# MSc Project Report

# Detection of hate speech written in Sinhala and Singlish language posted on social media by users in Sri Lanka using text analytics

Ramachandran Rajindra Sheran

2019

A report submitted as part of the requirements for the degree of
MSc in Big Data Analytics at Robert Gordon University, Aberdeen,
Scotland

# Abstract

Detection of hate speech is part of sentiment analysis and sentiment analysis refers to the task of natural language processing to determine whether a piece of text contains some subjective information and what subjective information it expresses, that is, whether the attitude behind this text is positive, negative or neutral. Understanding the opinions behind user-generated content automatically is of great help for commercial, social and political use, among others. The task can be conducted on different levels, classifying the polarity of words, sentences or entire documents. It is one of the most active research areas in natural language processing and text mining in recent years. Its popularity is mainly due to two reasons. First, it has a wide range of applications because opinions are central to almost all human activities and are key influencers of our behaviours. Whenever we need to make a decision, we want to hear others' opinions. Secondly, it presents many challenging research problems. Part of the reason for the lack of study before was that there was little opinionated text in digital forms. It is thus no surprise that the inception and the rapid growth of the field coincide with those of the social media on the Web. In fact, the research has also spread outside of computer science to management sciences and social sciences due to its importance to business and society as a whole.

The attempt done by this research is to detect hate content that in textual content uploaded to or exchanged through social media and are primarily in Sinhala or Singlish context. Here, the detection will not be restricted to one domain since the main aim of this research is to capture any hate content irrespective of the domain. Many previously used machine learning techniques used for mainly English and other native languages were researched in order to get the basic idea on how to perform sentiment analysis from the scratch.

Keywords: sentiment analysis, Singlish, Sinhala or Singlish context, social media, machine learning techniques, hate content, text mining