# Data Extraction from Social Media for Sentiment Analysis in order to Predict Sales and Correlated Items for Fashion Industry

C.S Kumanayaka

169315X

Faculty of Information Technology

University of Moratuwa

December 2018

# Data Extraction from Social Media for Sentiment Analysis in order to Predict Sales and Correlated Items for Fashion Industry

C.S Kumanayaka

169315X

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

**December  2018**

# Declaration

I declared that this thesis is written by myself and it has not submitted by other institution of education, degree or diploma of any other university. Information retrieved from unpublished and published work listed in the reference area.

Name of Student                                            Signature of Student

Date:

Supervised by

Name of Supervisor(s)                              Signature of Supervisor(s)

Date:

# Acknowledgment

Initially I am thankful to my supervisor Mr. Saminda Premaratne the valuable advice and supervision. I like to give my gratitude to him for his important time for long discussions throughout my research. His rightful guidance paved the way to achieve priceless achievements

Special thanks should also go to Mr Sankha Perera for the advice and guidance provided from a Statistical point of view. I am also thankful for Mr Prabhath Gunawardane for his valuable guidance and advice in Data Science. Gratitude should also go to Ms.Ayesha Kasthuriarachchi for the statistical advice and guidance.

Also, I should be thankful to all the lecturers about the assistance and guidance provided in carrying out the project. Finally, I should also be grateful to all my batch mates and all the other parties who helped us with their valuable ideas to proceed with my research successfully.

# Abstract

These days' Social media is being very popular for marketing. One of the most popular social media platforms is Facebook. Most of the Fashion brand have Facebook pages. Consumer expresses their feeling using comments and emotional buttons. The user interacts with the brand page using post, like, share comments. The analyzed data can give support to decision makers to the evaluation of the customer's feedback, identify potential customers and predict sales item for the upcoming month and predict correlated items. The purpose of this paper is to explain how to extract and prepare data collected on Facebook to perform sentiment Analysis.

# Table of Content

# List of Figures

# Abbreviations

FB      Facebook

ETL    Extraction transformation load

SSIS    SQL Server Integration service

SSRS   SQL Server Reporting service

BI       Business Intelligence