# SEMANTIC INFORMATION RETRIEVAL BASED ON TOPIC MODELING AND COMMUNITY INTERESTS MINING

R.P. Minuri Chathurika Rajapaksha

(158734M)

Degree of Master of Science in Artificial Intelligence

Department of Computational Mathematics

University of Moratuwa

Sri Lanka

March 2019

# SEMANTIC INFORMATION RETRIEVAL BASED ON TOPIC MODELING AND COMMUNITY INTERESTS MINING

R.P. Minuri Chathurika Rajapaksha

(158734M)

Dissertation submitted in partial fulfillment of the requirements for the degree Master of Science in Artificial Intelligence

Department of Computational Mathematics

University of Moratuwa

Sri Lanka

March 2019

# DECLARATION

I declare that this dissertation does not incorporate, without acknowledgment, any material previously submitted for a Degree or a Diploma in any University and to the best of my knowledge and belief, it does not contain any material previously published or written by another person or myself except where due reference is made in the text. Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my thesis/dissertation, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works (such as articles or books).

Name of Student                                        Signature of Student

Minuri Rajapaksha                                      Date:

The above candidate has carried out research for the Masters Dissertation under my supervision.

Supervised by                                          Signature of Supervisor

Dr. Thushari Silva                                     Date:

# Abstract

Search engines or localized software systems developed for information searching, play an important role in knowledge discovery. Proliferation of data in the web and social media has posed significant challenges in finding relevant information efficiently even using those search engines or other software systems. Moreover, those systems or engines tend to collect massive number of data, which could be useful for humans in various ways but overlook the meaning of the search phrases, hence generate irrelevant search results. A unit level searching i.e. searching information within a website or page is also not effective as they follow exact keyword matching techniques and ignore the semantic level matching of search phrases. In order to address those deficiencies, this research proposes a hybrid approach which use the semantics of data, community preferences as well as collaborative filtering techniques for semantic information retrieval. More specifically, Topic modeling based on Latent Dirichlet Allocation together with topic-driven based community detection methods are applied for identifying personalized search results and hence improve the relatedness of the research results. Based on the proposed hybrid approach a framework for semantic search that can easily be integrated to a software application has been implemented. The evaluation results confirm the effectiveness of search results which outperform benchmark approaches that follow traditional keyword search algorithms.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES