**MSc in Information Technology**

# Analyze Quality of Products in E-commerce Systems with Sentimental Analysis

Prepared by

P. M. A. U. Bandara

Index No : 158751L

Supervised by

Mr. SamindaPremaratne

Faculty of Information Technology

University of Moratuwa

**March 2019**

**MSc in Information Technology**


# Analyze Quality of Products in E-commerce Systems with Sentimental Analysis


Prepared by

P. M. A. U. Bandara

Index No: 158751L


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.


**March 2019**

# Declaration

I declare that this research is my own work and has not been 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 Information derived from the published or unpublished work of others has been acknowledged in the text and the list of references is given.

------------------------------

P. M. A. U. Bandara

(158751L)

2019/04/

I have supervised and accepted this thesis for the submission of the degree.

------------------------------

Mr. Saminda Premaratne

(Main Supervisor)

2019/04/

# Acknowledgements

I would like to express my genuine gratitude to my advisor Mr. S.C Premaratne for the continuous support of my research, for his patience, motivation, and immense knowledge. His direction helped me in all the period of research and writing of this thesis. I could not have imagined having a better advisor and mentor for my research. Also I thank my associates and staff in the University of Moratuwa who helped me carry out this research. Last but not the least, I would like to thank my family, my parents for supporting me spiritually throughout writing this thesis.

# Abstract

E-commerce websites getting extra significant and popular today since the vast differentiated and diversified information that is presented. Studies says that more than 80% of the world population is using these websites to purchase goods and services online. For these online customers, comments / feedbacks play a major role in decision making when buying the products from the market space. Hence the diversity and the popularity of the Online space, sales of these online products get increased with time. Therefore, it is not practical to review all the given product feedback and come to conclusion on purchasing the product for a consumer. Focusing on this, this study is urges to observe the success factors of online websites and how those aspects influence on online marketing to sales growth in any organization. Therefore, in this research a Analyze Quality of Products in E-commerce System is focused on analyzing the online consumers feedbacks or comments on various products using data mining techniques such as Sentimental and filtering analysis. The outcome from the study will show feature wise relativeness in the mobile phone domain. All procedures were based on the features extracted through a thorough literature review and existing apparatuses. This will aid to calculate a "Trust Score" for the online products and a general overview to achieve a higher trust score for e-commerce organization.

# Table of Contents

# List of Figures

# List of Tables