and David J Fleet. Adversarial manipulation of deep representations. In *ICLR (Poster)*, 2016.
- Hadi Salman, Jerry Li, Ilya P. Razenshteyn, Pengchuan Zhang, Huan Zhang, Sébastien Bubeck, and Greg Yang. Provably robust deep learning via adversarially trained smoothed classifiers. In *Advances in Neural Information Processing Systems 32, NeurIPS*, pages 11289–11300, 2019.
- Robert E Schapire and Yoav Freund. Boosting: Foundations and algorithms. *Kybernetes*, 2013.
- Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.
- Hans U Simon. An almost optimal pac algorithm. In *Conference on Learning Theory*, pages 1552–1563. PMLR, 2015.
- Yang Song, Rui Shu, Nate Kushman, and Stefano Ermon. Constructing unrestricted adversarial examples with generative models. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 8322–8333, 2018.
- 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.

- Leslie G. Valiant. A theory of the learnable. *Commun. ACM*, 27(11):1134–1142, 1984.
- V. N. Vapnik and A. Ya. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability & Its Applications*, 16(2):264–280, 1971.
- Chaowei Xiao, Jun-Yan Zhu, Bo Li, Warren He, Mingyan Liu, and Dawn Song. Spatially transformed adversarial examples. In *International Conference on Learning Representations*, 2018.
- Qiuling Xu, Guanhong Tao, Siyuan Cheng, and Xiangyu Zhang. Towards feature space adversarial attack. *arXiv preprint arXiv:2004.12385*, 2020.
- Larry Yaeger, Richard Lyon, and Brandyn Webb. Effective training of a neural network character classifier for word recognition. *Advances in neural information processing systems*, 9:807–816, 1996.
- Greg Yang, Tony Duan, J Edward Hu, Hadi Salman, Ilya Razenshteyn, and Jerry Li. Randomized smoothing of all shapes and sizes. In *Proceedings of the 37th International Conference on Machine Learning*, pages 693–10705, 2020a.
- Yao-Yuan Yang, Cyrus Rashtchian, Hongyang Zhang, Ruslan Salakhutdinov, and Kamalika Chaudhuri. A closer look at accuracy vs. robustness. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, pages 8588–8601, 2020b.
- Stephan Zheng, Yang Song, Thomas Leung, and Ian Goodfellow. Improving the robustness of deep neural networks via stability training. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4480–4488, 2016.



## Appendix A. Standard results from VC theory

Let  $X$  be a domain. For hypothesis  $h$  and  $B \subseteq X$  let  $h(B) = (h(b))_{b \in B}$ .

**Definition 17 (VC-dimension)** We say  $\mathcal{H}$  shatters  $B \subseteq X$  if  $|\{h(B) : h \in \mathcal{H}\}| = 2^{|B|}$ . The VC-dimension of  $\mathcal{H}$ , denoted by  $\text{VC}(\mathcal{H})$ , is defined to be the supremum of the size of the sets that are shattered by  $\mathcal{H}$ .

**Theorem 18 (Existence of Realizable PAC-learners Hanneke (2016); Simon (2015); Blumer et al. (1989))**

Let  $\mathcal{H}$  be a hypothesis class with bounded VC-dimension. Then  $\mathcal{H}$  is PAC-learnable in the realizable setting using  $O\left(\frac{\text{VC}(\mathcal{H}) + \log(1/\delta)}{\epsilon}\right)$  samples.

**Theorem 19 (Existence of Agnostic PAC-learners Haussler (1992))** Let  $\mathcal{H}$  be a hypothesis class with bounded VC-dimension. Then  $\mathcal{H}$  is PAC-learnable in the agnostic setting using  $O\left(\frac{\text{VC}(\mathcal{H}) + \log(1/\delta)}{\epsilon^2}\right)$  samples.

## Appendix B. Metric spaces

**Definition 20** A metric space  $(X, \text{dist})$  is called a doubling metric if there exists a constant  $M$  such that every ball of radius  $r$  in it can be covered by at most  $M$  balls of radius  $r/2$ . The quantity  $\log_2 M$  is called the doubling dimension.

**Definition 21** For a metric space  $(X, \text{dist})$ , a measure  $\mu$  defined on  $X$  is called a doubling measure if there exists a constant  $C$ , such that for all  $x \in X$  and  $r \in \mathbb{R}^+$ , we have that  $0 < \mu(\mathcal{B}_{2r}(x)) \leq C \cdot \mu(\mathcal{B}_r(x)) < \infty$ . In this case,  $\mu$  is called  $C$ -doubling.

It can be shown (Luukkainen and Saksman, 1998) that every complete metric with doubling dimension  $d$  has a  $C$ -doubling measure  $\mu$  for some  $C \leq 2^{cd}$  where  $c$  is a universal constant. For example, Euclidean spaces with an  $\ell_p$  distance metric are complete and the Lesbesgue measure is a doubling measure.

The following lemmas follow straightforwardly from the definitions doubling metric and measures.

**Lemma 22** Let  $(X, \text{dist})$  be a doubling metric equipped with a  $C$ -doubling measure  $\mu$ . Then for all  $x \in X$ ,  $r > 0$ , and  $\alpha > 1$ , we have that  $\mu(\mathcal{B}_{\alpha r}(x)) \leq C^{\lceil \log_2 \alpha \rceil} \cdot \mu(\mathcal{B}_r(x))$

**Proof** Since  $\mu$  is a measure, if  $B, B' \subseteq X$  such that  $B \subseteq B'$ , then  $\mu(B) \leq \mu(B')$ . Let  $R = 2^{\lceil \log_2 \alpha \rceil}$ . It's clear that  $R \geq \alpha$ . Therefore  $\mathcal{B}_{\alpha r}(x) \subseteq \mathcal{B}_R(x)$ . Expanding  $\mathcal{B}_r(x)$  by a factor of two  $\lceil \log_2 \alpha \rceil$  times, we get  $\mathcal{B}_R(x)$ , which means  $\mu(\mathcal{B}_R(x)) \leq C^{\lceil \log_2 \alpha \rceil} \cdot \mu(\mathcal{B}_r(x))$ . But since  $\mathcal{B}_{\alpha r}(x) \subseteq \mathcal{B}_R(x)$ , we get the desired result. ■

**Lemma 23** Let  $(X, \text{dist})$  be a doubling metric equipped with a  $C$ -doubling measure  $\mu$ . Let  $x, x' \in X$ ,  $r > 0$ , and  $\alpha > 1$  be such that  $\mathcal{B}_r(x') \subseteq \mathcal{B}_{\alpha r}(x)$ . Then  $\mu(\mathcal{B}_{\alpha r}(x)) \leq C^{\lceil \log_2(2\alpha) \rceil} \cdot \mu(\mathcal{B}_r(x'))$ .

**Proof** By Lemma 22, all we need to show is that  $B_{\alpha r}(x) \subseteq \mathcal{B}_{2\alpha r}(x')$ . Indeed, let  $y \in \mathcal{B}_{\alpha r}(x)$  be any point. Then, from triangle inequality, we have that

$$\begin{aligned} d(x', y) &\leq d(x, x') + d(x, y) \\ &\leq d(x, x') + \alpha r \end{aligned}$$

Moreover, since  $x' \in \mathcal{B}_{\alpha r}(x)$ , we have that  $d(x, x') \leq \alpha r$ . Substituting into the equation above, we get  $d(x', y) \leq 2\alpha r$ , which means  $y \in \mathcal{B}_{2\alpha r}(x')$ .  $\blacksquare$

Finally, we also get:

**Lemma 24** *For any family  $\mathcal{M}$  of complete, doubling metric spaces, there exist constants  $c_1, c_2 > 0$  such that for any metric space  $(X, \text{dist}) \in \mathcal{M}$  with doubling dimension  $d$ , there exists a measure  $\mu$  such that if a ball  $\mathcal{B}_r$  of radius  $r > 0$  is completely contained inside a ball  $\mathcal{B}_{\alpha r}$  of radius  $\alpha r$  (with potentially a different center) for any  $\alpha > 1$ , then  $0 < \mu(\mathcal{B}_{\alpha r}) \leq (c_1 \alpha)^{c_2 d} \mu(\mathcal{B}_r)$ .*

**Proof** We prove this when  $\mathcal{M}$  is the set of all complete, doubling metric spaces employing Lemmas 22 and 23, that can be found in Appendix, part B. We have that  $C^{\lceil \log_2(2\alpha) \rceil} \leq (2\alpha)^{2 \log_2 C}$ . Since  $\log_2 C \leq cd$ , we get  $(2\alpha)^{2 \log_2 C} \leq (2\alpha)^{cd}$ . Thus  $c_1 = 2$  and  $c_2 = c$ .  $\blacksquare$

**Corollary 25** *Suppose we have a constant  $\alpha_0 > 1$  such that we know that  $\alpha \geq \alpha_0$ . Then the bound in Lemma 24 can be further simplified to  $0 < \mu(\mathcal{B}_{\alpha r}) \leq \alpha^{\zeta d} \mu(\mathcal{B}_r)$ , where  $\zeta$  depends on  $\mathcal{M}$  and  $\alpha_0$ . Furthermore, if  $c_1 = 1$  then we can set  $\alpha_0 = 1$ .*

**Proof**  $(c_1 \alpha)^{c_2 d} = \alpha^{c_2 d (1 + \log_{\alpha} c_1)} \leq \alpha^{c_2 d (1 + \log_{\alpha_0} c_1)} = \alpha^{\zeta d}$  for  $\zeta = c_2 (1 + \log_{\alpha_0} c_1)$ . If  $c_1 = 1$ , then  $\zeta = c_2$  for all  $\alpha$ .  $\blacksquare$

## Appendix C. Proof of Lemma 11

Let  $X_{\text{err}} = \{z \in \mathcal{B}_{r(1+\gamma)}(x) \mid g(z) \neq y\}$ . Then, we have that  $\Sigma_{g,y}(x) = \mathbb{E}_{z \sim \mathcal{B}_{r(1+\gamma)}(x)} \mathbb{1}[g(z) \neq y] = \frac{\mu(X_{\text{err}})}{\mu(\mathcal{B}_{r(1+\gamma)}(x))}$ . Further, for all  $z \in \mathcal{B}_r(x)$ , we have  $\mathbb{E}_{z' \sim \mathcal{B}_{r\gamma}(z)} \mathbb{1}[g(z') \neq y] = \frac{\mu(X_{\text{err}} \cap \mathcal{B}_{r\gamma}(z))}{\mu(\mathcal{B}_{r\gamma}(z))}$ .

Let  $z \in \mathcal{B}_r(x)$ . Since this implies that  $\mathcal{B}_{r\gamma}(z) \subseteq \mathcal{B}_{r(1+\gamma)}(x)$ , the worst case happens when  $X_{\text{err}} \subseteq \mathcal{B}_{r\gamma}(z)$ . Therefore,

$$\begin{aligned} \sigma_{g,y}(x) &= \mathbb{E}_{z' \sim \mathcal{B}_{r\gamma}(z)} \mathbb{1}[g(z') \neq y] \\ &= \frac{\mu(X_{\text{err}} \cap \mathcal{B}_{r\gamma}(z))}{\mu(\mathcal{B}_{r\gamma}(z))} \\ &\leq \frac{\mu(X_{\text{err}})}{\mu(\mathcal{B}_{r\gamma}(z))} \\ &\leq \frac{\Sigma_{g,y}(x) \cdot \mu(\mathcal{B}_{r(1+\gamma)}(x))}{\mu(\mathcal{B}_{r\gamma}(z))} \\ &\leq \Sigma_{g,y}(x) \cdot \left(\frac{1+\gamma}{\gamma}\right)^{\zeta d}, \end{aligned} \tag{3}$$

where the last inequality is implied by Lemma 24. Thus,  $\Sigma(x) \leq \frac{1}{3} \cdot \left(\frac{1+\gamma}{\gamma}\right)^{-\zeta d}$  implies that  $\sigma(z) \leq 1/3$  as claimed.

## Appendix D. Compression based bounds

### D.1. Proof of Lemma 13

To prove the generalization bound for tolerant learning, we employ the following lemma that establishes generalization for compression schemes for adversarial losses:

**Lemma 26 (Lemma 11, (Montasser et al., 2019))** *For any  $k \in \mathbb{N}$  and fixed function  $\rho : \bigcup_{i=1}^k (X \times Y)^i \rightarrow Y^X$ , for any distribution  $P$  over  $X \times Y$  and any  $m \in \mathbb{N}$ , with probability at least  $(1 - \delta)$  over an i.i.d. sample  $S = ((x_1, y_1), (x_2, y_2), \dots, (x_m, y_m))$ : if there exist indices  $i_1, i_2, \dots, i_k$  such that*

$$\mathcal{L}_S^{\mathcal{U}}(\rho((x_{i_1}, y_{i_1}), (x_{i_2}, y_{i_2}), \dots, (x_{i_k}, y_{i_k}))) = 0$$

then the robust loss of the decompression with respect to  $P$  is bounded by

$$\begin{aligned} \mathcal{L}_P^{\mathcal{U}}(\rho((x_{i_1}, y_{i_1}), (x_{i_2}, y_{i_2}), \dots, (x_{i_k}, y_{i_k}))) \\ \leq \frac{1}{m - k} (k \ln(m) + \ln(1/\delta)) \end{aligned}$$

The above lemma implies that if  $(\kappa, \rho)$  is a compression scheme that compresses data-sets of size  $m$  to at most  $k \ln(m)$  data points, for class  $\mathcal{H}$  and robust loss  $\ell^{\mathcal{U}}$ , then the sample complexity (omitting logarithmic factors) of robustly learning  $\mathcal{H}$  in the realizable case is bounded by

$$m(\epsilon, \delta) = \tilde{O} \left( \frac{k + \ln(1/\delta)}{\epsilon} \right)$$

For the tolerant setting, since every sample that is realizable with respect to  $\ell^{\mathcal{V}}$  is also realizable with respect to  $\ell^{\mathcal{U}}$ , if a  $(\mathcal{U}, \mathcal{V})$ -tolerant compression scheme compresses to at most  $k \ln(m)$  data-points and decompresses all  $\ell^{\mathcal{V}}$ -realizable samples  $S$  to functions that have  $\ell^{\mathcal{U}}$ -loss 0 on  $S$ , then the lemma implies the above bound for the  $(\mathcal{U}, \mathcal{V})$ -tolerant sample complexity of learning  $\mathcal{H}$ .

### D.2. Proof of Lemma 14

The proof of this Lemma employs the notions of a sample being  $\epsilon$ -net or an  $\epsilon$ -approximation for a hypothesis class  $\mathcal{H}$ . A labeled data set  $S = ((x_1, y_1), (x_2, y_2), \dots, (x_m, y_m))$  is an  $\epsilon$ -net for class  $\mathcal{H}$  with respect to distribution  $P$  over  $X \times Y$  if for every hypothesis  $h \in \mathcal{H}$  with  $\mathcal{L}_P^{0/1}(h) \geq \epsilon$ , there exists an index  $j$  and  $(x_j, y_j) \in S$  with  $h(x_j) \neq y_j$ .  $S$  is an  $\epsilon$ -approximation for class  $\mathcal{H}$  with respect to distribution  $P$  over  $X \times Y$  if for every hypothesis  $h \in \mathcal{H}$  we have  $|\mathcal{L}_S^{0/1}(h) - \mathcal{L}_P^{0/1}(h)| \leq \epsilon$ . Standard VC-theory tells us that, for classes with bounded VC-dimension, sufficiently large samples from  $P$  are  $\epsilon$ -nets or  $\epsilon$ -approximations with high probability (Haussler and Welzl, 1987).

**Proof** We will employ a boosting-based approach to establish the claimed compression sizes. Let  $S = ((x_1, y_1), (x_2, y_2), \dots, (x_m, y_m))$  be a data-set that is  $\ell^{\mathcal{V}}$ -realizable with respect to  $\mathcal{H}$ . We let  $S_{\mathcal{V}}$  denote an “inflated data-set” that contains all domain points in the perturbation sets of the  $x_i \in S^X$ , that is

$$S_{\mathcal{V}}^X := \bigcup_{i=1}^m \mathcal{V}(x_i)$$

Every point  $z \in S_{\mathcal{V}}^X$  is assigned the label  $y = y_i$  of the minimally-indexed  $(x_i, y_i) \in S$  with  $z \in \mathcal{V}(x_i)$ , and we set  $S_{\mathcal{V}}$  to be the resulting collection of labeled data-points. (Note that since the

sample  $S$  is assumed to be  $\ell^{\mathcal{V}}$ -realizable, assigning it the label of some other corresponding data point in case  $z \in \mathcal{V}(x_i) \cap \mathcal{V}(x_j)$  for  $x_i \neq x_j$ , would not induce any inconsistencies). Now let  $D$  be the probability measure over  $S_{\mathcal{V}}^X$  defined by first sampling an index  $j$  uniformly from  $[j] = \{1, 2, \dots, j\}$  and then sampling a domain point  $z \sim \mathcal{V}(x_j)$  from the  $\mathcal{V}$ -perturbation set around the  $j$ -th sample point in  $S$ . Note that this implies that if  $D(B) \leq (\beta/m)$  for some subset  $B \subseteq S_{\mathcal{V}}^X$ , then

$$\mathbb{P}_{z \sim \mathcal{V}(x)}[z \in B] \leq \beta \quad (4)$$

for all  $x \in S^X$ .

We will now show that, by means of a compression scheme, we can encode a hypothesis  $g$  with binary loss

$$\mathcal{L}_D^{0/1}(g) \leq \beta/m. \quad (5)$$

Property 1 together with Equation 4 then implies that the resulting  $\mathcal{W}$ -smoothed function  $\bar{g}$  has  $\mathcal{U}$ -robust loss 0 on the sample  $S$ ,  $\mathcal{L}_S^{\mathcal{U}}(\bar{g}) = 0$ . Since the smoothing is a deterministic operation once  $g$  is fixed, this implies the existence of a  $(\mathcal{U}, \mathcal{V})$ -tolerant compression scheme.

Standard VC-theory tells us that, for a class  $G$  of bounded VC-dimension, for any distribution over  $X \times Y$ , and any  $\epsilon, \delta > 0$ , with probability at least  $(1 - \delta)$  an i.i.d. sample of size  $\Theta\left(\frac{\text{VC}(G) + \ln(1/\delta)}{\epsilon^2}\right)$  is an  $\epsilon$ -approximation for the class  $G$  (Haussler and Welzl, 1987). This implies in particular, that there exists a finite subset  $S_{\mathcal{V}}^f \subset S_{\mathcal{V}}$  of size at most  $\frac{4m^2 C \cdot \text{VC}(G)}{\beta^2}$  (for some constant  $C$ ) with the property that any classifier  $g \in G$  with empirical (binary) loss at most  $\beta/2m$  on  $S_{\mathcal{V}}^f$  has loss  $\mathcal{L}_D^{0/1}(g) \leq \beta/m$  with respect to the distribution  $D$ . We will choose such a set  $S_{\mathcal{V}}^f$  for the class  $G$  of  $T$ -majority votes over  $\mathcal{H}$  for  $T = 18 \ln(\frac{2m}{\beta})$ . That is

$$G = \{g \in Y^X \mid \exists h_1, h_2, \dots, h_T \in \mathcal{H} : \\ g(x) = \mathbb{1}[\sum_{i=1}^T h_i(x) \geq 1/2]\}$$

The VC-dimension of  $G$  is bounded by (Shalev-Shwartz and Ben-David, 2014)

$$\text{VC}(G) = \mathcal{O}(T \cdot \text{VC}(\mathcal{H}) \log(T \text{VC}(\mathcal{H}))) = \mathcal{O}(18 \ln(\frac{2m}{\beta}) \text{VC}(\mathcal{H}) \log(18 \ln(\frac{2m}{\beta}) \text{VC}(\mathcal{H}))).$$

We will now show how to obtain the classifier  $g$  by means of a boosting approach on the finite data-set  $S_{\mathcal{V}}^f$ . More specifically, we will use the boost-by-majority method. This method outputs a  $T$ -majority vote  $g(x) = \mathbb{1}[\sum_{i=1}^T h_i(x) \geq 1/2]$  over weak learners  $h_i$ , which in our case will be hypotheses from  $\mathcal{H}$ . After  $T$  iterations with  $\gamma$ -weak learners, the empirical loss over the sample  $S_{\mathcal{V}}^f$  is bounded by  $e^{-2\gamma^2 T}$  (see Section 13.1 in (Schapire and Freund, 2013)). Thus, with  $\gamma = 1/6$ , and  $T = 18 \ln(\frac{2m}{\beta})$ , we obtain

$$\mathcal{L}_{S_{\mathcal{V}}^f}^{0/1}(g) \leq \frac{\beta}{2m}$$

which, by the choice of  $S_{\mathcal{V}}^f$  implies

$$\mathcal{L}_D^{0/1}(g) \leq \beta/m$$

which is what we needed to show according to Equation 5.

It remains to argue that the weak learners to be employed in the boosting procedure can be encoded by a small number of sample points from the original sample  $S$ . For this part, we will

employ a technique introduced earlier for robust compression (Montasser et al., 2019). Recall that the set  $S$  is  $\mathcal{V}$ -robustly realizable, which implies that the set  $S_{\mathcal{V}}^f$  is (binary loss-) realizable by  $\mathcal{H}$ . By standard VC-theory, for every distribution  $D_i$  over  $S_{\mathcal{V}}^f$ , there exists an  $\epsilon$ -net of size  $\mathcal{O}(\text{VC}(\mathcal{H})/\epsilon)$  (Haussler and Welzl, 1987). Thus, for every distribution  $D_i$  over  $S_{\mathcal{V}}^f$  (that may occur during the boosting procedure), there exists a subsample  $S_i$  of  $S_{\mathcal{V}}^f$ , of size at most  $n = \mathcal{O}(3\text{VC}(\mathcal{H}))$  with the property that every hypothesis from  $\mathcal{H}$  that is consistent with  $S_i$  has binary loss at most  $1/3$  with respect to  $D_i$  (thus can serve as a weak learner for margin  $\gamma = 1/6$  in the above procedure). Now for every labeled point  $(x, y) \in S_i$ , there is a sample point  $(x_j, y_j) \in S$  in the original sample  $S$  such that  $x \in \mathcal{V}(x_j)$  and  $y = y_j$ . Let  $S'_i$  be the collection of these corresponding original sample points. Note that any hypothesis  $h \in \mathcal{H}$  that is  $\mathcal{V}$ -robustly consistent with  $S'_i$  is consistent with  $S_i$ . Therefore we can use the  $n$  original data-points in  $S'_i$  to encode the weak learner  $h_i$  (for the decoding any  $\mathcal{V}$ -robust ERM hypothesis can be chosen to obtain  $h_i$ ).

To summarize, we will compress the sample  $S$  to the sequence  $S'_1, S'_2, \dots, S'_T$  of  $n \cdot T = \mathcal{O}(\text{VC}(\mathcal{H}) \ln(\frac{m}{\beta}))$  sample points from  $S$ . To decode, we obtain the function  $g$  as a majority vote over the weak learner  $h_i$  and proceed to obtain the  $\mathcal{W}$ -smoothed function  $\bar{g}$ . This function  $\bar{g}$  satisfies  $\mathcal{L}_S^{\mathcal{U}}(\bar{g}) = 0$  and by this we have established the existence of a  $\mathcal{U}, \mathcal{V}$ -tolerant compression scheme of size  $\mathcal{O}(\text{VC}(\mathcal{H}) \ln(\frac{m}{\beta}))$  as claimed. ■