



PROJECT REPORT
MACHINE LEARNING BASED FRAUD
IDENTIFICATION ON E-TRANSACTION

LAUW, AGUNG WIJAYA ALAUY
16.K1.0050

Faculty of Computer Science
Soegijapranata Catholic University
2020

APPROVAL AND RATIFICATION PAGE

MACHINE LEARNING BASED FRAUD IDENTIFICATION ON E- TRANSACTION

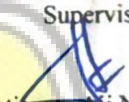
by

LAUW, AGUNG WIJAYA ALAUY – 16.K1.0050

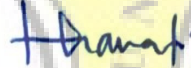
This project report has been approved and ratified
by the Faculty of Computer Science on January, 7, 2020


With approval,

Supervisor,


Robertus Setiawan Aji Nugroho, ST., Ph.D
NPP : 05812004264

Examiners,


1.) 
Rosita Herawati, ST., MIT
NPP : 05812004263

2.) 
Robertus Setiawan Aji Nugroho, ST., Ph.D
NPP : 05812004264

3.) 
Y.B. Dwi Setianto, ST., M.Cs
NPP : 05872017021



Dean of Faculty of Computer Science,


Robertus Setiawan Aji Nugroho, ST., MCompIT., PhD
NPP: 058.1.2004.264

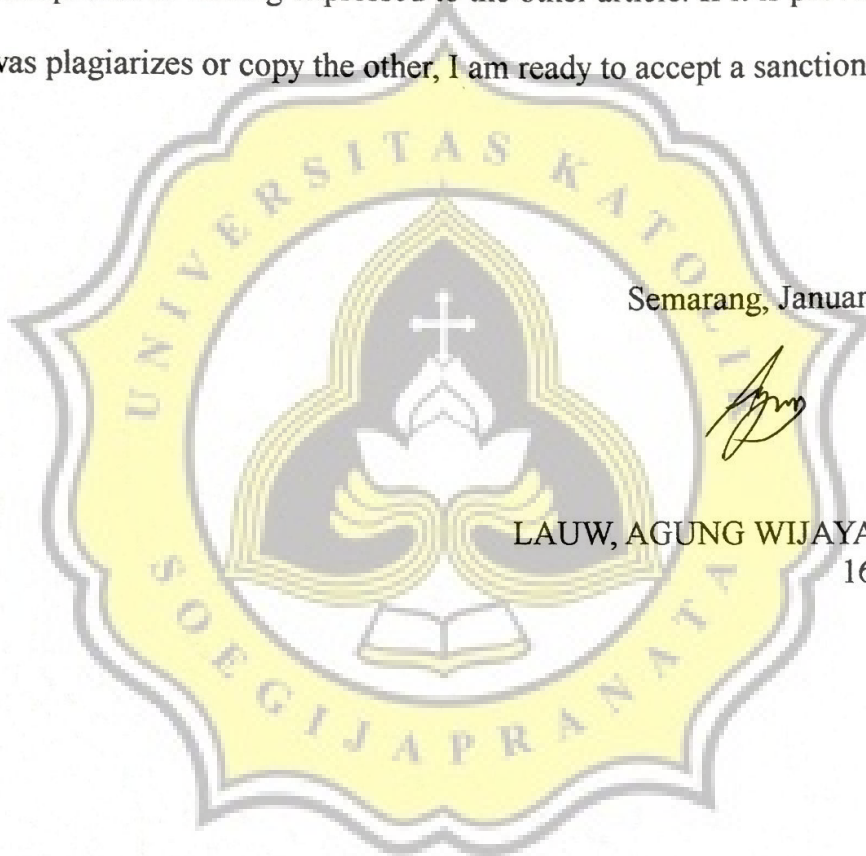
STATEMENT OF ORIGINALITY

I, the undersigned:

Name : LAUW, AGUNG WIJAYA ALAUY

ID : 16.K1.0050

Certify that this project was made by myself and not copy or plagiarize from other people, except that in writing expressed to the other article. If it is proven that this project was plagiarizes or copy the other, I am ready to accept a sanction.



Semarang, January, 7, 2020

LAUW, AGUNG WIJAYA ALAUY
16.K1.0050

ABSTRACT

Transaction fraud is fatal to any and all companies, costing them millions of moneys for the fraud cost and prevention. Machine learning to detect transaction fraud commonly uses classification method to determine the legitimacy of said transaction. It can either classify the transaction as a fraudulent or legitimate as a whole, or classify the types of fraud indication said transaction had done. For it to be acceptably fast, the training dataset should not be excessive in quantity but have high enough quality so as the model does not suffer accuracy issue instead.

Some studies include decision tree and random forest in which the random forest almost always yield higher accuracy; understandable considering random forest is an ensemble classifier consisting of many decision trees [1]. While other contenders are SVM and logistic regression, almost all of them boasts high accuracy rate for detection (above 80%) which indicates low complexity when determining a single transaction as far as testing goes. There are several types of fraud in Ecommerce, including but not limited to: Friendly Fraud where a customer (fraudster) complains and claims a refund or purchase, Clean Fraud where a fraudster uses a stolen credit card to make a purchase, Card Testing where the fraudster makes low purchase to validate stolen card information or randomly generated card number on a website with different specific notice like "Incorrect expiration date".

Some prevalent characteristic of a fraud includes but not limited to: customer is a first time customer, customer orders are bigger than average, customer is in an unusual location, customer orders same product but at high quantity, customer ships to multiple addresses, Several purchases with same IP but different card information, too many transaction in a short time span. The previous also includes potential false positives, so to get an accurate result requires a tally of points according to how suspicious or how many of the rules is broken.

Keyword: Classification, Fraud, Machine Learning

TABLE OF CONTENTS

| | |
|--|-----------|
| Cover..... | i |
| APPROVAL AND RATIFICATION PAGE..... | ii |
| STATEMENT OF ORIGINALITY..... | iv |
| ABSTRACT..... | v |
| TABLE OF CONTENTS..... | vi |
| ILLUSTRATION INDEX..... | vii |
| INDEX OF TABLES..... | viii |
| CHAPTER 1 INTRODUCTION..... | 1 |
| 1.1 Background..... | 1 |
| 1.2 Problem Formulation..... | 10 |
| 1.3 Scope..... | 11 |
| 1.4 Objective..... | 12 |
| CHAPTER 2 LITERATURE STUDY..... | 13 |
| CHAPTER 3 RESEARCH METHODOLOGY..... | 21 |
| CHAPTER 4 ANALYSIS AND DESIGN..... | 24 |
| 4.1 Analysis..... | 24 |
| 4.2 Desain..... | 25 |
| CHAPTER 5 IMPLEMENTATION AND TESTING..... | 26 |
| 5.1 Implementation..... | 26 |
| 5.2 Testing..... | 27 |
| CHAPTER 6 CONCLUSION..... | 33 |
| REFERENCES..... | A |
| APPENDIX..... | A |

ILLUSTRATION INDEX

| | |
|---|----|
| Illustration 1.1: Reported fraud cases occurred in U.S. in 2018, taken from https://shiftprocessing.com/credit-card-fraud-statistics/ | 2 |
| Illustration 1.2: Reported credit card fraud cases in U.S. from 2014 to 2018, taken from https://shiftprocessing.com/credit-card-fraud-statistics/ | 2 |
| Illustration 1.3: Data breaches by sector in 2018, taken from https://shiftprocessing.com/credit-card-fraud-statistics/ | 3 |
| Illustration 1.4: Losses suffered by method of contact, taken from https://shiftprocessing.com/credit-card-fraud-statistics/ | 4 |
| Illustration 1.5: Age of identity theft victims, taken from https://shiftprocessing.com/credit-card-fraud-statistics/ | 5 |
| Illustration 1.6: Physical impact to victims, taken from https://shiftprocessing.com/credit-card-fraud-statistics/ | 6 |
| Illustration 1.7: Emotional impact to victims, taken from https://shiftprocessing.com/credit-card-fraud-statistics/ | 6 |
| Illustration 1.8: General steps of phishing, taken from https://www.csoonline.com/article/2117843/what-is-phishing-how-this-cyber-attack-works-and-how-to-prevent-it.html | 8 |
| Illustration 2.1: How SVM works, taken from https://www.ibm.com/support/knowledgecenter/en/SS3RA7_15.0.0/com.ibm.spss.modeler.help/svm_howwork.htm | 15 |
| Illustration 2.2: Visualization of gaussian SVM. Area with dotted line is a collection of the data point's Gamma. The regular line is the line created by SVM to separate classes, taken from https://towardsdatascience.com/support-vector-machine-simply-explained-fee28eba5496 | 16 |
| Illustration 2.3: Different visualization of SVM kernels, taken from https://blog.easysol.net/machine-learning-algorithms-6/ | 16 |
| Illustration 2.4: Simple decision tree..... | 17 |

INDEX OF TABLES

| | |
|-----------------------------------|----|
| Table 5.1: Scenario 1 result..... | 27 |
| Table 5.2: Scenario 2 result..... | 28 |
| Table 5.3: Scenario 3 result..... | 28 |
| Table 5.4: Scenario 4 result..... | 28 |
| Table 5.5: Scenario 5 result..... | 29 |
| Table 5.6: Scenario 6 result..... | 29 |
| Table 5.7: Scenario 7 result..... | 29 |
| Table 5.8: Overall Result..... | 30 |

