



INFORMATICS  
INSTITUTE OF  
TECHNOLOGY

INFORMATICS INSTITUTE OF TECHNOLOGY

In Collaboration With

ROBERT GORDON UNIVERSITY

**Classification of Cyberbullying Romanized Sinhala Comments in  
Online Platforms**

A Dissertation By

Yasindi Hettiarachchi

20210303 | 2018109

Supervised By

Mr. Dinesh Asanka

Submitted in partial fulfillment of the requirements for the MSc in Business Analytics Degree  
at the Robert Gordon University

**April 2023**

# Abstract

This study investigates the detection of cyberbullying in Romanized Sinhala using various machine learning classifiers and feature extraction methods. The primary objective is to identify the most effective combination of classifier and feature extraction techniques for this task. We employ rule-based, Bag-of-Words (BoW), and Term Frequency-Inverse Document Frequency (TF-IDF) feature extraction methods, as well as additional features such as word count and gender. The classifiers studied include K-Nearest Neighbours (KNN), Voting, Random Forest, Support Vector Machines (SVM), Decision Tree, Naive Bayes, Multilayer Perceptron (MLP), AdaBoost, and Logistic Regression.

The results demonstrate that the performance of classifiers varies significantly across different feature extraction techniques. The rule-based method appears to be more effective than the TF-IDF and BoW methods for this task. The KNN, Voting, and Random Forest classifiers show the best performance when combined with rule-based, BoW, and TF-IDF feature extraction methods, respectively. Additionally, text pre-processing, such as emoticon removal, stop word removal, and tokenization, significantly impacts the performance of rule-based feature extraction.

Furthermore, the research explores the effect of incorporating gender and word count as independent variables on the performance of classifiers. The results reveal that the inclusion of word count significantly improves the performance of classifiers, while the impact of gender remains inconclusive due to insufficient data.

The study also tests several hypotheses, examining the performance differences among classifiers, the effectiveness of rule-based features, and the impact of text pre-processing on rule-based feature extraction. The results indicate that the selection of appropriate feature extraction methods and text pre-processing techniques play crucial roles in detecting cyberbullying content in Romanized Sinhala. The rule-based method, when combined with classifiers such as Random Forest, KNN, and Voting Classifier, demonstrates promising results for this task.

**Keywords:** Machine Learning, Cyberbullying, Text Mining, Romanized Sinhala, NLP, Feature Extraction