

LB/TH/43/2025
TH6012

**Enhancing the Robustness of Credit Card Fraud Detection
Systems Against Adversarial Attacks Using Machine
Learning**

Rathnayake R.M.C

219392J

MSc in Computer Science

Department of Computer Science and Engineering
Faculty of Engineering

University of Moratuwa
Sri Lanka

March 2025

**Enhancing the Robustness of Credit Card Fraud Detection
Systems Against Adversarial Attacks Using Machine
Learning**

Rathnayake R.M.C

219392J

Thesis/Dissertation submitted in partial fulfillment of the requirements for the degree
MSc in Computer Science

Department of Computer Science and Engineering
Faculty of Engineering

University of Moratuwa
Sri Lanka

March 2025

DECLARATION

I declare that this is my own work and this thesis/dissertation does not incorporate without acknowledgement any material previously submitted for a degree or diploma in any other University or Institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text. I retain the right to use this content in whole or part in future works (such as articles or books).

05/06/2025

Signature:

Date:

The above candidate has carried out research for the PhD/MPhil/Masters thesis/dissertation under my supervision. I confirm that the declaration made above by the student is true and correct.

Name of Supervisor: Prof. Sanath Jayasena

Signature of the Supervisor:

Date: 05th June 2025

DEDICATION

This research is dedicated to the countless individuals who have fallen victim to credit card fraud, whose experiences have fueled the urgency for robust and resilient fraud detection systems. It is also dedicated to the tireless efforts of researchers, practitioners, and financial institutions working diligently to protect consumers from these malicious activities. May this research contribute to the collective knowledge and efforts aimed at safeguarding the integrity of digital financial transactions.

ACKNOWLEDGEMENT

The successful completion of this research would not have been possible without the invaluable contributions and support of numerous individuals and organizations, I extend my deepest appreciation to my advisor, Prof. Sanath Jayasena, for providing unwavering guidance, mentorship, and encouragement throughout this research journey. Your expertise, patience, and belief in my abilities have been instrumental in shaping this work. Thank you for your invaluable contributions to this research.

I am profoundly humbled and grateful for the opportunity to have embarked on this journey, and I earnestly hope that the findings of this work will contribute to the advancement of knowledge and the betterment of society.

ABSTRACT

Credit card fraud detection systems are essential for protecting money transactions and stopping illegal use. The usefulness of these systems is seriously threatened by the increasing prevalence of adversarial assaults, as attackers try to alter input data in order to trick machine learning algorithms and evade detection. The problem of making credit card fraud detection systems more resilient to such hostile attacks is addressed in this study. It explores current defense tactics, analyses the status of adversarial attacks, and evaluates how they affect model performance. The study offers a comprehensive strategy that involves implementing several defense strategies, modelling hostile settings, and assessing their effectiveness using important performance indicators. The results show that the proposed approach greatly enhances the robustness of fraud detection models, effectively reducing the influence of adversarial manipulations

Credit card fraud detection systems face a significant challenge as adversarial attacks grow more complex, enabling attackers to alter data and avoid detection. By methodically evaluating defense tactics against these threats, this study seeks to improve fraud detection models. Various adversarial attack types, including hybrid, white-box, and black-box attacks, are simulated to identify weaknesses in current systems. Four defense mechanisms, namely, Neural Cleanse, Random Noise, General Adversarial Training, and Defensive Distillation are implemented and assessed. The results demonstrate that none of these methods offer total protection, although they improve model resilience to some degree. It backs up the claim that financial institutions do not yet have infallible defenses against hostile attacks. The study provides insightful information about enhancing fraud prevention by outlining the advantages and disadvantages of each strategy.

This study supports the development of more effective credit card fraud detection systems by offering actionable insights and a structured approach to designing and assessing protective strategies against adversarial threats. The findings have implications for financial institutions seeking to fortify their security posture and protect customers from fraudulent activities, fostering trust and confidence in digital financial transactions.

Keywords: Credit card fraud; Adversarial attacks; Robustness; Defense mechanisms; Machine learning; Security; Financial transactions; Fraud detection

TABLE OF CONTENTS

Declaration	i
Dedication	ii
Acknowledgement.....	iii
Abstract	iv
Table of Contents	v
List of Figures	viii
List of Tables.....	ix
List of Abbreviations.....	x
Chapter 1	1
Introduction	1
1.1 Credit Card fraud.....	1
1.1.1 Existing Approaches for the Credit Card Fraud Detection	3
1.1.2 Adversarial Attacks	4
1.2 Research Problem.....	6
1.3 Research aim and Objectives.....	6
1.3.1 Aim.....	6
1.3.2 Objectives	6
Chapter 2	8
Literature Review.....	8
2.1 Feature Selection for Unbalanced Data Set.....	9
2.2 Implement Credit Card Fraud Detection	11
2.3 Adversarial Attacks in Adversarial Environment.....	17
2.3.1 Zeroth Order Optimization.....	23
2.3.2 Deep Fool adversarial attack.....	26
2.3.3 Elastic Net adversarial attack.....	27
2.4 Robustness Predictions in Adversarial Environment	29
2.4.1 Defensive Distillation.....	29
2.4.2 General Adversarial Training.....	31
2.4.3 Random Noise.....	33

5.4.4 Neural Cleanse	35
Chapter 3	37
Proposed Method	37
3.1 Data Collection and Feature Selection	38
3.2 Training and Evaluating Fraud Detection Models	39
3.3 Simulation of Adversarial Environment	39
3.4 Implementation of Defense Mechanism	39
3.5 High Level Architecture.....	39
3.6 Summary	41
Chapter 4	43
Implementation	43
4.1 Dataset	44
4.2 Build a fraud detection models	44
4.2.1 Random Forest	45
4.2.2 Logistic Regression.....	45
4.2.3 Support Vector Machine (SVM).....	46
4.2.4 Decision Tree	46
4.3 Simulation of Adversarial Environment	46
4.3.1 Zeroth-Order Optimization (ZOO)	47
4.3.2 Deep Fool.....	47
4.3.3 Elastic Net attack.....	48
4.4 Implement of Defense Mechanism Technique	48
4.4.1 Defensive Distillation.....	48
4.4.2 General Adversarial Training.....	49
4.4.3 Random Noise.....	50
4.4.4 Neural Cleanse	50
4.6 Summary	50
Chapter 5	52
Evaluation	52
5.1 Evaluation Metrics	52
5.2 Evaluation of the Implementation.....	53
5.2.1 Blackbox Environment.....	53

5.2.1.1	Random Forest	53
5.2.1.2	Logistic Regression.....	57
5.2.1.3	Support Vector Machine (SVM).....	60
5.2.1.4	Decision Tree	63
5.2.1.5	Summary view in Blackbox environment.....	67
5.2.2	Whitebox Environment	68
5.2.3	Hybrid Environment	70
Chapter 6	72
Conclusion	72
References	73
APPENDIX A	76
Detail Evaluation results	76

LIST OF FIGURES

Figure	Description	Page
Figure 1.1	Example research for adversarial attack	5
Figure 2.1	Existing approaches for Credit Card fraud identifying.	8
Figure 2.2	Behavior of random forest algorithm.	15
Figure 2.3	Behavior of logistic regression algorithm.	16
Figure 2.4	Behavior of SVM algorithm	16
Figure 2.5	Behavior of decision tree algorithm	17
Figure 2.6	Adversarial Patch example	19
Figure 2.7	Overview of adversarial attack	27
Figure 3.1	Sample dataset	38
Figure 3.2	High level architecture of proposed solution	40
Figure 4.1	Activity flow of the implementation	43
Figure 4.2	Screenshot of final implementation	51

LIST OF TABLES

Table	Description	Page
Table 2.1	Summary of adversarial attack types	22
Table 5.1	Comparative analysis results of black box environments	67
Table 5.2	Comparative analysis results of white box environments	69
Table 5.3	Comparative analysis results of hybrid box environments	70

LIST OF ABBREVIATIONS

ACOFS	Ant Colony Optimization Feature Selection
AUC	Area Under the Curve
AUC-ROC	Area Under the ROC Curve
CNN	Convolutional Neural Network
DEFS	Differential Evolutionary Feature Selection
DNNs	Deep Neural Networks
ELM	Extreme Learning Method
FP	False Positive
FN	False Negative
GA	Genetic Algorithm
GAFS	Genetic Algorithm Feature Selection
GBM	Gradient Boosting Machine
GAT	Generative Adversarial Trainer
KNN	K-Nearest Neighbors
LR	Logistic Regression
LOF	Local Outlier Factor
LGBM	Light Gradient Boosting Machine
NB	Naive Bayes
PCA	Principal Component Analysis
PSOFS	Particle Swarm Optimization Feature Selection
RF	Random Forest
ROC	Receiver Operating Characteristic
RHSO	Rock Hyrax Swarm Optimization
SVM	Support Vector Machine
TN	True Negative
TP	True Positive