Adversarial VC-dimension and Sample Complexity of Neural Networks

Zetong Qi

Department of Electrical and Computer Engineering University of Wisconsin - Madison zqi38@wisc.edu

T.J. Wilder

Department of Computer Science University of Wisconsin - Madison twilder@wisc.edu

Abstract

Adversarial attacks during the testing phase of neural networks pose a challenge for the deployment of neural networks in security critical settings. These attacks can be performed by adding noise that is imperceptible to humans on top of the original data. By doing so, an attacker can create an adversarial sample, which will cause neural networks to misclassify. In this paper, we seek to understand the theoretical limits of what can be learned by neural networks in the presence of an adversary. We first defined the hypothesis space of a neural network, and showed the relationship between the growth number of the entire neural network and the growth number of each neuron. Combine that with the adversarial Vapnik-Chervonenkis(VC)-dimension of halfspace classifiers, we concluded the adversarial VC-dimension of the neural networks with sign activation functions.

1 Introduction

Machine learning has become the fastest growing area of computer science, and neural networks are among the most studied among all ML algorithms because of their impressive performance in areas like image recognition, natural language processing, etc. However, practical neural networks are often vulnerable to adversarial attacks: given an input x and any target label t, it is possible to find an x' that is very similar to x but which neural networks will misclassify as the target label t. This makes it difficult to apply neural networks in security critical areas like self-driving cars, malware detection systems, facial recognition systems to unlock devices, etc. Carlini et al. created the state of the art attacks on image recognition systems [1] and speech recognition systems [2]. Some have proposed methods, like adversarial training [3], for creating neural networks that are robust against adversarial attacks.

In this paper, we aim to understand more about neural networks in the presence of a test time adversary. In particular, we try to see if the sample complexity of neural networks is different in the presence of adversaries. We show that, for some simple networks, the sample complexity bound does not change with the introduction of an adversary. We prove this for neural networks with sign activation functions by combining known bounds on sample complexity with recent bounds on the adversarial VC-dimension of halfspace classifiers from Cullina et al. [4]

We first describe the adversarial PAC-learning framework. After that, we formally specify our neural network. Finally we prove that the VC-dimension bound is unchanged by the introduction of an adversary.

1.1 PAC-learning in the presence of an adversary

In this section, we setup the framework for PAC-learning in the presence of an evasion adversary, which is one that generates and presents the learner with adversarial examples during the test phase but does not interfere with the training.

In our setup, there is an unknown distribution $\mathcal{P}_{\mathcal{X} \times \mathcal{Y}}$. The learner receives training data $S = ((x_1, y_1), \dots, (x_n, y_n)) \sim \mathcal{P}^n_{\mathcal{X} \times \mathcal{Y}}$ and outputs $\hat{h} \in \mathcal{H}$. The adversary receives data $(x, y) \in \mathcal{P}_{\mathcal{X} \times \mathcal{Y}}$ and outputs some $x' \in N(x)$, where N(x) defines a neighborhood around x. The neighborhood of x is defined as $N(x) = \{x' \in \mathcal{X} : R(x, x') \leq \epsilon\}$ where R(x, x') defines a nearness relationship.

Under this new framework, the adversarial expected risk, L_P , is defined as the learner's risk under the true distribution $\mathcal{P}_{\mathcal{X} \times \mathcal{Y}}$ in the presence of an adversary constrained by the nearness relation R, and L_S is defined as the adversarial empirical risk with the same adversary

$$L_P(h, R) = \mathbb{E}_{(x,y) \sim \mathcal{P}_{\mathcal{X} \times \mathcal{Y}}} \left[\max_{x' \in N(x)} \ell(h(x'), y) \right]$$

$$L_S(h, R) = \frac{1}{n} \sum_{i=1}^{n} \left[\max_{x' \in N(x_i)} \ell(h(x'), y_i) \right]$$

The goal of the learner is to minimize the adversarial expected risk. Because calculating the expected risk over the distribution $\mathcal{P}_{\mathcal{X} \times \mathcal{Y}}$ is infeasible, the learner instead aims to minimize the adversarial empirical risk. The algorithm that selects the hypothesis \hat{h} that minimizes the adversarial empirical risk from the hypothesis space \mathcal{H} , with nearness relation R, is called Adversarial Empirical Risk Minimization and we define the objective as follows

$$AERM_{\mathcal{H},R}(S) = \operatorname{argmin}_{h \in \mathcal{H}} L_S(h,R)$$

1.2 Corrupted hypotheses

In this section, we describe the concept of a corrupted hypothesis class. The presence of an adversary forces the learner to learn with corrupted hypothesis. Instead of just predicting ± 1 , corrupted hypothesis also outputs \perp that means "always wrong".

For $\mathcal{Y} = \{-1, 1\}$, let the corrupted output space be $\widetilde{\mathcal{Y}} = \{-1, 1, \bot\}$. Now we can define the corruption function as a mapping from a hypothesis to the corrupted version $\kappa_R : (\mathcal{X} \mapsto \mathcal{Y}) \mapsto (\mathcal{X} \mapsto \widetilde{\mathcal{Y}})$:

$$\kappa_R(h) = x \to \begin{cases} -1 & \forall x' \in N(x) : h(x') = -1\\ 1 & \forall x' \in N(x) : h(x') = 1\\ \bot & \exists x'_0, x'_1 \in N(x) : h(x'_0) \neq h(x'_1) \end{cases}$$

Using this, we can define the set of corrupted hypotheses as $\widetilde{H} = \{\kappa_R(h) : h \in \mathcal{H}\}$. In other words, the hypothesis will predict "always wrong" when the test data x is in the neighborhood where different labels exist.

We also define the loss function λ , and the loss classes $\mathcal F$ and $\widetilde{\mathcal F}$ which are derived from $\mathcal H$ and $\widetilde H$ respectively

$$\lambda(h) = \mathcal{X} \times \widetilde{\mathcal{Y}} \to \{0, 1\}$$

$$\mathcal{F} = \{\lambda(h) : h \in \mathcal{H}\}$$

$$\widetilde{F} = \{\lambda(\widetilde{h}) : \widetilde{h} \in \widetilde{\mathcal{H}}\}$$

We can define the equivalent shattering coefficient in terms of the loss class as

$$\sigma'(\mathcal{F}, i) = \max_{(x', y) \in \mathcal{X}^i \times \mathcal{Y}^i} |\{(f(x_1', y_1), \dots, f(x_i', y_i)) : f \in \mathcal{F}\}|$$

Finally, the Adversarial VC-dimension is defined as

$$AVC(\mathcal{H}, R) = \sup\{n \in \mathbb{N} : \sigma'(\lambda(\widetilde{\mathcal{H}}), n) = 2^n\}$$

1.3 Adversarial VC-dimension of halfspace classifiers

For the sake of simplicity, we assume that the adversary has standard ℓ_p norm-based constraints that are usually imposed on evasion adversaries as described in the literature [1, 3]. As it turns out, the adversarial VC-dimension for halfspace classifiers corrupted by ℓ_p norm-constrained adversary is equal to the standard VC-dimension [4]. That is: Let \mathcal{H} be the family of halfspace classifiers of $\mathcal{X} \in \mathbb{R}^d$. Then the adversarial VC-dimension of \mathcal{H} in the presence of an adversary with ℓ_p norm-based constraints is $AVC(\mathcal{H},R)=d+1$.

2 Adversarial VC-dimension of Neural Networks with Sign Activation Function

2.1 Define the Neural Network

We define a general neural network described by a directed acyclic graph G=(V,E), where all neurons have the same activation function $\sigma(a)$. In a neural network of depth T, let V_0,\ldots,V_T be the layers of the neural network and let $E_{(t-1,t)}$ be the weights connecting the layers V_{t-1} and V_t . Any layer V_t has $|V_t|$ number of neurons. We can express our neural network's overall hypothesis space as a composition of the hypothesis spaces of each layer, $\mathcal{H} = \mathcal{H}^{(T)} \circ \mathcal{H}^{(T-1)} \circ \cdots \circ \mathcal{H}^{(1)}$ where $\mathcal{H}^{(t)} = \{f: \mathbb{R}^{|V_{t-1}|} \mapsto \mathbb{R}^{|V_t|}\}$. We analyze the adversarial VC-dimension and sample complexity for this family of hypotheses in the event of an evasion attack.

2.2 Neural Networks with Sign Activation Function

To simplify our hypothesis space, we choose the activation function to be $\sigma(a) = \mathbb{1}_{[a>0]}$, and let $\mathcal{H}^{(t)} = \{f : \mathbb{R}^{|V_{t-1}|} \mapsto \{\pm 1\}^{|V_t|}\}$. With this activation function, each neuron in each layer turns into a halfspace classifier:

$$\forall t \in [T], V_t = sign(\langle E_{(t-1,t)}, V_{t-1} \rangle)$$

Cullina et al. showed that the adversarial VC dimension of a halfspace classifier for $\mathcal{X}=\mathbb{R}^d$ is d+1 [4]. In the previous section, we showed that the hypothesis class of a neural network \mathcal{H} can be written as a composition of its layers, $\mathcal{H}=\mathcal{H}^{(T)}\circ\cdots\circ\mathcal{H}^{(1)}$. The following lemma shows that the growth function of a composition of hypothesis classes is bounded by the products of the growth functions of the individual classes.

Lemma 1 Let \mathcal{F}_1 be a set of functions from \mathcal{X} to \mathcal{Z} and let \mathcal{F}_2 be a set of functions from \mathcal{Z} to \mathcal{Y} . Let $\mathcal{H} = \mathcal{F}_2 \circ \mathcal{F}_1$ be the composition class. That is, $\forall f_1 \in \mathcal{F}_1$ and $f_2 \in \mathcal{F}_2 : \exists h \in \mathcal{H} \text{ s.t. } h(x) = f_2(f_1(x))$. The growth function of \mathcal{H} , $\tau_{\mathcal{H}}(m)$, is bounded by $\tau_{\mathcal{H}}(m) \leq \tau_{\mathcal{F}_1}(m)\tau_{\mathcal{F}_2}(m)$

$$\begin{aligned} \textit{Proof:} \\ \text{Let } C &= \{c_1, \dots, c_m\} \subseteq \mathcal{X}. \text{ Then} \\ &|\mathcal{H}_C| = |\{f_2(f_1(c_1)), \dots, f_2(f_1(c_m)) : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}| \\ &= \left| \bigcup_{f_1 \in \mathcal{F}_1} \{(f_2(f_1(c_1)), \dots, f_2(f_1(c_m))) : f_2 \in \mathcal{F}_2\} \right| \\ &\leq |\mathcal{F}_{1C}| \cdot \tau_{\mathcal{F}_2}(m) \\ &\leq \tau_{\mathcal{F}_1}(m) \tau_{\mathcal{F}_2}(m) \\ &\tau_{\mathcal{H}}(m) \leq \tau_{\mathcal{F}_1}(m) \tau_{\mathcal{F}_2}(m) \end{aligned}$$

Therefore, the growth function of the hypothesis space of neural networks is bounded by:

$$\tau_{\mathcal{H}}(m) \leq \prod_{t=1}^{T} \tau_{\mathcal{H}^{(t)}(m)}$$

In addition, each $\mathcal{H}^{(t)}$ can written as a product of individual neurons, $\mathcal{H}^{(t)} = \mathcal{H}^{(t,1)} \times \cdots \times \mathcal{H}^{(t,|V_t|)}$, where each $\mathcal{H}^{(t,j)}$ is a halfspace classifier: $\mathcal{H}^{(t,j)} = \{f : \mathbb{R}^{|V_{t-1}|} \mapsto \{\pm 1\}\}$. The following lemma shows that the growth function of the Cartesian product class is bounded by the products of the growth functions of the individual classes:

Lemma 2 For i=1,2, let \mathcal{F}_i be a set of functions from \mathcal{X} to \mathcal{Y}_i . Define $\mathcal{H}=\mathcal{F}_1\times\mathcal{F}_2$ to be the Cartesian product class. That is, $\forall f_1\in\mathcal{F}_1$ and $f_2\in\mathcal{F}_2:\exists h\in\mathcal{H}$ s.t. $h(x)=(f_1(x),f_2(x))$. The growth function of \mathcal{H} is bounded by: $\tau_{\mathcal{H}}(m)\leq\tau_{\mathcal{F}_1}(m)\tau_{\mathcal{F}_2}(m)$

$$\begin{split} & \textit{Proof:} \\ & \text{Let } C = \{c_1, \dots, c_m\} \subseteq \mathcal{X}. \text{ Then} \\ & |\mathcal{H}_C| = |\{((f_1(c_1), f_2(c_1)), \dots, (f_1(c_m), f_2(c_m))) : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}| \\ & = |\{((f_1(c_1), \dots, f_1(c_m)), (f_2(c_1), \dots, f_2(c_m))) : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}| \\ & = |\mathcal{F}_{1C} \times \mathcal{F}_{2C}| \\ & = |\mathcal{F}_{1C}| \cdot |\mathcal{F}_{2C}| \\ & \tau_{\mathcal{H}}(m) \leq \tau_{\mathcal{F}_1}(m)\tau_{\mathcal{F}_2}(m) \end{split}$$

Therefore, the growth function of each layer can be bounded by:

$$au_{\mathcal{H}^{(t)}}(m) \leq \prod_{i=1}^{|V_t|} au_{\mathcal{H}^{(t,i)}}(m)$$

Combining these two lemmas, the growth function of the entire neural network is bounded by:

$$\tau_{\mathcal{H}}(m) \leq \prod_{t=1}^{T} \prod_{i=1}^{|V_t|} \tau_{\mathcal{H}^{(t,i)}}(m)$$

Since the adversaries are only able to change the inputs, we only corrupt the first layer by introducing an adversary. However, since the adversarial VC-dimension of halfspace classifiers of dimension d is d+1, the same as the regular VC-dimension, we can say that the ith neuron in all hidden layers have an effective VC-dimension of $d_{t,i}$, where $d_{t,i}$ is the number of edges that are going into the ith neuron of the tth layer, assuming that one edge accounts for the bias term.

Using this alongside our growth function bound, we can use Sauer's lemma to show that:

$$\tau_{\mathcal{H}}(m) \le \prod_{t=1}^{T} \prod_{i=1}^{|V_t|} \left(\frac{em}{d_{t,i}}\right)^{d_{t,i}} \le \prod_{t=1}^{T} \prod_{i=1}^{|V_t|} (em)^{d_{t,i}} = (em)^{\left(\sum_{t=1}^{T} \sum_{i=1}^{|V_t|} d_{t,i}\right)}$$

Notice that $\sum_{t=1}^{T} \sum_{i=1}^{|V_t|} d_{t,i}$ is just the number of edges in the neural network, therefore we have:

$$\tau_{\mathcal{H}}(m) \le (em)^{|E|}$$

Let there be a set of size m that is shattered by the neural network. Therefore the growth number $\tau_{\mathcal{H}}(m) = 2^m$. Combining this with $\tau_{\mathcal{H}}(m) \leq (em)^{|E|}$, we have that:

$$2^m \le (em)^{|E|}$$

The following lemma shows that m must be $\mathcal{O}(|E|\log(|E|))$ in order to satisfy the inequality.

Lemma 3 If a neural network's hypothesis space \mathcal{H} has |E| number of parameters, let m be the size of the set \mathcal{H} shatters. If the inequality $2^m \leq (em)^{|E|}$ is satisfied, then \mathcal{H} has sample complexity of $\mathcal{O}(|E|\log(|E|))$.

Proof:

$$2^{m} \le (em)^{|E|}$$

$$m \log(2) \le |E| \log(em)$$

$$m \le \frac{|E|}{\log(2)} \log(em)$$

$$em \le \frac{e|E|}{\log(2)} \log(em)$$

Lemma A.1 in [5] states: "for $a \ge 0$, if $x \ge 2a \log(a)$ then $x \ge a \log(x)$ ". We can take the contrapositive to get: "if $x < a \log(x)$ then $x < 2a \log(a)$ " and apply it

$$em < \frac{2e|E|}{\log(2)} \log\left(\frac{e|E|}{\log(2)}\right)$$
$$m < \frac{2|E|}{\log(2)} \log\left(\frac{e|E|}{\log(2)}\right)$$

Ignoring the constants, we have:s

$$m \leq \mathcal{O}(|E|\log|E|)$$

Putting all of these pieces together, we now have a bound on the adversarial VC-dimension for our network.

Theorem 1 For a neural network with sign activation functions, the adversarial VC-dimension has the same bound as the regular VC-dimension.

Proof:

From Lemma 3, we can see that the AVC-dimension of neural networks with sign activation function is $\mathcal{O}(|E|\log|E|)$. We know from Theorem 20.6 in [5] that the regular VC-dimension of neural networks with sign activation function is also $\mathcal{O}(|E|\log|E|)$.

3 Extending corrupted hypotheses

One big issue with the corrupted hypothesis class as defined by Cullina et al. is that it limits the expressiveness of the hypotheses by requiring ± 1 output [4]. In order to overcome this issue, we introduce two generalizations of the corrupted hypothesis class. The first is a generalization to multi-class corrupted hypotheses.

$$\kappa_R(h) = x \to \begin{cases} h(x) & \forall x' \in N(x) : h(x) = h(x') \\ \bot & \exists x'_0, x'_1 \in N(x) : h(x'_0) \neq h(x'_1) \end{cases}$$

This is more of a generalization in notation than anything else because we simply output the normal label of x for the hypothesis if it agrees on all the neighbors of x, and reject it if it doesn't. In the ± 1 case, this version simplifies down to the original form of the equation.

However, this version is still impractical for problems without a reasonable number of outputs or for any problem which has more complex relationships between the outputs. To solve these problems, we introduce our second new form of the corrupted hypothesis class which we call the continuous corruption class.

$$\kappa_{R,d}(h) = x \rightarrow \begin{cases} h(x) & \forall x' \in N(x) : d(h(x),h(x')) \leq \delta \\ \bot & \exists x_0',x_1' \in N(x) : d(h(x_0'),h(x_1')) > \delta \end{cases}$$

Here, d defines a distance metric on \mathcal{Y} which is analogous to R, the distance metric on \mathcal{X} . Instead of looking for perfect agreement among all the neighbors, we instead allow normal predictions on x,

when any of the neighbors of x would have sufficiently similar outputs. This formulation now allows us to describe the corrupted version of any hypothesis class, instead of only binary ones. It should be possible to theoretically describe how easy it is to corrupt any hypothesis class by comparing the normal and corrupted versions.

This formulation directly relates to our own problem of discovering VC-dimension of neural networks in adversarial environments. The original framework restricted us to ± 1 output for the network and, for our proof of the bound, even for each node of the network. By using the continuous corruption class, we can instead look at adversarial elements of continuous functions, including many other common activation functions such as Sigmoids, Tanh, or ReLU. All of these are actually used in practice, while our version using sign activation is primarily a theoretical model.

Beyond accommodating for continuous functions, this generalization also allows for more interesting relationships between the output. For example, when classifying images for self-driving cars, it is a big problem if an adversary can make your model see a car instead of a human, but probably less of a problem if your model is made to see a car instead of a truck. Obviously there are may be problems either way, but when you have a lot of classes or continuous outputs, then intelligently picking your distance metric could help make the model robust without relying on perfect predictions.

4 Conclusion

We analyzed the adversarial VC-dimension and sample complexity for neural networks with sign activation function and showed that it could achieve the same bound as the non-adversarial case. Though Cullina et. al. did show that the AVC-dimension could be arbitrarily larger or smaller than the ordinary VC-dimension, we have shown that there exist learners, including at least some neural networks, which may be highly resilient to adversarial attacks given sufficient training data. In practice however, we find that neural networks are often highly susceptible to adversarial attacks. This may be because of the number of training samples used in practice is usually significantly less than the sample complexity, and the learned model is overfitting to the limited training set so it doesn't generalize well to the true distribution.

A natural extension to our work, which may help to show if that is true, is to find the adversarial VC-dimension of neural networks with some other activation functions such as Sigmoid, Tanh or ReLU. These are of particular interest because they are some of the most widely used in practice. The difficulty with the analysis of those activation functions arises from the fact that they are continuous and map $\mathbb{R} \mapsto \mathbb{R}$. Because each neuron has real output, it becomes hard to use the same proof techniques, because we cannot look at the AVC-dimension for a single neuron. That being said, there are many other methods which could potentially be used to prove these bounds. These include existing methods which bound the sample complexity for neural networks with sigmoid activation functions. For some of these functions, we may also be able to look at an adversarial Rademacher complexity and bound the sample complexity that way. We leave the exploration of these methods to future work.

References

- [1] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. *CoRR*, abs/1608.04644, 2016.
- [2] Hadi Abdullah, Washington Garcia, Christian Peeters, Patrick Traynor, Kevin R. B. Butler, and Joseph Wilson. Practical hidden voice attacks against speech and speaker recognition systems. *CoRR*, abs/1904.05734, 2019.
- [3] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples, 2014.
- [4] Daniel Cullina, Arjun Nitin Bhagoji, and Prateek Mittal. Pac-learning in the presence of evasion adversaries, 2018.
- [5] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, New York, NY, USA, 2014.