

# **IDENTIFYING HARMFUL COMMENTS FOR TAMIL LANGUAGE ON SOCIAL MEDIA**

Prepared by  
Ms. Disne Sivalingam  
198773R

Dissertation submitted to the Faculty of Information Technology, University of  
Moratuwa, Sri Lanka for the partial fulfillment of the requirements of the Degree of  
Master of Science in Information Technology

**July 2022**

## DECLARATION

We declare that this thesis titled “Identifying Harmful Comments for Tamil Language on Social Media” is my own work. Where any part of this thesis has not been submitted in any form for another degree or diploma at this university or other institution of tertiary education. Information derived from the published or unpublished work of others has been acknowledged. Due references have been provided on all supporting works of literature and resources.

Name of the Student:

Signature of Student:

Ms. Disne Sivalingam

***UOM Verified Signature***

Date: 30/07/2022

***UOM Verified Signature***

Supervised By:

Signature of Supervisor:

Dr. S. C. Premaratne

Senior Lecturer

Faculty of Information Technology

Date: 30/07/2022

University of Moratuwa

## **DEDICATION**

I dedicate this research report affectionately to the following:

My Research Supervisor, Dr. S. C Premaratne for dedicating his valuable time for guiding me throughout this.

My Husband and Parents for providing me with continuous encouragement throughout my years of study.

Finally for my friends who helped me on giving ideas and support.

Thank you.

## **ACKNOWLEDGEMENT**

I would like to acknowledge and convey my appreciation for those efforts and support in one way or another contributed to the successful completion of this research.

First, I am deeply grateful for the supervision throughout the one year of my research and the help received from my supervisor Dr. S. C. Premaratne, Department of Information Technology, Faculty of Information Technology, University of Moratuwa. I have learned so much from our project discussions. His willingness to motivate me contributed tremendously to my job.

Besides, I would like to thank all academic staff from the Faculty of Information and technology, who shared their vast knowledge throughout these two years by providing me with a good environment which influenced a lot to achieve this goal.

It is with great pleasure to thank the University of Moratuwa, Sri Lanka, for all the efforts and facilities that it has taken to contributions this postgraduate programme. Especially, my thanks and gratitude to my colleagues for their help and support to complete this research in many ways.

I am as ever, especially indebted to my parents, brother, and husband for their love and support throughout my life to improve my career.

## **ABSTRACT**

The era of social media, such as YouTube, Facebook, and Twitter adding comments to posts are being fun in the daily life of people. But this is also used to spread hate speech and organize hate based activities increasingly nowadays. Harmful and offensive text identification on social media platforms is being a trending research area over the last few years. In a country like Sri Lanka with multiple native languages, people like to comment on social media mostly in their native language. Tamil is one of the Languages commonly used and spoken in the North and East part of Sri Lanka. In recent years people like to comment not only in their native language they also comment in more than one language. In Sri Lanka, people use Singlish (Sinhala + English ) or Tanglish (Tamil + English).

Because of the rapid growth of hateful content on social media, there is an immediate need for an efficient and effective method to identify harmful content. A huge number of researches have been done and are being done for automated harmful content detection online. The complication of the Natural Language constructs builds this task very challenging.

A maximum of the research are done in the English Language. This research work aims to classify the code-mixed Tamil comments on social media by categorizing them as harmful and non-harmful by using machine learning models.

## TABLE OF CONTENTS

	<b>Page</b>
<b>DECLARATION</b>	<b>i</b>
<b>DECLARATION</b>	<b>ii</b>
<b>ACKNOWLEDGEMENT</b>	<b>iii</b>
<b>ABSTRACT</b>	<b>iv</b>
<b>TABLE OF CONTENTS</b>	<b>v</b>
<b>ABBREVIATIONS</b>	<b>viii</b>
<b>LIST OF FIGURES</b>	<b>ix</b>
<b>LIST OF TABLES</b>	<b>x</b>
<b>CHAPTER 01</b>	<b>1</b>
<b>Introduction</b>	<b>1</b>
<b>1.1. Prolegomena</b>	<b>1</b>
<b>1.2. Background &amp; Motivation</b>	<b>2</b>
<b>1.3. Research Problem Statement</b>	<b>3</b>
<b>1.4. Aim and Objective</b>	<b>3</b>
<b>1.4.1. Aim</b>	<b>3</b>
<b>1.4.2. Objectives</b>	<b>4</b>
<b>1.5. Scope of the Research</b>	<b>5</b>
<b>1.6. Proposed Solution</b>	<b>5</b>
<b>1.7. Overview of the Report</b>	<b>4</b>
<b>1.8. Summary</b>	<b>5</b>
<b>CHAPTER 02</b>	<b>6</b>
<b>Literature Review</b>	<b>6</b>
<b>2.1. Introduction</b>	<b>6</b>
<b>2.2. Related work in Text mining for Identifying Harmful Content on Social Media</b>	<b>6</b>
<b>Comments</b>	
<b>2.3. Harmful Content Identification</b>	<b>11</b>
<b>2.4. Available Tools for Tamil Language</b>	<b>12</b>
<b>2.5. Summary</b>	<b>12</b>
<b>CHAPTER 03</b>	<b>14</b>
<b>Technology adapted in Harmful Content Identification</b>	<b>14</b>

<b>3.1. Introduction</b>	<b>14</b>
<b>3.2. Text Mining Techniques</b>	<b>14</b>
<b>3.3. Machining Learning Classifiers</b>	<b>14</b>
<b>3.3.1. Logistic Regression( LR)</b>	<b>15</b>
<b>3.3.2. Support Vector Machine( SVM)</b>	<b>16</b>
<b>3.3.2. Naïve Bayes( NB)</b>	<b>16</b>
<b>3.4 Rapid Miner Studio</b>	<b>16</b>
<b>3.4.1. Weka</b>	<b>17</b>
<b>3.4.2. Text Processing</b>	<b>17</b>
<b>3.5. Summary</b>	<b>18</b>
<b>CHAPTER 04</b>	<b>19</b>
<b>A novel approach for Identifying Harmful Contents</b>	<b>19</b>
<b>4.1. Introduction</b>	<b>19</b>
<b>4.2. Hypothesis</b>	<b>19</b>
<b>4.3. Input</b>	<b>19</b>
<b>4.4. Output</b>	<b>19</b>
<b>4.5. Process</b>	<b>20</b>
<b>4.6. Summary</b>	<b>20</b>
<b>CHAPTER 05</b>	<b>21</b>
<b>Analysis and Design</b>	<b>21</b>
<b>5.1. Introduction</b>	<b>21</b>
<b>5.2. The High Level design of the Framework</b>	<b>21</b>
<b>5.3. Summary</b>	<b>22</b>
<b>CHAPTER 06</b>	<b>23</b>
<b>Implementation of the solution</b>	<b>23</b>
<b>6.1. Introduction</b>	<b>23</b>
<b>6.2. Data Corpus Construction</b>	<b>23</b>
<b>6.3. Text Preprocessing</b>	<b>25</b>
<b>6.4. FeatureExtraction</b>	<b>25</b>
<b>6.5. Feature Vectorization</b>	<b>26</b>
<b>6.6. Machine Learning based Classification</b>	<b>26</b>
<b>6.7. Performance Measurements</b>	<b>27</b>
<b>6.8. Summary</b>	<b>28</b>

<b>CHAPTER 07</b>	<b>29</b>
<b>Evaluation</b>	<b>29</b>
<b>7.1. Introduction</b>	<b>29</b>
<b>7.2. Evaluation for classification</b>	<b>29</b>
<b>7.3. Summary</b>	<b>32</b>
<b>CHAPTER 08</b>	<b>33</b>
<b>Conclusion and Future Work</b>	<b>33</b>
<b>8.1. Introduction</b>	<b>33</b>
<b>8.2. Conclusion</b>	<b>33</b>
<b>8.3. Limitations</b>	<b>34</b>
<b>8.4. Future Developments</b>	<b>34</b>
<b>8.5. Summary</b>	<b>34</b>
<b>CHAPTER 09</b>	<b>35</b>
<b>REFERENCES</b>	<b>35</b>



## **ABBREVIATIONS**

TF-IDF	Term Frequency-Inverse Document Frequency
WEKA	Waikato Environment for Knowledge Analysis
NB	Naïve Bayes
SVM	Support Vector Machine
LR	Logistic Regression
TN	True Negative
FP	False Positive
FN	False Negative
TP	True Positive

## **LIST OF FIGURES**

	<b>Page</b>
Figure 3.1- WEKA Extension	17
Figure 3.2- Text Processing Extension	17
Figure 5.1- The high level design of the Framework	21
Figure 6.1- Collected data corpus	23
Figure 6.2- Google form for labelling External data	24
Figure 6.3- Machine Learning based Classification Model	26
Figure 6.4- Performance Measurements	28

## LIST OF TABLES

	<b>Page</b>
Table 2.1- Summary of Feature Extraction Techniques	10
Table 2.2- Summary of Feature Vectorization Techniques	10
Table 2.3- Summary of Machine Learning Algorithms used for harmful content identification	10
Table 2.4- Summary of Performance values	11
Table 6.1- Details of Data Corpus	23
Table 7.1- Results with LR Classification Technique	29
Table 7.2- Results with SVM Classification Technique	30
Table 7.3- Results with Naïve Bayes Classification Technique	30
Table 7.4- Comparison between the test data results and external data results	31
Table 7.5- Details of the best fit model	31