# Sentiment indicator for e-mails

K.D.Chandima

I69304M

Faculty of Information Technology
University of Moratuwa
2019

# Sentiment indicator for e-mails

K.D.Chandima

169304M

Supervisor: Dr. Lochandaka Ranathunga

Dissertation submitted to the Faculty of Information Technology, University of Moratuwa, Sri Lanka for the partial fulfillment of the requirements of the Honours Degree of Bachelor of Science in Information Technology.
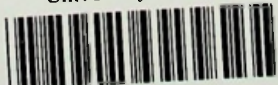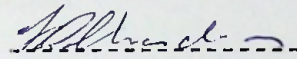
2019

## Declaration

I hereby declare that this project report entitled "Sentiment indicator for e-mails" contains my own work and has not been submitted and will not be submitted in any form for another degree or diploma at any university or other institution of tertiary education. Information derived from the published or unpublished work of others has been acknowledged in the text and a list of references is given.
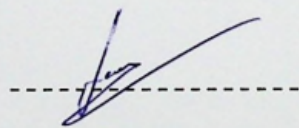
Name of Student: K. D. Chandima

Signature of Student

Date: 25 / 04 / 2019

Supervised By

Dr. Lochandaka Ranathunga

i

## Dedication

I would like to dedicate my project, "Sentiment indicator for e-mails";

To my project supervisor Dr. Lochandaka Ranathunga,

To my Father and sister who have given support for proof reading project documents,

To the Instructors of Faculty of Information Technology at University of Moratuwa who have supported me to find requirements for the project.

## Acknowledgments

I wish to express my sincere thanks to my supervisor Dr. Lochandaka Ranathunga for guiding me to have a successful outcome in this project work.

I would like to acknowledge and extend my heartfelt gratitude to the following people who have made the completion of this project possible –

My father and sister for the help given to proof read all the project documents by spending their valuable time,

Instructors of Faculty of Information Technology at University of Moratuwa who have supported me to find resources and requirements for the project.

Academic, non-academic staff members who helped me throughout this project,

Batch mates who motivated in doing the project, and friends.

iii

## Abstract

People receive large amount of e-mails from friends, relatives, companies, institutions and known and unknown people. It is not possible the user to read all the received e-mails during a busy day. Normally, people avoid reading most of the received e-mails by giving priority to the important e-mails. E-mails may contain happy news or sad news or defamatory contents or obscene contents. If there is a way for a user to get an idea of what kind of news would be there in just before he or she opens e-mails, then he or she can choose which e-mails should be opened first and find out which e-mails would make him happy. If the user can find out any obscene contents in an e-mail just before he opens it, then he or she can avoid the embarrassment of opening the email in front of a stranger and then the user will be able to delete it without simply opening it. A proper e-mail categorization is required and this attempt is to categorize e-mails according to the sentiment and the ingredient of the e-mail content, by indicating emoticons on the subject of e-mails. Text Preprocessing, Feature Extraction, Sentiment Classification methods are used to identify sentiments and proper emoticons are used to indicate sentiments on the subjects of e-mails. Machine learning and lexicon based approaches are used to predict the sentiment of emails. A better accuracy level is expected from machine learning in the process of sentiment extraction.

Table of Contents

## List of Figures