

**JADAVPUR UNIVERSITY**

**Bengali Text Classification Using SVM**

**BY**

Surajit Maity

EXAMINATION ROLL NO. : MCA2360012

**UNDER THE SUPERVISION OF**

**PROF.(DR.) KAMAL SARKAR**

**PROJECT SUBMMITION OF FULFILLMENT FOR THE DEGREE  
OF MASTER OF COMPUTER APPLICATION**

**IN THE**

**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING  
FACULTY OF ENGINEERING AND TECHNOLOGY**

**2023**

## **Certificate of Recommendation**

This is to certify that the thesis entitled “ Bengali Text Classification Using SVM” has been carried out by Surajit Maity (University Roll Number: 002010503002, Examination Roll No: MCA2360012, University Registration Number: 154210 (2020-2021)), under the guidance and supervision of Prof. (Dr.) Kamal Sarkar, Department of Computer Science and Technology, Jadavpur University, Kolkata, is being presented for the partial fulfilment of the Degree of Master of Computer Applications during the academic year 2022-2023. The research results presented in the thesis have not been included in any other paper submitted for the award of any degree in any other university or institute.

**Prof. (Dr.) Kamal Sarkar (Thesis Supervisor)**

**Department of Computer Science and Engineering**

**Jadavpur University, Kolkata-32**

**Countersigned**

**Prof. (Dr.) Nandini Mukhopadhyay**

**Head of Department**

**Computer Science and Engineering**

**Jadavpur University, Kolkata-32**

**Prof. Ardhendu Ghoshal**

**Dean**

**Faculty of Engineering and Technology**

**Jadavpur University, Kolkata-32**

**Jadavpur University**  
**Faculty of Engineering and Technology**  
**Department of Computer Science and Engineering**

**CERTIFICATE OF APPROVAL**

This is to clarify that the project entitled “**Bengali Text Classification Using SVM**” has been completed by Surajit Maity. This work is applied under the supervision of **Prof.(Dr.) Kamal Sarkar** in partial fulfilment for the award of the degree of **Master of Computer Applications** of the **Department of Computer Science and Engineering, Jadavpur University**, during the academic year **2022-2023**. The project report has been approved because it satisfies the tutorial requirements in respect of project work prescribed for the said degree.

.....  
**Signature of Examiner 1**

**Date:**

.....  
**Signature of Examiner 2**

**Date:**

Date:

**FACULTY OF ENGINEERING AND TECHNOLOGY  
JADAVPUR UNIVERSITY**

**Declaration of Originality and Compliance of Academic Ethics**

**I hereby declare that this thesis entitled “ Bengali Text Classification Using SVM” contains original research work by the undersigned candidate, as part of his degree of Master of Computer Applications.**

**All information in this document has been obtained and presented in accordance with academic rules and ethical conduct.**

**I also declare that, as required by these rules and conduct, I have fully cited and referenced all materials and results that are not original to this work.**

**Name : Surajit Maity**

**Class Roll No. : 002010503002**

**Examination Roll No. : MCA2360012**

**University Registration No. : 154210 (2020-2021)**

## **ACKNOWLEDGMENT**

I am delighted to convey my gratitude to my thesis supervisor Prof. (Dr.) Kamal Sarkar , Department of Computer Science & Engineering, my supervisor, for his overwhelming support and encouragement towards accomplishment throughout the duration of the project without which this work would not have been possible. His positive outlook and confidence inspired me and gave me confidence. I am grateful for his support, encouragement, and his valuable suggestions to complete this work. I feel deeply honored that I got this opportunity to work under him.

I would like to express my sincere thanks to all my teachers for providing a sound knowledge base and cooperation.

I would like to thank all the faculty members of the Department of Computer Science & Engineering of Jadavpur University for their continuous support. Last, but not the least, I would like to thank my batch mates for staying by my side when I need them.

---

**Surajit Maity**

**Roll No:** 002010503002

**Reg. No. :** 154210 (2020-2021)

**Examination Roll No. :** MCA2360012

## CONTENTS

<b>Abstract</b>	<b>7</b>
<b>1 Introduction</b>	<b>8</b>
<b>2 Literature Survey</b>	<b>11</b>
2.1 Bag of Words	11
2.2 TF-IDF	12
2.3 Machine Learning Algorithms	12
2.3.1 Naive Bayes Classifier	13
2.3.2 Neural Networks	14
<b>3 Dataset</b>	<b>16</b>
<b>4 Proposed Methodology</b>	<b>17</b>
4.1 Bengali Text Classification without keywords	17
4.2 Bengali Text Classification with keywords	21
<b>5 Evaluation and Results</b>	<b>26</b>
<b>6 Conclusion and Future Works</b>	<b>30</b>
<b>7 References</b>	<b>31</b>

## **Abstract**

In current scenario, lot of online news is available for different topics on Internet from which textual data is increasing rapidly. Due to this, text classification becomes essential to organize them properly so that important news can be searched easily as well as to avoid data loss. One effective solution for this problem is to classify the news into different classes or to extract most important and useful information. With the rapid growth of Text sentiment analysis, the demand for automatic classification of electronic documents has increased by leaps and bound. The paradigm of text classification or text mining has been the subject of many research works in recent time. In this paper we propose techniques for Bengali text classification. In one method we take all the words from document and in another method we extract keywords from each document. We found the Support vector machine (SVM) is the most appropriate to work with our proposed model. The achieved results show significant increase in accuracy compared to earlier methods.

# Chapter 1

## Introduction

The classification process in text classification involves assigning predefined categories or classes to text documents based on their content. This process is similar to function mapping in mathematics, where the goal is to map inputs to specific outputs.

In text classification, the texts to be classified are represented as a set  $D$ , where  $D = \{doc_1, doc_2, doc_3, \dots, doc_m\}$ , with  $m$  documents in total. Additionally, there is a set of predefined classes or categories  $C$ , where  $C = \{cls_1, cls_2, cls_3, \dots, cls_n\}$ , with  $n$  classes.

The classification process can be interpreted as mapping each document in set  $D$  to a specific class from set  $C$ . The goal is to accurately assign the most appropriate class to each document based on its content and characteristics. So the classification process can be interpreted as [1]

$$F : D \rightarrow C$$

It's worth mentioning that while your focus is on single-label classification, where each text is assigned to a single class, there are also techniques for multi-label classification, where texts can belong to multiple classes simultaneously. Multi-label classification [2] is used when texts can have multiple topics or attributes associated with them. Which will not be discussed here. In this study, all text is mapped to a single class.

To perform text classification, supervised machine learning methods are commonly used. These methods involve training a model on a labeled dataset, where each document is already assigned to a specific class. The model learns patterns and features from the labeled data and uses them to make predictions on new, unseen documents.

During the training phase, the model analyzes the textual content of the documents and extracts relevant features. These features can include word frequencies, word distributions, or more advanced linguistic features. The model then builds a representation of the documents and their corresponding classes.

Once the model is trained, it can be used to classify new, unlabeled documents. The model applies the learned patterns and features to the unseen text and predicts the most appropriate class based on its knowledge of the training data.

Text classification has various applications, including sentiment analysis, spam filtering, topic categorization, and document classification in many domains such as news articles, customer reviews, and social media posts.

The increasing amount of text information in various fields of life due to the rapid development of Internet and information technology, particularly in the era of big data. Manual text classification, which was previously employed to categorize information, is now deemed inadequate due to its time-consuming, labor-intensive, and costly nature. As a result, automatic text classification has emerged as a solution to efficiently and accurately summarize text from vast amounts of information. This automated approach has gained significant attention in both academic and industrial sectors, becoming a topic of considerable discussion.

Natural Language Processing (NLP) has gained prominence in recent years, leading to the development of numerous tools and techniques for text comprehension. One such tool is the Natural Language Processing Toolkit (NLTK) [3], which is specifically designed to statistically understand and process text data.

The NLTK offers several applications for text analysis, including:

- a. Word tokenization: This involves breaking down a sentence into individual words, enabling further analysis at the word level.
- b. Sentence tokenization: It entails splitting a paragraph into separate sentences, facilitating sentence-level analysis.
- c. Word count: This application involves determining the frequency of word occurrences in a particular document or corpus, providing insights into word usage.
- d. Stemming: This technique aims to reduce words to their root form, allowing for more effective processing in natural language

applications.

- e. Stop-word removal: Stop-words refer to common words in the English language, such as prepositions and conjunctions. Many text processing applications eliminate these stop-words as they typically carry little relevant information.

By utilizing these NLTK applications, text processing tasks become more efficient and enable the extraction of meaningful insights from text data.

Now-a-days the amount of information available on the web is tremendous and increasing at an exponential rate. To effectively manage this vast volume of data, automatic text classification has emerged as a crucial application and research topic. Since the advent of digital documents, automatic text classification has been employed to categorize and organize the abundance of web-based information.

Automatic text classification relies on machine learning techniques to develop classifiers that learn the distinguishing characteristics of different categories from a set of pre-classified documents. It plays a significant role in various tasks, including information extraction, summarization, text retrieval, and question-answering.

The data used for classification is typically heterogeneous, sourced from various platforms such as web pages, newsgroups, bulletin boards, as well as broadcast or printed news articles, scientific papers, movie reviews, and advertisements. These sources introduce diversity in formats, preferred vocabularies, and writing styles, even within a single genre. Consequently, automatic text classification becomes crucial for efficiently organizing and extracting information from these diverse sources.

## Chapter 2

### Literature Survey

Natural Language Processing (NLP) applications are commonly used for text processing, classification, and clustering tasks. One of the simplest text analytics applications is spam filtering, which is often employed in email systems. Spam filtering involves classifying incoming emails as either spam or legitimate based on the content within the email. Machine learning techniques are commonly utilized in spam filters to automatically classify emails as spam or non-spam based on the text present in the email. This application demonstrates how NLP can be used to automatically categorize and filter text data to enhance email management and reduce unwanted spam messages.

Different text processing techniques, are as follows.

#### 2.1 Bag of Words

Bag of words [4] is a technique where the sentence is represented by the occurrence or absence of a word in the given sentence. The sentence "টানা ব্যর্থ কেএল রাহুল" Each column represents a unique word from the sentence, and each row represents the binary occurrence (1) or absence (0) of that word in the sentence.

রবিবার	টানা	কেএল	ম্যাচে	সাফল্য	ব্যর্থ	বড়	গাড়ি	নিয়ন্ত্রণ	রাহুল
0	1	1	0	0	1	0	0	0	1

To address the disadvantages of the bag of words technique mentioned:

**Loss of Meaning:** Bag of words representation treats each word independently and disregards the order and context of the words. This can lead to a loss of meaning. For example, the sentences "Dog bites Man" and "Man bites Dog" would have the same bag of words representation, even though the meaning is different.

**Ignoring Word Frequency:** Bag of words representation does not consider the frequency of words in the document. It only indicates the presence or absence of a word. As a result, important information about the frequency or importance of certain words is not retained.

These limitations can be addressed by using more advanced techniques, such as n-grams or word embeddings, which consider the context and sequence of words to capture meaning more effectively.

## **2.2 TF-IDF**

TF-IDF, which stands for Term Frequency-Inverse Document Frequency [5], is a commonly used technique in natural language processing and information retrieval to overcome some of the limitations of the bag-of-words approach.

The bag-of-words model represents a text document as a collection of individual words, disregarding their order and context. It considers each word independently, assigning equal importance to all words in the document. This approach has limitations because it does not capture the semantic or contextual information of the words.

TF-IDF addresses this limitation by taking into account both the frequency of a word in a particular document (term frequency) and its rarity across all documents (inverse document frequency).

The term frequency (TF) component of TF-IDF calculates the frequency of a word in a specific document. It gives a higher weight to words that appear more frequently in the document, assuming that such words are more important.

The inverse document frequency (IDF) component measures the rarity of a word across all documents in the dataset. It assigns a higher weight to words that occur less frequently in other documents, assuming that rare words carry more significant meaning.

By combining the TF and IDF components, TF-IDF assigns higher weights to words that are both frequent in a document and rare across other documents. This approach helps to identify important words that are distinctive to a particular document and can provide valuable insights.

## **2.3 Machine Learning Algorithms**

Machine Learning algorithms are used to perform classification, regression, and forecasting tasks effectively. Machine learning algorithms are employed to accomplish tasks such as classification, regression, and forecasting, highlighting their effectiveness in these areas.

The algorithms used for classification tasks are Logistic Regression for binary classification, Decision Trees, Random Forest Classifier, Support Vector Machines, Naive Bayes Classifier, and Boosting Algorithms. There

are several algorithms that are commonly used for classification tasks, including Logistic Regression (specifically for binary classification), Decision Trees, Random Forest Classifier, Support Vector Machines, Naive Bayes Classifier, and Boosting Algorithms. These algorithms are known for their ability to classify data into different categories or classes.

Due to high sparsity and dimensions of the data, only Logistic Regression, Random Forest, and Naive Bayes classifiers are used in modeling the data. Because of the high sparsity (sparse data) and dimensions (large number of features) of the data, only three classifiers are utilized for modeling the data: Logistic Regression, Random Forest, and Naive Bayes. These algorithms are chosen likely because they can handle high-dimensional and sparse data effectively.

### 2.3.1 Naive Bayes Classifier

McCallum and Nigam [6] focuses on the Naive Bayes classifier for text classification. Naive Bayes classifiers are a collection of classification algorithms based on Bayes' Theorem. It is not a single algorithm but a family of algorithms where all of them share a common principle, i.e. every pair of features being classified is independent of each other.

Bayes' Theorem finds the probability of an event occurring given the probability of another event that has already occurred. Bayes' theorem is stated mathematically as the following equation:

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

Where A and B are events and  $P(B) \neq 0$

They apply Bayes' theorem in following way:

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

Where, y is class variable and X is a dependent feature vector (of size n) where:  $X = (x_1, x_2, x_3, \dots, x_n)$

Now, it's time to put a naive assumption to the Bayes' theorem, which is, independence among the features. If any two events A and B are independent, then,  $p(A,B) = p(A)p(B)$

Hence, we reach to the result:

$$P(y | x_1, x_2, \dots, x_n) = \frac{p(x_1|y)p(x_2|y)\dots p(x_n|y)p(y)}{p(x_1)p(x_2)\dots p(x_n)}$$

### 2.3.2 Neural Networks

The field of machine learning and data science has been greatly influenced by neural networks. These models, which are inspired by the human brain's structure and functioning, have become increasingly powerful with the availability of high computing engines. Neural networks can now consist of millions of parameters, enabling them to learn from vast amounts of data and make accurate predictions.

Prasanna and Rao [7] constructed a TFIDF matrix which given to the artificial neural network from which the neural networks learns. Training a neural network indeed involves two major steps, forward propagation and backward propagation.

**Forward Propagation:** In this step, the neural network takes a set of input data and computes the corresponding output. Each neuron in the network receives input signals, multiplies them by corresponding weights, sums them up, and applies an activation function to produce an output.

The inputs are multiplied by the weights, and the weighted sum is passed through an activation function. The activation function introduces non-linearity to the network, allowing it to learn complex relationships between inputs and outputs. Common activation functions include sigmoid, hyperbolic tangent (tanh), ReLU (Rectified Linear Unit), and others. The activation function maps the input to a specific range or threshold, which becomes the output of that neuron.

**Backward Propagation:** After the forward propagation step, the network compares the computed output with the desired output and calculates the error. Backward propagation, also known as backpropagation, aims to adjust the weights of the network to minimize this error.

The error is propagated backward through the network, starting from the output layer and moving towards the input layer. The weights of the neurons are updated based on the calculated error and the gradients of the activation function. This process is performed iteratively using optimization algorithms like gradient descent, which adjusts the weights in the direction that reduces the error.

The gradients are computed using the chain rule of calculus, which calculates the impact of each weight on the overall error. By adjusting the weights based on the gradients, the network gradually learns to make better predictions and minimize the error.

The process of forward and backward propagations is repeated iteratively, with the network adjusting the weights in each iteration to improve its

predictions. The goal is to find the optimal set of weights that minimize the error and make accurate predictions on new, unseen data.

During the feedforward propagation stage in neural network training, each input ( $x_i$ ) indeed receives an input signal.

**Input Signals:** Each input in the neural network receives a specific input signal ( $x_i$ ) associated with the input pattern or data point. These input signals represent the features or attributes of the input data.

**Hidden Units:** The input signal from each input is then sent to each hidden unit ( $z_1, z_2, z_3, \dots$ ). A hidden unit, also known as a neuron or node, is responsible for processing the input signals and generating an activation signal.

**Activation Calculation:** The hidden layer calculates its activation signal based on the received input signals. This calculation typically involves summing up the weighted input signals and applying an activation function to the result. The weights represent the connections between the inputs and the hidden units.

**Output Unit:** Once the activation signal is calculated for each hidden unit, it is sent to the output unit. The output unit uses the received activation signal(s) to compute the output or response for the given input pattern.

**Output Comparison:** After generating the output value, each output unit compares its activation value ( $y_k$ ) with the target value ( $t_k$ ). The target value represents the desired or expected output for the given input pattern. The output unit checks the activation value against the target value to determine the level of error or mismatch.

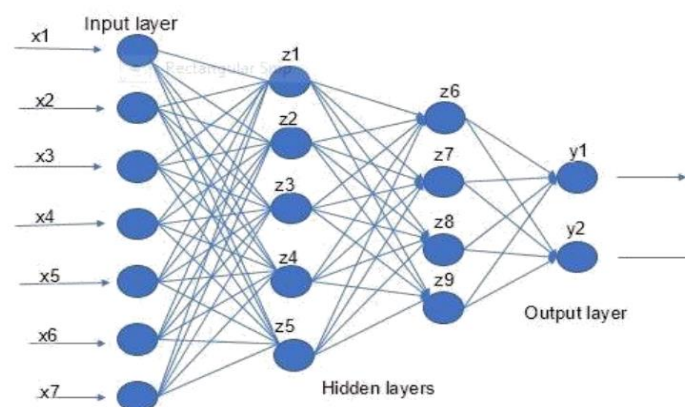


Figure 1: Neural Network Architecture

## Chapter 3

### Dataset

Our dataset is based on Bengali news. In recent years, news delivery websites have become a popular medium for individuals to access current events information. In our research the Bengali news have been gathered from the online news website 'Ei Samay'. Dataset contains almost 1756 news. News in these dataset belong to 38 different topics (labels). Each news record consists of several attributes from which we are using only 'Category' and 'Article'.

## Chapter 4

# Proposed Methodology

Our proposed Bengali text classification using SVM includes several methods: (1) Bengali Text Classification without keywords, (2) Bengali Text classification with keywords, (3) Evaluation and Results.

### 4.1 Bengali Text Classification without keywords

This method includes several steps: (1) Breaking into sentences, (2) Stop-Word Removal, (3) Document Representation using word embedding, (4) SVM Classifier, (5) Experiments, (6) Final Prediction. A block diagram of the proposed method is shown in Figure 2.

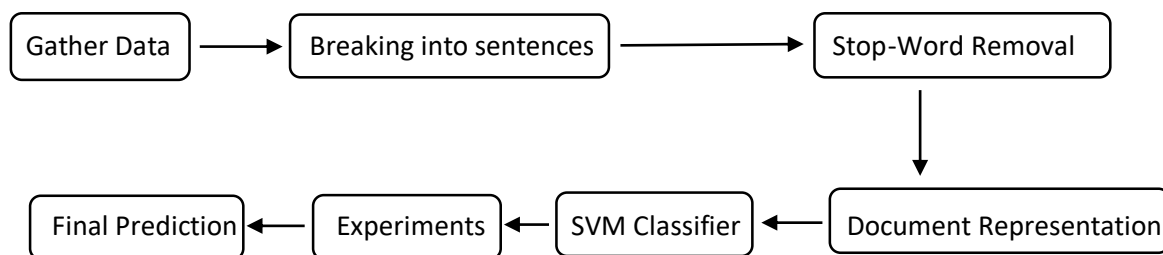


Figure 2: Bengali Text Classification without keywords

#### 4.1.1 Breaking into sentences

Sentence tokenization is performed using the sentence tokenize function of the NLTK library. This process involves breaking the input document set into individual sentences. Sentence tokenization helps in analyzing and processing text at a more granular level, focusing on the meaning and structure of individual sentences.

#### 4.1.2 Stop-Word Removal

Stop words are common words that often appear in a language but typically do not carry significant meaning or contribute to the understanding of the text. Bengali stop word file is used to remove stop words. Removing stop words helps in reducing noise and irrelevant

information, allowing the focus to be on more meaningful words and concepts.

### 4.1.3 Document Representation

After removing stop words, we get meaningful words for each document. For each word, we search for the corresponding word vector into the list of word vectors obtained from a Bengali pretrained word vector model called Fast Text. Then we calculate the average of the vectors for the words. Thus each document is represented as a vector. Figure 3 shows an example of a vector representing a document taken from our dataset.

```
[8.99625710e-02  1.12286594e-02  7.87209930e-03 -1.11322355e-02
-2.92406967e-05 -2.32450226e-02  1.28842127e-02  1.41462041e-02
-2.26229804e-02  2.44888794e-02 -4.75087343e-03 -6.12748647e-03
2.08181292e-02  9.99064092e-03  6.14736881e-03  2.23216501e-03
-7.41052488e-02  5.20292344e-03 -1.40514625e-02  2.13017528e-02
1.48526229e-02  8.22522294e-03 -1.2555592e-02 -4.24852899e-03
6.52046734e-03  1.41368508e-02  1.01988306e-02 -3.65309864e-02
-1.08608166e-02 -1.63836274e-02 -7.19999941e-03 -1.17701730e-02
-9.05848457e-04  1.02280616e-03  3.33391759e-03  1.71613928e-02
-3.51339992e-02 -2.38011635e-04 -1.2315759e-02 -1.69590656e-02
-7.10525969e-02  8.10926721e-03 -1.4128659e-03  2.36608074e-02
7.55321560e-03 -7.20117055e-03  1.71374213e-02 -4.52806987e-03
-6.37602527e-02  1.02543849e-02 -5.96199045e-03 -3.75087815e-03
-3.61403887e-04 -2.97602406e-03 -3.12099245e-02 -3.2323405e-03
-4.39583024e-03  1.87105257e-02 -3.11111071e-04  8.92952858e-04
-2.90409452e-02  1.56793665e-03 -2.69649229e-02 -1.72628724e-02
5.72801063e-03 -1.26625737e-02  1.40994135e-02  1.28304148e-02
-1.04157915e-02 -9.44737252e-03  2.39198841e-02  1.74029209e-02
-4.58011720e-02 -4.87269014e-02 -2.24888846e-02  1.41415223e-02
-3.87426978e-02 -2.47543849e-02  4.17584665e-02  5.62631665e-03
2.56140420e-02  2.92486057e-03 -1.14678370e-03  5.70643460e-03
-6.7239275e-03  5.60292415e-03 -2.33567205e-02  8.87602381e-03
2.02438496e-02  4.21947166e-02 -8.77836347e-03  2.97239628e-02
-5.32105193e-02 -9.82046500e-03 -7.44444595e-03  1.28654967e-04
2.35994746e-02  1.67245622e-02  1.69333406e-02  1.34736777e-03
-1.07785279e-02 -4.88005705e-03 -9.42685668e-04 -3.6594161e-03
-3.33918282e-02 -7.15204654e-03  3.67134507e-03  2.79081967e-02
9.61871352e-03 -3.56514752e-02 -5.59526421e-02 -2.28561442e-02
-6.35380344e-02  2.24397685e-02 -1.28312245e-01 -2.77493611e-02
-1.54514601e-02 -9.94152040e-04  1.31695922e-03  9.41982677e-03
2.11070161e-02  3.22573236e-03 -1.42503906e-02 -1.35122815e-03
-4.70175454e-02 -2.07725149e-02  9.83567163e-03  6.28473617e-03
6.27777958e-03 -4.62621509e-03  6.00877265e-03  4.21920244e-03
5.87549426e-02 -2.47826176e-03  1.92108304e-02 -5.55262947e-03
1.11391899e-02 -1.78578999e-02 -2.44385995e-03  1.43795320e-02
-4.65092702e-02 -1.42947270e-02  1.39982449e-02  4.16082144e-03
-1.61614046e-02  1.24052632e-02 -1.28690079e-02 -3.91754125e-03
4.12631501e-03 -1.64567232e-02  2.72052661e-02 -7.66023425e-03
-7.53801316e-02  1.97719336e-02 -1.22169638e-02 -1.22859655e-02
2.44594432e-02 -7.74362370e-03  1.05690039e-02 -1.5789574e-04
7.78888622e-02 -3.36081920e-03 -4.52046702e-02 -1.09627436e-02
5.12280385e-04  2.53485404e-02  1.71426926e-02 -5.54111004e-02
-1.64736854e-03  2.12514563e-03 -6.60818128e-04  3.24736885e-03
1.39239791e-02 -2.02608127e-02 -4.23134416e-02  5.48029249e-03
1.15204866e-02  2.22858465e-03  7.64861677e-03 -2.12269997e-02
-5.47134550e-03  3.26599928e-04  2.51111109e-03  1.21520506e-03
1.46058453e-02 -2.13906922e-02 -1.20198876e-02  1.46467807e-02
-4.07257341e-02  4.43806946e-02  1.01906406e-02 -6.64970418e-03
2.23497041e-02  1.16812848e-02 -7.74099331e-03  2.49413147e-03
-1.04726874e-02  2.94736889e-03  1.49941514e-03 -1.89742707e-02
3.49532068e-03  1.11140267e-02 -6.46198925e-04  6.98888907e-03
2.23567220e-02  2.37602291e-03  1.55771980e-02 -2.75555456e-03
1.06274849e-02 -3.01169627e-03  2.44894810e-02 -7.53216352e-03
9.61345620e-02 -2.69970912e-02 -4.94385976e-03  3.35436159e-03
1.20456125e-02  6.44099938e-03  1.18210567e-02  2.43684417e-03
8.47894792e-03  9.92163923e-03  1.78122837e-02  6.66356534e-02
8.60819884e-04 -7.02268779e-02  4.57602320e-03 -6.64251488e-02
-1.95684209e-02 -6.78128516e-03 -7.98712416e-03  2.14678366e-02
1.46321626e-02  1.68245705e-03  2.49838635e-02  7.50380766e-03
-2.42105150e-02  1.26964925e-02  1.15782624e-02  1.40128611e-02
1.41520507e-03  3.53801157e-03 -4.47602104e-03 -4.87193156e-03
-6.41344581e-02 -3.63304131e-02  3.76612736e-02  1.57485367e-03
3.64976823e-02 -6.04853779e-03 -3.00467553e-02 -2.74064410e-02
-6.20666929e-02  5.2239297e-03  1.25549679e-02 -7.73193017e-02
-3.15918103e-02  9.95087898e-03 -6.26987935e-02 -2.97824539e-02
2.94970698e-03  9.07602371e-04  1.90526212e-03  2.85140425e-02
-8.67309794e-03  1.22269010e-02  1.09970747e-02 -2.36140255e-03
-1.36658029e-02 -1.44856694e-02  1.94105245e-02  1.55847881e-03
-5.1380054e-02  1.26643266e-02  7.00525939e-03  1.89953148e-02
2.13809415e-02 -6.7000311e-03  1.77380938e-02  1.88350901e-02
-2.14157905e-02 -3.76783544e-03 -1.79198889e-02 -1.65023431e-02
-1.15204705e-02 -1.83726831e-02  1.38695920e-02  1.61421057e-02
-2.77649108e-02 -5.11988159e-02 -2.59304121e-02  7.25029642e-03
7.59298401e-03  1.19292279e-02 -2.94035061e-03 -7.16900500e-03]
```

Figure 3: document representation

### 4.1.4 Classifier

We import the SVM Classification Class from scikit-learn. Support vector machines (SVM) is powerful yet flexible supervised machine learning method used for classification. SVM seeks a decision surface to separate the training data points into two classes and makes decisions based on the

support vectors that are selected as the only effective elements in the training set. Support Vector Machine (SVM) have several advantages that make them stand out compared to other algorithms in certain scenarios. Here are some reasons why SVMs are often considered favourable:

Effective in high-dimensional spaces: SVM perform well in cases where the number of features is greater than the number of samples. This makes them suitable for tasks involving text classification or image recognition, where the dimensionality of the data is typically high.

Robust against overfitting: SVM aim to find the maximum margin, which promotes generalization and helps avoid overfitting. This means that SVM can handle noisy or overlapping data and still make accurate predictions on unseen data.

Versatility with kernel functions: SVM can handle non-linearly separable data by utilizing different kernel functions. The ability to transform the data into a higher-dimensional feature space allows SVM to find complex decision boundaries, making them capable of capturing intricate relationships in the data. The kernel used here is the “rbf” kernel which stands for Radial Basis Function.

Global optimization: The training of SVM involves solving a convex optimization problem, which means they can find the global optimum rather than getting stuck in local optima. This property contributes to their robustness and ensures that SVM consistently converge to a good solution.

Margin-based decision making: SVM focus on maximizing the margin between different classes. This approach leads to better generalization performance, as it encourages the model to learn decision boundaries that are less influenced by noisy or irrelevant data points.

Support for sparse data: SVM can handle datasets with a large number of features, even if most of them are zero (sparse data). This makes SVMs suitable for tasks involving text data or other high-dimensional sparse representations.

#### **4.1.5 Experiments**

News in these dataset belong to 37 different topics (labels). Each news

record consists of several attributes from which we are using only 'Category' and 'Article'.

The details of the dataset are available in the table 1.

Class	Number of Articles
Agriculture	50
Business	50
Health	50
Labor_and_Employment	50
Law	50
Miscellaneous	50
Music	50
Politics	50
Public_lands_and_water_management	24
Religion	50
Science	50
Social_welfare	11
Space	46
Sports_other_than_football_and_cricket	50
Caste	42
Technology	50
Transportation	44
Travel	50
Weather	50
World_and_international	50
Cinema	50
Computer	50
Cricket	50
Crime	50
Defence	50
Economy	50
Education	50
Banking	50
Election	50
Electronics	50
Energy	35
Entertainment	50
Environment	50
Family issues	50
Finance	50
Football	50
Government_Operations	50

Table 1: dataset description

We split the data into the training and the test set. 20% of the data kept as the test set and the remaining 80% used for training SVM model.

The training data is scaled, and its scaling parameters are determined by applying a `fit_transform()` to the training data. We do not want to be biased with our model. We want our test data to be a completely new and a surprise set for our model. The transform method helps us in this case.

Feature scaling is an additional step that can increase the speed of the program as we scale down the values of X to a smaller range. In this, we scale down both the train and the test to a small range of -2 to +2.

#### 4.1.5.1 Parameter tuning

Kernel parameters selects the type of hyperplane used to separate the data. Gamma: gamma is a parameter for non-linear hyperplanes. The higher the gamma value it tries to exactly fit the training data set. C: C is the penalty parameter of the error term. It controls the trade off between smooth decision boundary and classifying the training points correctly. How accuracy varies for different combination of gamma and C values shown in table 2.

gamma	C	Accuracy(in percentage)
0.001	1	52
0.002	1	67
0.003	1	72
0.001	10	85
0.002	10	89
0.003	10	89
0.001	100	92
0.002	100	92
<b>0.003</b>	<b>100</b>	<b>93</b>

Table 2: Parameter Tuning

## 4.2 Bengali Text Classification with keywords

This method includes several steps: (1) Keyword selection, (2) document representation and (3) SVM Classifier, (4) Experiments, (5) Final Prediction.

A block diagram of the proposed method is shown in Figure 4.

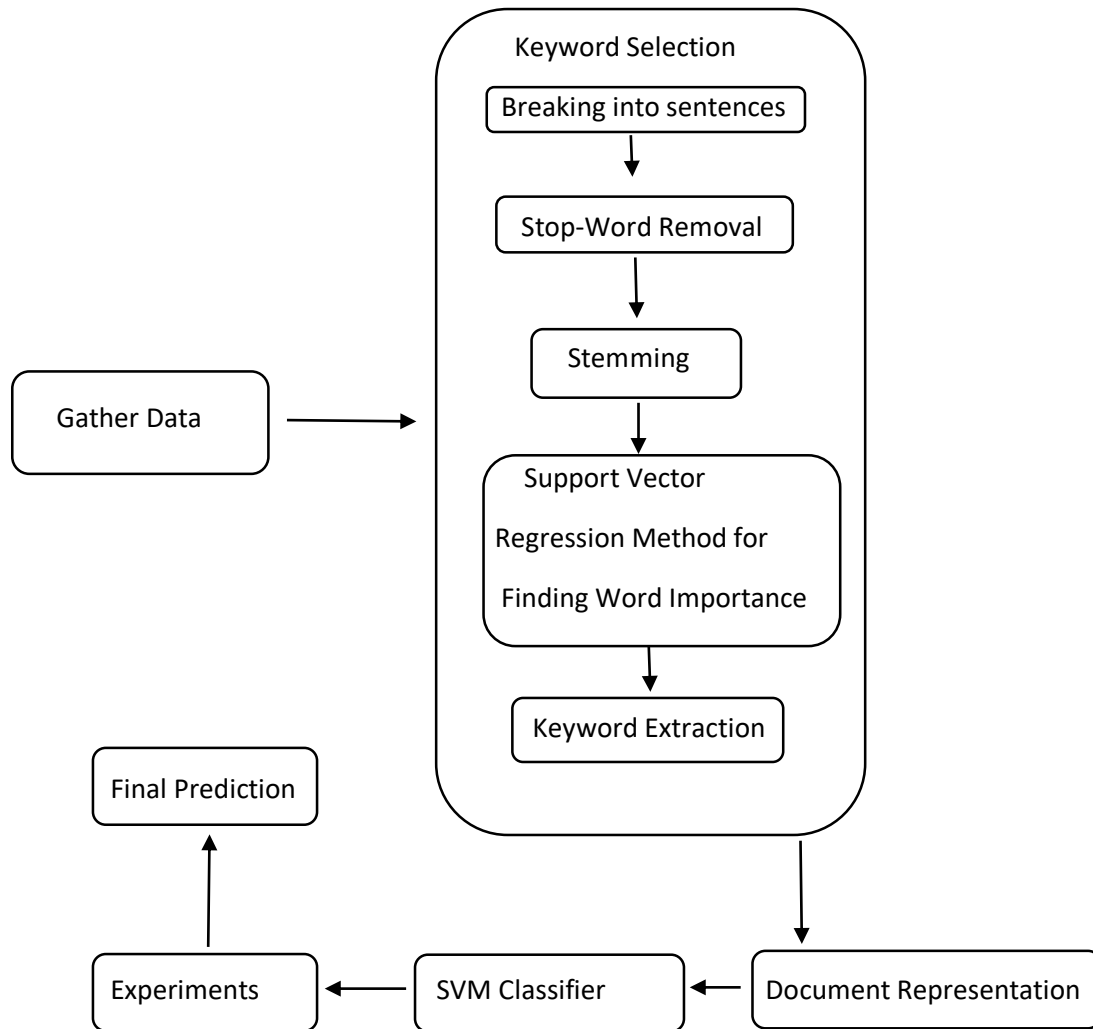


Figure 4: Bengali Text Classification with keywords

### 4.2.1 Keyword Selection

Chatterjee and Sarkar [8] use a method for predicting word importance (word weight) using the support vector regression model.

#### 4.2.1.1 Pre-processing

This step includes the process:

**Sentence Tokenization:** Sentence tokenization is performed using the sentence tokenize function of the NLTK library. This process involves breaking the input document set into individual sentences. Sentence

tokenization helps in analyzing and processing text at a more granular level, focusing on the meaning and structure of individual sentences.

**Stop Words Removal:** Stop words are common words that often appear in a language but typically do not carry significant meaning or contribute to the understanding of the text. Bengali stop word file is used to remove stop words. Removing stop words helps in reducing noise and irrelevant information, allowing the focus to be on more meaningful words and concepts.

**Stemming:** Stemming is a process that reduces words to their base or root form. It helps in standardizing different word variations to their common base, thereby reducing the dimensionality of the text data. In this step, Bengali steam word file is used. Stemming helps in consolidating related words and capturing their common essence.

#### **4.2.1.2 Word importance prediction using support vector regression model**

In the mentioned research, a supervised Support Vector Regression (SVR) model is utilized for predicting the importance of a term. While the frequency of a term is one factor in determining its importance, the research acknowledges that there are other features that contribute to the term's significance. We will discuss various features used for developing the SVR model.

**Word position in the document:** Since the word importance is dependent on the word's position in the document, we consider this as a feature. The sentences of the document are numbered from 1 to  $n$ . This feature is computed as follows:

$$\text{POS\_IN\_DOC} = \frac{\text{Position of the sentence in which the word occurs}}{\text{Total number of sentences in the document}} \quad (1)$$

**Global Term Frequency (GTF):** Since we work with multi-document, frequency of a word  $w$  in the entire input collection is considered as a separate feature. The average TF over the number of documents in the input cluster is taken as a feature.

$$\text{GTF}(w) = \frac{1}{|C|} \sum_{w \in C}^{|C|} \text{TF}(w) \quad (2)$$

Where  $\text{TF}(w)$  = total count of occurrences of  $w$  in the input collection.

$|C|$  = the size of the input collection

TF-IDF Local: Local term frequency (LTF) of a word is multiplied by its IDF value to define a new feature. The IDF value of the word  $w$  is computed over a corpus of  $N$  documents using equation (3):

$$ID(w) = \log \left[ 0.5 + \frac{N}{DF(w)} \right] \quad (3)$$

Where  $N$  = corpus size in terms of documents.

TF-IDF Global: This feature is defined by the product of GTF (defined in Equation 2) and IDF as follows:

$$TF - IDF\ Globa(w) = GTF(w) * IDF(w) \quad (4)$$

Proper Noun: The proper nouns like organization name, person name, etc. play an important role in terms selection. For this reason, we consider a feature that checks whether a word is a part of a proper noun or not. This is considered a binary feature. If  $w$  is the part of a proper noun then the feature value = 1, otherwise its value = 0.

Word length ( $w$ ): Word length is considered as a feature. It is observed that the words that are longer are highly informative. This is also considered a binary feature. If  $\text{length}(w) \geq 5$  then the value of this feature = 1, otherwise its value = 0

SVR model is used for predicting the degree of importance for each word. For training, the input to the regression model is represented in the form  $\langle x, y \rangle$ , where  $x$  is a vector of the values of features described above in this section and  $y$  is the target value.

While selecting keywords, the redundant words are removed from the keyword because the redundancy affects keyword quality. A word is selected for a keyword if its similarity with the previously selected word is less than a predefined threshold value.

#### 4.2.2 Document Representation

After keyword extraction, we get keywords for each document. For each keyword, we search for the corresponding word vector into the list of word vectors obtained from a Bengali pretrained word vector model called Fast Text. Then we calculate the average of the vectors for the

keywords. Thus each document is represented as a vector. Figure 5 shows an example of a vector representing a document taken from our dataset.

```
[ 1.54692316e-02 -5.40000014e-03 1.23692313e-02 -2.47846153e-02
-3.30769317e-03 2.41384618e-02 1.70384608e-02 -1.12076905e-02
2.33846158e-03 -3.67923044e-02 3.13615389e-02 1.97000001e-02
-3.38076912e-02 3.10153868e-02 1.12153841e-02 -1.51538476e-03
2.88538467e-02 8.23076991e-03 1.01692319e-02 -6.57692272e-03
-9.16923142e-03 7.91538414e-03 -5.73207734e-03 3.40538475e-02
1.43307690e-02 -3.83076840e-03 -1.13692321e-02 1.03846154e-02
2.70769116e-03 1.65769234e-02 -7.30769243e-03 -6.77153841e-02
-1.82207744e-03 -2.37923097e-02 -2.16923077e-02 -1.41769238e-02
-4.48461529e-03 7.90769234e-03 2.64307670e-02 1.22461533e-02
-3.03538442e-02 -3.84615315e-03 -1.17461532e-02 -2.93153841e-02
-1.05923070e-02 1.54384617e-02 3.80769162e-03 5.39230742e-02
2.00769212e-03 -4.83846139e-03 3.03384587e-02 -2.06307694e-02
2.88461545e-03 5.60769206e-03 -2.24153828e-02 8.55384581e-03
-2.43307687e-02 -1.01230777e-02 2.59615369e-02 -9.69230942e-03
-1.16538461e-02 2.30461489e-02 -1.74615451e-03 -6.46153872e-04
-2.45846175e-02 5.36153931e-03 -4.60153855e-02 -1.16230790e-02
-2.89999857e-03 -1.49076935e-02 3.08153816e-02 4.69230674e-03
-1.28692323e-02 4.16922895e-03 6.34615403e-03 4.17999998e-02
2.40138437e-02 -7.73615469e-02 -1.91307688e-02 1.19307702e-02
-1.39384614e-02 -5.98153882e-02 6.19615391e-02 1.43461535e-02
-5.18461596e-03 -2.43153851e-02 3.36922961e-03 -6.28461502e-03
-3.68461595e-03 -1.26384608e-02 -1.40615366e-02 1.09307701e-02
2.18330784e-02 4.58923057e-02 -2.35461540e-02 5.43923043e-02
-5.06153796e-02 -2.55615395e-02 -8.7307724e-02 1.46692330e-02
2.08692327e-02 2.90769245e-02 2.68307718e-02 1.66923064e-03
-1.02153849e-02 2.27692258e-03 2.71153841e-02 -1.99692287e-02
1.45307686e-02 -1.57999974e-02 2.01461539e-02 4.13461588e-02
2.58076945e-02 5.52207389e-02 -4.27307710e-02 2.61769220e-02
9.00000858e-04 3.93076949e-02 -1.71384449e-01 -4.41000015e-02
-1.15999989e-02 2.02076919e-02 -9.20769293e-03 7.46153854e-03
3.48923057e-02 7.39230814e-03 -4.28999998e-02 -2.67692301e-02
-9.20769320e-03 -3.05528476e-02 1.91230755e-02 1.89307686e-02
9.93846189e-03 -1.43846171e-02 -1.59923065e-02 -1.07384589e-02
6.63076853e-03 9.16153751e-03 -4.14615450e-03 1.89230789e-02
1.90769252e-03 -4.01307680e-02 -1.70769226e-02 1.80384610e-02
-6.53076917e-02 -6.50000060e-03 2.96153836e-02 1.11461533e-02
-1.19307684e-02 8.28461518e-03 1.33999996e-02 -4.50769200e-03
2.52461527e-02 7.66153860e-03 1.77461546e-02 -1.47923062e-02
1.99615378e-02 1.91000011e-02 -8.34615435e-03 -2.05846149e-02
4.36692275e-02 -1.86692309e-02 -1.25307692e-02 -5.76922949e-03
3.25615369e-02 -3.16307694e-02 3.84615408e-03 -8.38460808e-04
1.44153892e-02 4.48769206e-02 1.39307703e-02 -1.03307616e-01
-1.78153850e-02 4.38461575e-04 1.33076881e-03 1.21615399e-02
-8.40769056e-03 -3.15307714e-02 -3.18000018e-02 7.78461527e-03
1.94230769e-02 1.24307685e-02 3.22307670e-03 -5.08769192e-02
-1.19461529e-02 -1.00307669e-02 8.46154115e-04 1.24461539e-02
2.66938467e-02 -2.78000010e-02 -1.29769240e-02 3.82846147e-02
-4.84769270e-02 1.02515377e-01 2.47461535e-02 -6.22307695e-03
3.37846167e-02 3.28384601e-02 -3.92307800e-03 7.92307779e-03
-1.70846172e-02 1.58999991e-02 -1.45461550e-02 -2.13076938e-02
-2.34615314e-03 1.78076942e-02 5.78461541e-03 8.75384547e-03
-1.30615374e-02 2.56923051e-03 2.80769207e-02 5.65384608e-03
2.12769248e-02 -5.70769375e-03 3.55769247e-02 -1.89846177e-02
1.28538469e-02 -4.03923057e-02 2.46923137e-03 -9.34615452e-03
1.35923084e-02 -9.71538480e-03 1.11076944e-02 9.49230697e-03
-3.77692282e-02 2.90076919e-02 7.06923055e-03 8.19692214e-02
3.61418874e-04 -1.13507691e-01 1.23384604e-02 -6.69384589e-02
-2.19461527e-02 -1.18769221e-02 -1.41384620e-02 1.18923094e-02
1.22923087e-02 8.78461637e-03 4.80692312e-02 7.56154023e-03
1.90000003e-03 6.09230762e-03 8.71538464e-03 2.12153867e-02
2.16153706e-03 -1.66153710e-03 -9.65384673e-03 1.02538457e-02
-1.10461544e-02 -2.84769200e-02 7.86461532e-02 5.70000010e-03
4.69384603e-02 -2.94615375e-03 -1.24384612e-02 -4.39230800e-02
-7.08923116e-02 2.83076894e-03 1.65230762e-02 -1.07546158e-01
-4.15769219e-02 1.82307709e-03 -1.07015379e-01 -7.68461525e-02
9.37692262e-03 8.38461681e-04 -9.26153921e-03 3.90923098e-02
-3.06538455e-02 2.65923105e-02 8.67692381e-03 7.24615389e-03
-1.36923036e-02 -2.01769210e-02 -1.15383918e-04 1.56538468e-02
-1.03753850e-01 1.60153843e-02 2.38459848e-04 1.52615374e-02
1.50000001e-03 -5.57692349e-03 2.44307686e-02 4.36230712e-02
-2.51130728e-02 -1.78615393e-02 -4.48307730e-02 -5.91692328e-02
-1.57846175e-02 -3.17692310e-02 3.83461565e-02 -4.93846182e-03
-5.42692356e-02 -8.65384564e-03 -3.35307680e-02 1.45846158e-02
1.53846154e-02 1.90153848e-02 -1.10769272e-03 -2.11307704e-02]
```

Figure: 5

### 4.2.3 Classifier

Previously we have discussed in detail about SVM classifier in section 4.1.4. For this experiment we use the same classifier.

### 4.2.5 Experiments

Detailed dataset description is given in the section 4.1.5. We split the data into the training and the test set. 20% of the data kept as the test set and the remaining 80% used for training SVM model.

The training data is scaled, and its scaling parameters are determined by applying a `fit_transform()` to the training data. We do not want to be biased with our model. We want our test data to be a

completely new and a surprise set for our model. The transform method helps us in this case.

Feature scaling is an additional step that can increase the speed of the program as we scale down the values of X to a smaller range. In this, we scale down both the train and the test to a small range of -2 to +2.

#### 4.2.5.1 Parameter tuning

Kernel parameters selects the type of hyperplane used to separate the data. Gamma: gamma is a parameter for non-linear hyperplanes. The higher the gamma value it tries to exactly fit the training data set.

C: C is the penalty parameter of the error term. It controls the trade off between smooth decision boundary and classifying the training points correctly.

How accuracy varies for different combination of gamma and C values shown in table 3.

gamma	C	Accuracy(in percentage)
0.001	1	54
0.002	1	64
0.003	1	73
0.001	10	80
0.002	10	82
0.003	10	82
0.001	100	84
0.002	100	84
<b>0.003</b>	<b>100</b>	<b>84</b>

Table 3: Parameter Tuning

## Chapter 5

# Evaluation and Results

This section shows the experimental results based on the proposed algorithms. The results for the category classification will be given. The evaluation indicators include: Precision, Recall, f-1 score and Accuracy.

Precision measures the proportion of correctly classified positive instances (documents assigned to a particular class) out of all instances classified as positive. It indicates the ability of the SVM classifier to accurately identify documents belonging to a specific category.

$$\text{Precision} = \frac{TP}{(TP + FP)}$$

Recall also known as sensitivity or true positive rate, measures the proportion of correctly classified positive instances out of all actual positive instances. It reflects the ability of the SVM classifier to capture all relevant documents of a particular category

$$\text{Recall} = \frac{TP}{(TP + FN)}$$

Accuracy score that indicates the percentage of correctly classified documents when classified by SVM. Higher accuracy values indicate a more successful classification model.

$$\text{Accuracy} = \frac{(TP + TN)}{(TP + TN + FP + FN)}$$

F1-score is the harmonic mean of precision and recall. It provides a balanced measure of the classifier's performance, considering both precision and recall. Higher F1-scores indicate better overall performance.

$$\text{F1 - score} = 2 \frac{(\text{Precision} * \text{Recall})}{(\text{Precision} + \text{Recall})}$$

TP: True Positive means that news which is classified to its correct class.

FN: False Negative means that news is classified to a wrong class.

TN: True Negative means that news which does not belong to that class and is misclassified.

FP: false positive is an outcome where the model incorrectly predicts the *positive* class.

Results are obtained from Bengali Text Classification without keywords are shown in table 4 and Bengali Text classification with keywords are shown in table 5.

class	precision	recall	f1-score
Agriculture	1.00	1.00	1.00
Banking	1.00	0.90	0.95
Business	0.91	1.00	0.95
Caste	0.89	1.00	0.94
Cinema	1.00	0.92	0.96
Computer	1.00	1.00	1.00
Crime	1.00	1.00	1.00
Defence	1.00	1.00	1.00
Economy	0.88	1.00	0.93
Education	1.00	0.92	0.96
Electronics	0.75	1.00	0.86
Energy	1.00	0.88	0.93
Entertainment	1.00	1.00	1.00
Environment	1.00	1.00	1.00
Family issues	1.00	1.00	1.00
Finance	1.00	1.00	1.00
Football	1.00	1.00	1.00
Government Operatios	0.90	1.00	0.95
Health	1.00	0.90	0.95
Labor_and_Employment	1.00	1.00	1.00
Law	1.00	1.00	1.00
Miscellaneous	1.00	1.00	1.00
Music	0.92	1.00	0.96
Politics	0.89	0.89	0.89
Public_lands and_water_managemen t	1.00	0.40	0.57
Science	0.89	1.00	0.94
Space	0.67	1.00	0.80
Sports_other_than_f ootball_and_cricket	1.00	1.00	1.00
Technology	1.00	1.00	1.00
Transportation	1.00	1.00	1.00
Weather	1.00	0.73	0.84
World_and_Internati onal	0.78	1.00	0.88
Cricket	1.00	1.00	1.00
Election	1.00	0.75	0.86
Religion	0.69	1.00	0.81
Travel	0.91	1.00	0.95

Accuracy: 0.94

Table 4: results of Bengali Text Classification without keywords

Cross Validation Scores:

[0.96875 0.94886364 0.93465909 0.94034091 0.93181818]

Average CV Score: 0.9448863636363637

Number of CV Scores used in Average: 5

class	precision	recall	f1-score
Agriculture	0.92	0.86	0.89
Banking	0.76	0.62	0.68
Business	0.91	1.00	0.95
Caste	1.00	0.36	0.53
Cinema	1.00	1.00	1.00
Computer	1.00	1.00	1.00
Crime	0.92	1.00	0.96
Defence	0.93	1.00	0.96
Economy	0.87	1.00	0.93
Education	0.82	1.00	0.90
Electronics	0.56	1.00	0.71
Energy	0.75	0.60	0.67
Entertainment	0.56	1.00	0.71
Environment	0.85	1.00	0.92
Family issues	0.88	1.00	0.94
Finance	0.89	1.00	0.94
Football	0.86	0.92	0.89
Government Operatios	0.64	1.00	0.78
Health	0.80	0.89	0.84
Labor_and_Employment	0.88	0.88	0.88
Law	0.91	1.00	0.95
Miscellaneous	0.89	1.00	0.94
Music	1.00	1.00	1.00
Politics	0.93	1.00	0.96
Public_lands and_water_managemen t	0.67	0.67	0.67
Science	0.67	0.86	0.75
Social_welfare	1.00	1.00	1.00
Space	0.83	1.00	0.91
Sports_other_than_f ootball_and_cricket	0.91	0.91	0.91
Technology	0.89	1.00	0.94
Transportation	0.83	0.83	0.83
Weather	0.87	1.00	0.93
World_and_Internati onal	0.85	1.00	0.92
Cricket	0.67	1.00	0.80
Election	1.00	1.00	1.00
Religion	0.70	1.00	0.82
Travel	0.82	1.00	0.90

Accuracy: 0.85

Table 4: results of Bengali Text Classification with keywords

Cross Validation Scores:

[0.85795455 0.84090909 0.84659091 0.85511364 0.86647727]

Average CV Score: 0.8534090909090908

Number of CV Scores used in Average: 5

## **Chapter 6**

### **Conclusion and Future Works**

With the development of machine learning and deep learning, text classification comes into a new generation, the accuracy is higher and the time consuming is lower. In this thesis, the process of text classification is introduced, extracting keyword from document. Text classification not only can be used in news classification but also in other areas such as spam detection and sentiment analysis. In future these algorithms can be tested on larger corpora. Moreover these algorithms can be improved so that efficiency of categorizations could be improved. A combination of algorithm can be used in order to achieve clustering in a faster way.

## References

- [1] Li, Saihan. "Text classification Based on Machine Learning Methods." (2019).
- [2] Gibaja, Eva, and Sebastián Ventura. "A tutorial on multilabel learning." *ACM Computing Surveys (CSUR)* 47.3 (2015): 1-38.
- [3] Loper, Edward, and Steven Bird. "Nltk: The natural language toolkit." *arXiv preprint cs/0205028* (2002).
- [4] Vajrala, Ajith. "Text Classification." (2019).
- [5] Das, Bijoyan, and Sarit Chakraborty. "An improved text sentiment classification model using TF-IDF and next word negation." *arXiv preprint arXiv:1806.06407* (2018).
- [6] McCallum, Andrew, and Kamal Nigam. "A comparison of event models for naive bayes text classification." *AAAI-98 workshop on learning for text categorization*. Vol. 752. No. 1. 1998.
- [7] Prasanna, P. Lakshmi, and D. Rajeswara Rao. "Text classification using artificial neural networks." *International Journal of Engineering & Technology* 7.1.1 (2018): 603-606.
- [8] Soma Chatterjee & Kamal Sarkar (2022). Predicting Word Importance Using a Support Vector Regression Model for Multi Document Text Summarization. In: International conference on Advance in Data-driven Computing and Intelligent System (ADCIS-2022). (Accepted)