

**WORD LEVEL LANGUAGE IDENTIFICATION OF CODE-  
MIXING TEXT IN SOCIAL MEDIA USING NLP**

Kasthuri Shanmugalingam

168287D

Degree of Master of Science in Artificial Intelligence

Department of Computational Mathematics

University of Moratuwa

Sri Lanka

April 2019

# **WORD LEVEL LANGUAGE IDENTIFICATION OF CODE- MIXING TEXT IN SOCIAL MEDIA USING NLP**

Kasthuri Shanmugalingam

168287D

Thesis submitted in partial fulfilment of the requirements for the degree of Master of Science  
in Artificial Intelligence

Department of Computational Mathematics

University of Moratuwa

Sri Lanka

April 2019

## **Declaration**

I declare that this is my own work and this thesis does not incorporate without acknowledgement any material previously submitted for a Degree or Diploma in any other University or institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text.

Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my thesis/dissertation, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works (such as articles or books).

Name of Student

Kasthuri Shanmugalingam

Signature of Student:

Date:

The above candidate has carried out research for the Master's Dissertation under my supervision.

Name of Supervisor

Dr. Sagara Sumathipala

Signature of Supervisor:

Date:

## **Acknowledgements**

I would like to make this opportunity to express my sincere gratitude to my supervisor Dr. Sagara Sumathipala for his valuable guidance extended throughout the research. This research would not have been completed to success without his immense support and guidance. Further, I would like to thank prof. Asoka Karunananda to gave valuable guidance to the documentation of the work during the lectures. And I would like to thank all academic staff of the Department of Computational Mathematics to gave sufficient knowledge to complete this research.

I wish to express my sincere thanks to my family members and colleagues who stood beside me whenever I need and helped me always with support, advice, and encouragement to complete my research work successfully.

## Abstract

Automatic analyzing and extracting useful information from the noisy social media content are currently getting more attention from the research community. Recent days people easily mixing their native language along with the English language together to express their thoughts in social media, using the Unicode characters written in Roman Scripts. Thus these types of noisy code-mixed text are characterized by a high percentage of spelling mistakes with phonetic typing, wordplay, creative spelling, abbreviations, Meta tags, and so on. Identification of languages at word level become as necessary part for analyzing the noisy content in social media. It would be used as an intimate language identifier for chatbot application by using the native languages.

For this study used Tamil-English and Sinhala-English code-mixed text from social media. Natural Language Processing (NLP) and Machine Learning (ML) technologies used to identify the language tags at the word level. A novel approach proposed for this system implemented as machine learning classifier based on features such as Tamil Unicode characters in Roman scripts, dictionaries, double consonant, and term frequency used for Tamil-English code-mixed text and features such as Sinhala Unicode characters written in Roman scripts, dictionaries, and term frequency used for Sinhala-English code-mixed text.

Different machine learning classifiers such as Support Vector Machines (SVM), Naive Bayes, Logistic Regression, Random Forest and Decision Trees used in the model evaluation process. Ten-fold cross-validation used to evaluate the performance based on language tags at the word level. Among that the highest accuracy of 89.46% was obtained in SVM classifier and 90.5% was obtained in Random Forest classifier for Tamil-English (Tanglish) and Sinhala-English (Singlish) code-mixed text respectively.

In the testing process of Tanglish model with SVM and Singlish model with Random Forest gave accuracy as 93.87% and 95.83% respectively for the testing unseen data. Tanglish model with SVM gave F-Measure for 'tam' and 'eng' tags were 0.965 and 0.894 respectively. Singlish model with Random Forest gave F-Measure for 'sin' and 'eng' tags were 0.975 and 0.929 respectively. So this the evidence that most of the times the Tanglish model with SVM and Singlish model with Random Forest predict the language labels correctly at word level.

# Table of Contents

Declaration	i
Acknowledgements	ii
Abstract	iii
List of Figures	vii
List of Tables	viii
<b>Chapter 1 Introduction</b>	<b>1</b>
1.1 Prolegomena	1
1.2 Objectives	2
1.3 Background and Motivation	2
1.4 Code Mixing Problem in Brief	3
1.5 Proposed Solution	4
1.6 Resource Requirements	4
1.7 Structure of the Thesis	4
1.8 Summary	5
<b>Chapter 2 Code Mixing in Social Media – Practices and Issues</b>	<b>6</b>
2.1 Introduction	6
2.2 State of the art of language identification of code-mixed text in social media	6
2.2.1 Code Mixing	6
2.2.2 Language Identification	8
2.3 Future Trends	12
2.4 Summary	13
<b>Chapter 3 Natural Language Processing and Machine Learning</b>	<b>14</b>
3.1 Introduction	14
3.2 Artificial Intelligence	14
3.3 Natural Language Processing	15
3.3.1 Natural Language Tool Kit (NLTK)	16
3.3.2 Pandas	16
3.3.3 Numpy	17
3.4 Machine Learning	17
3.4.1 Supervised Learning	17
3.4.2 Weka and Classifiers used in development of models	18

3.4.2.1 Support Vector Machine (SVM)	19
3.4.2.2 Logistic Regression	20
3.4.2.3 Naïve Bayes	20
3.4.2.4 Decision Tree	22
3.4.2.5 Random Forest	23
3.5 Summary	24
<b>Chapter 4 Novel Approach to Language Identification of Code-Mixing Text</b>	<b>25</b>
4.1 Introduction	25
4.2 Hypothesis	25
4.3 Process	25
4.4 Input	26
4.5 Output	27
4.6 Users	27
4.7 Summary	27
<b>Chapter 5 Design</b>	<b>28</b>
5.1 Introduction	28
5.2 Language Identification System for Tamil-English Code-Mixed Text	28
5.2.1 Dataset Description	29
5.2.2 Preprocessing	29
5.2.3 Feature identification of Tamil-English code-mixed text	29
5.2.3.1 Tamil Unicode characters in Roman scripts	29
5.2.3.2 Language-specific dictionaries	31
5.2.3.3 Double consonants	31
5.2.3.4 Term Frequency	31
5.3 Language Identification System for Sinhala-English Code-Mixed Text	31
5.3.1 Dataset Description	32
5.3.2 Preprocessing	33
5.3.3 Feature identification of Sinhala-English code-mixed text	33
5.3.3.1 Sinhala Unicode Characters in Roman Scripts	33
5.3.3.2 Language-Specific Dictionaries	35
5.3.3.3 Term Frequency	35
5.4 Summary	35

<b>Chapter 6 Implementation</b>	<b>36</b>
6.1 Introduction	36
6.2 Language Identification System for Tamil-English Code-Mixed Text	36
6.2.1 Preprocessing	36
6.2.2 Feature identification of Tamil-English code-mixed text	38
6.2.3 Model Development for Tamil-English code-mixed text	40
6.3 Language Identification System for Sinhala-English Code-Mixed Text	41
6.3.1 Preprocessing	41
6.3.2 Feature identification of Sinhala-English code-mixed text	43
6.3.3 Model Development for Sinhala-English code-mixed text	45
6.4 Summary	45
<b>Chapter 7 Evaluation</b>	<b>46</b>
7.1 Introduction	46
7.2 Experimental design	46
7.2.1 Experimental design for Model Evaluation	46
7.2.1.1 Evaluation Strategy for predictive Models	47
7.2.2 Experimental design for Testing of Models	48
7.2.2.1 Evaluation Strategy for Testing Models	49
7.3 Experimental Results	49
7.3.1 Experiment Results for Model Evaluation	49
7.3.2 Experiment Results for Model Testing	53
7.4 Summary	53
<b>Chapter 8 Conclusion and Future Work</b>	<b>54</b>
8.1 Introduction	54
8.2 Concluding remarks	54
8.3 Limitation and Future work	55
8.4 Summary	56
<b>References</b>	<b>57</b>
<b>Appendix A</b>	<b>60</b>
<b>Appendix B</b>	<b>67</b>