

**ADOPTION OF MACHINE LEARNING
ALGORITHM FOR ANALYSING SUPPORTERS
AND NON-SUPPORTERS FEEDBACK ON
POLITICAL POSTS**

OGUNFOLAJIN MARUFF TUNDE

**FACULTY OF COMPUTER SCIENCE AND
INFORMATION TECHNOLOGY
UNIVERSITI OF MALAYA
KUALA LUMPUR**

2022

**ADOPTION OF MACHINE LEARNING ALGORITHM FOR
ANALYSING SUPPORTERS AND NON-SUPPORTERS
FEEDBACK ON POLITICAL POSTS**

OGUNFOLAJIN MARUFF TUNDE

**DISSERTATION SUBMITTED IN FULFILMENT
OF THE REQUIREMENTS FOR THE DEGREE
OF MASTER COMPUTER SCIENCE (APPLIED
COMPUTING)**

**FACULTY OF COMPUTER SCIENCE AND INFORMATION
TECHNOLOGY
UNIVERSITI OF MALAYA
KUALA LUMPUR**

2022

UNIVERSITY OF MALAYA
ORIGINAL LITERARY WORK DECLARATION

Name of Candidate: **Ogunfolajin Maruff Tunde**

Matric No: **WOA160026 /17020404/1**

Name of Degree: **Master of Computer Science (Applied Computing)**

Title of Project Paper/Research Report/Dissertation/Thesis (“this Work”):
Adoption of Machine Learning Algorithm for Analysing Supporters and Non-Supporters Feedback on Political Posts.

Field of Study: **Sentiment Analysis.**

I do solemnly and sincerely declare that:

- (1) I am the sole author/writer of this Work;
- (2) This Work is original;
- (3) Any use of any work in which copyright exists was done by way of fair dealing and for permitted purposes and any excerpt or extract from, or reference to or reproduction of any copyright work has been disclosed expressly and sufficiently and the title of the Work and its authorship have been acknowledged in this Work;
- (4) I do not have any actual knowledge nor do I ought reasonably to know that the making of this work constitutes an infringement of any copyright work;
- (5) I hereby assign all and every right in the copyright to this Work to the University of Malaya (“UM”), who henceforth shall be owner of the copyright in this Work and that any reproduction or use in any form or by any means whatsoever is prohibited without the written consent of UM having been first had and obtained;
- (6) I am fully aware that if in the course of making this work, I have infringed any copyright whether intentionally or otherwise, I may be subject to legal action or any other action as may be determined by UM.

Candidate Signature:

Date:

Subscribed and solemnly declared before,

Witness’s Signature:

Date:

Name:

Designation:

ADOPTION OF MACHINE LEARNING ALGORITHM FOR ANALYSING SUPPORTERS AND NON-SUPPORTERS FEEDBACK ON POLITICAL POSTS ABSTRACT

Sentiment Analysis is a field that deals with the problem of identifying and extracting sentiment (or opinion) from data (particularly textual data). Studies have shown how user perception can have a strong influence on policies and decision-making processes in a place, society, and nation. This thesis is based on the application of sentiment classification algorithm to tweet data with the goal of classifying messages based on the polarity of sentiment towards a particular topic (or subject matter). Political analysts often communicate with the public and exchange information through the social media platform. Their activities (otherwise termed cyber-trooping) could have either positive, negative, or neutral feedbacks (perceptions) in the public space. Thus, there is a need to automate the process of identifying and predicting (positive, negative, or neutral class) these cyber-trooping data. This work employed the use of machine learning approach. Four conventional classification algorithms: naïve bayes (NB), support vector machines (SVM), nearest neighbor (k -NN), and decision trees (J48) classifiers are implemented in identifying and categorizing tweet data of three political figures in Malaysia: Dato Seri Anwar, Dato Hadi Awang, and Lim Guang Eng, as either positive, negative, or neutral perceptions. The method was implemented using Java and the results of the simulation were evaluated using five standard performance metrics: accuracy, AUC, precision, recall, and f -Measure. The support vector machines

(SVM) algorithm obtained the overall best results of 94.5% accuracy, 91.8% precision, 91.7% recall, and 91.1% *f*-Measure while the naïve bayes (NB) algorithm obtained the best AUC score of 0.944 with the tweet data of Dato Seri Anwar.

Keywords: Cyber-trooper, Perception, Sentiment, Twitter, Algorithms, Naïve bayes, Support vector machine, nearest neighbor, decision trees.

Universiti Malaya

ABSTRAK

Analisis Sentimen ialah bidang yang menangani masalah mengenal pasti dan mengekstrak sentimen (atau pendapat) daripada data (terutamanya data teks). Kajian telah menunjukkan bagaimana persepsi pengguna boleh mempunyai pengaruh yang kuat ke atas dasar dan proses membuat keputusan di sesuatu tempat, masyarakat dan negara. Tesis ini adalah berdasarkan aplikasi teknik klasifikasi sentimen untuk data tweet dengan matlamat untuk mengklasifikasikan mesej berdasarkan kekutuban sentimen terhadap topik tertentu (atau subjek). Penganalisa politik sering berkomunikasi dengan orang ramai dan bertukar maklumat melalui platform media sosial. Aktiviti mereka (atau disebut cyber-trooping) boleh mempunyai maklum balas (persepsi) sama ada positif, negatif atau neutral di ruang awam. Oleh itu, terdapat keperluan untuk mengautomatiskan proses mengenal pasti dan meramalkan (kelas positif, negatif atau neutral) data siber-trooping. Kerja ini menggunakan pendekatan pembelajaran mesin. Empat algoritma pengelasan konvensional: pengelas naïve bayes (NB), mesin vektor sokongan (SVM), jiran terdekat (k-NN), dan pokok keputusan (J48) dilaksanakan dalam mengenal pasti dan mengkategorikan data tweet tiga tokoh politik di Malaysia: Dato Seri Anwar, Dato Hadi Awang, dan Lim Guang Eng, sama ada persepsi positif, negatif atau neutral. Kaedah ini dilaksanakan menggunakan Java dan keputusan simulasi dinilai menggunakan lima metrik prestasi standard: ketepatan, AUC, ketepatan, ingat semula dan f-Measure. Algoritma mesin vektor sokongan

(SVM) memperoleh keputusan terbaik keseluruhan iaitu 94.5% ketepatan, 91.8% ketepatan, 91.7% ingat semula, dan 91.1% f-Measure manakala algoritma naïve bayes (NB) memperoleh skor AUC terbaik 0.944 dengan data tweet Dato Seri Anwar.

Kata kunci: Cyber-trooper, Persepsi, Sentimen, Twitter, Algoritma, Naïve bayes, Mesin vektor sokongan, jiran terdekat, pokok keputusan.

Universiti Malaya

ACKNOWLEDGEMENT

First and foremost, my sincere appreciation goes to my supervisor, Dr. Raja Jamilah Yusof, Faculty of Computer Science and Information Technology, University of Malaya. I thank her for being good to me during my thesis period whenever I ran into a trouble or had a question about my research or writing. She consistently allowed this paper to be my own work but steered me in the right direction whenever she thought I needed it.

Finally, I must express my profound gratitude to my parents and to my hall-mates for providing me with unfailing support and continuous encouragement throughout my years of study and through the process of researching and writing this thesis. This accomplishment would not have been possible without them. Thank you.

TABLE OF CONTENTS

TITLE	i
DECLARATION	iii
ABSTRACT	iv
ABSTRAK	v
ACKNOWLEDGEMENT	vii
TABLE OF CONTENTS	viii
LIST OF TABLES	xi
LIST OF FIGURES	xii
LIST OF SYMBOLS AND ABBREVIATIONS	xiii
CHAPTER 1 INTRODUCTION	1
1.1 Background	1
1.2 Motivation	4
1.3 Problem Statement	5
1.4 Research Aim and Objectives	6
1.5 Research Questions	6
1.6 Scope of Research	7
1.7 Research Methodology	7
1.7 Significance of Research	8
1.8 Overall Structure of the Thesis	8
CHAPTER 2 LITERATURE REVIEW	10
2.1 Introduction	10

2.2	Machine Learning: Overview	10
2.2.1	Definition	10
2.2.2	Concepts	11
2.2.3	Machine Learning: Grouping & Categories	13
2.2.4	Machine Learning Applications: Review	14
2.3	Data Classification Algorithm	19
2.3.1	Data Classification Algorithms	20
2.3.2	Research Works in Data Classification	22
2.4	Sentiment Analysis: Overview	26
2.5	Sentiment Analysis and Cyber-Trooping: Review	27
2.6	Chapter Summary	32
CHAPTER 3 METHODOLOGY		33
3.1	Introduction	33
3.2	Research Design	34
3.2.1	Pre-Research Stage	34
3.2.2	Data Collection	35
3.2.3	Data Pre-Processing	40
3.2.4	String To Word Vector	36
3.2.5	TF-IDF	36
3.3	Sentiment Classification	42
3.4	Performance Measures	42
3.5	Expected Outcome	45
3.6	Environment	45
3.7	Chapter Summary	46
CHAPTER 4 RESULTS AND DISCUSSION		47
4.1	Introduction	47
4.2	Classification Results Analysis	47
4.3	Chapter Summary	51
CHAPTER 5 CONCLUSION AND RECOMMENDATION		53
5.1	Introduction	53
5.2	Research Contributions	53

5.3	Research Limitations and Future Recommendation	54
5.4	Concluding Remarks	55
REFERENCES		56
APPENDICES		73

Universiti Malaya

LIST OF TABLES

2.1	Summary of Related Works	31
3.1	Confusion Matrixes in terms of TP, TN, FP, FN	43
4.1	Sentiment Classification Performance of the Classifiers with the Tweet Data of Dato Seri Anwar in terms of Accuracy, Precision, Recall, F-Measure, and AUC	47
4.2	Sentiment Classification Performance of the Classifiers with the Tweet Data of Dato Hadi Awang in terms of Accuracy, Precision, Recall, F-Measure, and AUC	48
4.3	Sentiment Classification Performance of the Classifiers with the Tweet Data of Lim Guang in terms of Accuracy, Precision, Recall, F-Measure, and AUC	49

LIST OF FIGURES

2.1	Data Classification Process	20
2.2	Sentiment Classification Algorithms	27
3.1	Proposed Research Framework	33
3.2	Research Design	36
3.3	Dato Seri Anwar's Tweets on Social	36
3.4	Dato Seri Anwar's Tweets on Politics	37
3.5	Dato Seri Anwar's Tweets on Religion	37
3.6	Dato Hadi Awang's Tweets on Social	38
3.7	Dato Hadi Awang's Tweets on Politics	38
3.8	Dato Hadi Awang's Tweets on Religious	39
3.9	YB LIM GUANG ENG's Tweets on Politics	40
3.10	YB LIM GUANG ENG's Tweets on Religion	40

LIST OF SYMBOLS AND ABBREVIATIONS

AI	-	Artificial Intelligence
ML	-	Machine Learning
ACC	-	Accuracy
TF-IDF	-	Term Frequency Inverse Document Frequency
AUC	-	Area Under (ROC) Curve
DT	-	Decision Trees
FN	-	False Negative
FP	-	False Positive
TP	-	True Positive
TN	-	True Negative
kNN	-	k Nearest Neighbor
NB	-	Naïve Bayes
SVM	-	Support Vector Machines

CHAPTER 1

INTRODUCTION

1.1 Background

The massive technological growth over the years has made the field of machine learning one of the mainstays of information technology (Ivanovic and Radovanovic, 2015). Machine learning as defined by Arthur Samuel is a field of study that gives computers the ability to learn without being explicitly programmed (Das *et al.*, 2015). The field of machine learning received much attention in recent years with the development of many successful machine learning applications such as data mining programs, information retrieval systems and autonomous vehicles.

In the field of Artificial Intelligence (AI), Machine Learning (ML) can be defined as the capability of an AI system to improve its performance over a time period through acquiring new knowledge and skills as well as its ability to reorganize the existing knowledge based on the newly acquired knowledge. Learning is considered as a parameter for intelligent machines to be able to make decisions in a more optimized form as well as work smoothly. The concept of ML is based on training machines to be able to detect patterns and adapt to new circumstances.

An emerging field in machine learning is Sentiment Analysis (SA) which is a study involving recognizing and interpreting opinions, sentiments, emotions, attitudes, and feelings related to subject matters, issues, or events. With the proliferation of information technology and the internet, there is continual increase in data, particularly exchange of communication among internet users on several contents ranging from politics to business, sports to education among others. Due to this massive exploration of the social media space, political enthusiasts (sometimes referred to as cyber-troopers) often convey their political world-views, ideas, emotions etc., via public social media platforms like twitter. Cyber-trooping is a social phenomenon among political parties in Malaysia. It is a well-known activity, conducted online to counter what they felt were unbalanced or ill-informed opinions circulating in cyberspace.

As earlier mentioned, political parties in Malaysia use social media as a way of communicating to the people. They believe through the internet information either positive, negative or neutral can be conveyed to the citizens. There are many issues pertaining to the use of the internet for the purpose of communicating opinion especially in social media with the desired goal of changing users' perception either positively, negatively, or neutral. Statistics from 2017 Malaysians Communications and Multimedia Commission (MCMC 2017) report on internet users, stated that, from 32 million people in Malaysia, 24.5 million are users, which is approximately (76.9%) have access to the internet while others 7.5 million approximately (23.1%) do not have access to the internet. Figures also showed that Malaysian households with computer and mobile access increased to 74.1 percent and 98.1 percent respectively, compared with 67.6 percent and 97.9 percent in 2015. At the same time, the share of internet users aged 15 years and older increased by nine percent to 80.1 percent in 2017, from 71.1 percent in 2015. The percentage of people using computers increased to 69.8

percent by 1.1 percentage points compared to 68.7 percent in 2015, while smartphones used for internet access increased to 97.7 percent compared to 97.5 percent in 2015. Therefore, there is a need to manage users' perceptions in relation to cyber-trooping activities.

Government statistics continued to show that the highest rates of internet penetration were found in the highly developed Klang Valley area, comprising the capital city of Kuala Lumpur (99.9%) and the state of Selangor (99.7%). In many urban areas, including shopping malls, restaurants, hotels and tourist destinations, free Wi-Fi is often provided. Penetration rates remained low in the under-developed, less populated eastern Malaysian states of Sabah (43.3 percent) and Sarawak (51.8 percent), where the majority of residents belong to indigenous groups. Government figures showed a slight gender imbalance in access rates, with men accounting for 59.4% of mobile and internet users. The age group with the largest share of users was between the ages of 20 and 24 (22 per cent). The average age of internet users (32.4 years) and non-users (50.7 years) showed an increase over the average of 2014 indicating that older age groups constitute an increasing share of the online community. Based on these statistics, it shows that more than the average Malaysians are on the Internet which makes it easy for cyber-troopers to perform their activities through social media channels.

In this work, sentiment analysis (SA) will be used to identify and categorize user opinion in reaction to different issues on social media. It determines a lot of social media decisions most especially when it comes to perception and decision. Sentiment Analysis grows exponentially because of the importance of automation in the extraction and processing of information to determine a person's general opinion. The

research work is based on employing machine learning algorithms for identifying cyber-trooping activities.

1.2 Motivation

According to the theory postulated by Albert Bandura which states that “when people observe a model performing a behavior and the consequences of that behavior, they remember the sequence of events and use this information to guide subsequent behaviors.” Perception is very important and relatively stronger than reality. It is a way of apprehending reality and experience through the senses, making it possible to discern figure, form, language, behavior and action. Individual perception influences the opinion, judgment, understanding of a person or a situation, meaning of an experience and how one responds. Perceptions are subjective, for example, when individuals or groups engage in an intersubjective dialogue or dialect where the potential for different interpretations exist, they often “see” entities quite differently based on different life contexts and contingencies. This “perceptual disparity” where two different subjective perceptions of the same event or experience are contradictory, occurs in the intersubjective space between two people or groups and can be a source of misunderstanding, injustice, and human conflict. The significance of human perceptions and their impacts are the motivation drives of this research project. Specifically, this study is motivated by the political activity related to Red Shirts leader Datuk Seri Jamal Yunus where he purported made an allegation against former Bersih chairman Maria Chin Abdullah, including linking her to terrorist group Islamic State (IS), based on report he read on social media. Jamal assumed based on online reports

that Maria's late husband Yunos Lebai Ali was involved in terrorist activities (Newspaper dated on 2016). He made the allegation without proper verification of the fact. With the help of sentiment analysis algorithm which could have been employed to check if the report of cyber-troopers was based on positive, negative, or neutral perception, possibly he might not have been a victim of such.

1.3 Problem Statement

The problem investigated in this research work is based on employing sentiment analysis algorithm and algorithms to identify cyber-trooping activities to manage users' perception on social media. There are many sentiment analysis algorithms available but not used for alerting system for cyber trooping activities in order to manage the situation of user perception on social media.

However, to do so, it is imperative for users to be aware of the various cyber-trooping activities that can change their perception. Based on the Malaysian internet users' statistics, facts established that most Malaysians have access to the internet. This makes it easy for cyber-troopers to perform their activities through several online (social media) platforms. Notable of such is the twitter platform which over the years has become the 'second home' and most desired choice of people around the world (in particular Malaysia) to convey, exchange, and communicate among one another on topics, issues ranging from politics to economy, businesses to sports etc.

In this age and time, with the advancement in information technology, machine learning, and data analytics, there is a great need to develop (model) predictive algorithms capable of self-learning to automate the prediction of users' perception on matters/issues on social media (particularly twitter). This will facilitate and support

government decision making policies and other stakeholders including academics, jurists, and the general public in having a broader and clearer understanding of matters, especially when controversies, disputes erupt among people (as a result of cyber-troopers), which if not properly managed could become chaotic and damaging. Thus, this study attempts to address the cyber-trooping problem in Malaysia political space by employing the use of sentiment analysis algorithm. To the best of our knowledge, there is little/no existing work in this direction of automating the prediction/identification of users' perceptions on data associated with Malaysian political figures.

1.4 Research Aim and Objectives

The aim of this research is to employ sentiment analysis for the adoption of machine learning algorithm for analysing supporters and non-supporters feedback on political posts. To achieve this aim, the study proposed the following stated objectives:

1. To employ a sentiment analysis (SA) algorithm to manage users' perception against cyber-trooping activities in Malaysia.
2. To apply the algorithm using four machine learning algorithms: naïve bayes (NB), support vector machines (SVM), nearest neighbor (k-NN), and decision trees (J48) on selected political tweets.
3. To analyze the sentiment analysis results using the five conventional performance measures available.

1.5 Research Questions

This study addressed these research questions.

RQ1 – What are the sentiment analysis used in literature to manage users' perception against cyber-trooping activities in social media? RO1

RQ2 – How to implement naïve bayes (NB), support vector machines (SVM), nearest neighbor (k-NN), and decision trees (J48) on selected political tweets. RO2

RQ3- What are the political tweets among political parties in Malaysia to be used for sentiment analysis? RO2

RQ4- How to evaluate and analyze the simulation results using five conventional performance measures: accuracy, precision, recall, *f*-Measure, and AUC methods.? RO3

1.6 Scope of Research

The research work is about employing machine learning algorithms to identify cyber-trooping activities in Malaysia political space. The study implements four of the conventional machine learning algorithms: naïve bayes (NB), nearest neighbor (k-NN), support vector machines (SVM), and decision trees (J48), to perform the sentiment analysis task of the tweet activities of three political figures in Malaysia namely Dato Seri Anwar, Dato Hadi Awang and Lim Guang. However, neutralizing the perception is not the scope of this study.

1.7 Research Methodology

The methodology to conduct this research will focus on three modules:

- Data preprocessing module for preprocessing data
- Feature representation module which will be used to extract data from tweet
- Sentiment classification which will be used for classify sentiment analysis

1.7.1 Data processing

This is responsible to decrease the size of the feature set to make it suitable for learning algorithms. It is required because the proposed methodology of tweets selected may contain several features. In order to be able to determine high quality data, it is crucial to preprocess the data. For the process to be easier, data preprocessing is divided into four categories: data cleaning, data integration, data reduction, and data transformation.

1.7.2 Features Representation

This is responsible for extracting features from preprocessed tweets. Bag-of-Words algorithm deployed to convert training tweets into numeric representation. Bag-of-Word will enable to learn a vocabulary of known words from all of the tweets. Once learning vocabulary is done, Bag-of-Word will describe the presence of known words within a tweet.

1.7.3 Sentiment Classification

Conventional classification algorithms such as naïve bayes, decision trees etc., will be used for classifying users' perceptions.

1.8 Significance of Research

The proposed study will be useful in the task of identifying elements of cyber-trooping on social media. The novelty in this work is the practical applicability of sentiment analysis algorithms and algorithms for identifying cyber-trooping activities in Malaysia political space.

1.9 Overall Structure of the Thesis

The thesis covers introductory chapter, review of literature, research methodology, algorithm design and implementation, algorithms evaluation and validation, results, discussion and conclusion. Chapter one contains introductory research background, motivation, problem statement, objectives of the research, research questions, scope of the research, summary of methodology, target group, thesis structure, and the significance of the thesis. Chapter two elaborates the literature review of the research findings in related works, identifying research gaps, and concept of cyber-trooping. The methodology of this research has been elaborated in chapter three. The experimental results are presented and analyzed in chapter four while the study concludes in chapter five with direction to future works.

CHAPTER 2

LITERATURE REVIEW

2.1 Introduction

A literature review seeks to describe, summarize, evaluate, clarify and/or integrate the content of primary reports. This chapter covers the overview of important fields related to the research study found in literature as well as the different algorithms and algorithms commonly used in the study domain. The existing related works are reviewed, analyzed and summarized.

2.2 Machine Learning: Overview

This section presents an overview of the field of machine learning with focus on important concepts and methods. Furthermore, recent machine learning works and applications found in literature are reviewed. Generally, the social media has provided the internet users with the platforms to publish their written and multimedia contents, express their feelings and emotions about particular subjects via the internet. However, some users, including the politicians have abused these platforms by performing various acts such as pushing some ill-informed opinions to the public through the social media. Therefore, cyber-trooping is used to create a balanced opinion within the cyberspace. Since the internet and the virtual community have caused severe negative

consequences to the welfare of society, and creating social problems such as cyber-aggression, or in some cases called cyber-bullying. The accurate classification of some of these comments using machine learning can help to take measures to diminish this phenomenon. Although multiple works have focused on detecting these phenomenon, for other applications (Sanur et al. 2019; Ghulam et al. 2021; Chong19; Sintaha et al. 2018; Balakrishnan et al. 2020). Moreover, these works used limited number of methods and datasets. Furthermore, there is a lack of datasets that are concerned with this topic. We propose the use of Machine Learning to detect and analyze the phenomenon. We apply various classification algorithms to the dataset, and we use various sentiment analysis. To evaluate the performance of the study, linear and probabilistic classifiers are used.

2.2.1 Machine Learning Definition

Machine learning is a field that focuses on the design and development of algorithms (or models) for the purpose of making predictions (from data). Machine learning is often highly regarded as a multi-disciplinary field covering important subject and application areas including (*but not limited to*) computer science, mathematics, statistics, robotics, information security, and cognitive science. Conventional among learning algorithms include artificial neural networks (ANN), decision trees (J48), and nearest neighbors (*k*-NN). The term '*learning*' as often used in machine learning (Arunakrathi et al., 2020) refers to the process where a system (or machine) learns from experience, and then utilize the knowledge to make decisions (predictions) from input data. Machine learning algorithm is suitable for recognizing and storing word

patterns easily for a certain sentiment class in an opinion sentence. Therefore, it is appropriate to classify sentiment data properly using this method (Ghulam et al. 2021).

2.2.2 Machine Learning Concepts

There are some standard concepts that are mostly adopted in machine learning application works. These include observation, data representation, attributes (features), label, model, and evaluation measures.

A. Observation: In machine learning, an observation (also known as data point) is an instance used either for training and/or testing (evaluating) a machine learning model. It could be in synthetic form (e.g., iris data), or real-world scenario (e.g., a cancerous patient).

B. Data Representation: It is a standard procedure in machine learning to transform, process, and analyze experimental datasets. The input data must be machine readable and could be of different types:

- Numeric, such as: integers [e.g., -16,0,1,198, ...]; real [e.g., 17.65,199.99, ...]
- Categorical, [e.g., race, age group, sex, grade, ...]
- Binary, [000_s and 111_s]
- Text data

C. Attributes (or Features): A data sample is often associated with one or more variables. These variables are otherwise termed attributes or features. For example, working with real-world data involving cancerous patients, common attributes that

could be associated with the data are variables such as: age, gender, blood group, etc. Another example is the well-known iris dataset from UCI machine learning repository (<https://archive.ics.uci.edu/ml/datasets/Iris>). It has 4 attributes: sepal length, sepal width, petal length, and petal width. Attributes could be of any type, however in most machine learning models, the numerical data type is often required. The attributes could exist in discrete form (e.g., integer values), or continuous form (e.g., real numbers).

D. Label (or Class): A good experimental dataset for a machine learning problem is recommended to have one special variable/attribute known as the class or label. A class attribute is the value to be predicted using a machine learning model. Depending on the simulation environment such as Matlab, Weka, or Python, the class label is often the last column (attribute value) of a particular observation (data instance). Likewise, a class label could either be an integer or real value.

E. Model: A machine learning model is a function that takes attributes vector as input and uses an algorithm to predict the output, which is the label for a particular data instance. As mentioned earlier, some of the conventional learning algorithms often used to construct machine learning models include decision trees (Ghiasi & Zendejboudi, 2020; H. Lu & Ma, 2020; Vu et al., 2021), neural networks (Christou et al., 2019; DeLatta et al., 2019; L. Li, 2021; Y. Park & Yang, 2019; Rezaei-Ravari et al., 2021), logistic regression (Belciug, 2021; Dedetürk & Akay, 2020; Dzhamtyrova & Kalnishkan, 2020; Fan et al., 2020), random forests (Arora & Kaur, 2020; Matsuo et al., 2020; Simsekler et al., 2020; Z. Xu et al., 2020), support vector machines (Q. Gu et al., 2021; H. Jiang et al., 2020; Moreira et al., 2019; Nanglia et al., 2020; Rakhmetulayeva et al., 2018; Ryabtseva & Skomorokhov, 2020; Soumaya et al.,

2021), and nearest neighbors (Arian et al., 2020; Hamed et al., 2020; Jianyun Lu et al., 2018; Z. Pan et al., 2020; Talavera-Llames et al., 2019).

F. Evaluation Measures: To validate effectively a machine learning model's performance, some evaluating measures (also referred to as performance metrics) are employed. This is an important critical stage of a machine learning project, evaluating the performance of the experimenting methods justifies how efficient and accurate the methods predict the class labels. In addition, performance metrics help to determine if a model needs optimization. There are several standard performance measures commonly applied in machine learning projects such as accuracy, precision, recall, AUC, ROC curve, PR curve, F-measure, and error rate.

2.2.3 Machine Learning: Grouping & Categories

Machine learning projects are commonly grouped into three categories. These are: supervised ML, unsupervised ML, and semi-supervised ML.

- **Supervised Machine Learning:** The task of supervised machine learning (D'hooge et al., 2020; Hajji et al., 2020; Le et al., 2020; Puranik et al., 2020; Reyes et al., 2020) is to predict the output of an observation whose class membership is known (predefined). The class membership (label) could be of real-value (such as regression), or discrete value (such as classification). Input variables could be of any type such numerical or categorical.

- ***Unsupervised Machine Learning:*** In unsupervised machine learning problems (Hyun et al., 2020; Reis et al., 2019; Watson, 2020; Yassine et al., 2021), the class memberships to be predicted are unlabeled (unknown). Given an observation, an unsupervised machine learning model makes predictions by finding hidden structures (or similarities) from the sample data. A typical unsupervised machine learning application is clustering.
- ***Semi-supervised Machine Learning:*** In this type of machine learning, a hybridization approach is adopted. It involves learning from few labeled samples combined with large unlabeled samples. The goal in semi-supervised ML is to make effective use of all available data (labeled or unlabeled), hence improving the accuracy of supervised learning (which is based only on labeled data). Semi-supervised based applications (Dunham et al., 2020; Mukherjee & Prasad, 2020) include natural language processing (NLP), automation speech recognition, computer vision etc.

2.2.4 Machine Learning Applications: Review

As previously noted, the goal of machine learning is to train a system for decision making purposes without being explicitly programmed. Over the years, there have been continual increasing growth and developments in the field of artificial intelligence (AI) and machine learning (ML). Several algorithms and methods have been developed to solve both simple and complex problems. Machine learning algorithms and models have been applied extensively and exhaustively in numerous real-world applications such as: *image recognition and classification, medical*

diagnosis, online fraud detection, automatic natural language processing, speech recognition, traffic alerts and prediction, weather forecasting, autonomous driving, stock market trading, virtual personal assistant, web spam filtering, robotics, big data analytics, sports and gaming etc.

A new classification approach based on low-dimensional feature vector (LDFV) was proposed by Amiri, et al., (2020) for classification of HIS images. The dimensionality reduction method is based on reducing the feature space while retaining relevant spectral information. The proposed method was combined with conventional classifiers: SVM, k -NN. It showed significant reduction in processing time. In the work of Adegun and Vadapalli (2020), an extreme learning machine (ELM) method was applied along with a support vector machine (SVM) for micro-expression recognition. Furthermore, feature extraction methods: local binary pattern (LBP) and local binary pattern on three orthogonal planes (LBP-TOP) were applied on apex micro-expression frames and micro-expression videos respectively. The work concluded showing that the application of machine learning algorithms in automating micro-expressions recognition is faster and obtained better results. Also, the work of Geng, et al. (2020) proposed a multi-scale deep feature learning network with bilateral filtering (MDFLN-BF) for the synthesis of aperture radar (SAR) image classification. The proposed method was used to extract discriminative features and reduce labeled samples requirements. In the results, the classification method performed better compared to related methods. In the work of Guo and Yuan (2020), a semi-supervised learning method with adaptive aggregated attention (AAA) was proposed for automating wireless capsule endoscopy (WCE) image classification. The proposed method was used in minimizing discriminative angular (DA) and Jensen-Shannon divergence (JS) losses of both labeled and unlabeled data. The method achieved a high

accuracy of 93.17%. In addition, a deep learning model based on mobileNet was proposed in the work of Pan et al. (2020). The proposed model functioned as a welding defect feature extractor. It achieved about 97% prediction accuracy result. Other recent machine learning applications for image classification and recognition problems include (Bulat et al., 2019; Cen et al., 2021; Hénaff et al., 2019; F. Luo et al., 2020; Qin et al., 2020).

In speech recognition, Zhang (2020) proposed a hybrid system based on logistic regression and WBCS algorithms. The proposed system was used in automating speech emotion recognition (SER). Subsequently, a random forest learning model was further applied for feature preprocessing. The experimental results showed the proposed hybrid system achieved satisfactory results. In the work of Yao et al., (2020), a framework was proposed based on three classifiers: deep neural network (DNN), convolution neural network (CNN), and recurrent neural network (RNN). Three predictive models were constructed: low-level descriptors recurrent neural network (LLDs-RNN), mel-spectrograms convolution neural network (MS-CNN), and high-level statistical functions deep neural network (HSF-DNN). The models were used for the categorical recognition of four discrete emotions: angry, happy, neutral, and sad. Tuncer et al., (2021) proposed a nonlinear multi-level feature generational model based on cryptographic structure. The proposed method was applied for feature generation and selection. The work further implements three sub-methods: a multi-level feature generation using Tunable Q wavelet transform (TQWT), twine shuffle pattern for feature generation, and an iterative neighborhood component analysis (INCA). The TQWT method was used to generate high-level, medium-level, and low-level wavelet coefficients. Twine shuffle pattern method was used to extract features from the decomposed wavelet coefficients. Finally, INCA was used to select

significant features from the experimental datasets. The results showed good predictive performance with the proposed methods. In the work of Langari et al., (2020), an autonomous speech emotion recognition system (SER) was designed for identifying different emotional classes through the extraction and selection of effective features from speech signals. Results showed significant performance across the datasets experimented. For further study, some of the recent machine learning applications in speech recognition could be found in (Bandela & Kumar, 2021; He & Dong, 2020; Issa et al., 2020; Kwon, 2020; Linhui Sun et al., 2019).

In medical diagnosis, Jabir et al., (2020) proposed a new similarity-based measure for picture fuzzy sets (PFSs). The proposed method, comprising of two parameters: level of uncertainty (t) and L_p norm, was used to diagnose several predefined medical problems such as: malaria, chest pain, cough, typhoid, heart problem among others. Jiang et al., (2020) proposed a recursive neural knowledge network (RNKN) for multi-disease diagnosis. It combined medical knowledge based on first order logic with a recursive neural network. The results showed the proposed method had good diagnosis accuracy. Other recent works in medical diagnosis include (M. Chen et al., 2020; A. Das et al., 2019; Ker et al., 2019; Z. Li et al., 2020).

Also, machine learning algorithms and algorithms have been applied in financial fraud detection with a growing number of works. Bagga et al., (2020) used conventional machine learning algorithms in detecting credit card fraud data. The ML algorithms are logistic regression, k -nearest neighbors, naïve bayes, random forest, ada boost, multilayer perceptron, ensemble learning, pipelining, and quadrant discriminative analysis. The results obtained from the work showed the pipelining method had the best performance among the classification models. A cascaded forest

and XGBoost (CFXGB) model was proposed in the work of Thejas et al., (2020) for detecting click fraud in commercial industries. The model combined two learning methods for feature transformation and classification. Results showed that the proposed CFXGB model achieved better performance across multiple click-fraud datasets. Wang and Chen (2020) work was on detecting fraud campaigns in product reviews. To achieve this, an online monitoring algorithm was proposed. The method works in two phases: monitored online reviews to generate most abnormal review subsequences (MARSSs), then conditional random fields are employed to label each review as fake or genuine. The fraud detection method proved to be effective and efficient. Among other recent works in this category include (Askari & Hussain, 2020; Patil et al., 2018; Rtayli & Enneya, 2020).

In natural language processing (NLP), Ochieng (2020) proposed a framework for converting natural language to a semantic query language for databases (SPARQL). The work used a query tool called PAROT to handle user's queries. It employed a heuristic-based algorithm for converting user's queries to user's triples. The user's triples are further processed into ontology triples, and finally into SPARQL queries. The experiment showed the proposed method achieved satisfactory results. Also, Nawaz et al., (2020) proposed a framework based on local weights (LW) and global weights (GW) for modeling Urdu language. Vector space model (VSM) was used as the baseline framework for sentence weighting. Results showed satisfactory results with the local weight model. Further recent works in NLP are found in (Brouer & Benabbou, 2019; Cattoni et al., 2021; Shengxue, 2020; Tang et al., 2020).

In sports, Daud and Abbasi (2021) employed the use of machine learning algorithms for rising star players' prediction in the game of basketball. The work applied support vector machines (SVMs), naïve bayes (NB), classification and

regression trees (CART), Bayesian network (BN), and maximum entropy Markov model (MEMM) machine learning models. The results showed that the application of machine learning was very effective and useful in predicting rising star players. Also, a deep neural network model (VGG-16) was proposed in Rangasamy et al. (2020) for the automation of activity recognition in hockey sport. Likewise, Lei, and Yuenian (2020) proposed the use of particle swarm optimization (PSO) in automating sports image detection. For further study, some of other recent machine learning applications to real-world problems could be found in (Fialho et al., 2019; Soltani & Morice, 2020; D. Yang, 2021; Zhou et al., 2021).

2.3 Data Classification Algorithm

Classification is a classical data mining algorithm based on machine learning concepts. This algorithm is used in classifying items in a set of data into one of predefined sets of classes or groups. Classification algorithm is a supervised learning approach that is used to predict the target class for each data point in a data sample.

There are two main steps (phases) in data classification as shown in Figure 2.1. The first step is the model construction used in analyzing the training dataset of a database and the second step is model usage where the constructed model is used for prediction. In a classification task such as text classification, the accuracy of the classification algorithm is estimated based on the percentage of test samples in the dataset correctly classified.

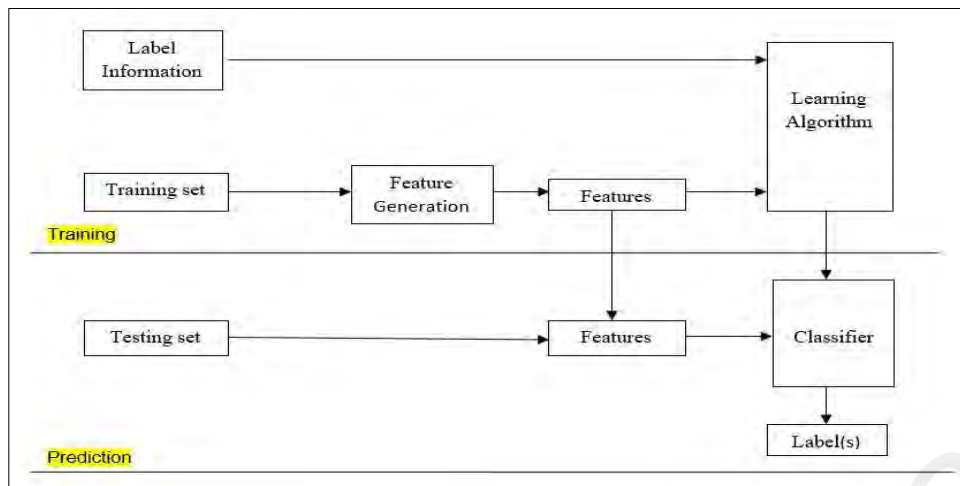


Figure 2.1: Data Classification Process (Reyes et al., 2020)

2.3.1 Data Classification Algorithms

A growing number of data mining algorithms have been applied to text classification problems, including the *Bayes probabilistic approach* (Tang et al., 2020), *decision trees* (Ghiasi and Zendehboudi 2020), *neural networks* (Christou et al. 2019), *support vector machines* (SVM) (Q. Gu et al. 2021; H. Jiang et al. 2020; Moreira et al. 2019; Soumaya et al. 2021), and *k-nearest neighbor* (Arian et al. 2020; Hamed et al. 2020).

The *k*-NN classifier is an instance-based learning algorithm that has shown to be very simple but effective for text classification problems (Arian et al., 2020). It is a non-parametric method used in classification and works by calculating the Euclidean distance between points. In classifying a new document x , the algorithm ranks the document's neighbor in the training set, and then uses the class of k most similar neighbors to predict the class of a new document (also known as majority vote). The Euclidean distance is given as:

$$d(x, x_i) = \sqrt{\sum_{i=1}^n (x_j - x_{ij})^2} \quad (2.1)$$

where x is the new point, x_i is the existing point across all input attributes j .

The naïve bayes classifier greatly simplifies learning by assuming that features are independent given class and has proven effective in many practical applications, including text classification. The classifier is a simple probabilistic model based on the Bayes rule. Given a class C , the probability of a particular document d to belong to C is given as:

$$P(C_i | d) = \frac{P(d | C_i) * P(C_i)}{P(d)} \quad (2.2)$$

SVM is one of the most widely used and applied classification methods. It has been successfully applied to many application domains. SVMs are typically used for learning classification, regression, or ranking functions. The algorithm works by searching a separating hyperplane to separate between samples with a maximal margin (Jiang et al., 2020, Moreira et al., 2019). The separating hyperplane is:

$$w^T x + b = 0 \quad (2.3)$$

To classify an unseen document d , the sign of $w^T x + b$ must be known. This is further shown as:

$$w^T x_i + b \geq 1 \text{ or } w^T x_i + b \leq -1 \quad (2.4)$$

Decision tree is one of the most popular and powerful approaches in data mining used to extract knowledge by making decision rules from a large amount of available information. The algorithm is a tree-like structure which classifies an input sample into one of its possible classes (Vu et al., 2021). Decision tree is a way of

representing a sequence of rules that leads to a class or value. It consists of three fundamentals: root node, internal node, and leaf node. In decision tree classification algorithm, each node specifies a test to be performed on a single attribute. The goal is to create a model that predicts the value of a target variable based on several input variables.

The algorithm can further be described as the combination of mathematical and computational algorithms to aid the description or categorization of a given set of data. The data generally takes the form:

$$(x, Y) = (x_1, x_2, x_3, \dots, x_k, Y) \quad (2.5)$$

where Y is a dependent variable to be classified or predict, x is a vector with input variables $x_1, x_2, x_3, \dots, x_k$ to be used for the classification of Y .

2.3.2 Recent Works in Data Classification

The existence of textual data, be it in offline or online forms, have provided us with a mass amount of information. Due to the increasing rise in information, the study in automated text classification has drawn researchers from many artificial intelligence areas.

The work of Chen *et al.*, (2021) was based on a modified term weighting scheme with distinguishing feature vector (DFS) for text classification problems. The study identified an issue with the existing term weighting scheme which is lack of effectiveness in assessing the distribution of information (terms) across the training documents. The proposed TF-MDFS method was evaluated using 19 multiclass text datasets. Results showed the proposed method achieved the best results in terms of classification accuracy.

A three-stage hybrid text classification framework was proposed by Li *et al.*, (2021) for medical text classification problems. The hybridized method combined gated attention-based bi-directional long short-term memory (ABLSTM) and regular expression-based classifier. The combined classification framework allows the recurrent neural network (RNN) to weigh words effectively in order of importance. Results on real-world medical query data showed that the proposed method proved to be effective in selecting domain-specific and topic-related features.

Luo (2021) worked on the implementation of support vector machines (SVM) for classifying English textual documents. The results showed the SVM classifier outperformed other machine learning methods implemented. Godavarthi and Sowjanya (2020) used kNN, MLP, and XGBoost machine learning classifiers for automating the classification of abstract texts related to covid 19 pandemic. Similarly, nearest neighbor (kNN) classification algorithm was implemented in Chen *et al.*, (2020) for real-world classification of Lao news text. Also, a filter-based feature selection method was proposed by Cekik and Uysal (2020) for short text data classification using rough set theory. The results showed the feature selection method proved to be competitive with other dimensionality reduction methods.

A lazy fine-tuning naïve bayes (LFTNB) method was proposed by ElHindi *et al.*, (2020) to address classification performance issues associated with the standard naïve bayes (NB) algorithm. The proposed method uses the nearest neighbors of an instance query to get the probability estimation used by NB classifier. UCI datasets were used to evaluate the method and the results showed the proposed LFTNB outperformed the classical NB algorithm. Similarly, in the work of Kolluri and Razia (2020), naïve bayes classification algorithm was employed for text classification problems.

The work of Tang *et al.*, (2020) was based on exploring term weighting methods for text representation and classification. The study identified an existing problem associated with text representation. The authors experimented with four term weighting methods. Results showed the proposed methods have the potential of improving the performance of text classification algorithms. In addition, a novel term weighting scheme was proposed by Dogan and Uysal (2020) for text classification problems. The proposed strategy was based on assessing non-occurrence information of terms. The method performs intra-class document scaling in order to have better terms representation across documents. Two standard classifiers SVM and kNN were used to test the proposed scheme and the results showed the proposed weighting scheme outperformed existing schemes across the experimental datasets.

Li *et al.*, (2020) worked on improving convolutional neural network (CNN) for text classification using recursive data pruning. The study identified that CNNs ignore filtering some irrelevant words which may lead to unsatisfactory classification results. The proposed pruning method was used to evaluate the features generated at the pooling layer. The experimental results showed the proposed model outperformed the baseline CNN model. Also, the research work of Kim *et al.*, (2020) was based on exploring the use of capsule networks for text classification problem. Capsule based networks have shown to be effective over the traditional convolutional neural networks. In their work, a simple routing method was recommended to effectively reduce the computation complexity of dynamic routing.

Rohidin *et al.*, (2020) worked on associates rules of fuzzy softset-based classification. The study identified two setbacks associated with classifiers including more processing time and lower accuracy performance. The authors proposed a new classification model called class-based fuzzy soft associates (CBFSA). The proposed

CBFSA model is a hybrid of associates rules and fuzzy softset model. Results showed the proposed model is more accurate and efficient compared to other classifiers.

Zhan *et al.*, (2020) proposed an improved sandwich neural network (SNN) model for extracting semantic and structural representations of textual documents. In addition, a knowledge attention method was proposed to fuse external semantic and structural knowledge in order to improve the performance of text classification algorithms. The results showed the proposed method significantly reduced the structural complexity of attention models.

An end-to-end entity classification system was proposed in the work of Li *et al.*, (2019). The proposed system was based on a neural network classification method. In addition to the entity-based classification system, a fusion model was proposed to fuse the entity types found in the sentence documents. The results showed the proposed system proved to be effective for entity-based classification tasks.

Hu *et al.*, (2016) worked on active learning for text classification with reusability. The research work focused on exploring the reusability problem in text classification by measuring the impact of using different classification algorithms in the active learning process and also in the classification applications that make use of the outcomes of the process. The authors experimented four text classification problem scenarios using three of the common text classification algorithms namely: Naïve Bayes, Support Vector Machine, and Nearest Neighbor classifiers.

Saikrishna *et al.*, (2016) worked on statistical compression-based models for text classification. In their research, two approaches were proposed for text classification using the Minimum Message Length (MML) and Probabilistic Finite State Automata (PFSAs). The proposed approaches were experimented on the Enron Spam dataset and the result showed a good classification accuracy. Other researches

in text classification include (Brindha *et al.*, 2016; Sharma and Singh, 2016; Zewen, 2016).

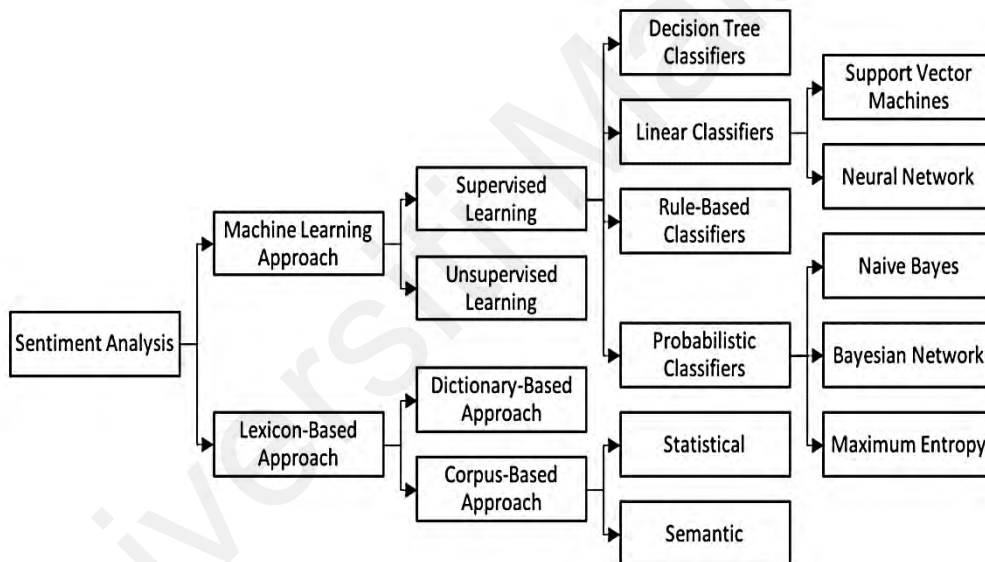
In addition, text classification has been studied to minimize the problem of information overload phenomenon in various problem domains, including news categorization (Mulahuwaish *et al.*, 2020; Huang and Chen, 2020), web searching (Correa *et al.*, 2020; Makkar and Kumar, 2020; Buber and Diri, 2019), medical document indexing (An *et al.*, 2021; Zhao *et al.*, 2021; Faris *et al.*, 2021; Timsina *et al.*, 2016), and sentiment analysis (Sazzed and Jayarathna, 2021; Vashishtha and Susan, 2020; Ullah *et al.*, 2020; Goyal, 2016). The information overload phenomenon is also present in understanding content of the Holy Quran (Adeleke *et al.*, 2018).

2.4 Sentiment Analysis: Overview

Sentiment Analysis (SA) is an emerging technology that deals with the study of peoples' emotions, opinions, relationships, and behaviors (Sazzed and Jayarathna, 2021; Vashishtha and Susan, 2020; Ullah *et al.*, 2020). In practice, Psychologists (experts trained to study human behavioral patterns) apply sentiments process through the concept of hypothesis, whereas data scientists could extract useful and meaningful patterns from people through data. Thus, the algorithm of sentiment analysis can be referred to as the computational process of identifying and categorizing opinions, thoughts, and ideas through textual data.

Sentiment analysis are expressed in two different categories: polarity and subjectivity. The polarity-based sentiment analysis measures textual data such as twitter data and classify them (using classification algorithms) into either positive, negative, or neutral. On the other hand, the subjectivity-based SA approach is the

classification of a sentence either subjective or objective. It measures from 0.0 to 1.0, where 0.0 index means very objective and 1.0 means very subjective. Sentiment classification is basically a text classification problem. There are two strategies that could be applied in sentiment classification: (a) applying a standard supervised machine learning algorithm and (b) using a classification method designed specifically for sentiment classification. This study work employed the polarity-based sentiment analysis approach in classifying twitter data of selected Malaysia politicians as either positive, negative, or neutral. To do so, standard supervised machine learning algorithms are applied to manage users' perception.



2.5 Sentiment Analysis and Cyber-Trooping: Review on Users' Perception

Yadigar N. Imamverdiyev (2015) reported on the issues involved in developing cyber-troops. The work studied the experience of developed countries and the international military organizations in this field using publicly available data. In the

study, the main aspects of cyber-troops' formation, objectives and functions; the structural and organizational models of cyber-command; and cyber-troops' weapons arsenals and human potential. The study only incorporated the international cooperation issues in the development of cyber-troops to manage users' perception to ensure the sovereignty of each country in cyberspace.

In some previous reports, when the activities of cyber-troopers were reported in 2004, some teams of cyber activists were mobilized by the youths of a political group known as UMNO, which is at that time the largest component of the Barisan National party. They were believed to have operated on various social media platforms, including blogs, mailing lists etc (Cheong, 2020). It was observed by the UMNO that cyber-troopers were hired to react to what they perceived to be lies, slanders and false allegations (Tan 2013); in short, they were mobilized to balance out the opposition playground.

Ana Valdiva, M.Victoria Luzon, Eric Cambria and Francisco Herrera (2015) worked on models for detecting neutrality guided by consensus voting among sentiment analysis methods prior to users' opinion classification step. In their research paper, they included the neutrality proximity function that assigns weights to polarities according to its proximity to the neutral point. They conclude by introducing two polarity aggregation models based on a Weighting Average using the proximity function as well as Induced Ordered Weighted Averaging guided by linguistic quantifiers to represent the majority concept respectively.

Dimitrios E. Pournarakis, Dionisios N. Sotiropoulos, George M. Giaglis (2017) introduced a new method for eliciting influential factors that govern brand equity assessment by mining and analyzing consumer perceptions from online social network

data. They designed science-based research approach and described the design and evaluation of the computational model that lays out a proposed method on how consumer perceptions could be detected, starting from a marketing perspective.

Ankit, Nabizath Saleena (2018) discussed the issue of identification of the most suitable sentiment classifier that can correctly classify the tweets in their articles “An Ensemble Classification System for Twitter Sentiment Analysis”. They proposed ensemble classifiers and the method performed better than stand-alone classifiers as well as majority voting ensemble classifiers on different types of data.

Marouane Birjali, Abderrahim Beni-Hssane, Mohammed Erritali (2017) reported on Machine Learning and Semantic Sentiment Analysis based Algorithms for suicide Sentiment Prediction in Social Network to address the lack of terminological resources related to suicide. The research work experimented with Machine learning algorithms, semantic sentiment analysis and Weka as a tool. Their method based on machine learning algorithms and semantic analysis can extract predictions of suicidal idea using Twitter data. This work verifies the effectiveness of performance in terms of accuracy and precision on semantic sentiment analysis that could thinking of suicide.

Ana Valdivia, M. Victoria Luzon, Erik Cambria Francisco Herrera (2018) introduced consensus vote models for detecting and filtering neutrality in sentiment analysis on empower neutrality by characterizing the boundary between positive and negative review with the goal of improving the model’s performance. They experimented using Bing, Vader, CoreNLP, Microsoft Azure, SentiStrength. They concluded in their result that detecting neutrality based on a consensus improves classification precision. The ALH-Pron model get the best result on average. It weighs the polarity of 3 out of 6 less extreme SAMS.

Aggregation methods outperformed single models in most cases, which led us to conclude that neutrality is key for distinguishing between positive and negative for improving sentiment classification. Demitrios E. Pournarakis, Dionisios N. Sotiropoulos, George M. Giaglis (2017) introduced a computational model for mining consumer perceptions in social media by using computational model that combined topic and sentiment classification to elicit influential subject from consumer perception in social media. Sample uber transportation through twitter is used as a case study. Their methods are based on prerequisites, data collection and preparation, VSM-based corpus vectorization, sentiment classification using support vector machines (SVMs), LDA, topic modelling, GA clustering sentiment per topic, output IT metric, and consumer perceptions. The results obtained presented consumer perceptions and produce insights for two fundamental brand equity dimensions: brand awareness and brand meaning. Simultaneously, they improve clustering results, in comparison to the K-means approach.

Raja-Jamilah Raja Yusof, Azah-Anir Norman, Siti- Soraya Abdul-Rahman, Nurul'addilah Nazri, Zulkifli Mohd-Yusoff (2016) studied cyber-volunteering, social media affordance in fulfilling NGO social missions by analyzing the issue on voluntary behavior through social media usage based on affordance theory. They used thematic coding based on identified key social media affordance: Visibility, Editability, Persistence, Virtual collaboration, Synthetic representation, Individual, and Collective affordance by interview on the basis of convenience sampling. The findings demonstrated that social media affordance related to cyber-volunteering are achieved through promoting, training, fundraising, knowledge sharing, and problem-solving activities. The affordance is highly influenced by cyber-volunteering behavior through work culture and personal privacy. The collective, individualized and visibility

affordance are more associated with cyber volunteering behavior than persistence, virtual collaboration and editability.

Gillian Warner-Soderholm, Andy Bertsch, Everlyn Sawe, Dwight Lee, Trina Wolfe, Josh Meyer, Josh Engel, Uepati Normann Fatilua. (2018) carried out a study on Computer in Human Behavior: Who trust Social Media? Sampling LinkedIn, and Instagram on What degree a users' perception of trust varies depending on gender age or amount of time spent using social media. Tools on Convenience population sampling (n=214) comes from university student and staff because they tend to be more frequent internet and social media users. Five different validated scale to measure: Benevolence, Competence, Integrity, Identification and Concern using Age, News and Social media preference. The findings demonstrated that Woman and younger users have the highest expectations for integrity and trusting others and expecting others to show empathy.

Rosyidah Muhamad (2010) studied Political Blogging and the Public Sphere in Malaysia and of deliberative democracy through a case study of Najib's blog. Her findings showed that the willingness of the government to publish information related to policy would be one element of deliberative democracy. However, Najib avoided critical and political debate in his blog and also selects comments that are not critical by readers which contradict the principle of deliberative democracy. If the use of the internet is to contribute to the public sphere and promote the principle of deliberate democracy, political blogging must be substantive not as in the case of Malaysia Prime Minister: Najib.

This research aims to investigate how the public is responding to the threat of digital technologies being used to undermine democratic processes. While propaganda and the influencing of political decisions are nothing new, the ever-increasing use of

digital technologies has made the issue of misinformation and the subverting of political discourse a pressing issue (Grigsby 2017).

The existence of false and falsified information has become a societal issue beyond the narrow topic of elections. About 15% of Twitter users in 2017 were estimated to be bots (Prier 2017), while the use of internet technologies to propagate false claims and narratives have been used as a political tool for altering narratives around incidents as diverse as the Salisbury poisoning (EU vs Disinfo 2019), and the Hong Kong protests (Facebook 2019). The ‘normality’ of misinformation has been accepted to such an extent that a UK parliamentary committee described it as an expected part of modern political life (Digital, Culture, Media and Sport Committee 2019).

False and misleading information, and its ability to spread rapidly online, has been described as a societal vulnerability and a threat to democracy. While some uncertainties remain, especially the extent to which misinformation is effective (Benkler, Faris, and Roberts 2018), the possible consequences are indeed grave. Countries struggling with this challenge are presently adopting policies that can prevent misinformation from propagating too widely in their societies (Niels et al. 2020). Table 2.1 presents some previous studies that offered some backgrounds on the issue of misinformation, cyber-trooping and methods of analyzing them for proper checks and balances.

Table 2.1: Summary of related works

Authors (Year)	Research Work	Methods of classification and Analysis	Limitations
Yadigar et al. (2015)	Review of cyber-trooping activities in developed countries.	-	Inability to experiment with machine learning algorithms.
Valdiva et al. (2015)	Proposed polarity models for detecting neutrality of sentiment analysis methods.	Sentiment Analysis	Based on comparative analysis of sentiment analysis methods.
Pournarakis et al. (2017)	A new method to determine factors influencing brand equity assessments based on customers' perceptions.	Sentiment Analysis and Machine learning	Limited to analysis of customers' perceptions for marketing purposes.
Ankit (2018)	An ensemble classification system for twitter sentiment analysis	Sentiment Analysis	Based on ensemble of classifiers
Birjali et al. (2017)	Application of machine learning and semantic sentiment analysis-based algorithms for suicide sentiment prediction	Semantic Sentiment Analysis and Machine Learning	Based on machine learning algorithms limited to semantic analysis.
Valdivia et al. (2018)	A consensus vote model for detecting and filtering neutrality in sentiment analysis	Sentiment Analysis	Based on consensus voting strategy.
Demitrios et al. (2017)	A computational model for mining consumer perceptions on social media	Machine learning through genetic algorithm	Based on opinion mining
Jamilah et al. (2016)	A thematic based coding to identify key social media affordance.	qualitative methodology using semi-structured interviews and analysis	Based on natural language processing

Gillian et al. (2018)	A computer-based study of human behavior	-	Based on analysis of human perceptions
Rosyidah (2010)	A study of political blogging and public sphere In Malaysia: Case study of Najib's b;ogh	-	Based on review of political blogging. Inability of algorithmic implementation

2.6 Chapter Summary

This chapter presented overview of the concepts and algorithms of machine learning and data classification. Also, this chapter explained several of the algorithms commonly used in data classification tasks. The steps and methods for the research study are carefully and systematically presented in Chapter three of this thesis.

CHAPTER 3

METHODOLOGY

3.1 Introduction

Research methodology is defined as the analysis of the principle of methods, the systematic study of methods that can or have been applied within a discipline for a particular procedure or set of procedures. This research methodology is very important for producing decisions that can fulfill the purpose and objective of the research.

This chapter discusses the methodology employed in the research study with respect to the steps taken in actualizing the set objectives. It includes the framework of the proposed sentiment analysis approach. Each of these is systematically placed in the rest of the sections that follow. In section 3.2, the experimental data, data preprocessing, classification, and the output results are presented. The experimental data are made up of tweets from three Malaysian political figures: Dato Seri Anwar, Dato Hadi Awang and Lim Guang twitter. The preprocessing phase involved the use of standard StringToWordVector and Term Frequency Inverse Document Frequency (TF-IDF) methods. The classification process adopted the polarity-based sentiment classification approach using standard classification algorithms. The work is based on applying sentiment analysis in identifying cyber-trooping activities as positive,

negative, or neutral perspectives as illustrated in Figure 3.2 of the proposed research framework.

3.2 Research Design

The research process as shown in Figure 3.1 consists of four phases; which include researching and identifying the problem gap, review recent works in literature, proposing the research objectives, designing the methodology, preprocessing and analyzing the data, and writing the report.

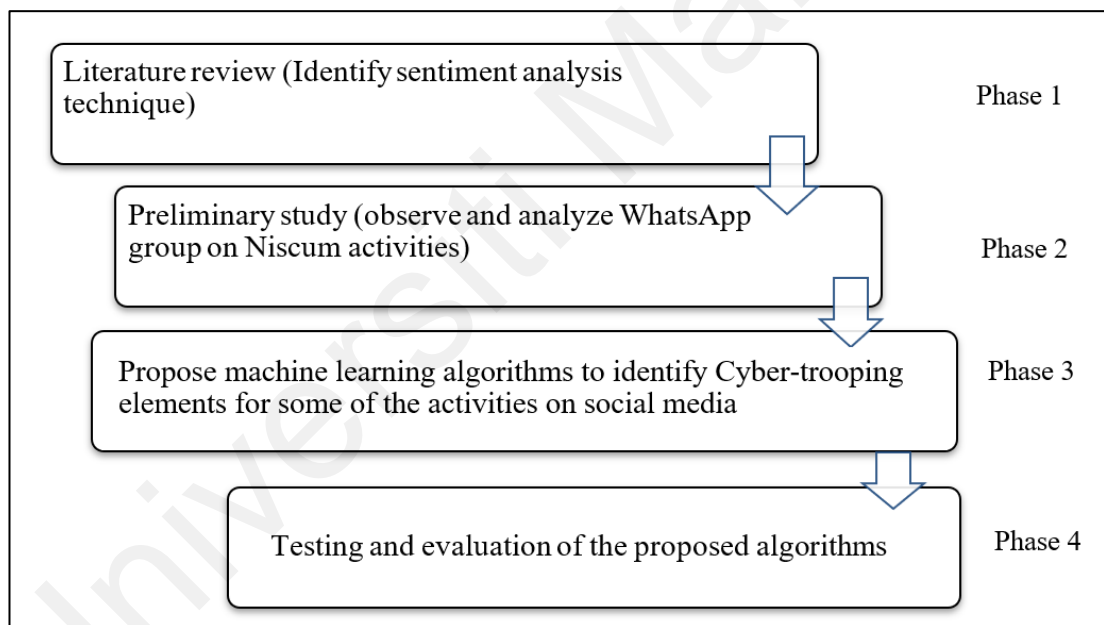


Figure 3.1: Research Process

3.2.1 Identification of Problem Gap

scope and conducting a literature review. These have been carefully established in chapter one and two of this report. In order to provide a step by step process of identifying the sentiment analysis and evaluating the algorithm, it was suggested by (Shi et al.,2016) that each emoji in the tweets or comment should be matched with the sentiment that each emotion expresses. As for a microblog, the frequency of occurrence of the emotions are counted and checked whether the emotion is negative or positive. If the frequency of occurrence of negative emotions is more than positive emotions, the microblog would be labeled as negative and vice versa. However, if a comment does not have emojis or there are equal number of emojis, then they will be sent to a Naive Bayes Classifier.

3.2.2 Data Collection

Data collection is a process of collecting data from different sources. These data can be obtained from various sources of data. Among the data sources, include reading, observation of the situation, exchange ideas and opinions and self-study. In data collection stage, the focus is more on primary data. Primary data is data which are generated by a researcher who is responsible for the design of the study and the collection, analysis and reporting of the data. Some of the samples of the tweet data experimented in this study can be found in the appendices.

Dato Seri Anwar Ibrahim's tweets (observations) on social related issues (refer to Figure 3.3) generated much feedbacks (responses) from the public targeted audience. Specifically, the tweets received over 80000 responses. In addition, the tweets garnered lots of likes of about 26,000 in numbers. Often, a tweet would be classified as positive if it attracted a high number of feedbacks accompanied with likes

(acceptance). On his political observations (refer to Figure 3.4), the highest number of responses were above 250 with over a thousand likes garnered likewise. The Dato's view on religious matters (as shown in Figure 3.5) had the highest number of responses with number of likes exceeding a thousand.

Dato Hadi Awang's tweets collections on several of the listed categories were graphed in Figures 3.6 to 3.8 in terms of social, political, and religious categories respectively. In Figure 3.6, Dato Hadi's tweets on social matter gathered average responses with the peak at 469 responds while a total number of likes total 1600 were achieved. On his political observations in Figure 3.7, it garnered at the peak 690 responses and the highest peak likes at 1900. This is an improvement over his social thoughts. The Dato's view on religion (Figure 3.8) had at the highest 1400 responses which is the overall highest in the three categories of his tweets. Subsequently, his observation on religious matter had at peak 577 likes.

Finally, YB LIM GUANG ENG tweets on politics and religion are graphed in Figures 3.9 and 3.10 respectively. His tweets on politics (Figure 3.9) generated a low response of 52 with 240 likes while on religion (Figure 3.10) garnered around 69 responses and 409 likes. Therefore, there is a possibility of mostly negative comments from the tweets due to the low response.

The data were analyzed in order to fetch appropriate information regarding the research topic. Data was collected through social media platform such as twitter. In this study, to determine the element of cyber-troopers, the collection of tweets for three public figures in political parties in Malaysia will be made and categorized based on public reactions with either positive, negative or neutral perspectives.

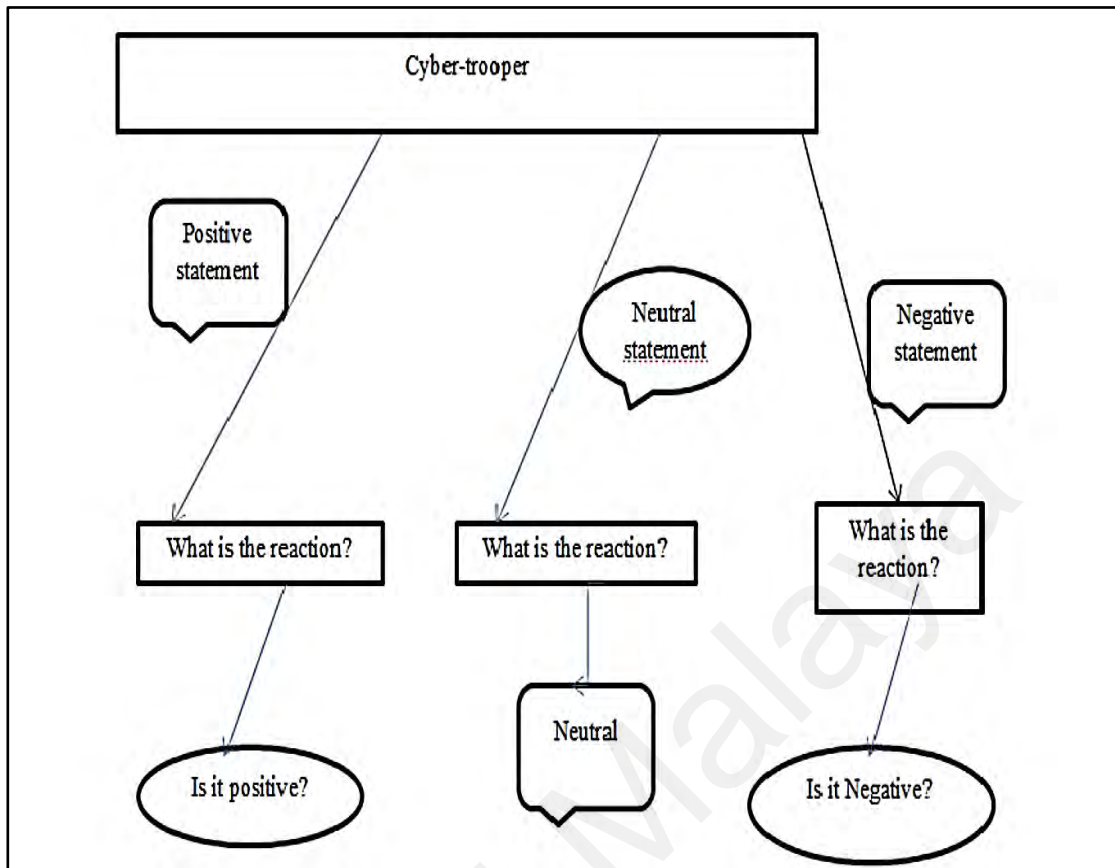


Figure 3.3: Dato Seri Anwar's Tweets on Social

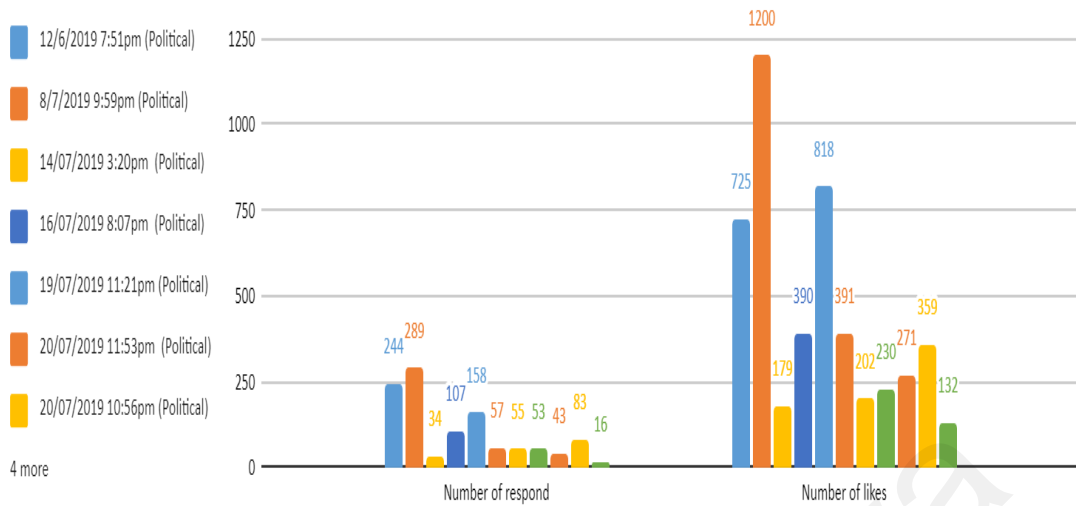


Figure 3.4: Dato Seri Anwar's Tweets on Politics

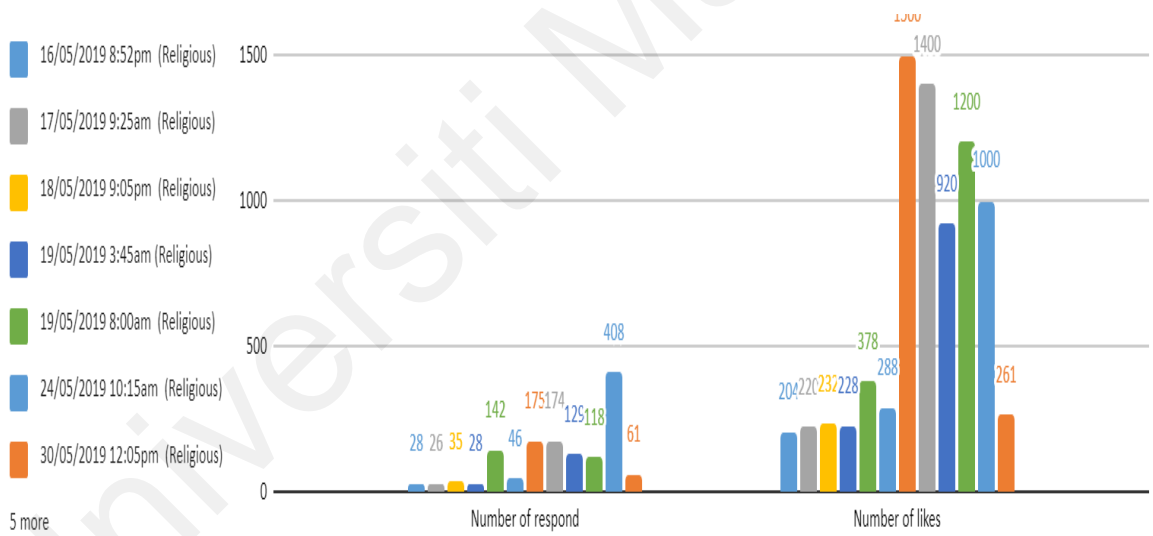


Figure 3.5: Dato Seri Anwar's Tweets on Religion

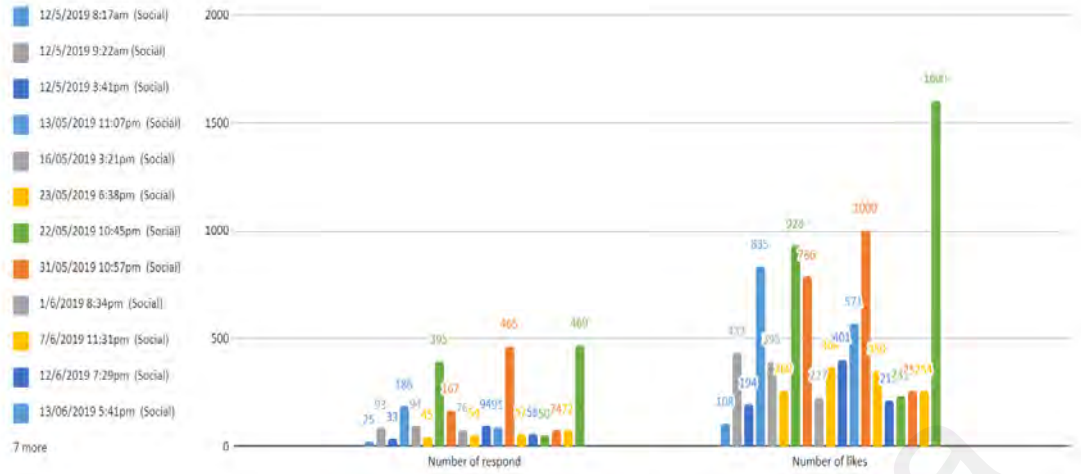


Figure 3.6: Dato Hadi Awang's Tweets on Social

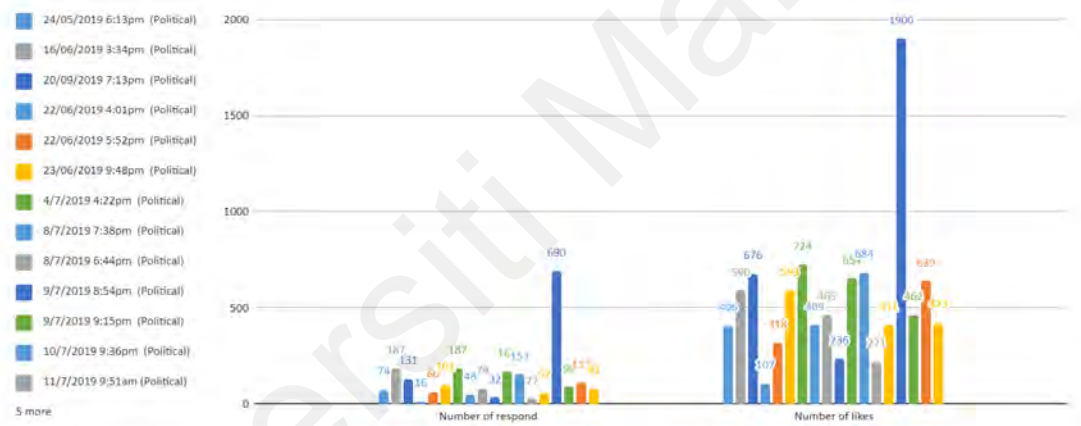


Figure 3.7: Dato Hadi Awang's Tweets on Politics

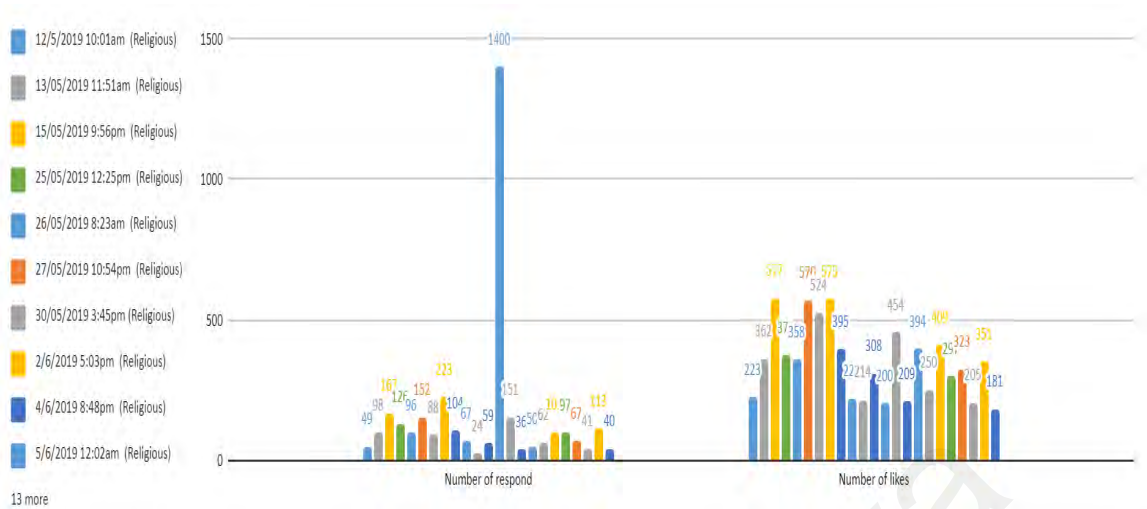


Figure 3.8: Dato Hadi Awang's Tweets on Religious

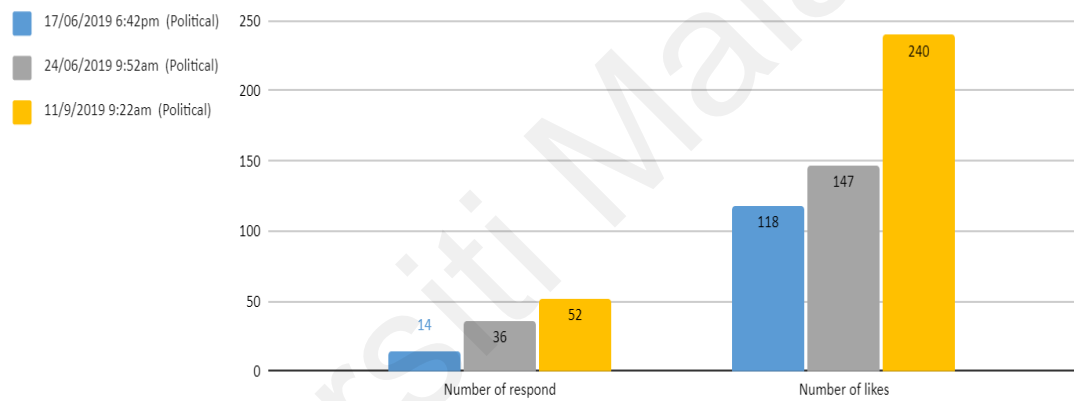


Figure 3.9: YB LIM GUANG ENG's Tweets on Politics

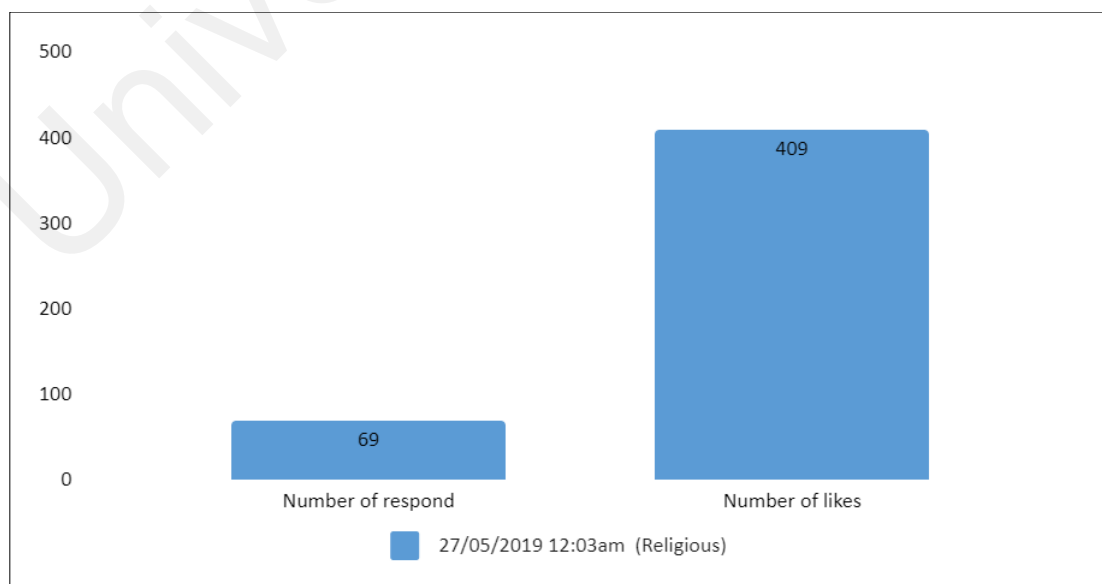


Figure 3.10: YB LIM GUANG ENG's Tweets on Religious

3.2.3 Data Pre-processing

Feature generation is an important step in data preprocessing which involves the process of extracting features (or keywords) from a normalized data. For the experimental work, the data is first converted to Attribute-Relation File Format (ARFF) which is the standard file format for machine learning using WEKA (Waikato Environment for Knowledge Analysis). The following are the steps taken in generating features from the textual data:

3.2.4 String to Word Vector

String To Word Vector is an unsupervised filter standard tool used for converting string attributes into a set of attributes representing word occurrence information from the text contained in the strings. This process can also be termed as Tokenizing (Subramaniaswamy *et al.*, 2018).

3.2.5 Term Frequency-Inverse Document Frequency (TF-IDF)

TF-IDF is a standard term weighting scheme used in accessing and measuring the significance of a word to a document in a collection. Its value increases proportionally to the number of times a word appears in a document. In text classification, TF-IDF is one of the methods used for stop-words filtering. It is an important step in the dimensionality reduction process involving the removal of redundant, irrelevant words

from text. TF-IDF is a product of two statistical weighting methods, Term Frequency (TF) and Inverse Document Frequency (IDF).

A) *Term Frequency*

Term frequency $Tf(t, d)$ is defined as the number of times a given term t (word/token) appears in a document d . Mathematically, Term Frequency ($Tf(t, d)$) is defined as:

$$Tf(t, d) = 0.5 + \frac{0.5 \times f(t, d)}{\text{Maximum Occurrences of words}} \quad (3.1)$$

where *Maximum Occurrences of words* is denoted with: $\text{Max} \{f^t, d: t \in d\}$.

B) *Inverse-Document Frequency (IDF)*

The Inverse-Document Frequency (IDF) is a measure of how much information a word provides. In other words, IDF is a method of evaluating if a term is common or rare across all documents in a collection (Chen *et al.*, 2016).

Mathematically, IDF is given as:

$$idf(t, D) = \log \frac{N}{|\{d \in D: t \in d\}|} \quad (3.2)$$

Where N is the total number of documents in D . $|\{d \in D: t \in d\}|$ is the number of documents where the term t appears.

Term Frequency-Inverse Document Frequency (TF-IDF) could then be given as:

$$tfidf(t, d, D) = tf(t, d) \cdot idf(t, D) \quad (3.3)$$

A high weight in TF-IDF is reached by a high term frequency in a given document with a low document frequency of the term in the whole document collection.

3.3 Sentiment Classification

In classification problems, there are several standard algorithms also known as classifiers implemented for classification. In this study, four standard classification algorithms were employed for the implementation of the classification task. It includes Naïve Bayes (NB), k Nearest Neighbor (kNN), Support Vector Machines (SVM), and Decision Trees (J48). These are the most commonly used classifiers in text classification problems (Hu *et al.*, 2016). The detailed explanation of each had been presented in Chapter 2 of the literature review.

3.4 Performance Measures

The study employed standard performance metrics commonly used in text classification problems for evaluating the sentiment classification methods.

The performance evaluators in the study include: Accuracy, Precision, Recall, F-Measures, and Area Under (ROC) Curve (AUC). The specific target of the study is to achieve high classification accuracy at less computational runtime and high AUC values (closer to 1) across the classifiers implemented.

➤ Confusion Matrix

A confusion matrix contains information about actual and predicted classification done by a classification algorithm. The correctly classified instances are called True Positive (*TP*) and True Negative (*TN*). Incorrectly classified instances are called False

Positive (*FP*) and False Negative (*FN*). Table 3.1 shows a sample of a typical confusion matrixes of the three classes in terms of *TP*, *TN*, *FP*, and *FN*.

Table 3.1: Confusion Matrixes in terms of *TP*, *TN*, *FP*, *FN*

<p>Class <i>a</i></p> <table border="1" data-bbox="480 454 676 629"> <tr> <td><i>TP</i> 311</td> <td><i>FN</i> 34</td> </tr> <tr> <td><i>FP</i> 6</td> <td><i>TN</i> 100</td> </tr> </table>	<i>TP</i> 311	<i>FN</i> 34	<i>FP</i> 6	<i>TN</i> 100	<p>Class <i>b</i></p> <table border="1" data-bbox="1031 454 1227 629"> <tr> <td><i>TP</i> 30</td> <td><i>FN</i> 12</td> </tr> <tr> <td><i>FP</i> 22</td> <td><i>TN</i> 387</td> </tr> </table>	<i>TP</i> 30	<i>FN</i> 12	<i>FP</i> 22	<i>TN</i> 387
<i>TP</i> 311	<i>FN</i> 34								
<i>FP</i> 6	<i>TN</i> 100								
<i>TP</i> 30	<i>FN</i> 12								
<i>FP</i> 22	<i>TN</i> 387								
<p>Class <i>c</i></p> <table border="1" data-bbox="480 719 676 893"> <tr> <td><i>TP</i> 59</td> <td><i>FN</i> 5</td> </tr> <tr> <td><i>FP</i> 23</td> <td><i>TN</i> 364</td> </tr> </table>	<i>TP</i> 59	<i>FN</i> 5	<i>FP</i> 23	<i>TN</i> 364	<p>Summary of Classes (<i>a, b, c</i>)</p> <table border="1" data-bbox="1031 719 1227 893"> <tr> <td><i>TP</i> 400</td> <td><i>FN</i> 51</td> </tr> <tr> <td><i>FP</i> 51</td> <td><i>TN</i> 851</td> </tr> </table>	<i>TP</i> 400	<i>FN</i> 51	<i>FP</i> 51	<i>TN</i> 851
<i>TP</i> 59	<i>FN</i> 5								
<i>FP</i> 23	<i>TN</i> 364								
<i>TP</i> 400	<i>FN</i> 51								
<i>FP</i> 51	<i>TN</i> 851								

The performance of a classification algorithm can be calculated based on the figures (numbers) in the confusion matrix.

➤ Accuracy

Classification Accuracy is the percentage or proportion of the total number of predictions that are correctly classified. Accuracy can be calculated as:

$$accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (3.13)$$

TP is True Positive (instances correctly classified as Positive)

TN is True Negative (instances correctly classified as Negative)

FP is False Positive (instances incorrectly classified as Positive)

FN is False Negative (instances incorrectly classified as Negative)

Working example:

Using the summarized confusion matrix in Table 3.3:

$$\begin{aligned} \text{accuracy } (A) &= \frac{(400 + 851)}{(400 + 51 + 851 + 51)} \\ A &= 93\% \end{aligned}$$

The higher the accuracy rate, the better a classifier performed. Here, the Naïve Bayes classifier had 93% accuracy.

➤ Precision and Recall

Precision is the number of positive instances correctly classified over all samples.

Recall is the number of positive instances classified over all the positive instances.

Recall is also known as sensitivity. Precision and Recall are calculated as:

$$\text{precision} = \frac{TP}{TP + FP} \quad (3.14)$$

$$\text{recall} = \frac{TP}{TP + FN} \quad (3.15)$$

➤ F-Measures

F-Measure is the harmonic mean of Precision and Recall. It is calculated as:

$$f - \text{measure} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (3.16)$$

Working example:

From the summarized confusion matrix above, precision, recall, and *f*-Measures can be calculated as:

$$\begin{aligned} \text{precision} &= \frac{400}{(400 + 51)} \\ \text{precision} &= 0.887 \end{aligned}$$

$$recall = \frac{400}{(400 + 51)}$$

$$recall = 0.887$$

$$f - measure = \frac{(2 \times 0.887 \times 0.887)}{(0.887 + 0.887)}$$

$$f - measure = 0.902$$

➤ Area Under (ROC) Curve (AUC)

AUC is one of the most popularly used performance metrics in classification problems. AUC values reflect the overall ranking performance of a classifier. The performance metric over the years has been proven to be better than classification accuracy metric for evaluating the classifier performance. Its value ranges from 0 to 1, where AUC = 1 corresponds to perfectly correct classification, AUC = 0.5 corresponds to classification by chance, and AUC = 0 corresponds to an inverted classification. As AUC value tends to 1, the most perfect classification algorithm.

3.5 Analysis of the Results

The ultimate goal of data classification is predicting the target class (or labels) for each data point in a data sample. This study has successfully identified cyber-trooping activities and also evaluate users' perception on social media using sentiment analysis approaches. Standard machine learning algorithms have been used to classify which class users' perceptions belong to either positive, negative, or neutral. The experimental results of the proposed sentiment classification approach are evaluated using conventional performance metrics in classification problems: Accuracy, Precision, Recall, F-Measure, and Area under the (ROC) curve (AUC).

When applying the classification algorithms to all datasets with all classifier, the Naïve Bayes (NB), Nearest Neighbor (kNN), Support Vector Machines (SVM), and Decision Trees (J48) were evaluated based on the aforementioned performance metrics. However, when the classification algorithms were applied to the tweet dataset with performance metrics, the results were compared with respect to the tweet data of Dato Seri Anwar, Dato Hadi Awang and Lim Guang in terms of Accuracy, Precision, Recall, F-Measure, and AUC. The results showed that SVM recorded the best accuracy followed by nearest neighbor, naïve bayes, and J48 algorithms respectively in that particular order. It was also shown that the SVM algorithm gives the highest values for the F-Measure followed by the RF, NB, J48, and KNN respectively. From the experimental results, it can be observed that SVM gives better results as the size of datasets increases. It was noted also that the results are better when applying classification algorithms with the application of machine learning. The next chapter discusses the experimental classification results.

3.6 Environment

The experimental setup was carried out on WEKA (v3.8.0) obtained from (<http://www.waikato.ac.nz/ml/weka/>). WEKA is an open source data mining tool developed at the University of Waikato, New Zealand. The software implements data mining algorithms using Java language. The state-of-the-art tool is used for developing machine learning (ML) algorithms with applications for data preprocessing, feature selection, classification, regression, clustering, and association rules. Furthermore, the machine learning software has tools for visualizing results. It allows users to quickly try out and compare different machine learning methods on new data sets. Its modular,

extensible architecture allows sophisticated data mining processes to be built up from the wide collection of base learning algorithms and tools provided. Extending the toolkit is easy thanks to a simple API, plugin mechanisms and facilities that automate the integration of new learning algorithms with WEKA's graphical user interfaces. The workbench includes algorithms for regression, classification, clustering, association rule mining and attribute selection.

Weka enables implementation of learning that can easily apply to datasets in classifying sentiment analysis. The essence of the weka library is to use a learned model to generate predictions on new instances.

WEKA supports various data mining tasks like data preprocessing, binning, clustering, regression and feature selection. It is not capable of multi-relational data mining, but there is separate software for converting a collection of linked database tables into a single table that is used for processing using WEKA.

The classification problem is a supervised learning task that consists of assigning a class label to an unclassified tuple according to an already classified instance set, that is used as a training set for the algorithm.

Steps adopted for performing Experiments are as below:

Step 1: Importing the dataset in WEKA. The first step performed was to import the dataset into the WEKA tool. To perform this step a simple import procedure for textual datasets called TextDirectoryLoader component was used.

Step 2: After importing the dataset it is converted and saved in the ARFF format.

Step 3: After that, a relation is created by containing 2000 instances and two attributes "Text" and "Class". The figure shows the uniform distribution of the attribute Class. Blue color represents reviews of negative polarity and Red color represents reviews of positive polarity

Step 4: Then the StringToWordVector filter is applied.

Step 5: Then the AttributeSelection filter is applied.

Step 6: After applying the AttributeSelection filter, the results are obtained.

Step 7: Three algorithms are performed on the data generated from above steps. The three algorithms are Naïve Bayes, K Nearest Neighbour, and Random Forest.

3.7 Chapter Summary

This chapter presented the proposed framework and design of the research study. It consists of data collection, transformation and pre-processing, sentiment classification, and classification results (output). The application of machine learning for the detection of users' perception against cyber-trooping and analysis of the data was presented. Several classification algorithms of the dataset were applied, and several sentiment analysis was used. Performance evaluation of the study was also done using linear and probabilistic classifiers. This includes Accuracy, Precision, Recall, F-Measures, and Area Under (ROC) Curve (AUC). For the implementation of the classification task, four standard classification algorithms, which are the most commonly used classifiers in text classification problems were used. In the study, Naïve Bayes (NB), k Nearest Neighbor (kNN), Support Vector Machines (SVM), and Decision Trees (J48) were employed. The next chapter discusses the experimental results.

CHAPTER 4

RESULTS AND DISCUSSION

4.1 Introduction

This chapter detailed the experimental results of the study following carefully the application of machine learning to management of information that are pushed into the virtual community. The purpose of the research work is to apply sentiment analysis algorithm to identify cyber-trooping activities in the Malaysia political space. To achieve the set purpose, the study employed the use of four standard machine learning algorithms. These include: Naïve Bayes (NB), k Nearest Neighbor (kNN), Support Vector Machines (SVM), and Decision Trees (J48). The proposed sentiment analysis model was evaluated on selected twitter data of three political figures in Malaysia: Dato Seri Anwar, Dato Hadi Awang and Lim Guang. In addition, 5 standard performance metrics namely: Accuracy, Precision, Recall, F-Measures, and Area under (ROC) Curve (AUC) were employed for validating the classification results.

4.2 Classification Results Analysis

In this section, the results obtained from the experimental work are analyzed. For this study, the selected classification algorithms were all implemented on the textual data. A standard 10-fold cross validation method was used for implementing the

classification algorithms. The resulting tables consist of classification results obtained with the classifiers on the tweet data of the three selected politicians. The detailed experimental results analysis for the classification algorithms is as follows:

Table 4.1: Sentiment Classification Performance of the classifiers with the tweet data of Dato Seri Anwar in terms of Accuracy, Precision, Recall, F-Measure, and AUC

Metrics	Classification Algorithms			
	NB	SVM	k-NN	J48
Accuracy (↑)	92.5%	94.5%	90.7%	87.1%
Precision (↑)	0.899	0.918	0.850	0.784
Recall (↑)	0.887	0.917	0.860	0.807
F-Measures (↑)	0.892	0.911	0.845	0.786
AUC (↑)	0.944	0.846	0.768	0.695

Table 4.2: Sentiment Classification Performance of the classifiers with the tweet data of Dato Hadi Awang in terms of Accuracy, Precision, Recall, F-Measure, and AUC

Metrics	Classification Algorithms			
	NB	SVM	k-NN	J48
Accuracy (↑)	88.3%	93.3%	92.8%	86.7%
Precision (↑)	0.871	0.902	0.891	0.801
Recall (↑)	0.828	0.9	0.891	0.8
F-Measures (↑)	0.838	0.892	0.891	0.801
AUC (↑)	0.93	0.809	0.846	0.760

Table 4.3: Positive Sentiment Classification Performance of the classifiers with the tweet data of Lim Guang in terms of Accuracy, Precision, Recall, F-Measure, and AUC

Metrics	Classification Algorithms			
	NB	SVM	k-NN	J48
Accuracy (↑)	91%	90%	87.3%	85.2%
Precision (↑)	0.856	0.856	0.791	0.730
Recall (↑)	0.860	0.856	0.809	0.778
F-Measures (↑)	0.858	0.857	0.797	0.733
AUC (↑)	0.894	0.729	0.745	0.619

From the results, the competitiveness of the classification algorithms is observed. Implementing naïve bayes algorithm had the least classification accuracy of 88.3% with the tweet data of Dato Hadi Awang while decision trees (J48) algorithm had the least AUC value of 0.619 with that of Lim Guang. This is mostly due to the nature of textual data i.e., the presence of a high level of dimensionality which could have influence on classifier's performance. Nevertheless, the highest classification accuracy of 94.5% was achieved using support vector machines (SVM) algorithm on tweet data of Dato Seri Anwar. Similarly, the naïve bayes (NB) algorithm performed best, obtained the highest AUC value of 0.944 with Anwar's tweet data.

Consistently, three of the four classification algorithms implemented in this work namely: naïve bayes (NB), support vector machines (SVM), and nearest neighbor (k-NN) obtained above 90% accuracy results with the tweet data of Dato Anwar, while the decision trees (J48) classifier is placed least with 87.1% accuracy performance. As earlier mentioned, classification algorithms are sensitive to the nature of the data experimented. Often, decision trees (J48) algorithms perform well with large size of data compared with SVM which often perform excellently with medium

to relatively large data. Working with Dato Hadi Awang's tweets, the four algorithms obtained accuracy results of 88.3%, 93.3%, 92.8%, 86.7% using naïve bayes, SVM, nearest neighbor, and J48 algorithms respectively. Again here, similar to Dato Seri Anwar's results, the support vector machines (SVM) algorithm achieved the best accuracy result of 93.3%. Finally, assessing the accuracy result of the classifiers with Lim Guang data, two of the classification algorithms: naïve bayes (NB) and SVM jointly obtained 90% mark of accuracy result with the NB classifier placed first with the best accuracy of 91%. Nearest neighbor (k-NN) classifier obtained 87.3% accuracy value while at the least position is the decision trees (J48) algorithm with the lowest accuracy value of 85.2%. Overall, the support vector machines (SVM) algorithm achieved the best accuracy result of 94.5% with Dato Seri Anwar's data.

In addition, assessing the performance of the classification algorithms in terms of precision, recall, and *f*-Measure obtained mixed results. The best precision value of 91.8% was obtained using SVM algorithm with Dato Seri Anwar's data while the least precision score of 73% was obtained using decision trees (J48) algorithm with Lim Guang's data. Similarly, the best recall value of 91.7% was obtained using SVM algorithm with Dato Anwar's data while the least value of 77.8% was obtained using decision trees (J48) algorithm with Lim Guang's data. Likewise, evaluating the performance of the classification algorithms based on *f*-Measure metric achieved the best value of 91.1% using SVM algorithm with Dato Anwar's data while the least value of 73.3% was obtained using decision trees (J48) algorithm with Lim Guang's data.

Combining different evaluation measures is a good strategy for reflecting the effectiveness of the experimenting classification algorithms. Accuracy is the most widely applied metric for classification problems. A good classifier is expected to have

a high accuracy value. However, a good accuracy might not actually correlate with good predictions due to the presence of false predictions. But, when accuracy is combined with other metrics such as: precision and recall, then the effectiveness of the classification methods could be clearer. A high precision and recall show that the classifier achieved accurate predictions (high precision) and majority are positive results (high recall). A classifier with high recall but low precision returns many results (predictions), but most (predicted labels) are incorrect. On the other hand, a classifier with high precision but low recall achieved few results (predictions), but most (predicted labels) are correct. However, the best classifier has high precision and recall, with many results (predictions) are positive (correct).

4.3 Chapter Summary

In this chapter, depth discussion of the experimental work has been done. The experimental work comprises of tweet data of three political figures in Malaysia namely: Dato Seri Anwar, Dato Hadi Awang, and Lim Guang. Furthermore, four standard classifiers were implemented independently across the three experimental data using standard 10-fold cross validation method. The classifiers are naïve bayes (NB), support vector machines (SVM), nearest neighbour (k -NN), and decision trees (J48). Detailed analysis of the experimental results was presented including the performance metrics used for results evaluation and comparison. The results were tabulated and comparisons were made between the experimental data in terms of classification accuracy, AUC value, precision, recall, and f -Measure. and computational runtime. The results showed that combining multiple validating algorithms such as accuracy, precision, recall etc., proved to be the best approach to

assess the quality of classification algorithms. The overall highest accuracy result of **94.5%** was achieved by support vector machines (SVM) with Dato Seri Anwar's data. Similarly, the overall highest AUC value of **0.944** was achieved by SVM algorithm on same data. Managing users Perception could be better using naive bayes and J48 tree. The next chapter concludes the study with recommendations and directions to future work.

Universiti Malaya

CHAPTER 5

CONCLUSION AND RECOMMENDATION

5.1 Introduction

This chapter summarizes the research study, noting the study's significance contributions and directions to future works. The study is based on proposing a sentiment analysis algorithm for identifying and categorizing cyber-trooping activities using machine learning algorithms. Section 5.2 highlights the research contributions while section 5.3 noted the research limitations and future recommendations. Finally, Section 5.4 presents the research conclusion.

5.2 Research Contributions

As earlier established, perception is a way of understanding reality (in its broader scale) and reacts accordingly based on one's opinion, judgement, and reasoning. Naturally, humans tend to view things in different perspectives. A group of people will perceive matters differently based on several varying (individual) characteristics. This study attempts to aid the identification and categorization of these human perspectives into three key classes: positive, negative, and neutral perception. To do so, the study carefully proposed the following objectives:

1. Determination of a sentiment analysis (SA) algorithm for cyber-trooping activities in Malaysia.
2. Application of the SA algorithm using four machine learning algorithms: naïve bayes (NB), support vector machines (SVM), nearest neighbor (k -NN), and decision trees (J48) on selected political tweets.
3. Evaluation and analysis of the simulation results using five conventional performance measures: accuracy, precision, recall, f -Measure, and AUC methods.

The first objective is utilization of a sentiment analysis (SA) algorithm for cyber-trooping activities in Malaysia. The study achieved the stated objective by applying a sentiment analysis algorithm using machine learning-based approach for the classification task. In the second objective, the research work achieved the stated purpose by implementing successfully the sentiment analysis algorithm using four conventional machine learning algorithms. These include: naïve bayes (NB), nearest neighbor (k -NN), support vector machines (SVM), and decision trees (J48) algorithms. The algorithm was successfully implemented and achieved satisfactory results. Finally, the research work achieved the third objectives which is on evaluating and comparing the classification performance of the standard machine learning algorithms. This was carried out using five standard performance metrics: accuracy, AUC, precision, recall, and f -Measure. The classification algorithm achieved competitive results, support vector machine (SVM) algorithm obtained the overall best results of 94.5% accuracy, 0.918 precision, 0.917 recall, and 0.911 f -Measure, while the naïve bayes (NB) algorithm achieved the best AUC value of 0.944 with the tweet data experimented.

5.3 Research Limitations and Future Recommendation

This research work is limited to testing the tweet data of three popular political figures in Malaysia: Dato Seri Anwar, Dato Hadi Awang, and Lim Guang, using four conventional classification algorithms: naïve bayes, SVM, k -NN, and decision trees (J48). The proposed sentiment analysis algorithm was implemented in Java and WEKA (an open source library for machine learning projects). The classifiers implemented produced competitive and satisfactory results. The study recommends validating the proposed sentiment analysis algorithm with more datasets and classification algorithms. In addition, the study recommends evaluating the classifiers with more performance methods for more convincing and satisfactory assessments.

5.4 Concluding Remarks

This study proposed an adoption of machine learning algorithm for analysing supporters and non-supporters' feedback on political posts. Four conventional classification algorithms: naïve bayes, SVM, k -NN, and J48 were experimented with tweet data of three political figures in Malaysia: Dato Seri Anwar, Dato Hadi Awang, and Lim Guang. Results were obtained and analyzed, support vector machines (SVM) algorithm achieved the overall best results with the tweet data of Dato Seri Anwar.

REFERENCES

- Adegun, I. P., & Vadapalli, H. B. (2020). Facial micro-expression recognition: A machine learning approach. *Scientific African*, 8.
- Amiri, K., Boualleg, Y., & Farah, M. (2020). Radiometric indices-based spectro-spatial approach for hyperspectral image classification. *Egyptian Journal of Remote Sensing and Space Science*.
- Arian, R., Hariri, A., Mehridehnavi, A., Fassihi, A., & Ghasemi, F. (2020). Protein kinase inhibitors' classification using K-Nearest neighbor algorithm. *Computational Biology and Chemistry*, 86, 107269.
- Arora, N., & Kaur, P. D. (2020). A Bolasso based consistent feature selection enabled random forest classification algorithm: An application to credit risk assessment. *Applied Soft Computing Journal*, 86, 105936.
- Arunakranthi, G., Rajkumar, B., Shekhar, V. C., & Harshavardhan, A. (2020). Materials Today : Proceedings Advanced patterns of predictions and cavernous data analytics using quantum machine learning. *Materials Today: Proceedings*.
- Askari, S. M. S., & Hussain, M. A. (2020). IFDTC4.5: Intuitionistic fuzzy logic based decision tree for E-transactional fraud detection. *Journal of Information Security and Applications*, 52, 102469.
- Bagga, S., Goyal, A., Gupta, N., & Goyal, A. (2020). ScienceDirect Credit Card Fraud Detection ICITETM2020 using Pipeling and Ensemble Learning Credit Card Fraud Detection using Ensemble a Pipeling and Goyal c Learning. *Procedia Computer Science*, 173(2019), 104–112.
- Bandela, S. R., & Kumar, T. K. (2021). Unsupervised feature selection and NMF denoising for robust Speech Emotion Recognition. *Applied Acoustics*, 172, 107645.
- Belciug, S. (2021). Parallel versus cascaded logistic regression trained single-hidden feedforward neural network for medical data. *Expert Systems With Applications*, 170(April 2020), 114538.

- Brindha, S., Prabha, K., and Sukumaran, S. (2016). Pattern Document Weight Discovery For Text Classification Mining. *2016 International Conference on Communication and Electronic Systems, IEEE, 2016*, pp. 2–6.
- Brouer, M., & Benabbou, A. (2019). ATLASLang MTS 1: Arabic Text Language into Arabic Sign Language Machine Translation System. *Procedia Computer Science, 148(Icids 2018)*, 236–245. <https://doi.org/10.1016/j.procs.2019.01.066>
- Bulat, A., Tzimiropoulos, G., Kossaifi, J., & Pantic, M. (2019). Improved training of binary networks for human pose estimation and image recognition. *ArXiv*
- Cattoni, R., Di Gangi, M. A., Bentivogli, L., Negri, M., & Turchi, M. (2021). MuST-C: A multilingual corpus for end-to-end speech translation. *Computer Speech and Language, 66*, 101155.
- Cekik, R. & Uysal, A. (2020). A novel filter feature selection method using rough set for short text data. *Expert Systems with Applications, 160*, 113691.
- Cen, F., Zhao, X., Li, W., & Wang, G. (2021). Deep feature augmentation for occluded image classification. *Pattern Recognition, 111*, 107737.
- Chen, Z., Zhou, L.J., Li, X., Zhang, J., & Huo, W. J. (2020). The lao text classification method based on KNN. *Procedia Computer Science, 166*, 523-528.
- Chen, L., Jiang, L., & Li, C. (2021). Modified DFS-based term weighting scheme for text classification. *Expert Systems with Applications, 168*, 114438.
- Christou, V., Tsipouras, M. G., Giannakeas, N., Tzallas, A. T., & Brown, G. (2019). Hybrid extreme learning machine approach for heterogeneous neural networks. *Neurocomputing, 361*, 137–150.
- D’hooge, L., Wauters, T., Volckaert, B., & De Turck, F. (2020). Inter-dataset generalization strength of supervised machine learning methods for intrusion detection. *Journal of Information Security and Applications, 54*.
- Das, S., Dey, A., Pal, A., and Roy, N. (2015). Applications of Artificial Intelligence in Machine Learning: Review and Prospect. *International Journal of Computer Applications, 115(9)*, pp. 31-41.
- Das, A., Acharya, U. R., Panda, S. S., & Sabut, S. (2019). ScienceDirect Deep learning based liver cancer detection using watershed transform and Gaussian mixture model algorithms. *Cognitive Systems Research, 54*, 165–175.
- Dedetürk, B. K., & Akay, B. (2020). Spam filtering using a logistic regression model trained by an artificial bee colony algorithm. *Applied Soft Computing Journal*,

91, 106229.

- DeLatte, D. M., Crites, S. T., Guttenberg, N., & Yairi, T. (2019). Automated crater detection algorithms from a machine learning perspective in the convolutional neural network era. *Advances in Space Research*, 64(8), 1615–1628.
- Dogan, T., & Uysal, A. K. (2020). A novel term weighting scheme for text classification: TF-MONO. *Journal of Informetrics*, 14(4), 101076.
- Dunham, M. W., Malcolm, A., & Welford, J. K. (2020). Improved well log classification using semisupervised Gaussian mixture models and a new hyperparameter selection strategy. *Computers and Geosciences*, 140, 104501.
- Dzhamtyrova, R., & Kalnishkan, Y. (2020). Universal algorithms for multinomial logistic regression under Kullback–Leibler game. *Neurocomputing*, 397, 369–380.
- El Hindi, K. M., Aljulaidan, R. R., & AlSalman, H. (2020). Lazy fine-tuning algorithms for naïve Bayesian text classification. *Applied Soft Computing Journal*, 96, 106652.
- Fan, Y., Bai, J., Lei, X., Zhang, Y., Zhang, B., Li, K. C., & Tan, G. (2020). Privacy preserving based logistic regression on big data. *Journal of Network and Computer Applications*, 171, 102769.
- Fialho, G., Manhães, A., & Teixeira, J. P. (2019). Predicting Sports Results with Artificial Intelligence - A Proposal Framework for Soccer Games. *Procedia Computer Science*, 164, 131–136.
- Geng, J., Jiang, W., & Deng, X. (2020). Multi-scale deep feature learning network with bilateral filtering for SAR image classification. *ISPRS Journal of Photogrammetry and Remote Sensing*, 167, 201–213.
- Ghiasi, M. M., & Zendehboudi, S. (2020). Application of Decision Tree-Based Ensemble Learning in the Classification of Breast Cancer. *Computers in Biology and Medicine*, 128, 104089.
- Godavarthi, O., & Sowjanya, M. A. (2021). Classification of covid related articles using machine learning. *Materials Today: Proceedings*.
- Gu, Q., Chang, Y., Li, X., Chang, Z., & Feng, Z. (2021). A novel F-SVM based on FOA for improving SVM performance. *Expert Systems with Applications*, 165, 113713.
- Guo, X., & Yuan, Y. (2020). Semi-supervised WCE image classification with adaptive

- aggregated attention. *Medical Image Analysis*, 64, 101733.
- Hajji, M., Harkat, M. F., Kouadri, A., Abodayeh, K., Mansouri, M., Nounou, H., & Nounou, M. (2020). Multivariate feature extraction based supervised machine learning for fault detection and diagnosis in photovoltaic systems. *European Journal of Control*.
- Hamed, Y., Ibrahim Alzahrani, A., Shafie, A., Mustaffa, Z., Che Ismail, M., & Kok Eng, K. (2020). Two steps hybrid calibration algorithm of support vector regression and K-nearest neighbors. *Alexandria Engineering Journal*, 59(3), 1181–1190.
- Hénaff, O. J., Razavi, A., Doersch, C., Ali Eslami, S. M., & Van Den Oord, A. (2019). Data-efficient image recognition with contrastive predictive coding. *ArXiv*.
- He, Y., & Dong, X. (2020). Microprocessors and Microsystems Real time speech recognition algorithm on embedded system based on continuous Markov model. *Microprocessors and Microsystems*, 75, 103058.
- Hu, R., Mac, B., and Delany, S. J. (2016). Active learning for text classification with reusability R. *Expert Systems With Applications*, 45, pp. 438–449.
- Hyun, S., Kaewprag, P., Cooper, C., Hixon, B., & Moffatt-Bruce, S. (2020). Exploration of critical care data by using unsupervised machine learning. *Computer Methods and Programs in Biomedicine*, 194.
- Issa, D., Demirci, M. F., & Yazici, A. (2020). Biomedical Signal Processing and Control Speech emotion recognition with deep convolutional neural networks. *Biomedical Signal Processing and Control*, 59, 101894.
- Ivanovic, M., and Radovanovic, M. (2015). Modern Machine Learning Algorithms and their Applications. *International Conference on Electronics, Communications and Networks*.
- Jabir, M., Kumam, P., Deebani, W., Kumam, W., & Shah, Z. (2020). Bi-parametric distance and similarity measures of picture fuzzy sets and their applications in medical diagnosis. *Egyptian Informatics Journal*.
- Jiang, H., Zou, B., Xu, C., Xu, J., & Tang, Y. Y. (2020). SVM-Boosting based on Markov resampling: Theory and algorithm. *Neural Networks*, 131, 276–290.
- Ker, J., Bai, Y., Yee, H., Rao, J., & Wang, L. (2019). Automated brain histology classification using machine learning. *Journal of Clinical Neuroscience*, 66, 239–245.

- Kim, J., Jang, S., Park, E., & Choi, S. (2020). Text classification using capsules. *Neurocomputing*, 376, 214–221.
- Kolluri, J., & Razia, S. (2020). Text classification using Naïve Bayes classifier. *Materials Today: Proceedings*.
- Kwon, S. (2020). MLT-DNet : Speech Emotion Recognition Using 1D Dilated CNN Based on Multi-Learning Trick Approach Interaction Technology Laboratory, Department of Software , Sejong University, *Expert Systems With Applications*, 114177.
- Langari, S., Marvi, H., & Zahedi, M. (2020). Informatics in Medicine Unlocked Efficient speech emotion recognition using modified feature extraction. *Informatics in Medicine Unlocked*, 20, 100424.
- Le, S., Pellegrini, E., Green-Saxena, A., Summers, C., Hoffman, J., Calvert, J., & Das, R. (2020). Supervised machine learning for the early prediction of acute respiratory distress syndrome (ARDS). *Journal of Critical Care*, 60, 96–102.
- Lei, H. A. N., Lei, T., & Yuenian, T. (2020). SPORTS IMAGE DETECTION BASED ON PARTICLE SWARM OPTIMIZATION ALGORITHM. *Microprocessors and Microsystems*, 103345.
- Li, Q., Dong, J., Zhong, J., Li, Q., & Wang, C. (2019). A neural model for type classification of entities for text. *Knowledge-Based Systems*, 176, 127-132.
- Li, Q., Li, P., Mao, K., & Lo, E. Y. M. (2020). Improving convolutional neural network for text classification by recursive data pruning. *Neurocomputing*, 414, 143–152.
- Li, L. (2021). Software Reliability Growth Fault Correction Model Based on Machine Learning and Neural Network Algorithm. *Microprocessors and Microsystems*, 80, 103538.
- Lu, H., & Ma, X. (2020). Hybrid decision tree-based machine learning models for short-term water quality prediction. *Chemosphere*, 249, 126169.
- Luo, F., Huang, Y., Tu, W., & Liu, J. (2020). Local manifold sparse model for image classification. *Neurocomputing*, 382, 162–173.
- Matsuo, R., Yamazaki, T., Suzuki, M., Toyama, H., & Araki, K. (2020). A random forest algorithm-based approach to capture latent decision variables and their cutoff values. *Journal of Biomedical Informatics*, 110, 103548.
- Meher, S. K. (2019). Semisupervised self-learning granular neural networks for remote sensing image classification. *Applied Soft Computing Journal*, 83, 105655.

- Moreira, C. A., Philot, E. A., Lima, A. N., & Scott, A. L. (2019). Predicting regions prone to protein aggregation based on SVM algorithm. *Applied Mathematics and Computation*, 359, 502–511.
- Mukherjee, S., & Prasad, S. (2020). A spatial–spectral semisupervised deep learning framework using siamese networks and angular loss. *Computer Vision and Image Understanding*, 194, 102943.
- Nanglia, P., Kumar, S., Mahajan, A. N., Singh, P., & Rathee, D. (2020). A hybrid algorithm for lung cancer classification using SVM and Neural Networks. *ICT Express*. <https://doi.org/10.1016/j.ict.2020.06.007>
- Nawaz, A., Bakhtyar, M., Baber, J., Ullah, I., Noor, W., & Basit, A. (2020). Extractive Text Summarization Models for Urdu Language. *Information Processing and Management*, 57(6), 102383. <https://doi.org/10.1016/j.ipm.2020.102383>
- Ochieng, P. (2020). PAROT: Translating natural language to SPARQL. *Expert Systems with Applications: X*, 5, 100024.
- Pan, H., Pang, Z., Wang, Y., Wang, Y., & Chen, L. (2020). A New Image Recognition and Classification Method Combining Transfer Learning Algorithm and MobileNet Model for Welding Defects. *IEEE Access*, 8, 119951–119960.
- Pan, Z., Wang, Y., & Pan, Y. (2020). A new locally adaptive k-nearest neighbor algorithm based on discrimination class. *Knowledge-Based Systems*, 204, 106185.
- Park, Y., & Yang, H. S. (2019). Convolutional neural network based on an extreme learning machine for image classification. *Neurocomputing*, 339, 66–76.
- Patil, S., Nemade, V., & Soni, P. K. (2018). Predictive Modelling for Credit Card Fraud Detection Using Data Analytics. *Procedia Computer Science*, 132, 385–395.
- Puranik, T. G., Rodriguez, N., & Mavris, D. N. (2020). Towards nline prediction of safety-critical landing metrics in aviation using supervised machine learning. *Transportation Research Part C: Emerging Technologies*, 120, 102819.
- Qin, J., Pan, W., Xiang, X., Tan, Y., & Hou, G. (2020). A biological image classification method based on improved CNN. *Ecological Informatics*, 58, 101093.
- Rakhmetulayeva, S. B., Duisebekova, K. S., Mamyrbekov, A. M., Kozhamzharova, D. K., Astaubayeva, G. N., & Stamkulova, K. (2018). Application of

- Classification Algorithm Based on SVM for Determining the Effectiveness of Treatment of Tuberculosis. *Procedia Computer Science*, 130, 231–238.
- Rangasamy, K., Amir, M., Azmina, N., & Fathiah, N. (2020). Hockey activity recognition using pre-trained deep learning model. *ICT Express*, 6(3), 170–174.
- Reis, I., Rotman, M., Poznanski, D., Xavier Prochaska, J., & Wolf, L. (2019). Effectively using unsupervised machine learning in next generation astronomical surveys. *ArXiv*, 34, 100437.
- Reyes, O., Pérez, E., Luque, R. M., Castaño, J., & Ventura, S. (2020). A supervised machine learning-based methodology for analyzing dysregulation in splicing machinery: An application in cancer diagnosis. *Artificial Intelligence in Medicine*, 108.
- Rezaei-Ravari, M., Eftekhari, M., & Saberi-Movahed, F. (2021). Regularizing extreme learning machine by dual locally linear embedding manifold learning for training multi-label neural network classifiers. *Engineering Applications of Artificial Intelligence*, 97, 104062.
- Rtayli, N., & Enneya, N. (2020). Enhanced credit card fraud detection based on SVM-recursive feature elimination and hyper-parameters optimization. *Journal of Information Security and Applications*, 55(September), 102596.
- Rohidin, D., Samsudin, N. A., & Deris, M. M. (2020). Association rules of fuzzy soft set based classification for text classification problem. *Journal of King Saud University - Computer and Information Sciences*,.
- Ryabtseva, V., & Skomorokhov, A. (2020). Critical power prediction using SVM algorithms. *Procedia Computer Science*, 169(2019), 198–202.
- Saikrishna, V., Dowe, D. L., and Ray, S. (2016). Statistical Compression-Based Models for Text Classification. *Fifth International Conference on Eco-Friendly Computing Communication and Systems, IEEE*, pp. 1-6.
- Sharma, N., and Singh, M. (2016). Modifying Naive Bayes Classifier for Multinomial Text Classification. *IEEE International Conference on Recent Advances and Innovations in Engineering, 2016*, pp. 1-7.
- Shengxue, Z. (2020). English corpus translation system based on FPGA and machine learning. *Microprocessors and Microsystems*, 103464.
- Simsekler, M. C. E., Qazi, A., Alalami, M. A., Ellahham, S., & Ozonoff, A. (2020). Evaluation of patient safety culture using a random forest algorithm. *Reliability*

Engineering and System Safety, 204, 107186.

- Soltani, P., & Morice, A. H. P. (2020). Augmented reality tools for sports education and training. *Computers and Education*, 155, 103923.
- Soumaya, Z., Drissi Taoufiq, B., Benayad, N., Yunus, K., & Abdelkrim, A. (2021). The detection of Parkinson disease using the genetic algorithm and SVM classifier. *Applied Acoustics*, 171, 107528.
- Talavera-Llames, R., Pérez-Chacón, R., Troncoso, A., & Martínez-Álvarez, F. (2019). MV-kWNN: A novel multivariate and multi-output weighted nearest neighbours algorithm for big data time series forecasting. *Neurocomputing*, 353, 56–73.
- Tang, Z., Li, W., Li, Y., Zhao, W., & Li, S. (2020). Several alternative term weighting methods for text representation and classification. *Knowledge-Based System*, 207, 106399.
- Thejas, G. S., Dheeshjith, S., Iyengar, S. S., Sunitha, N. R., & Badrinath, P. (2020). *Jou. Machine Learning with Applications*, 100016.
- Tuncer, T., Dogan, S., & Acharya, U. R. (2021). Knowledge-Based Systems Automated accurate speech emotion recognition system using twine shuffle pattern and iterative neighborhood component analysis algorithms. *Knowledge-Based Systems*, 211, 106547.
- Vu, Q. V., Truong, V. H., & Thai, H. T. (2021). Machine learning-based prediction of CFST columns using gradient tree boosting algorithm. *Composite Structures*, 259, 113505.
- Wang, Z., & Chen, Q. (2020). Monitoring online reviews for reputation fraud campaigns. *Knowledge-Based Systems*, 195, 105685.
- Watson, L. M. (2020). Using unsupervised machine learning to identify changes in eruptive behavior at Mount Etna, Italy. *Journal of Volcanology and Geothermal Research*, 405, 107042.
- Xu, Z., Shen, D., Nie, T., & Kou, Y. (2020). A hybrid sampling algorithm combining M-SMOTE and ENN based on Random forest for medical imbalanced data. *Journal of Biomedical Informatics*, 107, 103465.
- Yang, D. (2021). Online sports tourism platform based on FPGA and machine learning. *Microprocessors and Microsystems*, 80, 103584.
- Yao, Z., Wang, Z., Liu, W., Liu, Y., & Pan, J. (2020). Speech emotion recognition using fusion of three multi-task learning-based. *Speech Communication*, 120, 11–

19.

- Yassine, A., Mohamed, L., & Al Achhab, M. (2021). Intelligent recommender system based on unsupervised machine learning and demographic attributes. *Simulation Modelling Practice and Theory*, 107, 102198.
- Zewen, C. (2016). Short Text Classification Based on Wikipedia and Word2vec. *2016 2nd IEEE International Conference on Computer and Communications, IEEE, 2016*, pp. 1195–1200.
- Zhan, Z., Hou, Z., Yang, Q., Zhao, J., Zhang, Y., & Hu, C. (2020). Knowledge attention sandwich neural network for text classification. *Neurocomputing*, 406, 1-11.
- Zhang, Z. (2020). Speech feature selection and emotion recognition based on weighted binary cuckoo search. *Alexandria Engineering Journal*.
- Zhou, Y., Wu, X., & Li, X. (2021). Prediction model of sports injury based on dynamic sampling and transfer learning. *Microprocessors and Microsystems*, 80, 103583.
- N. Saleena, "ScienceDirect An Ensemble Classification System for Twitter Sentiment Analysis An Ensemble Classification System for Twitter Sentiment Analysis," *Procedia Comput. Sci.*, vol. 132, no. Iccids, pp. 937–946, 2018.