



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: August 31, 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	4
2.1.3	CNNs: Convolutional Neural Networks	5
2.1.4	From Neural Networks to Deep Neural Networks	5
2.2	Adversarial Examples Theory	5
2.3	Defenses	5
2.4	Kernel Based Activation Functions	5
3	Related Works	7
3.1	K-Winners Take All	7
3.2	Smooth Adversarial Training	7
4	Solution Approach	9
4.1	Lipschitz Constant Approach	9
4.2	Fast is Better than Free Adversarial Training	9
5	Evaluation	11
5.1	VGG Inspired Architectures Results	11
5.2	Explofing Gradients with KafNets	11
5.3	ResNet20 Inspired Architectures Results	11
6	Future Works	13
7	Conclusions	15

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

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*.

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

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

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

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.

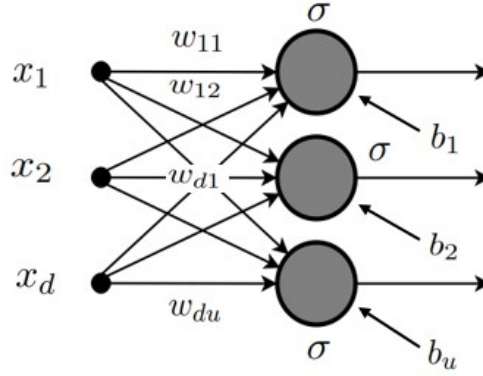


Figure 2.1. Graphic representation of a one layer NN also know as MLP(ref.)

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 $z_l = \sigma_l(z_{l-1}) = \sigma_l(W_l x + b_l)$ we write a L -layered $f_{\mathbf{W}, \mathbf{b}}$:

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

. This general but still basic form of Neural Network is known as *Feed Forward Neural Network*

2.1.2 Training

There are still a couple of important pieces left to define to have a proper working Neural Network. How are parameters learned?

2.1.3 CNNs: Convolutional Neural Networks

2.1.4 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

