



NAIF ARAB UNIVERSITY FOR SECURITY SCIENCE
COLLEGE OF COMPUTER INFORMATION SECURITY
DEPARTMENT

A NEW MODEL FOR DETECTING BLACK ACCOUNTS IN SOCIAL MEDIA

BY:

YOUSEF SAMIR ALHARBI

A THESIS

SUBMITTED IN PARTIAL FULFILMENT OF THE REQUIREMENTS FOR
THE DEGREE OF:

MASTER OF SCIENCE IN INFORMATION SECURITY

TO:

INFORMATION SECURITY DEPARTMENT
THE COLLEGE OF COMPUTER AND INFORMATION SECURITY
NAIF ARAB UNIVERSITY FOR SECURITY SCIENCES

2017/2018

SUPERVISED BY:

DR. FAHAD ALHARBY

NAIF ARAB UNIVERSITY FOR SECURITY SCIENCES

Table of Contents

Abstract	II
المستخلص	III
Dedication	IV
Acknowledgment	V
Table of Figures	VIII
Table of Tables	IX
List of Acronyms	X
Chapter One: Introduction	1
1.1 Introduction	1
1.2 Problem statement	3
1.3 Research Questions.....	5
1.4 Research objectives	5
1.5 Research Significance.....	6
1.6 Research Methodology	7
1.7 Research Background	10
1.8 Research General Structure	14
1.9 Thesis contributions.....	16
1.10 Ethical considerations.....	17
1.11 Thesis outline.....	18
Chapter Two: Background and Literature Review	19
2.1 Introduction	19
2.2 Twitter Online Social Network.....	20
2.3 Data Mining.....	24
2.3.1 Textual Mining	27
2.3.2 Text Classification	31
2.3.3 Text Classification Techniques.....	36
2.4 Features.....	43
2.4.1 Feature Extraction.....	44

2.4.2 Feature Selection	46
2.5 Related Works in social networks mining:	49
2.5.1 Related work in non-Arabic corpus	49
2.5.2 Related work in Arabic corpus	54
Chapter Three: The research proposed model	57
3.1 Introduction	57
3.2 Data Set.....	59
3.3 Classification Technique	63
Chapter Four Results and discussion.....	65
4.1 Introduction	65
4.2 Initial considerations.....	65
4.3 Evaluation Measures.....	65
4.4 Dataset	66
4.5 K-Fold.....	67
4.6 TF-IDF.....	67
4.7 Hyperparameters Optimization.....	68
4.8 Training the Model	73
4.9 Results	73
Chapter Five: Conclusion.....	75
5.1 Finding and Suggestions.....	75
5.2 Limitation of Work.....	75
5.3 Summary of Results.....	76
5.4 Future Work.....	76
References	78
Appendix A: Python Script	89
Appendix B: JAVA Program	98

Table of Figures

Figure 1. Most popular Social Networks by country [4].	2
Figure 2. Hyper plane created for data sets point (mathematical principle of the SVM) [35]	37
Figure 3. Optimal margin between closest data set points [36].	38
Figure 4. Decision Tree Nodes[41].	40
Figure 5. Nuarel Netwrok Layers [46].	43
Figure 6. Percentage of Black Tweets and Good Tweets in Dataset.	66
Figure 7. K-Fold.	67
Figure 8. Accuracy during K-Fold.	71
Figure 9. Recall during K-Fold.	72
Figure 10. Precision during K-Folds.	72

Table of Tables

Table1. AUtomatic Feature selection for the four purposes in the [61]study	51
Table 2. Data set Twitter accounts.....	60

List of Acronyms

Acronyms	Details
ANN	Artificial Neural Networks.
URL	Uniform Resource Locator.
KDD	Knowledge Discovery and Data Mining
AI	Artificial Intelligent.
ML	Machine Learning.
DL	Deep Learning.
DM	Data Mining.
KNN	K Nearest Neighbors.
SVM	Support Vector Machines
SVN	Support Vector Network.
LIBSVM	Library for Support Vector Machine.
LP	Label PowerSet.
CLR	Calibrated Label Ranking.
WEKA	Waikato Environment for Knowledge Analysis.
LC	Label Calibration.
RT	Ranking and Threshold.
TC	Text Classification.
NB	Naïve Bays
MNB	Multinomial naïve bayes
MFCC	Mel-Frequency Cepstral Coefficients.
SFM	Spectral Flatness Measure.
RBFKernel	Radial Basis Function Kernel.
LIBINAR	Library for Large Linear Classification
TF-IDF	Term Frequency-Inverse Document Frequency
HTML	Hyper Text Markup Language
Ws	Slang Word.

LIWC	Words Linguistic Inquiry and Word Count dictionary.
BoW	Bag of words.
TF	Term Frequency.
TO	Term Occurrences.
BTO	Binary Term Occurrences.
VSM	Vector Space Model.
NLP	Natural Languages Processing
LSA	Latent Semantic analysis.
IE	Information extraction
TR	Training set
TE	Test set
QP	Quadratic programming
PCA	Principle component analysis
IG	Information Gain
SFS	Sequential learning selection
RL	Reinforcement learning algorithm
GBDT	Gradient boosted decision trees
KTEA	Korean data set
API	Application programming interface
EM	Expectation maximization
LDA	Latent dirichlet allocation
TLC	Lightweight clustering
XML	Extensible Markup Language
XLS	Microsoft Excel
PTT or Tweepy	Python Twitter tool
Twitter4J	Java class named
LSTM	Long short-term memory
CSV	comma-separated values