

**PRE-TRAINED LANGUAGE MODEL-BASED SEMI-  
SUPERVISED LEARNING APPROACH FOR  
CONTENT-BASED EMAIL CATEGORIZATION**

Nuwan Dinusha Kankanamge

189383 U

Degree of Master of Science in Artificial Intelligence

Department of Computational Mathematics

University of Moratuwa

Sri Lanka

October 2020

# **PRETRAINED LANGUAGE MODEL-BASED SEMI-SUPERVISED LEARNING APPROACH FOR CONTENT-BASED EMAIL CATEGORIZATION**

Nuwan Dinusha Kankanamge

189383 U

This thesis submitted in partial fulfilment of the requirements for the degree of Master of Science in Artificial Intelligence

Department of Computational Mathematics

University of Moratuwa

Sri Lanka

October 2020

# Declaration

I declare that this dissertation does not incorporate, without acknowledgement, 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. I also hereby give consent for my dissertation, if accepted, to be made available for photocopying and for interlibrary loans, and for the title and summary to be made available to outside organizations.

## *UOM Verified Signature*

Name of Student

Signature of Student

N.D. Kankanamge

Date: 22/05/2021

Supervised by

Name of Supervisor(s)

Signature of Supervisor(s) *UOM Verified Signature*

Dr. (Ms.) A.T.P. Silva

Date: 24/05/2021

## **Dedication**

I dedicate this thesis to my parents for the support they have given me to fulfil my dreams. I dedicate this work to people at the University of Moratuwa for their help and guidance throughout this research.

## **Acknowledgements**

I would like to express my sincere gratitude to my supervisor Dr. A.T.P. Silva for guiding me throughout the project. She has given his fullest cooperation to me whenever I sought advice.

I must thank Prof. Asoka Karunananda for the advice given, which had a significant impact on making my mind to think differently and his advice and techniques helped me a lot in the process of developing this solution.

I must also thank all members of the lecture panel. During one year and three months, these valuable lectures helped me to think differently.

For this research, I had to refer to many books and research papers as reference. I would like to thank all the authors of those publications.

Also, my batch mates helped me a lot in many ways. I would like to pay my gratitude for them. They made my life enjoyable during the course period.

I would like to place my gratitude to my loving parents and wife for always encouraging me on higher studies.

Last but not least, I would like to thank all my colleagues and others who are not mentioned, for all the support extended to me. Without their dedication, the project would not have been successful.

## Abstract

Recent developments in neuroscience have revolutionized modern trends in artificial intelligence. Artificial neural networks (ANN), which is the artificial model of the human brain, have started to dominate in the field of artificial intelligence. The major usage of ANN is for data classification and prediction. There are numerous applications of ANN, ranging from health, education, entertainment, and business.

Email classification has been an issue for many of the large organizations as it needs human interaction. There are many artificial intelligence-based solutions have been proposed. When it comes to content-based email filtering, many recent researchers have identified that the use of ANN-based approaches are much more useful than conventional natural language modelling methods, as the volume of data increased. One reason for this is ANN has been able to capture some of the hidden styles of writing which have not been captured by conventional natural language processing. However conventional ANN has been suffering from lack of labeled data for training. This has been the major drawback of conventional ANN approach as generating labeled data needs human interaction and therefore making it a costly process. This has limited ANN solutions from providing a generic approach for email classification in any domain since to succeed, it needs large a number of labeled data from each of these domains to train the particular ANN.

This thesis report on our research on content-based email classification using semi-supervised learning which will address the issues with conventional ANN. Semi-supervised learning was introduced around 15 years back but came to play an important role in the field of artificial intelligence recently. Semi-supervised learning provides a solution to this issue as it needs a minimum amount of labeled data for training and it can use unlabeled data to increase its' accuracy. Proposed solution is a multi-view core-training approach that takes labeled emails, unlabeled emails and the names of the different categories as inputs. Output of the project is a trained model that can classify emails to given categories. We have tested our solution with 10000 training samples where only 10% to 20% were given to the system as labeled data and others were used as unlabeled data. We managed to achieve around 0.888 accuracy which is more than 5% accuracy improvement from the total system.

# Contents

<b>1. Introduction</b> .....	1
<b>1.1. Prolegomena</b> .....	1
<b>1.2. Importance of Email Categorization</b> .....	1
<b>1.3. Aim and Objectives</b> .....	2
<b>1.4. Summary</b> .....	3
<b>2. Developments and Challenges in Email Categorization</b> .....	4
<b>2.1. Introduction</b> .....	4
<b>2.2. Early Developments in Email Categorization</b> .....	4
<b>2.3. Modern Developments in Email Categorization</b> .....	5
<b>2.4. Challenges and New Trends in Email Categorization</b> .....	6
<b>2.5. Problem Definition</b> .....	8
<b>2.6. Summary</b> .....	8
<b>3. Technologies Adapted</b> .....	9
<b>3.1. Introduction</b> .....	9
<b>3.2. Semi-Supervised Learning</b> .....	9
<b>3.2.1. Multi-View Core-Training</b> .....	11
<b>3.3. Pre-Trained Language Models</b> .....	12
<b>3.3.1. BERT</b> .....	12
<b>3.3.2. ELMO</b> .....	13
<b>3.3.3. Why BERT and ELMO</b> .....	14
<b>3.4. Summary</b> .....	15
<b>4. Approach</b> .....	16
<b>4.1. Introduction</b> .....	16
<b>4.2. Hypothesis</b> .....	16
<b>4.3. Input</b> .....	16
<b>4.4. Output</b> .....	17
<b>4.5. Process</b> .....	17
<b>4.6. Features</b> .....	18
<b>4.7. Users</b> .....	18
<b>4.8. Summary</b> .....	18
<b>5. Design of the Semi-Supervised Content-Based Email Categorization Solution</b> .....	19
<b>5.1. Introduction</b> .....	19
<b>5.2. Data Pre-Processing Module</b> .....	20
<b>5.3. Feature Extractor Module</b> .....	20

5.4.	<b>Semi-Supervised Training Module</b> .....	21
5.5.	<b>Summary</b> .....	23
6.	<b>Implementation of the Semi-Supervised Content-Based Email Categorization Solution</b> ....	24
6.1.	<b>Introduction</b> .....	24
6.2.	<b>Dataset Selection</b> .....	24
6.3.	<b>Data Pre-Processing Module</b> .....	24
6.4.	<b>Feature Extractor Module</b> .....	26
6.4.1.	<b>BERT Feature Extractor</b> .....	26
6.4.2.	<b>ELMO Feature Extractor</b> .....	27
6.5.	<b>Semi-Supervised Training Module</b> .....	28
6.6.	<b>Summary</b> .....	33
7.	<b>Evaluation</b> .....	34
7.1.	<b>Introduction</b> .....	34
7.2.	<b>Evaluation of Text Pre-Processing Module</b> .....	34
7.3.	<b>Evaluation of Semi-Supervised Learning Process</b> .....	35
7.3.1.	<b>Results by Changing the Amount of Labelled and Unlabelled Data</b> .....	37
7.3.2.	<b>Results by Changing the Threshold Value</b> .....	38
7.3.3.	<b>Results by Increasing Iterations</b> .....	40
7.4.	<b>Evaluation with Benchmark Dataset</b> .....	41
7.5.	<b>Summary</b> .....	42
8.	<b>Conclusion &amp; Further Work</b> .....	43
8.1.	<b>Introduction</b> .....	43
8.2.	<b>Conclusion</b> .....	43
8.3.	<b>Limitations &amp; Further Work</b> .....	43
8.4.	<b>Summary</b> .....	44
	<b>References</b> .....	45
	<b>Appendix A</b> .....	47



## List of Figures

Figure 3.1. Semi-supervised learning types.....	10
Figure 3.2. ELMO Architecture.....	14
Figure 4.1. High-level approach of the diagram.....	17
Figure 5.1. High-level design diagram.....	19
Figure 5.2. Classifier Architecture.....	21
Figure 5.3. Multi-view core-training Architecture.....	22
Figure 7.1. Impact of data preprocessing.....	35
Figure 7.2. Impact of data balancing – BERT.....	36
Figure 7.3. Impact of data balancing – ELMO.....	37
Figure 7.4. Impact of labeled & unlabeled data amounts.....	38
Figure 7.5. Results for different threshold values.....	39
Figure 7.6. Results for 15 iterations.....	40
Figure 7.7. Results for 8 iterations for AG News Classification Dataset.....	42

## **List of tables**

Table 2.1 - Summary of literature review.....	07
---	----

## **Abbreviations**

ANN	-	Artificial Neural Network
SVM	-	Support Vector Machine
NLP	-	Natural Language Processing
LSTM	-	Long Short-Term Memory