



Tribhuvan University
Institute of Science and Technology

**A Comparative Study of Naive Bayes and Support
Vector Machine Classifier for Nepali News
Classification**

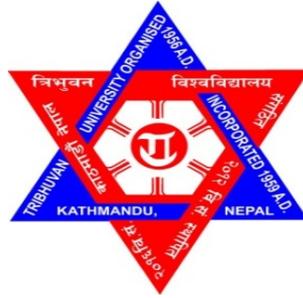
Dissertation
Submitted To

Central Department of Computer Science & Information Technology
Kirtipur, Kathmandu, Nepal

In partial fulfillment of the requirements
for the Degree of Master of Science in Computer Science & Information
Technology

Submitted By
Ganesh Bahadur Ayer
Roll No. : 18 (2010-2012)
Date: 28 Dec, 2015

Supervisor
Mr. Nawaraj Paudel



Tribhuvan University
Institute of Science and Technology
Central Department of Computer Science and Information Technology

Date: 28 Dec, 2015

Student's Declaration

I hereby declare that I am the only author of this work and that no sources other than the listed here have been used in this work.

.....

Ganesh Bahadur Ayer

M.Sc. CSIT (2010-2012)

Central Department of Computer Science

And Information Technology,

Institute of Science and Technology,

Kirtipur, Kathmandu, Nepal



Tribhuvan University
Institute of Science and Technology
Central Department of Computer Science and Information Technology

Date: 28 Dec, 2015

Supervisor's Recommendation

I hereby recommend that the dissertation prepared under my supervision by **Mr. Ganesh Bahadur Ayer** entitled “**A Comparative Study of Naive Bayes and Support Vector Machine Classifier for Nepali News Classification**” be accepted as in fulfilling partial requirement for the completion of Master's Degree of Science in Computer Science & Information Technology. In my best knowledge this is an original work in computer science.

Mr. Nawaraj Paudel

Head,
Central Department of Computer Science and Information Technology,
Institute of Science and Technology,
Kirtipur, Kathmandu, Nepal



Tribhuvan University
Institute of Science and Technology
Central Department of Computer Science and Information Technology

Date: 28 Dec 2015

LETTER OF APPROVAL

We certify that we have read this dissertation work and in our opinion it is appreciable for the scope and quality as a dissertation in the partial fulfillment of the requirements of Master's Degree of Science in Computer Science & Information Technology.

Evaluation Committee

Mr. Nawaraj Paudel
Head of Department
Central Department of Computer Science
& Information Technology
Tribhuvan University
Kirtipur

(External Examiner)

Mr. Nawaraj Paudel
Head of Department
Central Department of Computer Science
and Information Technology (T.U)
(Supervisor)

(Internal Examiner)

ACKNOWLEDGEMENTS

Above all, I thank God for his blessing and submit my graduate to my almighty for providing strength and confidence in me to complete this work.

Secondly, I would like to extent my, gratitude and sincerest thanks to my respected supervisor **Mr. Nawaraj Poudel** Head of Computer Science & IT Department Kirtipur, TU (Kathmandu, Nepal) for his impressive tutelage, constructive criticism and intellectual support best owed for me sacrificing his invaluable time. His crucial role to make this report culminate is indescribable.

I would like to express my gratitude to respected teachers Prof. Dr. Shashidhar Ram Joshi, Prof. Dr. Subarna Shakya, Prof. Sudarshan, Karanjeet, Mr. Min Bahadur Khati, Mr. Bishnu Gautam, Mr. Jagdish Bhatt, Mr. Bikash Balami, Mr. Sarbin Syami, Mr. Dhiraj Pandey, Mr. Arjun Singh Saud and others staffs of CDCSIT for granting me broad knowledge and inspirations within the time of period of two years.

I wish to thank to all my colleagues and friends especially **Mr. Tej Bahadur Sahi**, Mr. Ashok Kumar Pant, Mr. Niranjana Kathayat, Mr. Dipak Prasad Bhatt and Mr. Rajendra Prasad Joshi for supporting me directly and indirectly in this research work.

From the beginning to end in all count, I would like to thank my family members for their love, support and encouragement.

As we know that, there won't be 100% accuracy and efficiency in any work done by both machine any human, so there may be some errors in my project. But, I have done my best to complete this dissertation, so any suggestion regarding the mistakes of this work will be always welcomed.

Ganesh Bahadur Ayer

ABSTRACT

Automated document classification is the task of assigning the given document into some class of interest. Text classification is the subset of document classification as document can be text, image, music, etc. Document classification has many applications in library science, information science, computer science and others. It can be used for intellectual categorization of documents, indexing of documents, filtering of spams, routing of emails, identification of language, classification of genre, etc.

The problem of automated document classification can be solved in supervised, unsupervised or semi-supervised way. Most of the learning and classification algorithms use document attributes and human inference to learn and classify given documents. In this dissertation work, many Natural Language Processing (NLP) techniques are used for document processing and attribute selection. And, two learning based classification techniques are used namely, Support vector machine (SVM) and Naive Bayes Classifier.

For the evaluation of the system, we have created Nepali text datasets for five classes of documents: Business, Crime, Education, Health and Sports. There are two separate datasets for training and testing of the system. SVM classification system has the average system accuracy rate of 86.34%, precision rate of 84% recall rate of 94.4%. Similarly, Naive Bayes classification system has the average system accuracy rate of 88.8%, precision rate of 92.23% and recall rate of 88.87%.

Keywords:

Automated Document Categorization, Text Classification, Natural language processing, Nepali language, Preprocessing, Feature extraction, Artificial Neural Networks, Support vector machine, Naive Bayes Classifier

Table of Contents

Acknowledgement

Abstract

List of Figures

List of Tables

Abbreviation

INTRODUCTION	1
1.1 Introduction.....	1
1.1.1 Document Classification	1
1.1.2 Text Classification	2
1.1.3 Principles of Text Classification	4
1.1.4 Methods of Text classification	4
1.1.4.1 Manual text classification	4
1.1.4.2 Automatic text classification.....	4
1.1.5 Nepali Text Classification.....	5
1.2 Applications of Document Classification	5
1.3 Motivation.....	6
1.4 Problem Definition.....	6
1.5 Objective.....	7
1.6 Contribution of this Dissertation.....	7
1.7 Outline of the Document.....	8
BACKGROUND AND LITERTURE REVIEW	9
2.1 Literature review	9
2.2 Machine Learning	13
2.2.1 Types of Machine Learning Algorithms	13
2.3 Overview of Data Mining Concept	14
2.4 Data Mining	14
2.4.1 Purpose of Data mining	15

2.4.2 Technique in Data mining.....	15
2.4.2.1 Association rule.....	15
2.4.2.2 Clustering.....	15
2.4.2.3 Decision Tree.....	15
2.4.2.4 Neural Network.....	16
2.4.2.5 Classification and prediction.....	16
2.5 Automatic Text Classification	16
2.5.1 Supervised Text Classification.....	16
2.5.2 Unsupervised Text Classification	17
2.5.3 Semi-supervised text Classification.....	17
2.6 Automatic Text Classification Techniques	17
2.6.1 Support Vector Machines.....	17
2.6.2 Artificial neural networks	17
2.6.3 K-Nearest Neighbor	18
2.6.4 Decision Tree	19
2.6.5 Naive Bayes Probabilistic Model.....	20
2.7 Document Representation	20
2.7.1 Document Term Matrix	20
2.7.2 TF-IDF Method.....	21
RESEARCH METHODOLOGY	22
3.1 Automatic Text Classification System Overview	22
3.2 Data acquisition	23
3.3 Preprocessing	23
3.3.1 Stop word removal.....	24
3.3.2 Symbols removal	24
3.3.3 Stemming	25
3.4 Feature Extraction	26
3.4.1 Term frequency-inverse document frequency (TF-IDF)	26
3.5 Bayes Theorem	27
3.5.1 Naive Bayesian Classifier	27
3.5.2 The Naive Bayes probabilistic model	28
3.6 Support Vector Machine	29

3.6.1 Multi-Class Classification.....	30
3.7 Categories of classification for Nepali documents.....	31
3.8 System Evaluation Measures	31
3.8.1 Average System Accuracy	31
3.8.2 Precision.....	31
3.8.2 Recall	32
PROGRAM DEVELOPMENT AND IMPLEMENTATION	33
4.1 Development Methodology and Tools.....	33
4.1.1 Programming Language Used.....	33
4.1.2 Eclipse IDE	33
4.2 Machine Learning Library and Plug-ins	34
4.2.1 Weka	34
EXPERIMENTATIONS AND RESULTS.....	36
5.1 Training and Testing Datasets.....	36
5.1.1 Training Datasets	39
5.1.1.1 Dataset I	39
5.1.1.2 Dataset II.....	39
5.1.2 Testing Dataset.....	40
5.2 Data Dictionaries	40
5.2.1 Stop Word Dictionary	40
5.2.2 Symbol Dictionary	41
5.3 Experimentation Results	41
5.3.1 Experiment 1	41
5.3.2 Experiment 2.....	42
5.4 Result Analysis	43
CONCLUSION.....	44
6.1 Conclusion	44
6.2 Limitations and Future Scope	44
Appendix A.....	49
Appendix B.....	50

List of Tables

Table 5.1: Dataset I.....	39
Table 5.2: Dataset II.....	39
Table 5.3: Test Dataset	40
Table 5.4: Experimentation Results (Experiment 1).....	41
Table 5.5: Experimentation Results (Experiment 2).....	42
Table 5.6: Aggregate System Results	43

List of Figures

Figure 1.1: Flat classification.....	2
Figure 1.2: Hierarchical classification	2
Figure 1.3: Binary classification	3
Figure 1.4: Multiple class single level classification	3
Figure 3.1: Top Level DFD of Automated Text Classification System.....	22
Figure 3.2: The Detail Architecture of Automatic Text Classification System for Nepali Language.	23
Figure 3.3: Sample Stop Words.	24
Figure 3.4: Support Vector Machine.....	29
Figure: 5.1: Sample Crime Documents.....	36
Figure: 5.2: Sample Education Documents.....	37
Figure: 5.3: Sample Health Documents	37
Figure: 5.4: Sample Health Documents	38
Figure: 5.5: Sample Business Documents.....	38
Figure 5.6: Stop Word Dictionary.....	40
Figure 5.7: Symbol Dictionary.....	41
Figure 5.8: Graph of Experiment 1.	42
Figure 5.9: Graph of Experiment 2.	43

LIST OF ABBREVIATIONS

AI Artificial Intelligence

ANN Artificial Neural Network

API Application Programming Interface

BPN Back Propagation Network

CDT C/C++ Development Components

IDE Integrated Development Environment

KNN K-Nearest Neighbor

MT Machine Translation

MAP Maximum A Posteriori

NB Naive Bayes

NN Neural Network

NLP Natural Language Processing

PCA Principal Component Analysis

SVM Support Vector Machine

SDK Software Development Kit

TC Text Classification

TF-IDF Term Frequency - Inverse Document Frequency