

# **HANDLING ADVERSARIES IN IMAGE RECOGNITION DEEP NEURAL NETWORKS**

Pivithuru Thejan Amarasinghe

(209307X)

Degree of Master of Science

Department of Computer Science and Engineering

Faculty of Engineering

University of Moratuwa

Sri Lanka

March 2022

# **HANDLING ADVERSARIES IN IMAGE RECOGNITION DEEP NEURAL NETWORKS**

Pivithuru Thejan Amarasinghe

(209307X)

Thesis/Dissertation submitted in partial fulfilment of the requirements for the degree  
Master of Science in Data Science

Department of Computer Science and Engineering

Faculty of Engineering

University of Moratuwa  
Sri Lanka

March 2022

## **Declaration**

I declare that the thesis as my own work and it doesn't contain any sort of external material without an acknowledgement. Also, according to my knowledge I haven't used any previously published materials in order to prepare the thesis. All the work carried out to prepare the thesis are individual work and do not copy any existing work without an acknowledgement.

I like to grant permission to the University of Moratuwa to reuse or share my work in any sort of format with external parties. Also, I like to keep my right to use my work in future.

Signature:

Date:

I supervised the research conducted by the above candidate.

Signature of the supervisor:

Date:

## **Acknowledgments**

I wish to sincerely thank my supervisor, Dr. Charith Chitraranjan for providing me with the necessary guidance throughout the project. He provided me with ideas from the beginning of the project and helped me with finding the research idea. Also, special gratitude should be given to our MSc coordinator for encouraging me to full fill my research under dedicated timelines. Also, I would like to extend my gratitude to all the academic staff of the University of Moratuwa for the great contribution they made for me during study. Also, I am grateful to my family members and all the friends who helped me throughout the project.

## **Abstract**

Deep neural networks play a vital role in image recognition. There are so many mission-critical applications that use deep neural networks for image recognition. With the popularization of deep neural networks, attackers have identified their downsides of them when it comes to image recognition. Some ways can create images that can fool even deep neural networks. These images are commonly known as adversarial images. So attackers use these adversarial images to fool image recognition neural networks to develop a negative picture about using neural networks for image recognition. And even sometimes, attackers use these loopholes to conduct criminal activities as well. Keeping all these aspects in mind the idea of the research is to develop a viable solution that can tackle the main two attack techniques. The research will focus on developing adversarial images using main attacking techniques and developing a defense mechanism for those attacks. The defense technique used in the research is a combination of two techniques called adversarial training and defense distillation. As the outcome of the project accuracy of the proposed solution is measured against a typical deep neural network-based image recognition system using data samples containing adversarial images.

# Table of Contents

Declaration.....	i
Acknowledgments.....	ii
Abstract.....	iii
Table of Contents.....	iv
List of Figures.....	vi
List of Tables.....	vii
1 INTRODUCTION.....	1
1.1 Type I Error Based Attacks.....	1
1.2 Type II Error Based Attacks.....	1
1.3 Defense Methods.....	2
1.3.1 Proactive Defense.....	2
1.3.2 Reactive Defense.....	2
1.4 Contribution of Research.....	2
2 RESEARCH PROBLEM.....	3
3 RESEARCH OBJECTIVES.....	4
4 LITERATURE REVIEW.....	5
4.1 Adversarial Sample Generation.....	8
4.1.1 White-box Attacks.....	9
4.1.2 Black-box Attacks.....	15
4.1.3 Grey-box Attacks.....	15
4.1.4 Poisoning Attacks.....	16
4.1.5 Physical Attacks.....	16
4.2 Countermeasures.....	18
4.2.1 Gradient Masking.....	18
4.2.2 Robust Optimization.....	19
4.2.3 Adversarial Examples Detection.....	22
5 METHODOLOGY.....	25
5.1 Dataset.....	27
5.1.1 Type I.....	28
5.1.2 Type II.....	29
5.2 New Dataset.....	30
5.3 Adversarial Training.....	32
5.3.1 Defense Distillation.....	34
5.3.2 Evaluation.....	35
6 Results and Evaluation.....	38
6.1 Adversarial Images.....	38
6.2 Model Training.....	41
6.3 Accuracy.....	42

7	Conclusion .....	52
8	REFERENCES .....	53

## List of Figures

Figure 1: Hyper-planes .....	10
Figure 2: Adversarial image from the fast gradient .....	11
Figure 3: Adversarial image from Biggio’s attack.....	12
Figure 4: Adversarial image from the spatial transformation .....	13
Figure 5: Physical attacks .....	17
Figure 6: 3D printed physical attacks .....	18
Figure 7: Color depth variation.....	23
Figure 8: Physical attack from 3D printed glass frames .....	24
Figure 9: High-level methodology .....	26
Figure 10: Image Classifying Neural Network’s Behaviour for A Type I Image .....	28
Figure 11: Image Classifying Neural Network’s Behaviour for A Type II Image .....	29
Figure 12: New Data Sets for Type I Error.....	31
Figure 13: New Data Sets for Type II Error .....	31
Figure 14: New Data Set with Type I & Type II Errors .....	32
Figure 15: The LeNet Architecture [2] .....	32
Figure 16: Adversarial Training.....	34
Figure 17: Defense Distillation .....	34
Figure 18: Architecture of the simple neural network .....	35
Figure 19: Model Testing.....	36
Figure 20: Stratified Sampling .....	36
Figure 21: Training & Testing Data Sets .....	37
Figure 22: Model accuracies with 6% adversarial data .....	45
Figure 23 : F1-Score in the base model compared to the upgraded model.....	47



## List of Tables

Table 1: Features of each layer of the upgraded LeNet model .....	33
Table 2: Features of each layer of the base LeNet model .....	34
Table 3: Features of each layer of the simple neural network .....	35
Table 4: Type I Error .....	39
Table 5: Type II Error .....	41
Table 6: Number of data points in data sets .....	42
Table 7: Number of classes for data sets .....	42
Table 8: Model accuracies with 2% adversarial data .....	43
Table 9: Model accuracies with 4% adversarial data .....	44
Table 10: Model accuracies with 6% adversarial data for MNIST fashion data .....	46
Table 11: Special Scenarios with the Upgraded Model .....	50