# Adversarial Machine Learning

*A Project Report Submitted
in Partial Fulfillment of the Requirements
for the Degree of*

**Bachelor of Technology**

*by*

**Vaibhav Nagrale**
(112001046)



**COMPUTER SCIENCE AND ENGINEERING**

**INDIAN INSTITUTE OF TECHNOLOGY PALAKKAD**

# CERTIFICATE

This is to certify that the work contained in the project entitled "**Adversarial Machine Learning**" is a bonafide work of **Vaibhav Nagrale (Roll No. 112001046)**, carried out in the Department of Computer Science and Engineering, Indian Institute of Technology Palakkad under my guidance and that it has not been submitted elsewhere for a degree.

**Vivek Chaturvedi**

Assistant/Associate Professor

Department of Computer Science & Engineering

Indian Institute of Technology Palakkad

# Acknowledgements

I would like to take this opportunity to express my deepest gratitude to my mentor, **Vivek Chaturvedi**.

**Vaibhav Nagrale**

IIT Palakkad

Date: 27/10/2023

# Abstract

In the pursuit of my Bachelor's Thesis Project (BTP) in Adversarial Machine Learning, my primary objective was to explore the development of novel attack or defense strategies.

I conducted an extensive review of research papers on various attack methods and applied pre-existing attacks to different datasets and models. Through rigorous experimentation, I systematically tested a range of attacks, each tailored to specific perturbations, and meticulously documented their outcomes. Additionally, I employed graphical representations and the GradCAM visualization technique to provide insightful and visually compelling observations.

This report encapsulates the comprehensive findings and insights gained from these investigations, contributing to the ever-evolving landscape of adversarial machine learning research.

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

## 1.1 Introduction to Adversarial Machine Learning

Adversarial machine learning is the study of the attacks on machine learning algorithms, and of the defenses against such attacks. It focuses on exploiting vulnerabilities in machine learning models by crafting adversarial inputs to deceive them.

## 1.2 Significance and motivation for the BTP

The motivation for this BTP stems from the critical issue of adversarial attacks on artificial intelligence systems. As highlighted in the research paper "Review of Artificial Intelligence Adversarial Attack and Defense Technologies," [1] such attacks limit the potential of AI in crucial security domains. Enhancing AI robustness against adversarial threats is paramount for advancing the field, making our BTP a significant endeavor.

## 1.3 Research objectives

The primary objective of this study is to explore the development of novel attack or defense strategies, aiming to advance the understanding and capabilities in the field of adversarial machine learning. Additionally, it seeks to contribute to the ongoing efforts to improve the security and robustness of machine learning models in the face of adversarial attacks.

# Chapter 2

# Literature Review

## 2.1 Insights into Adversarial Attacks

To gain valuable insights into the realm of adversarial attacks and their implications within the domain of security, I extensively reviewed a selection of research papers on BlackBoxAttacks and WhiteBoxAttacks.

### 2.1.1 BlackBoxAttacks

- **Decision Based** : Focuses on altering model outputs without access to model internals. [2]

- **Gen Attack** : Evolves adversarial examples using genetic algorithms. [3]

- **Swarm Optimization** : Utilizes swarm intelligence for crafting adversarial inputs. [4]

### 2.1.2 WhiteBoxAttacks

- **Carlini Wagner** : Generates adversarial examples using an optimization-based approach. [5]

- **FGSM** : Perturbs input data based on gradient information. [6]

3

- **JSMA** : Manipulates features based on Jacobian saliency maps.

[7]

## 2.2 Conclusion

I got an overview on ways of attacking and types of attack under them. It has shed light on different strategies employed by researchers and attackers to manipulate machine learning models, exposing vulnerabilities in their decision-making processes.

# Chapter 3

# Methodology

After getting idea on adversarial attacks by reading research papers I implemented some pre-existing attacks on model. This section has overview of attacks implementation.

## 3.1 Description of Datasets:

- **ImageNet** : I utilized the ImageNet dataset, an extensive online dataset, for my research. ImageNet comprises a diverse collection of images spanning various categories, making it a widely used resource in computer vision and machine learning.

- **X-ray Pneumonia** : In addition, I incorporated the X-ray Pneumonia dataset, which exclusively contains X-ray images of patients diagnosed with pneumonia. This dataset served as a specific and targeted subset relevant to the medical domain.

## 3.2 Overview of Model:

For my experiments, I selected the ResNet-18 model as the backbone. ResNet-18 is a well-established convolutional neural network architecture known for its robust performance in image classification tasks. It offers a suitable balance between accuracy and

computational efficiency, making it an ideal choice for this study.

## 3.3 Attack Strategies Applied:

In this research, I implemented a series of adversarial attack strategies to evaluate the model's robustness against adversarial perturbations. These included:

- **Fast Gradient Sign Method (FGSM)** : FGSM is a gradient-based attack method that generates adversarial examples by perturbing input data based on the sign of the gradient with respect to the loss function.

- **Projected Gradient Descent (PGD)** : PGD is an iterative variant of FGSM that applies multiple iterations to produce more potent adversarial examples.

- **Basic Iterative Attack** : This method extends the iterative approach to craft adversarial examples by perturbing the input multiple times.

- **Fast Minimum Norm Attack (FMN)** : The goal of FMNAttack is to generate minimal perturbations that, when applied to input data, lead to misclassifications by the model.

- **DeepFool Attack** : The DeepFool Attack is an iterative method that computes the smallest perturbation needed to mislead the model.

I executed these attacks using the Foolbox library in Python, a versatile and widely-used library for crafting and evaluating adversarial attacks.
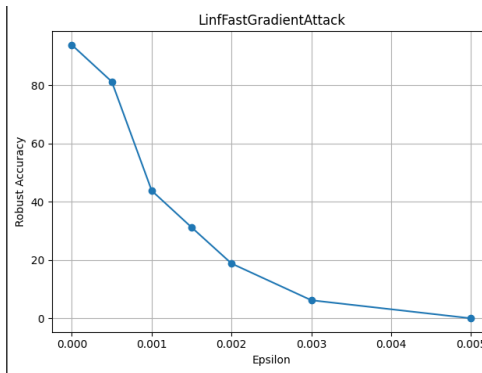
# Chapter 4

# Results and Analysis

This section has the results of the analysis by applying a series of adversarial attacks to the model. The attacks are executed with varying perturbation magnitudes (epsilon) to assess the model's performance under different levels of adversarial influence.

Here are graphical results for different attacks :

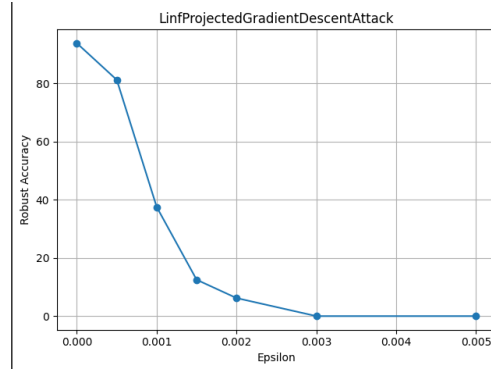## 4.1 Robustness Analysis using Adversarial Attacks:

In figures below we see decrease in robust accuracy with increase in epsilon depicting success of attack for decreasing accuracy of model.
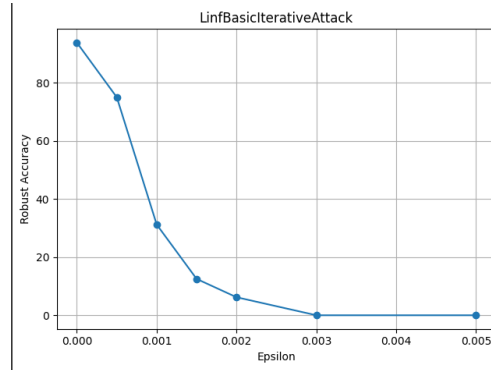
- **Fast Gradient Sign Method:** (Fig. 4.1)



**Fig. 4.1**: Fast Gradient Sign Attack Graph
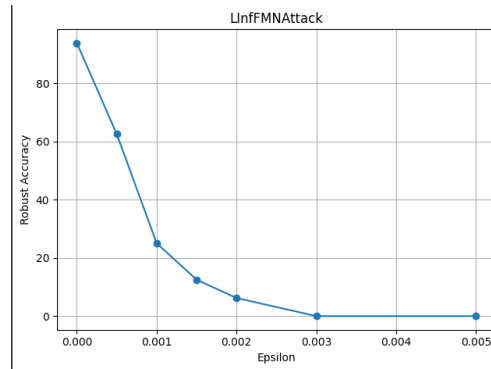
- **Projected Gradient Descent:** (Fig. 4.2)



**Fig. 4.2**: Projected Gradient Descent Attack Graph

- **Basic Iterative Attack:** (Fig. 4.3)



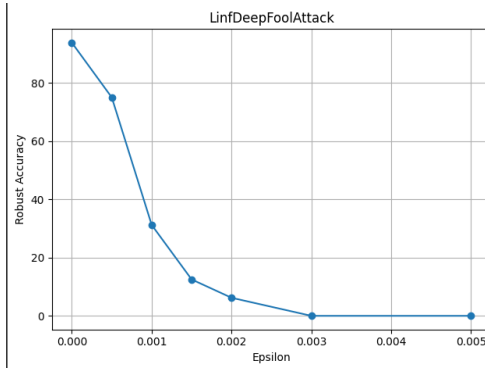**Fig. 4.3**: Basic Iterative Attack Graph

- **Fast Minimum Norm Attack:** (Fig. 4.4)



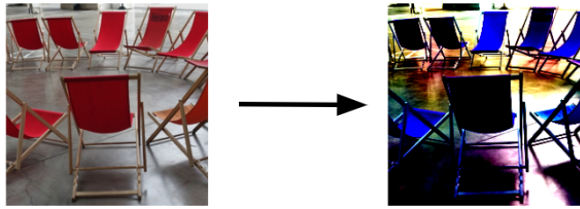**Fig. 4.4**: Fast Minimum Norm Attack Graph
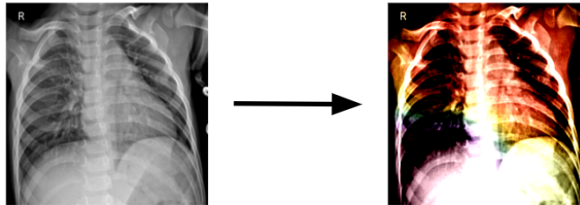
- **DeepFool Attack:** (Fig. 4.5)



**Fig. 4.5**: DeepFool Attack Graph

## 4.2 Visual Insights using GradCAM:

In addition to quantitative analysis, I employed GradCAM (Gradient-weighted Class Activation Mapping) to provide a qualitative perspective on model behavior. GradCAM images for select attack scenarios are included to reveal the regions of the image that the model focuses on during classification decisions.



**Fig. 4.6**: Chair GradCAM image (imagenet dataset)



**Fig. 4.7**: Lungs GradCAM image (xray pneumonia dataset)

In Figure (Fig. 4.6), we observe the GradCAM of an adversarial image depicting a chair.

Notably, the chair appears in blue, signifying a shift in the model's attention away from the chair and consequently leading to incorrect predictions. Similarly, in Figure (Fig. 4.7), the presence of the red region in the upper portion of the X-ray image, due to adversarial perturbations, results in the model focusing on the wrong area, leading to inaccurate predictions.

This structured approach allows for a clear and organized presentation of the results and analysis section, combining quantitative metrics with qualitative insights for a comprehensive evaluation of the model's performance under adversarial conditions.

# Chapter 5

# Conclusion and Future Work

## 5.1 Summary of key findings

In conclusion, this research has provided valuable insights into the robustness of the ResNet-18 model against various adversarial attacks. Our analysis, supported by quantitative results and qualitative GradCAM visualizations, sheds light on the model's performance under different attack scenarios. We oberved for epsilon greater than 0.005, attacks preformed well in attacking RestNet-18 model.

## 5.2 Contributions to the field of Adversarial Machine Learning

These findings contribute to our understanding of adversarial machine learning and highlight the importance of security aspects.

## 5.3 Implications for future research

Having gained insights into various adversarial attacks and their practical applications, the subsequent focus of this study lies in the development and implementation of Resource Efficient Decision-based Imperceptible Attack for Machine Learning [8]. This section outlines the research direction and approach for the upcoming investigation.

# References

[1] Shilin Qiu, Qihe Liu, Shijie Zhou, and Chunjiang Wu, "Review of artificial intelligence adversarial attack and defense technologies," 2019. [Online]. Available: https://www.mdpi.com/2076-3417/9/5/909

[2] Wieland Brendel, Jonas Rauber, and Matthias Bethge, "Decision-based adversarial attacks: Reliable attacks against black-box machine learning models," 2018. [Online]. Available: https://arxiv.org/abs/1712.04248

[3] Moustafa Alzantot, Yash Sharma, Supriyo Chakraborty, Huan Zhang, Cho-Jui Hsieh, and Mani Srivastava, "Genattack: Practical black-box attacks with gradient-free optimization," 2019. [Online]. Available: https://arxiv.org/abs/1805.11090

[4] Quanxin Zhang, Kunqing Wang, Wenjiao Zhang, and Jingjing Hu, "Attacking black-box image classifiers with particle swarm optimization," vol. 7, 2019. [Online]. Available: https://ieeexplore.ieee.org/document/8876844

[5] Nicholas Carlini and David Wagner, "Towards evaluating the robustness of neural networks," 2017. [Online]. Available: https://arxiv.org/abs/1608.04644

[6] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy, "Explaining and harnessing adversarial examples," 2015. [Online]. Available: https://arxiv.org/abs/1412.6572

[7] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami, "The limitations of deep learning in adversarial settings," 2015. [Online]. Available: https://arxiv.org/abs/1511.07528

[8] Faiq Khalid, Hassan Ali, Muhammad Abdullah Hanif, Semeen Rehman, Rehan Ahmed, and Muhammad Shafique, "Red-attack: Resource efficient decision based attack for machine learning," 2019. [Online]. Available: https://arxiv.org/abs/1901.10258