



SAPIENZA
UNIVERSITÀ DI ROMA

Robustness Of Deep Neural Networks Using Trainable Activation Functions

Computer Science - Informatica LM-18

Corso di Laurea Magistrale in Computer Science

Candidate

Federico Peconi

ID number 1823570

Thesis Advisor

Prof. Simone Scardapane

Academic Year 2019/2020

Thesis defended on Something October 2020
in front of a Board of Examiners composed by:
Prof. Nome Cognome (chairman)
Prof. Nome Cognome
Dr. Nome Cognome

Robustness Of Deep Neural Networks Using Trainable Activation Functions
Master's thesis. Sapienza – University of Rome

© 2020 Federico Peconi. All rights reserved

This thesis has been typeset by L^AT_EX and the Sapthesis class.

Version: September 4, 2020

Author's email: peconi.1823570@studenti.uniroma1.it

*Dedicated to
Donald Knuth*

Abstract

This document is an example which shows the main features of the $\text{\LaTeX} 2_{\epsilon}$ class `sapthesis.cls` developed by with the help of GuIT (Gruppo Utilizzatori Italiani di \TeX).

Acknowledgments

Ho deciso di scrivere i ringraziamenti in italiano per dimostrare la mia gratitudine verso i membri del GuIT, il Gruppo Utilizzatori Italiani di T_EX, e, in particolare, verso il prof. Enrico Gregorio.

Contents

1	Introduction	1
1.1	Intriguing Properties of Neural Networks	1
1.2	Smooth Activation Functions and Robustness	1
1.3	Structure of the Thesis	1
2	Fundamentals	3
2.1	Deep Neural Networks	3
2.1.1	Definition	3
2.1.2	Training	5
2.1.3	Activation Functions	9
2.1.4	CNNs: Convolutional Neural Networks	10
2.1.5	From Neural Networks to Deep Neural Networks	10
2.2	Adversarial Examples Theory	10
2.3	Defenses	10
2.4	Kernel Based Activation Functions	10
3	Related Works	11
3.1	K-Winners Take All	11
3.2	Smooth Adversarial Training	11
4	Solution Approach	13
4.1	Lipschitz Constant Approach	13
4.2	Fast is Better than Free Adversarial Training	13
5	Evaluation	15
5.1	VGG Inspired Architectures Results	15
5.2	Explofing Gradients with KafNets	15
5.3	ResNet20 Inspired Architectures Results	15
6	Future Works	17
7	Conclusions	19

Chapter 1

Introduction

1.1 Intriguing Properties of Neural Networks

Here we informally state the problem of adversarial attacks in ML models especially wrt to Neural Networks. Why is it of fundamental importance for the progress of the field from both practical (nns cant yet be deployable in critical scenarios for such reasons) and theoretical (Madry arguments around interperatability and robustness) perspectives

1.2 Smooth Activation Functions and Robustness

Recently a link has been proposed between activation functions and the robustness of Neural Networks (Smooth Adversarial Training). In particular, authors showed how they managed to improve the robustness by replacing the traditional Rectified Linear Units activation functions with smoother alternatives such as ELUs, SWISH, PReLU

Building up from this result we thought we could find benefits by leveraging recently proposed smooth trainable activation functions called Kernel Based Activation Functions (Scardapane et al.), which already showed great results in standard tasks, in the context of adversarial attacks.

1.3 Structure of the Thesis

Description of the remaining chapters

Chapter 2

Fundamentals

In this chapter, the basic concepts needed to understand the main arguments for the thesis are introduced. Pointers to more appropriate and detailed resources on the topics are given throughout

2.1 Deep Neural Networks

Broadly speaking, the field of Machine Learning is the summa of any algorithmic methodology whose aim is to automatically find meaningful patterns inside data without being explicitly programmed on how to do it. Well known examples are: Search Trees (ref.), Support Vector Machines (ref.), Clustering (ref.) and, more recently, Neural Networks (ref.). During the last two decades Neural Networks gained a lot of attention for their outstanding performances in different tasks like image classification (ref. ImageNet, over human level), speech and audio processing(ref).

2.1.1 Definition

Neural Networks (NNs) are often used in the context of Supervised Learning where the objective is to model a parametric function $f_\theta: \mathcal{X} \rightarrow \mathcal{Y}$ given n input-output pairs $S = \{(\mathbf{x}_i, y_i)_{i=1}^n\}$ with $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \mathcal{Y}$ such that

$$f_\theta \sim f$$

where f is assumed to be the real input-output distribution that we want to learn. In plain words, this means that we want to find the best set of parameters θ^* for the model such that, for any unseen input x_{new} we have that $f_{\theta^*}(x_{new})$ is as close as possible to $f(x_{new})$

For the sake of explanation, assume the input, which in practice can be very complex and unstructured e.g. made of: graphs, text, sounds, ecc, to be embedded in an input space $\mathcal{X} = \mathbb{R}^d$. The simplest form of a neural network is then given by

$$f_{W,b}(x) = \sigma(Wx + b)$$

where the parameters of the network are the elements of a $u \times d$ matrix W and a u -dimensional vector called bias. The last element applied at the end is the

σ function which consists of a non-linear function acting element-wise and is the key component to introduce non linearity in NNs allowing them to model highly non-linear functions. We call it *activation function*.

$$f_{W,b}(x) = [\sigma(W_1^T x + b_1), \sigma(W_2^T x + b_2), \dots, \sigma(W_u^T x + b_u)]$$

Where W_i and b_i are respectively the i -th row of W and the i -th element of the bias.

Historically, the whole picture was somehow biologically inspired and had an intuitive explanation. Indeed, if we think at W as weights i.e. w_{ij} as the importance the model gives to the input x_i for how much it contributes to the $f_W(x)_j$ -th component and define σ to be

$$\sigma(W_j^T x + b_j) = \begin{cases} 1 & W_j^T x \geq -b_j \\ 0 & W_j^T x < -b_j \end{cases} \quad (2.1)$$

then it is easy to see that here the bias is acting like a threshold which discriminates between *activating* or not the j -th component depending on how much importance was given. Due to this analogy with the behaviour of neurons in the brain we call each component *neuron*, non-linearities activation functions and the whole model neural network.

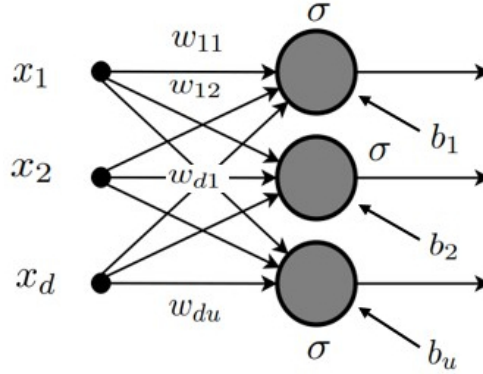


Figure 2.1. Graphic representation of a one layer NN also know as MLP(ref.)

In general, the idea of a layer of neurons can be recursively extended by stacking more layers together, all of which are described by a matrix of weights, a bias and an activation function and letting the output of one becoming the input of the subsequent. The resulting model is the mathematical composition of the layers, thus if we let L be the number of layers, $z_0 = \mathbf{x}$ and $z_l = \sigma_l(W_l z_{l-1} + b_l)$ we write a L -layered $f_{\mathbf{W},\mathbf{b}}$:

$$f_{\mathbf{W},\mathbf{b}} = z_L = \sigma_L(W_L \sigma_{L-1}(\dots(W_2 \sigma_1(W_1 z_0 + b_1) + b_2) \dots) + b_L)$$

With $\mathbf{W} = \{W_1, \dots, W_L\}$ and $\mathbf{b} = \{b_1, \dots, b_L\}$. We will call the first layer *input layer*, the middle layers *hidden layers* and the last layer *output layer*. This general but still basic form of Neural Network is known as *Feed Forward Neural Network* Fig. 2.2.

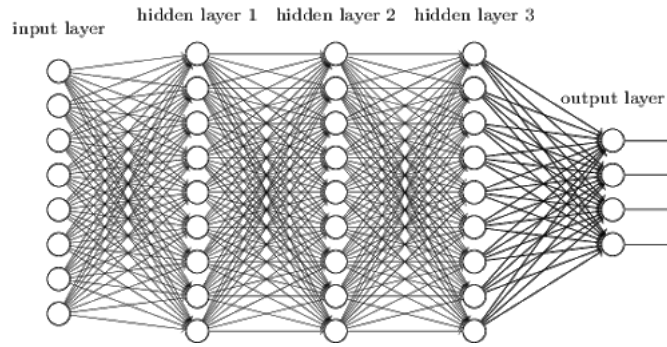


Figure 2.2. A Feedforward Neural Network with 3 hidden layers(ref. Michael A. Nielsen, Neural Networks and Deep Learning, Determination Press', 2015)

2.1.2 Training

There are still a couple of important pieces left to define to develop a properly working Neural Network. For instance, how are parameters computed? And in particular, with respect to what we compute them?

Loss Function

Before we realized that our goal is to maximize the approximation of the ground-truth input-output relation that lies under the data, therefore there is a need to introduce some metric to quantify this approximation. Call *loss function* $L(f_\theta(x), y): \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ such metric. Common choices are (ref.):

- Least Square: $\|f_\theta(x) - y\|^2$ for regression tasks
- Binary Cross-Entropy: $y \log(f_\theta(x)) + (1 - y) \log(1 - f_\theta(x))$ for binary classification
- Categorical Loss Function: $-\sum_{c=1}^C y_c \log(f_\theta(x)_c)$ for multicategory classification with C classes.

Intuitively, a good loss function will map bad approximations to high values and good approximations to smaller ones. Nevertheless, those are only point-wise estimates of the error, hence the best empirical solution learnable from the training set S would be

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n L(f_\theta(x_i), y_i) \quad (2.2)$$

which is called *empirical risk minimization*.

Gradient Descent

So far we have described, given an input \mathbf{x} , how we can compute the output of a feedforward neural network by means of compositions of dot products and non-linear transformations between matrices, starting from the first to the very last of the layers in what is called a *forward pass*. As it turns out, to attempt to solve 2.2 we need to follow the exact opposite path. Indeed, every optimization algorithm used in

practice makes use of the same subroutine called Backpropagation (ref.) introduced in the 1970s, which allows to compute, starting from the output layer and going backwards, the partial derivative of the loss function with respect to each weight in the network. Moreover it does so efficiently requiring only one *backward pass*.

Let $i^l = W_l z_{l-1} + b_l$ be the weighted input to the l -th layer, then the key observation is that the only way W_l can affect the loss function is by affecting linearly the next layer which in turn affects its next layer and so on. In particular assume we add a little change Δi_j^l to the j -th element of i^l so that the neuron will output $\sigma(i_j^l + \Delta i_j^l)$, this change will eventually propagate in the network causing the overall loss LS to change by an amount $\frac{\partial LS}{\partial i_j^l} \Delta i_j^l$. For brevity, denote the gradient of the weighted input on the j -th neuron $\delta_j^l = \frac{\partial LS}{\partial i_j^l}$, then the following holds:

$$\delta_j^L = \frac{\partial LS}{\partial i_j^L} = \frac{\partial LS}{\partial z_j^L} \frac{\partial z_j^L}{\partial i_j^L} = \frac{\partial LS}{\partial z_j^L} \sigma'_L(i_j^L) \quad (2.3)$$

and equally, taking into account the whole output layer

$$\delta^L = \nabla_{z^L} LS \odot \sigma'(i^L) \quad (2.4)$$

where \odot is the element-wise product and ∇_x the vector of the partial derivatives $\partial LS / \partial x$. That is, the gradient with respect to the weighted input to the last layer is given, using the chain rule, by the gradient with respect to the activation of the last layer times the derivative of the last activation function. Similarly, for any hidden layer l we note that:

$$\delta^l = ((W^{l+1})^T \delta^{l+1}) \odot \sigma'(i^l) \quad (2.5)$$

When we apply the transpose weight matrix, $(W^{l+1})^T$, think intuitively of this as moving the previous layer's gradient backward, giving a measure of the gradient at the output of the l -th layer. We then take the product $\sigma'(i^l)$ which again moves the gradient backward through the activation function in layer l , giving the gradient of the weighted input to layer l .

By combining 2.4 with 2.5 we can compute the gradient δ^l for any layer in the network. We start by using 2.4 to compute on the last layer, then apply equation 2.5 to compute δ^{L-1} , then the same equation again to compute δ^{L-2} , and so forth, all the way back until the input layer. Since our intent is to retrieve the gradients for every weights of the network, we are left to show how δ^l relates to them, here we provide such relation without giving an explicit proof which instead can be found in many texts like (ref Nielsen chapter 2).

$$\frac{\partial LS}{\partial b_j^l} = \delta_j^l \quad (2.6)$$

$$\frac{\partial LS}{\partial w_{i,j}^l} = z_i^{l-1} \delta_j^l. \quad (2.7)$$

Remark how we already know how to compute each element on the right sides of these equations, moreover, given that the activation function and its derivative is efficiently

computable, we will be able to efficiently get the sought gradients in just one pass. It is worth mention that Backpropagation is actually a special case of a more generic set of programming techniques that go under the name of *Automatic Differentiation* (Ref.) to numerically evaluate the derivative of a function specified by a computer program. Such techniques are usually implemented in modern numerical libraries building variations of a data structure called *computational graph*. Well known examples are *Autograd* in *Pytorch* (Ref.) or *GradientTape* in *TensorFlow* (Ref.).

What does it mean to be able to compute partial derivatives of the loss? It means being able to understand where and how the loss decreases and thus we can exploit such information to find better and better weights solutions. This is the idea behind the Gradient Descent algorithm (Ref.). In particular, the gradient of a weight is nothing but the direction inside the weight-space where the loss function is increasing, therefore what we want to do is to follow the opposite direction. Formally, this translates in the following weight update rules:

$$w_l^t \rightarrow w_l^{t+1} = w_l^t - \frac{\eta}{n} \sum_j \frac{\partial LS_{\mathbf{x}_j}}{\partial w_l^t} \quad (2.8)$$

$$b_l^t \rightarrow b_l^{t+1} = b_l^t - \frac{\eta}{n} \sum_j \frac{\partial LS_{\mathbf{x}_j}}{\partial b_l^t} \quad (2.9)$$

where w_l^t are the values of the weights for layer l -th during the t -th pass, η is a small positive constant called *learning rate* chosen by the user accordingly and the gradients are averaged among all samples in the training set. Most importantly,

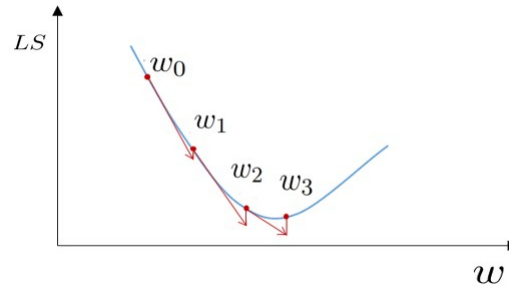


Figure 2.3. A 1-D representation of 4 gradient descent steps

note how we take the negative of the gradients, meaning that we are following the direction in which the loss decreases Fig. 2.3. The distance between two consecutive weights Δ_{w^t} will be directly proportional to both the learning rate picked and the averaged gradient.

In practice, however, when very often we are dealing with thousands or millions of data points, becomes unfeasible to compute every pass over the entire training set, thus what is done is to split the data into so called *mini-batches* and then apply the update rules on each mini-batch until we scanned the whole data. The entire scan is called *epoch* and the resulting algorithm Stochastic Gradient Descent (Ref.). Lastly, the size of a mini-batch is another hyperparameter that should be tuned by the developer, keeping in mind that the bigger the size the more stable will be our training the smaller the size the faster the training.

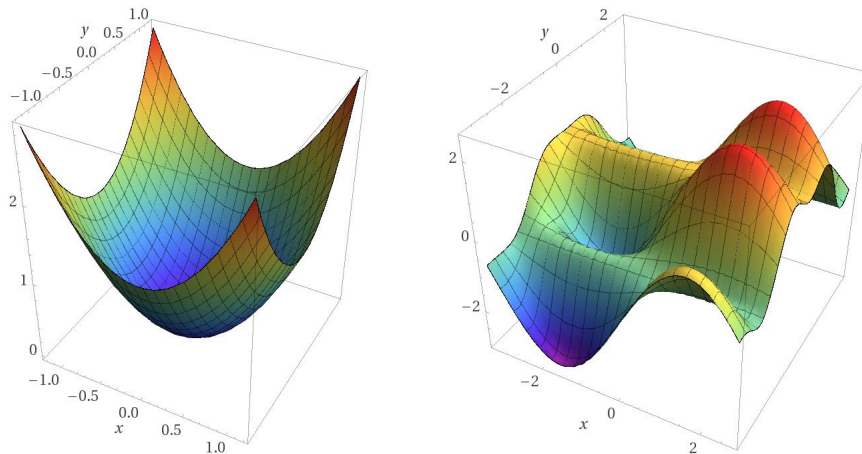


Figure 2.4. Convex and Non-convex optimization landscapes

As stated earlier, neural networks can model extremely non-convex input-output functions therefore in principle there is no guarantee that Gradient Descent will find the optimal solution to 2.2 Fig. 2.4. Indeed, the optimal solution would be the global minimum of our weighted loss function but there is no apparent way for the algorithm to distinguish between global, local minimum or saddle points Fig. 2.5. However, it turns out that in practice Gradient Descent works fairly well once we correctly tune hyperparameters and run the algorithm from different initial values. (Ref .) Moreover, lately authors have been proposed different methods to improve the convergence and the efficiency by smart changes of the learning rate during the training process (ref . cyclical learning rates)

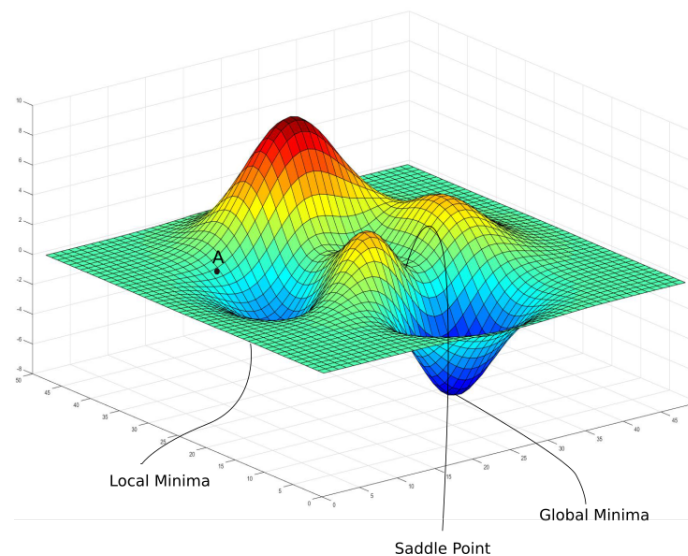


Figure 2.5. Visualization of global, local and saddle points. How can A reach the global minimum?

2.1.3 Activation Functions

We have seen how the learning process works optimizing the empirical risk by means of gradients computation. Therefore to be able to optimize anything, a neural network needs to have only differentiable components. However, before in 2.1 we discussed the so called *step function*, an activation function that, even if biologically inspired and easier to justify, is not differentiable at the origin and the derivative is 0 elsewhere. Thus if we think again about how Backpropagation works, we see that employing such activation function would make the weight updates impossible since already δ^L would be either undefined or 0.

To deal with this issue, one of the first proposed activation functions was an approximation of the step function known as *sigmoid* (Ref.)

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}} \quad (2.10)$$

which is differentiable everywhere with continuous derivatives (property that we will refer as *smoothness* through the chapters) and maps to $[0, 1]$ values. Nevertheless,

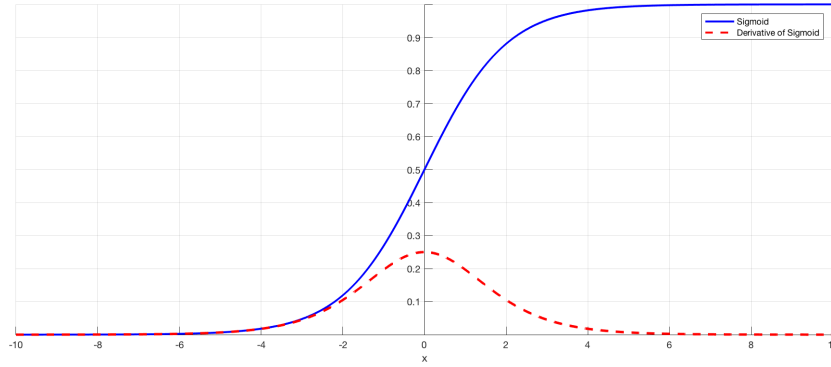


Figure 2.6. Plot of the Sigmoid function and its first derivative.

as the number of layers increases and the network becomes sufficiently deep, the sigmoid suffers from the *vanishing gradient* problem (Ref.) Fig. 2.6 due to its derivative being $[0, 0.25]$ bounded. Mainly for this reason it is not widely adopted in practice.

More in general, the sigmoid lies inside a class of activation functions known as *squashing* i.e. monotonically non-decreasing functions Σ that satisfy

$$\lim_{x \rightarrow -\infty} \sigma(x) = c, \quad \lim_{x \rightarrow \infty} \sigma(x) = 1. \quad (2.11)$$

For example, another kind of activation function of this type is the *hyperbolic tangent*, defined as

$$\tanh(x) = \frac{\exp\{x\} - \exp\{-x\}}{\exp\{x\} + \exp\{-x\}}, \quad (2.12)$$

which was found to allow for universal expressiveness of nets (Ref.). However, as for the sigmoid, squashing functions tend to be prone to vanishing and exploding gradients (Ref.).

Nowadays, the most used activation function in neural networks for different applications is the *rectifier linear unit* (ReLU), first introduced in (Ref. Hahnloser et al. in 2000) and defined as the positive part of its argument

$$\text{ReLU}(x) = \max(0, x), \quad (2.13)$$

allows for efficient training and alleviates the exploding gradient problem (having derivative either 0 or 1), introducing only one point of non-differentiability. Moreover, it promotes *sparseness* in the network, which is usually beneficial (Ref.). One problem with ReLUs though is that the neuron's value, that get pushed to a big negative number, might stay stucked in 0 for essentially all inputs, in a so called *dead state*. If many neurons in the network die this can afflict the model capacity and can be seen as a form of vanishing gradient problem. To overcome this problem, a slightly different activation functions can be used, called *leaky ReLU* (Ref.):

$$\text{LeakyReLU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha x & \text{otherwise,} \end{cases} \quad (2.14)$$

where $\alpha > 0$ is a user-defined constant usually set to small values such as 0.01. Even if this solutions solves the dying neurons problem, it does affect the sparseness property of ReLUs.

By definition, the mean of output values from a ReLU is always positive. *Exponential linear unit* (ELU) try to normalize their inputs:

$$\text{ELU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha(\exp\{x\} - 1) & \text{otherwise,} \end{cases} \quad (2.15)$$

saturating negative values at a user-defined value $-\alpha$ which is usually set to 1. Conversely to ReLU and LeakyReLUs, the derivative is continuous therefore the function is smooth and for the negative values is defined as $\text{ELU}(x) + \alpha$.

One more recent activation function

2.1.4 CNNs: Convolutional Neural Networks

2.1.5 From Neural Networks to Deep Neural Networks

2.2 Adversarial Examples Theory

Formal definition of what is an adversarial attacks plus currently well known attacks

2.3 Defenses

Review of the literature on defenses to improve robustness: provable robustness, adversarial training

2.4 Kernel Based Activation Functions

Chapter 3

Related Works

3.1 K-Winners Take All

3.2 Smooth Adversarial Training

Chapter 4

Solution Approach

Comparing the activations's distributions for different activation functions (ReLU, KWTa, Kafs) seem to suggest Kafs might be good candidates to improve model robustness

4.1 Lipschitz Constant Approach

On the limitations of current Lipschitz-Constant based approaches especially when involving Kafs

4.2 Fast is Better than Free Adversarial Training

Adversarial training (Madry et al.) and current methods to improve the efficiency (Fast is better than free)

Chapter 5

Evaluation

5.1 VGG Inspired Architectures Results

5.2 Explofing Gradients with KafNets

The exploding gradients problem with KafResNet, why is it happening? (still to clarify)

5.3 ResNet20 Inspired Architectures Results

Chapter 6

Future Works

Different Kernels, resolve the exploding gradient problem and scale to ImageNet
Perform more adaptive attacks to assess the robustness of kafresnets as is the current standard (Carlini et al.)

Chapter 7

Conclusions

This thesis tries to add to the bag of evidences in literature that smoother architectures might benefit improvements in adversarial resiliency

