# TAMIL NEWS CLUSTERING

M.S. Faathima Fayaza

179318T

Degree of Master of Science in Computer Science

Department of Computer Science and Engineering

University of Moratuwa
Sri Lanka

May 2020

# TAMIL NEWS CLUSTERING

M.S.Faathima Fayaza

179318T

Thesis/Dissertation submitted in partial fulfillment of the requirements for the degree Master of Science in Computer Science specializing in Data Science Engineering and Analytics

Department of Computer Science and Engineering

University of Moratuwa
Sri Lanka

May 2020

# DECLARATION

I declare that this is my own work and this dissertation does not incorporate without acknowledgement any material previously submitted for degree or diploma in any other university or institute of higher learning and to the best of my knowledge and belief, it does not contain any material previously published or written by another person except where the acknowledgement is made in the text.

Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my 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).


Signature: ………………                      Date: …………….

Name: M.S. Faathima Fayaza


The above candidate has carried out research for the Master of Science thesis under my supervision.


Signature of the supervisor: …………………         Date: ……………..

Name: Dr. Surangika Ranathunga

# ABSTRACT

The web has an abundance of online news articles that are updated frequently. Readers face difficulty in discovering content of interest from the overwhelming news sources and feel tired browsing various websites. This situation is valid in the case of Tamil online news as well, and the number of online news articles published in the Tamil language is on the rise. To address this issue, news aggregators and clustering techniques come into play. Even though there are many news aggregators available for languages like English, the only news aggregator that supports Tamil is Google news, which is a noticeable shortage. Google news mainly covers the Indian news and gives high weightage to the words that appear on the headline rather than those appearing in the body of the news when searching for the news [1].

This research focuses on clustering Tamil online news articles into related topics. There are several clustering techniques and similarity measures used to cluster the documents in the literature for other languages. Tamil is an agglutinative language, meaning that the techniques used for English documents might not readily work for Tamil. The purpose of this research is to study the techniques available for other languages and develop a mechanism to cluster the Tamil online news articles according to their content similarity.

As the first step of this study, ten different datasets were created by collecting news from nine different news providers. Data was collected on nonadjacent days to get diversified data. TF-IDF and word embedding techniques were used to create vector representations of data. One pass algorithm and affinity propagation algorithm were used to cluster the news articles, since the number of clusters cannot be predefined and there is a high number of single news clusters. We achieved the best solution when applying word embedding with one pass algorithm. As another contribution of this research, we were able to create a Tamil word embedding model with 21,077,843 words.

**Keywords:** Clustering, TF-IDF, Word embedding, One pass algorithm, Affinity propagation, Cosine similarity, Crawler

# ACKNOWLEDGEMENTS

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| Abbreviation | Description |
| --- | --- |
| URL | Universal Resource Locator |
| WWW | World Wide Web |
| HTML | Hyper Text Markup Language |
| ANN | Artificial Neural Networks |
| VSM | Vector Space Model |