

# Machine and Deep Learning Methods with Manual and Automatic Labelling for News Classification in Bangla Language

Istiaq Ahmad<sup>1</sup>, Fahad AlQurashi<sup>1</sup>, and Rashid Mehmood<sup>2,\*</sup>

<sup>1</sup>Department of Computer Science, Faculty of Computing and Information Technology, King Abdulaziz University, Jeddah 21589, Saudi Arabia

<sup>2</sup>High Performance Computing Center, King Abdulaziz University, Jeddah 21589, Saudi Arabia

\*Corresponding author: RMehmood@kau.edu.sa

## ABSTRACT

Research in Natural Language Processing (NLP) has increasingly become important due to applications such as text classification, text mining, sentiment analysis, POS tagging, named entity recognition, textual entailment, and many others. This paper introduces several machine and deep learning methods with manual and automatic labelling for news classification in the Bangla language. We implemented several machine (ML) and deep learning (DL) algorithms. The ML algorithms are Logistic Regression (LR), Stochastic Gradient Descent (SGD), Support Vector Machine (SVM), Random Forest (RF), and K-Nearest Neighbour (KNN), used with Bag of Words (BoW), Term Frequency-Inverse Document Frequency (TF-IDF), and Doc2Vec embedding models. The DL algorithms are Long Short-Term Memory (LSTM), Bidirectional LSTM (BiLSTM), Gated Recurrent Unit (GRU), and Convolutional Neural Network (CNN), used with Word2vec, Glove, and FastText word embedding models. We develop automatic labelling methods using Latent Dirichlet Allocation (LDA) and investigate the performance of single-label and multi-label article classification methods. To investigate performance, we developed from scratch Potrika, the largest and the most extensive dataset for news classification in the Bangla language, comprising 185.51 million words and 12.57 million sentences contained in 664,880 news articles in eight distinct categories, curated from six popular online news portals in Bangladesh for the period 2014-2020. GRU and Fasttext with 91.83% achieve the highest accuracy for manually-labelled data. For the automatic labelling case, KNN and Doc2Vec at 57.72% and 75% achieve the highest accuracy for single-label and multi-label data, respectively. The methods developed in this paper are expected to advance research in Bangla and other languages.

**Keywords** Natural Language Processing · news classification · Bangla language · word embedding · machine learning · deep learning · automatic labelling · single label classification · multi-label classification

## 1 Introduction

The primary objective of text classification is to determine the class or sentiment of the unknown texts. We can define the problem as follows. Assume, we have  $n$  texts,  $x = x_1, x_2, \dots, x_n$ , and each of them is assigned a category from a set of categorical values  $l$ , where  $l = \{l_1, l_2, \dots\}$ . The training dataset is applied for generating a classification model, which relates the features to one of the class labels. The trained classification model can ascertain the unknown class from the text. Typically, texts are not tagged; we have to do so manually, which is the most time-consuming and challenging task. Additionally, without tagged texts, it's complicated to develop a classification model. Text classification has made continuous success in many applications such as sentiment analysis [1], information retrieval [2], information filtering, knowledge management, document summarization [3], spam mail detection [4], recommended systems, and many others, which have become immense and boundless.

About 238 million people speak Bangla natively or as a second language throughout the world (2021) [5]. As a result, this language has carved out a niche for itself in different information-exchanging media. Bangla, with a large number of online newspapers, blogs, Wikipedia, eBooks, literature, and so on, may be considered to be following the NLP's action ground contest in the imminent future. Each day, a lot of events happening around the world, and some of those events become more trendy discussion topics for a certain time. Most of the news media are engaged in presenting the most popular events every time. Everyone desires to follow the most influential and frequently discussed events among

a large number of events happening around us at a specific time. To get the most contemporary discussion topics and events, text analysis can automatically detect them more precisely and speedily. The research on text analysis.

This paper introduces several machine and deep learning methods with manual and automatic labelling for news classification in the Bangla language. In the case of manual labelling, we implemented several machine (ML) and deep learning (DL) algorithms. The ML algorithms are Logistic Regression (LR), Stochastic Gradient Descent (SGD), Support Vector Machine (SVM), Random Forest (RF), and K-Nearest Neighbour (KNN), used with Bag of Words (BoW), Term Frequency-Inverse Document Frequency (TF-IDF), and Doc2Vec embedding models. The DL algorithms are Long Short-Term Memory (LSTM), Bidirectional LSTM (BiLSTM), Gated Recurrent Unit (GRU), and Convolutional Neural Network (CNN), used with Word2vec, Glove, and FastText word embedding models. To address the challenges related to the arduous task of manual labelling, we develop automatic labelling methods using Latent Dirichlet Allocation (LDA), an unsupervised topic modelling algorithm and investigate the performance of single-label and multi-label article classification methods.

We developed Potrika – the largest and the most extensive dataset for news classification in the Bangla language – comprising 185.51 million words and 12.57 million sentences contained in 664,880 news articles, and used it to investigate the proposed ML and DL methods [6]. Potrika is a single-label news article textual dataset in the Bangla language curated for NLP research from six popular online news portals in Bangladesh (Jugantor, Jaijaidin, Ittefaq, Kaler Kontho, Inqilab, and Somoyer Alo) for the period 2014-2020. The articles are classified into eight distinct categories (National, Sports, International, Entertainment, Economy, Education, Politics, and Science & Technology). GRU and Fasttext with 91.83% achieve the highest accuracy for manually-labelled data. For the automatic labelling case, KNN and Doc2Vec at 57.72% and 75% achieve the highest accuracy for single-label and multi-label data, respectively. The lower performance for automatic-labelling-based classification is because it uses ML algorithms compared to the case of classification with manually-labelled data where the best performance was obtained using a DL algorithm. We will extend our work in the future to include DL methods for automatic labelling.

The NLP methods developed in this paper and the techniques for their extensive analyses are expected to advance research in Bangla and other languages.

**Hardware and Software:** We use the Quadro RTX-6000 GPU, which has 4608 CUDA Parallel-Processing Cores, 576 tensor cores, and 72 RT Cores. The GPU memory is 24 GB of GDDR6. We use Python as the programming language along with machine and deep learning libraries like Tensorflow, Keras, Scikit-Learn, Gensim, etc. We use data visualization libraries like Seaborn and Matplotlib to visualize the evaluation results.

The following is how the paper is structured: Section 2 describes the literature review of text classification, followed by the Bangla text classification (Sections 2.1 to 2.3) and other language text classification (Section 2.4). Section 3 discusses the proposed methodologies of our research including, the overview of methodology and framework architecture 3.1, dataset 3.2, preprocessing 3.3, feature extraction 3.4, and word embedding techniques 3.5. Machine and deep learning methods for manual labelling are described in the Sections 3.6 and 3.7. Section 3.8 explains the methodology for creating the automatically labelled dataset using the unsupervised topic modeling method, and Section 3.9 discusses the methodology of multi-label news article classification. Subsequently, the results of all proposed methods are discussed in Sections 4 and 5, which depict the manual labeling, and automatic single labeling with multi-labeling news article classification results, respectively. Section 6 describe the discussion of the paper. Finally, in Section 7, we conclude with recommendations for further research.

## 2 Literature Review

To address text classification [7], several machine and deep learning-based approaches have been introduced. In this part, we will go through how to classify Bangla text using sentiment analysis, multi-domain, and topic modeling methods. We also go through several methods for classifying other language-related texts.

### 2.1 Sentiment Classification

The core idea of sentiment analysis, or opinion mining, is to analyses the addressed text, if the text expression holds positive, negative, or neutral meaning. For sentiment analysis in Bangla, TF-IDF was applied to a small dataset using machine learning algorithms (see [8, 9, 10]). The word embedding method named word2vec was proposed by [11] for Bangla sentiment analysis based on the Bangla comments. Bangla tweet data is also used for sentiment analysis. For example, Asimuzzaman et al. [12] used an adaptive neuro-fuzzy system for Bangla tweet classification. For sentiment detection, Hasan et al. [13] proposed WordNet and SentiWordNet as tools but the major limitation of this research was proposed tools were developed specifically for English. Tuhin et al. [14] predicted six individual emotions using ML

algorithms such as SVM and NB. Further, NB, DT, KNN, SVM, and K-means clustering were also used by Rahman et al. [15] to predict some basic emotions from the text. In addition, mutual information-based feature selection methods and the multi NB algorithm proposed by Paul et al. [16] for predicting sentiment polarity. N-gram and SVM based Bangla sentiment analysis proposed by Taher et al. [17]. A popular English tool called VADER was proposed by Amin et al. [18] to predict Bangla sentiment.

A deep learning-based algorithm named LSTM was proposed by Hassan et al. [19] for sentiment analysis, where they used 10k Bangla and romanized Bangla text (BRBT) dataset with binary and a categorical cross-entropy loss function. Furthermore, the CNN-based method was proposed by Alam et al. [20].

## 2.2 Multi-domain Text Classification

Alam et al. [21] presented a new dataset for Bengali news articles which contains about 350K articles in five categories (State, Economy, International, Entertainment, and Sports). In their dataset, 65% of the data is labelled as State and 13.5%, 8.5%, 8% and 5% are labelled as Sports, International, Entertainment, and Economy respectively. They have applied machine learning algorithms with two word embedding techniques such as Word2Vec and TFIDF. In another study, the Word2vec embedding model was implemented with KNN and SVM classification algorithms by Ahmed et al. [22] for news document classification. A classification technique based on cosine similarity and Euclidean distance based on a set of 1000 documents was recommended by Dhar et al. [23]. They measure the  $\beta 0$  threshold using the 90th percentile formula for both the distance measures and calculate the score based on the distance. In another research, the dimensional reduction technique with TFIDF (40% of TF) was developed by Dhar et al. [24] where they used a total of 1960 Bangla text documents from five categories (Sports, Business, Science, Medical, and State) with 632,924 tokens and applied machine learning algorithms. The classification algorithm LIBLINEAR achieved the highest accuracy. For 40 thousand news samples divided into 12 categories, Mojumder et al. [25] suggested DL algorithms such as BiLSTM, CNN, and convolutional BiLSTM, and fastText as word embedding techniques. The Bangla article classification based on transformers was proposed by Alam et al. [26]. In this study, they used multilingual transformer models to classify Bangla text in several areas.

## 2.3 Topic Modeling-based Text Classification

Scarce research has been performed to classify Bangla text using topic modeling. Helal and Mouhoub [27] find the key topics in the Bangla news corpus using LDA with a bigram model and classify them by applying similarity measures. They evaluated the proposed model using the LDA and Doc2Vec models and compared the similarity scores. They point out that in some specific news articles, the LDA performance is better than the Doc2Vec model. Alam et al. [28] also proposed an LDA-based topic modeling algorithm using 70k Bangla news articles. They detect 5 distinct news article topics (National, Sports, International, Technology, and Economy) and another topic called 'others' which exclude the following distinct topics.

Most of the above research work has been done on machine learning algorithms with small datasets, but there has been remarkably little works on deep learning algorithms for Bangla article classification because there is no comprehensive dataset for Bangla article classification.

## 2.4 Text Classification for Other Languages

This section provides an overview of different text classification methods for other languages. Shaw et al. [29] implemented ML techniques including random forest, KNN, and logistic regression to classify the news into five categories (Entertainment, Business, Politics, Sports and Technology) based on the BBC news dataset. In terms of efficiency among these algorithms, it turned out that logistic regression has better performance for all the categories. Another research on three machine learning algorithms, namely SVM, neural network, and decision tree, has been done by Raychaudhuri et al. [30] for text classification. The authors used the UCI dataset on US congressional voting that consists of 16 features, 435 instance examples, 335 examples of the training dataset, 50 examples of the testing dataset, and 50 examples of the validation dataset. They used variable C, which controls the training error. When C=1, SVM performed better than the neural network and when C=1000, the neural network performed better. The outcome also revealed that a fully grown decision tree produced better results than a smaller decision tree.

The data augmentation technique is most popular in computer vision research when the amount of data is small or imbalanced. Recently, the text data augmentation technique is noted for small text datasets. Wei and Zou [31] proposed this technique to increase the text classification performance. The following operations are proposed for data augmentation: synonym replacements, random insertion of synonyms of a word, randomly swapping two words positions in a sentence, and randomly removing words in a sentence.

Recently, the attention mechanism has become an efficient approach to determine the important erudition to achieve excellent outcomes. Numerous studies have been carried out on attention mechanisms and architecture. For text classification, several novel methods are also proposed [32, 33, 34, 35]. An attention-based LSTM network was proposed by Zhou et al. [32] to classify cross-lingual sentiments, where they used English and Chinese as the source and target languages, respectively. A Convolutional-Recurrent Attention Network (CRAN) was proposed by Du et al. [33]. Their proposed architecture includes a text encoder using RNN, and an attention extractor using CNN. The experimental result shows that the model effectively extracts the salient parts from sentences along with improving the sentence classification performance. Liu et al. [34] proffered attention-based convolution layer and BiLSTM architecture, where the attention mechanism provides a focus for the hidden layers output. The BiLSTM is used to extract both previous and following context, while the convolutional layer retrieves the higher-level phrase from the word embedding vectors. Their experimental results get comparable results for all the benchmark datasets.

The state-of-the-art graph-based neural network methods for text classification have been gaining increasing attention recently. A text graph convolutional network (TextGCN) was proposed by Yao et al. [36], which is more notable for its small training corpus for text classification. To learn the TextGCN for the corpus, word co-occurrence and the relation between the word document based single text graph was developed. Another tensor graph convolutional network (TensorGCN) has been proposed by Liu et al. [37]. They develop the text graph tensor based on semantic, syntactic, and sequential contextual information. After that, two types of propagation learning are performed on the text graph tensor called intra-graph propagation to aggregate information from neighboring nodes and inter-graph propagation to tune heterogeneous information between graphs.

Capsule network is another state-of-the-art method for text classification that is inherent to CNNs. Several studies based on the capsule network have been conducted [38, 39, 40]. In capsule networks, capsules are locally invariant groups that learn to recognize the presence of visual entities and encode their characteristics into vectors. It also requires a nonlinear function called squashing, whereas neurons in a CNN act independently. However, equivariance and dynamic routing are the two most essential characteristics of Capsule Networks that distinguish them from standard Neural Networks. A Capsule network with dynamic and static routing based text classification methods was proposed by Kim et al. [39]. Static routing achieved higher accuracy than dynamic routing. Yang et al. [38] introduced a cross-domain capsule network and illustrated the transfer learning applications for single-label to multi-label text classification and cross-domain sentiment classification. An attention mechanism-based capsule network system called Deep Refinement was suggested by Jain et al. [40]. Their proposed method achieved 96% accuracy for text classification compared to BiLSTM, SVM, and C-BiLSTM for the Quora insincere question dataset.

Traditional text classification techniques use manually labelled datasets that are monotonous and time-consuming. Recently, a few dataless text classification techniques, for example, the Laplacian seed word topic model (Lap-SWTM) [41], and seed-guided multi-label topic model (SMTM) [42] have recently been proposed to solve this challenge. Anantharaman et al. [43] proposed large and short text classification non-negative matrix factorization, LDA, and LSA (latent semantic analysis). LSA with TFIDF was proposed by Neogi et al. [44] for text classification. To increase the accuracy, they used entropy. A self-training LDA based semi-supervised text classification method was proposed by pavlinek et al. [45] for text classification.

## 2.5 Research Gap, Novelty, and Contributions

Text datasets, often known as corpora, are used to study linguistic phenomena including text classification, morphological structure, word sense disambiguation, language evolution over time, and spelling checking. The quality and amount of the corpus have a big impact on the research output. A well-structured, comprehensive corpus can yield far superior study results. In comparison to the English language, there has been inadequate study done due to the paucity of the Bangla corpus and the complicated grammatical structure. In this paper, our contributions are as follows:

- We are the first to use a comprehensive Bangla newspaper article dataset called Potrika [6, 46] to classify eight distinct news article classes, including Education, Entertainment, Sports, Politics, National, International, Economy, and Science & Technology.
- We implement both machine learning (ML) including logistic regression, SGD, SVM, RF and KNN algorithms, and deep learning (DL) including CNN, LSTM, BiLSTM, and GRU algorithms for single label news article classification. We perform BOW, TFIDF, and Doc2Vec word embedding models for ML algorithms. For DL algorithms, we apply word embedding models such as word2vec, glove, and fasttext that were developed based on the Potrika dataset. These word embedding models are not only valuable for news article classification but also for other NLP tasks like text summarization, named entity recognition, Bangla automatic word prediction, question-answering systems, etc. Further, we evaluate and scrutinise the results for both cases.



- Manual labeling is the most difficult and time-consuming task for classification datasets. In the following paper, we investigate the possibility of using the topic modeling algorithm to automatically label the news article dataset and compare the performance of the automatically labelled dataset with that of the manually labelled dataset. Additionally, we also develop another multi-label dataset based on the automatic label dataset and evaluate the multi-label news article classification’s performance.

The NLP work proposed in this paper builds on our earlier NLP works applied to several sectors and multiple languages including transportation [47, 48, 49], healthcare [50, 51], education [52, 53], and smart cities [54, 55, 56, 57]. We expect that this paper will significantly increase the impact of our work particularly in the Bangla language.

### 3 Methodology and Design

In this section, we describe our methodology for the research presented in this paper. We begin with an overview of our methodology in Section 3.1 followed by a description of the dataset in Section 3.2. Data preprocessing and feature extraction are explained in Sections 3.3 and 3.4. Word Embedding Models are discussed in Section 3.5. Machine learning and deep learning techniques are discussed in Section 3.6 and 3.7, respectively. Subsequently, we explain our methodology for automatically creating labels for the news items. The details of automatic labeling with single labels are provided in Section 3.8 and the details of marking news items with multiple labels are given in Section 3.9.

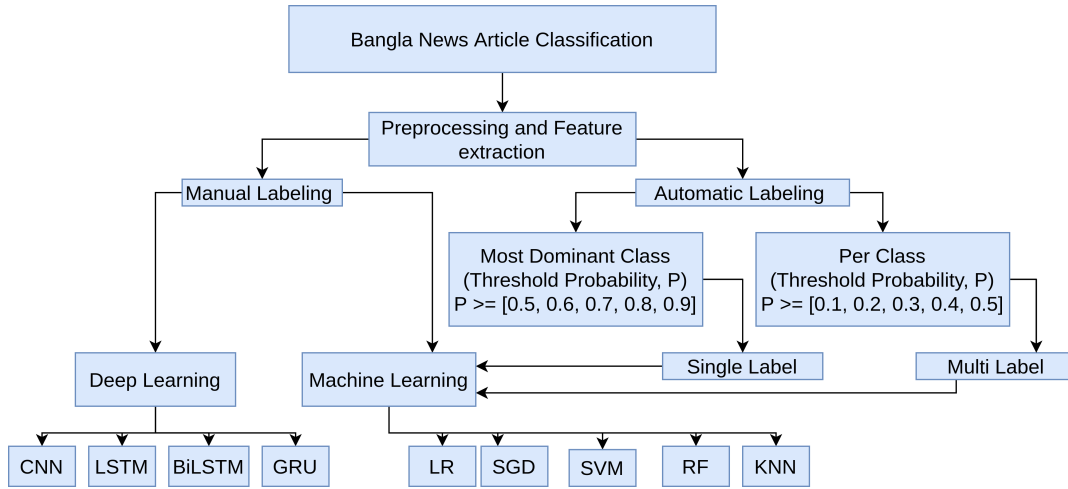


Figure 1: System Process for Article Classification Overview

#### 3.1 Methodology Overview

As mentioned earlier, the aim of the paper is to investigate the performance of machine and deep learning-based news classification in the Bangla language using manual and automatic labeling of news documents. Towards this end, we explore, firstly, the classification of manually labelled news data using machine and deep learning algorithms and, secondly, the classification of automatically labelled news items using single label and multi-label approaches. The automatic labeling is done using topic modeling. The overview and detailed architecture of our methodology are depicted in Figure 1 and Figure 2 and its algorithmic flow is provided in Algorithm 1.

---

#### Algorithm 1 Master Algorithm

---

**Input:** *ReadtrainDF, testDF, potrikaDF*

**Output:** *Evaluationofnewsarticleclassification*

---

- 1: *cleanTrText, cleanTsText, cleanText*  $\leftarrow$  *preprocessing (trainDF, testDF, potrikaDF)*
  - 2: *evaluation*  $\leftarrow$  *man\_ML(cleanTrText, trainDF.class, cleanTsText, testDF.class)*
  - 3: *w2vmodel, gloveModel, fasttextmodel*  $\leftarrow$  *getWordEmbeddingModels(cleanText)*
  - 4: *evaluation*  $\leftarrow$  *man\_DL(w2vmodel, gloveModel, fasttextmodel)*
  - 5: *evaluation, autoLabelingDF*  $\leftarrow$  *auto\_singleLabel(trainDF, testDF)*
  - 6: *evaluation*  $\leftarrow$  *auto\_multiLabel(autoLabelingDF)*
-

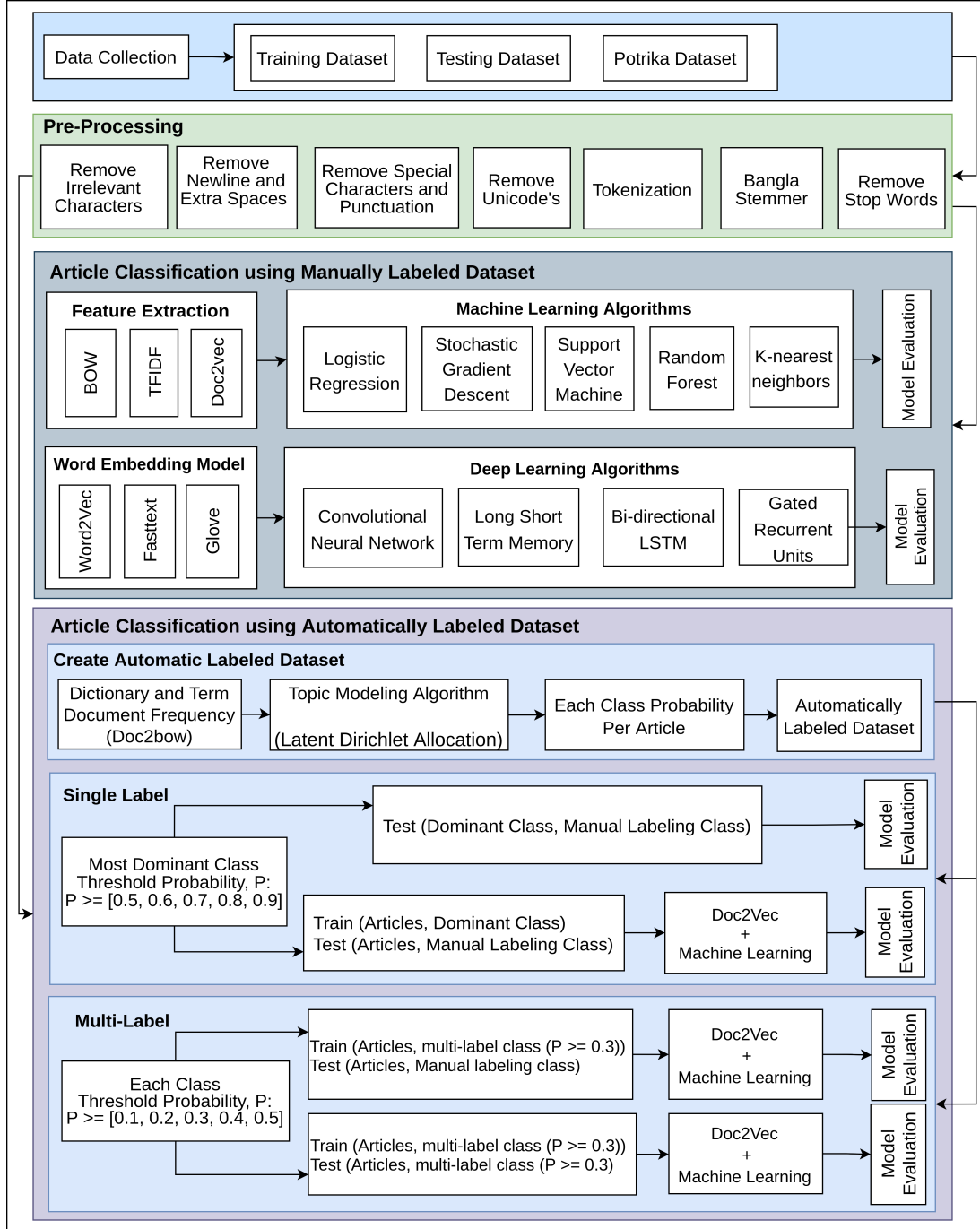


Figure 2: A Detailed System Process for Article Classification

### 3.2 Dataset

We used Potrika (for details, see [6, 46]), a large single-labelled Bangla News article dataset derived from six popular online news portals, including Jugantor, Jaijaidin, Ittefaq, Kaler Kontho, Inqilab, and Somoyer Alo, and divided into eight classes: National, Sports, International, Economy, Education, Politics, and Science & Technology. The dataset has five columns for each class: news article, class, headline, publish date, and news source. Over 665K news articles, 12.5 million sentences, and 185.5 million words are included in the Potrika dataset. We used the Potrika dataset for word embedding models, and a total of 120K articles for text classification, including 100K for training and 20K for

testing. Each class has 12.5K and 2.5K articles in the training and testing sets, respectively. The total amount of words each document/article vs the total number of documents/articles is depicted in Figure 3.

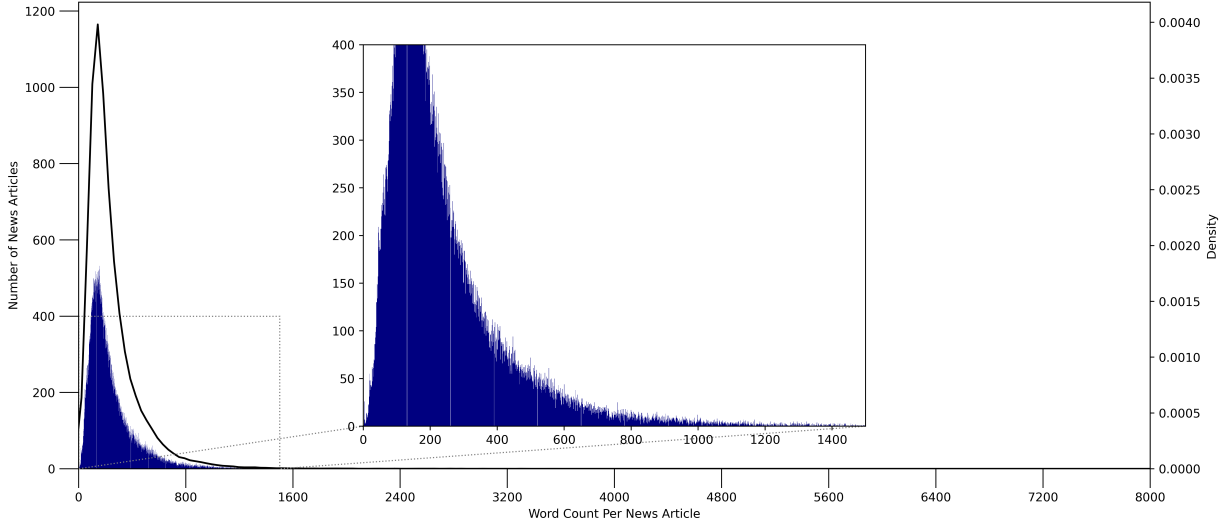


Figure 3: Distribution of Document Word Counts

Table 1 depicts the short form of each category used in the result section.

Table 1: Short Form of Each Category

Sports	SP
National	NA
Economy	EC
Entertainment	EN
Politics	PO
International	IN
Education	ED
Science&Technology	ST

### 3.3 Pre-Processing

Most of the texts include several irrelevant words, such as stop words, misspellings, unrecognized characters, Unicode, and so on, which might have inimical impacts on research. This section briefly explains the pre-processing technique to clean the texts. We have used similar pre-processing techniques for all the proposed methods. Moreover, we also provide an algorithm 2 to illustrate the pre-processing techniques.

**Algorithm 2** Pre-processing**Input:** *trainDF*, *testDF*, and *potrikaDF***Output:** *cleanTrText*, *cleanTsText*, *cleanText*


---

```

1: cleandataset  $\leftarrow$  initialize with empty list
2: for  $i = 0$  to datasetLength do
3:   cleanArticle  $\leftarrow$  article[ $i$ ]
4:   cleanArticle  $\leftarrow$  re.sub('[a-zA-Z0-9]', '', cleanArticle)
5:   cleanArticle  $\leftarrow$  cleanArticle.replace('\n', '')
6:   for  $x = 0$  to length of bangladigitlist do
7:     cleanArticle  $\leftarrow$  cleanArticle.replace(x, ' ')
8:   end for
9:   cleanArticle  $\leftarrow$  translate(str.maketrans(", ", string.punctuation))
10:  cleanArticle  $\leftarrow$  cleanArticle.replace(u'\uf06c', '')
11:  cleanArticle  $\leftarrow$  re.sub(r'\x9d', r'', cleanArticle)
12:  cleanArticle  $\leftarrow$  re.sub(' *', ' ', cleanArticle)
13:  cleanArticle  $\leftarrow$  cleanArticle.split()
14:  bs  $\leftarrow$  stemmer.BanglaStemmer()
15:  cleanArticle  $\leftarrow$  [bs.stem(word) for word in cleanArticle]
16:  cleanArticle  $\leftarrow$  [word for word in cleanArticle if not word in bangla_stopwords]
17:  cleanArticle  $\leftarrow$  ' '.join(cleanArticle)
18:  cleandataset.append(cleanArticle)
19: end for

```

---

There are numerous optional Bangla words that do not have major consequences in classification algorithms. We manually generate a list of 542 stop words based on the significance of these words and remove from the texts (see Figure 4).

An enormous number of special and punctuation characters (i.e., ‘!#\$%&’()\*+?@’) are included in the texts. This punctuation and special characters are removed.

In the texts, we found some Unicode characters including ‘\uf06c’, ‘\u200c’, ‘\u09e5’, ‘\x9d’, etc., which we removed from the text to improve the accuracy of classification models.

A single word might appear in diverse forms in a document, all of which have the same semantic meaning. Stemming refers to consolidating different versions of words into the same feature location. In our research, we used python Bangla stemmer (version 1.0) to margin the same words. Figure 4 shows some examples before and after applying the Bangla stemmer.

Tokenization is a preprocessing technique that distributes texts into tokens, which can be words, phrases, symbols, or other significant components. To process the tokenization of the texts, text classification requires a parser. Figure 4 shows some examples before and after applying tokenization.

The feature extraction approach has been used by a number of researchers to deal with the problem of losing syntactic and semantic correlations inside words. We used the n-gram technique (1-gram, 2-gram, 3-gram) to address the syntactic problem. We used the 1-gram technique for manual labeling and the 3-gram technique for automatic labeling. Figure 4 shows some examples of n-gram.

### 3.4 Feature Extraction

Feature extraction or document representation is the most prominent method for enhancing the performance of text classification, since it enables us to determine which features are most relevant for the intended classification job. It should be done more accurately because of the vast dimensional of text features and the presence of irrelevant or noisy data. Usually, texts are unstructured and unorganized data that needs to be transformed into structured data using mathematical modeling as part of a classifier. Word embedding techniques are well-known feature extraction approaches that are covered in Section 3.5.

### 3.5 Word Embedding Models

Several word embedding methods are discussed in this section, including BOW, TFIDF, Doc2Vec, word2vec, fasttext, and glove. BOW counts the total occurrence of a word in a document, and TFIDF assesses the significance of a word

Technique	Example
Stop Words	গোটা, আমরা, আবার, উনি, এবং, গেল, কোন
Stemmer	before: ডলারের, হাইকমিশনের, রপ্তানিকারকদের, করেছিলেন, প্রতিষ্ঠানের after: ডলার, হাইকমিশন, রপ্তানিকারক, করেছি, প্রতিষ্ঠান
Tokenization	before: প্রধান নির্বাহী কর্মকর্তা আনিসুর রহমান after: প্রধান, নির্বাহী, কর্মকর্তা, আনিসুর, রহমান
N-gram	sentence: সোহেল মিয়া দেশে টাকা পাঠাননি 1-gram: সোহেল, মিয়া, দেশে, টাকা, পাঠাননি 2-gram: সোহেল মিয়া, মিয়া দেশে, দেশে টাকা, টাকা পাঠাননি 3-gram: সোহেল মিয়া দেশে, মিয়া দেশে টাকা, দেশে টাকা পাঠাননি

Figure 4: Text Pre-Processing Techniques Example

in a document, not its frequency. Word2Vec trains neural networks on documents and outputs a vector for each word. The embedding result catches whether words appear in similar contexts. The word2vec skip-gram model seeks to obtain co-occurrence one window at a time, whereas the glove seeks to obtain the counts of overall statistics on how often it appears. Both word2vec and glove train on the smallest units of words. Fasttext is a word2vec extension that treats each word as a set of n-grams. The following approach enables embedding to be trained on a smaller corpus and then generalized to unknown or rare words. A brief discussion of these word embedding methods is given below. We developed word embedding models for deep learning-based news article classification using the Potrika dataset. The procedure is depicted by the algorithm 3.

---

**Algorithm 3** Word Embedding Model

---

**Input:** *cleanText***Output:** *Word Embedding Models*

- 1:  $w2vmodel \leftarrow word2vecmodel(cleanText)$
  - 2:  $glovemodel \leftarrow glovemodel(cleanText)$
  - 3:  $fasttextmodel \leftarrow fasttextmodel(cleanText)$
- 

**3.5.1 Bag of Words (BOW)**

The most basic type of numerical text description is the BOW model. To produce an unordered list of words without syntactic, semantic, or POS labeling, extract just the uni-gram words from the BOW. The document is represented by these groupings of words. There are various disadvantages to this strategy. If the new sentences contain new words, then the vocabulary dimension and vector length will grow. Moreover, multiple 0s may occur in the vectors, resulting in a sparse matrix that must be avoided.

**3.5.2 Term Frequency Inverse Document Frequency (TFIDF)**

A statistical metric for determining the importance of a word in a corpus or collection of documents is the TFIDF model. The significance of a word increases proportionally to its occurrence in the document, which remains offset by its recurrence in the corpus. We can divide TFIDF into two segments: term frequency (TF) and inverse document frequency (IDF). TF denotes a metric that measures the occurrence of a term in the current document. Since each document varies in length, a term may occur more frequently in large documents than in diminutive documents. For normalization, the TF is often divided by the length of the document. IDF denotes a metric that measures how rarely the word occurs in all documents. Equation 1, 2, and 3 describe the formula of TFIDF method, where  $\zeta$  = number of documents in total,  $\beta$  = total amount of words in the document,  $\alpha$  = number of times a word  $w$  occurs in a document, and  $\delta$  = number of documents in which the word  $w$  occurs.

$$TF_w = \frac{\alpha}{\beta} \quad (1)$$

$$IDF_w = \log \frac{\zeta}{\delta} \quad (2)$$

$$TF\_IDF_w = TF_w \times IDF_w \quad (3)$$

### 3.5.3 Doc2Vec

Doc2Vec is an extension of the word2vec (CBOW) model introduced by Mikilov and Le [58]. The word2vec model runs similarity queries to predict the subsequent words, whereas in the Doc2Vec model, we tag the document with extra tag vectors called document unique. When the word vectors  $W$  are trained, the document vector  $D$  is also trained, and it retains a numeric representation of the document at the conclusion of the training. The model is performing as memory that recalls what is lacking from the present context or the paragraph's topic. This model is also known as paragraph vector with distributed memory. The execution time for Doc2Vec is about 2.42 hours. There is a total of 664,883 tokens and 368,894 vocabularies in the Doc2Vec model.

### 3.5.4 Word2vec

The word2vec was mentioned as an enhanced word embedding design by T. Mikolov et al. [59, 60]. This approach used deep neural networks with two hidden layers, continuous bag-of-words (CBOW), and the continuous skip-gram model to create a high-dimensional vector per word and retain syntactic and semantic information of sentences. For the target phrase, the CBOW model is represented by numerous words. For instance, the context terms "aeroplane" and "military" for the target phrase "air-force." The continuous Skip-gram model, on the other hand, aims to maximize a word's classification depending on another word in the same phrase. The application of word2vec is enormous in deep learning, such as text classification, language modeling, question and answer systems, machine translations, image captioning, speech recognition, and so on. The execution time for word2vec is about 4 hours. There is a total of 664,883 tokens and 368,894 vocabularies in the word2vec model.

### 3.5.5 Glove

Global Vectors (GloVe) [61] is a word embedding technique that is very comparable to the word2vec method. Here, each word is presented by a high-dimensional vector and trained on the surrounding words in a large corpus. Glove predicts surrounding words by using dynamic logistic regression to maximize the likelihood of occurrence of a context word given a core word. Other pre-trained word vectorizations with 100, 200, and 300 dimensions are available from the glove model, which has been trained on larger corpora. However, this pre-trained model is not suitable for the Bangla language. We trained the Glove model for the Potrika dataset for 300 dimensions. In the vocabulary building, 123,133,606 tokens are processed and total unique words 1,165,377 with 3 minimum occurrences accepted, and the total vocabulary size is 368,894. There are a total of 30 epochs and the execution time is 2.3 hours.

### 3.5.6 Fasttext

FastText [62] is an extension of the word2vec model that encodes each word as an n-gram of characters rather than learning vectors for words directly. It supports the training of CBOW or Skip-gram models using negative sampling, softmax or hierarchical softmax loss functions, and in terms of word representations and sentence classification, it performs admirably, especially for rare words, by utilizing character-level information. We trained a fasttext model with 300 dimensions, and the execution time was about 8.43 hours. There is a total of 664,883 tokens and 368,894 vocabularies in the fasttext model.

## 3.6 Machine Learning Algorithms

Before applying ML algorithms, we took the training and testing datasets separately and preprocessed both datasets to create a clean dataset. We discussed the preprocessing techniques in Section 3.3. After preprocessing the dataset, we applied three-word embedding techniques 3.5 including BOW, TFIDF, and Doc2Vec for feature extraction. We used ML algorithms named RF, LR, KNN, SGD, and SVM for the news article classification. Algorithm 4 shows the process for news article classification using machine learning algorithms.

**Algorithm 4** News Article Classification with Manual Labeling and Machine Learning**Input:** *cleanTrText, trainDF.class, cleanTsText, testDF.class***Output:** *Evaluation of news article classification*

- 1: *train*  $\leftarrow$  *cleanTrText, trainDF.class*
- 2: *test*  $\leftarrow$  *cleanTsText, testDF.class*
- 3: *xtrain, ytrain, xtest, ytest*  $\leftarrow$  *bow(train, test, maxfeature=300)*
- 4: *xtrain, ytrain, xtest, ytest*  $\leftarrow$  *tfidf(train, test, maxfeature=300)*
- 5: *d2vmodel, traintag, testtag*  $\leftarrow$  *Doc2VecModel(train, test, maxfeature=300)*
- 6: *xtrain, ytrain, xtest, ytest*  $\leftarrow$  *Doc2Vec(d2vmodel, traintag, testtag)*
- 7: *evaluation*  $\leftarrow$  *mlAlgorithms(xtrain, ytrain, xtest, ytest)*

**3.7 Deep Learning Algorithms**

With three-word embedding techniques such as Word2vec, Glove, and Fasttext, we applied DL algorithms such as CNN (Convolutional Neural Network), LSTM (Long Short Term Memory), biLSTM (Bi-Directional LSTM), and GRU (Gated recurrent units). We classified the news articles using the manually labelled dataset. Algorithm 5 shows the process for news article classification using deep learning algorithms.

**Algorithm 5** News Article Classification with Manual Labeling and Deep Learning**Input:** *cleanTrText, cleanTsText, w2vmodel, glovemodel, fasttextmodel***Output:** *Evaluation of news article classification*

- 1: *train*  $\leftarrow$  *cleanTrText, trainDF.class*
- 2: *test*  $\leftarrow$  *cleanTsText, testDF.class*
- 3: *xtrain, ytrain, xtest, ytest*  $\leftarrow$  *train, test* ▷ apply tokenizer, pad\_sequences
- 4: *w2v*  $\leftarrow$  *embeddingMatrix(w2vmodel, cleanTrText)*
- 5: *fasttext*  $\leftarrow$  *embeddingMatrix(fasttextmodel, cleanTrText)*
- 6: *glove*  $\leftarrow$  *embeddingMatrix(glovemodel, cleanTrText)*
- 7: *evaluation*  $\leftarrow$  *dlAlgorithms(embeddingMatrix, train, test)*

The CNN architectures is shown in Figure 5. For LSTM, BiLSTM, and GRU models, we used a similar architecture with a batch size of 64 and an epoch number of 8. The architecture of these three algorithms are shown in Figure 6.

**3.8 Automatic Single labelled News Article Classification**

Topic modeling is an unsupervised machine learning method for analyzing and exploring latent information and expression patterns in multiple documents. In this study, we explore the behaviour of news article classification for the unlabelled dataset using the Latent Dirichlet Allocation (LDA) topic modeling technique. Each news article is determined by the probability distribution over multiple topics or classes while a given topic or class is described as a probability distribution across words.

Algorithm 6 describes pseudo code to create an automatically labelled dataset. In our case, the number of classes is 8. We merge training and testing datasets and preprocess the merged articles. Doc2bow was used for the feature extraction, which includes the news article id and its frequency in every article. We use various parameter tunings, including n-gram, chunk size, passes, and iterations, to find the best LDA model. Based on the keywords and perplexity value, we chose the finest LDA model. After finding the best model, we have 8 lists of keywords based on the trained model, and each list defines a class. We manually labelled the names of the lists as news article class names based on the keywords. Table 7 demonstrates the top 5 keywords for each of the classes. We retrieve an array of the probability distribution for each news article, which reflects how much the news article belongs to each class. To find the dominant class for each news article, we choose the class with the highest contribution probability. Finally, we create an automatically labelled dataset that contains all 8 class probabilities and the most dominant class, the manually labelled class, and the preprocessed article.

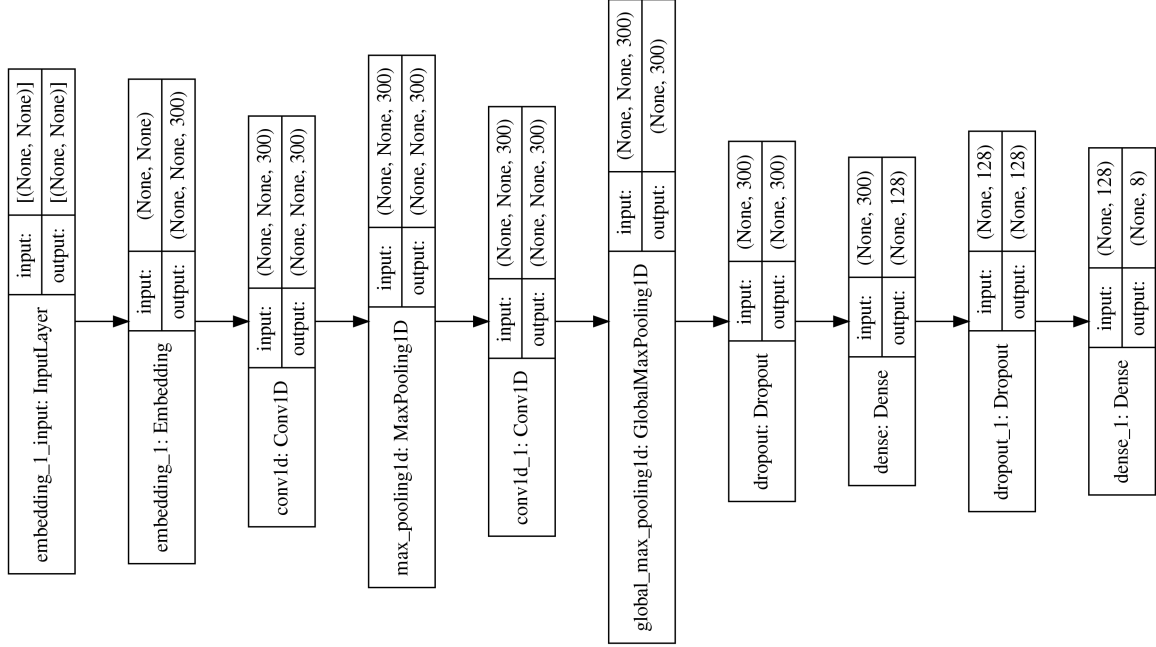


Figure 5: CNN Architecture

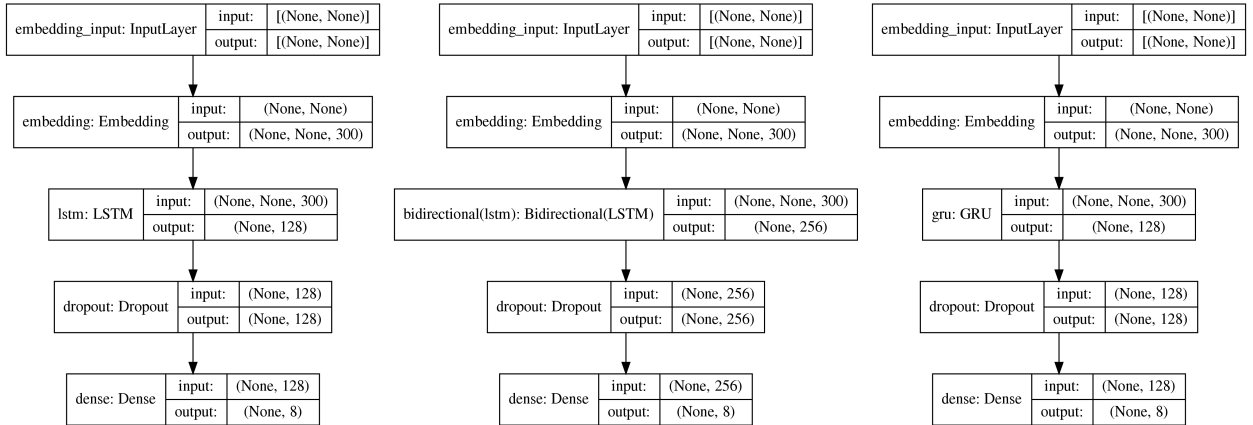


Figure 6: Architecture of LSTM, BiLSTM, and GRU

**Algorithm 6** News Article Classification with Automatic Labeling for Single Label**Input:** *trainDF*, *testDF*, *ngram*, *[chunksize]*, *[passes]*, *[iterations]*, *[th]***Output:** *Evaluation of single label news article classification*

- 1: *cleanArticlesDF*  $\leftarrow$  *trainDF.append(testDF)*
- 2: *cleanArticles*  $\leftarrow$  *preprocessing(cleanArticlesDF)*
- 3: *ngram*  $\leftarrow$  *make\_ngram(cleanArticles, [ngram])*
- 4: *corpus*  $\leftarrow$  *doc2bow(ngram)*
- 5: *ldaModel*  $\leftarrow$  *getBestModel(corpus, [chunksize], [passes], [iterations])*
- 6: *domClass*  $\leftarrow$  *getMostDominantClass(ldaModel)*
- 7: *classProb*  $\leftarrow$  *getEachClassProbPerArticle(ldaModel)*
- 8: *autoLabelDF*  $\leftarrow$  *getAutoLabelDF(classProb, domClass, cleanArticlesDF.class, ngram)*
- 9: *evaluation*  $\leftarrow$  *autoLabelDF.class, autoLabelDF.domClass*
- 10: *xtrain*, *ytrain*, *xtest*, *ytest*  $\leftarrow$  *getTrainTestSet(autoLabelDF, [th])*
- 11: *evaluation*  $\leftarrow$  *mlAlgorithms(xtrain, ytrain, xtest, ytest)*



Figures 8 and 9 describe the distribution of word counts by original topics and dominant topics respectively. Figure 10 shows the number of documents for each topic probability contribution.

Topic	Keywords
Sports	দল, ম্যাচ , ক্রিকেট , খেলা, রান Team, Match, Cricket, Play, Run
National	পুলিশ, গ্রাম, উপজেলা, স্থানীয়, এলাকা Police, Village, Sub-District, Local, Area
Economy	টাকা, লেনদেন, দাম, ব্যাংক, আয় BDT, Transaction, Cost, Bank, Income
Entertainment	ছবি, অভিনয়, নাটক, গান, সিনেমা Picture, Acting, Drama, Song, Movie
Politics	লীগ, প্রার্থী, নেতা, বিএনপি, নির্বাচন League, Candidate, Leader, BNP, Election
International	দেশ, বাংলাদেশ, ভারত, প্রেসিডেন্ট, যুক্তরাষ্ট্র Country, Bangladesh, India, President, USA
Science & Technology	ফোন, মোবাইল, ফেসবুক, প্রযুক্তি, ডিজিটাল Phone, Mobile, Facebook, Technology, Digital
Education	বিশ্ববিদ্যালয়, শিক্ষার্থী, পরীক্ষা, ভর্তি, অধ্যাপক University, Student, Exam, Admission, Professor

Figure 7: Top 5 Keywords of Each Cluster using LDA

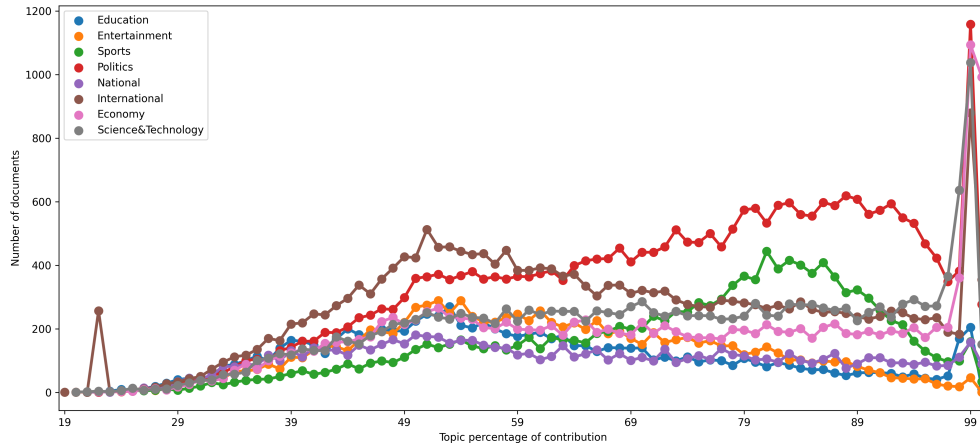


Figure 10: Probability Map

### 3.9 Automatic Multi-labelled News Article Classification

A set of labels  $L$ , a set of examples  $X$ , where each example  $x$  is connected with a subset of the relevant  $L$  labels, can be used to represent the multi-label classification problem. The primary purpose of this is to create an  $L$ -dimensional target vector  $y \in \{0, 1\}^L$ , where  $y_i = 1$  represents the relevant  $i$ -th label, whereas  $y_i = 0$  represents the irrelevant label for the current case. To address the multi-label text classification issue, there are primarily three types of approaches: problem transformation methods (Label Powerset, Binary Relevance, and Classifier Chains), Algorithm adaptation methods (KNN, decision trees), and neural network models. In this paper, we perform binary relevance learning that

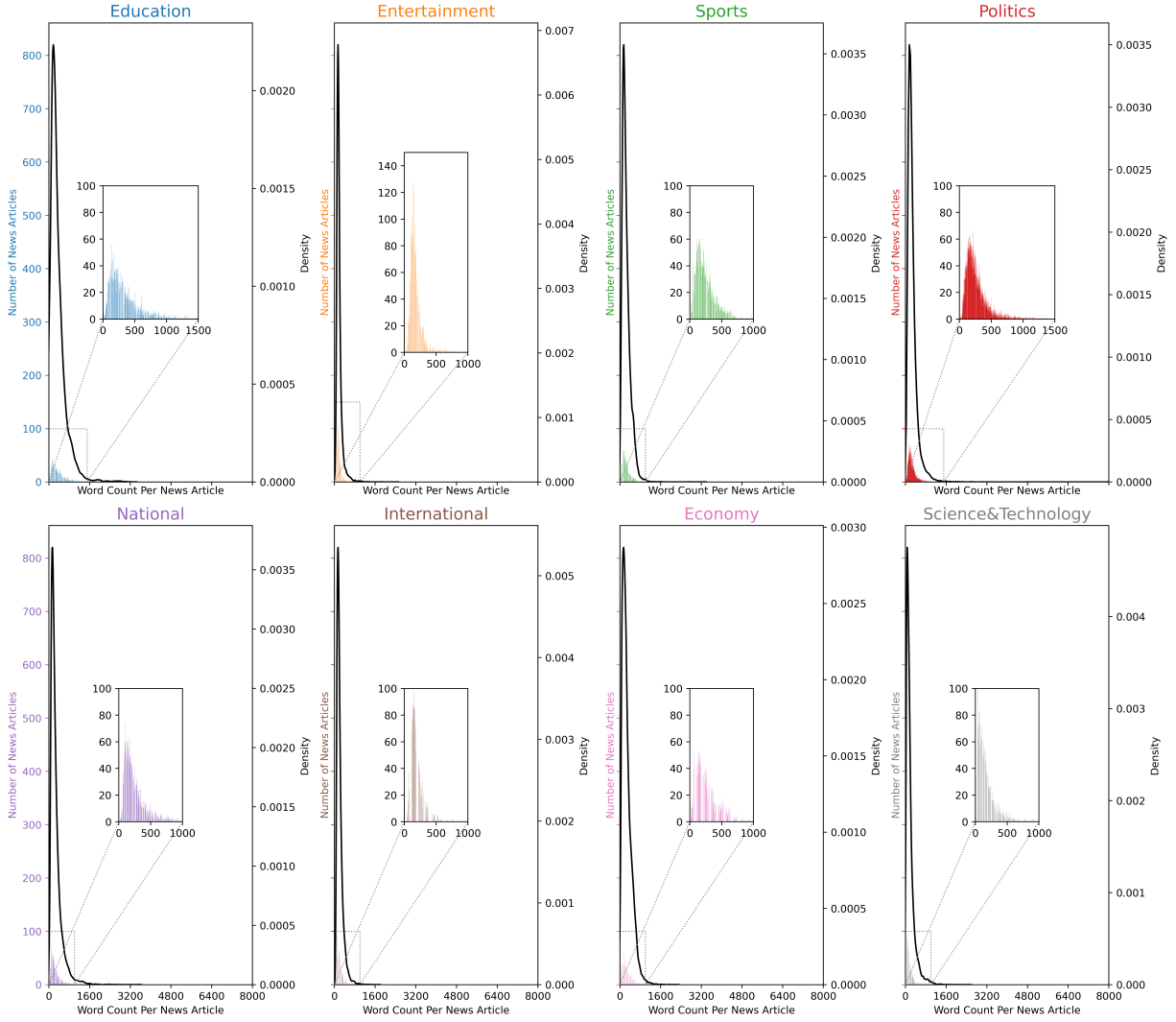


Figure 8: Distribution of Document Word Counts (Original Topic)

constructs  $L$  binary classifiers by training on the  $L$  labels individually and combining all the classifiers results into a multi-label prediction by overlooking the associations between labels.

In the previous section, we created an automatically labelled dataset containing the 8 class probabilities for each item. Algorithm 7 describes the pseudo code for creating a multi-labelled dataset using this dataset. We use MultiLabelBinarizer to generate the multi-label classes. We apply the same ML techniques with the Doc2Vec feature extraction methodology to train the model. As the prediction for an example or instance consists of a collection of labels that may be entirely correct, partially correct (with varying degrees of accuracy), or completely incorrect, evaluating a multi-label classifier is more difficult than evaluating a single-label classifier. We will use example-based matrices to evaluate the multi-label classifiers performance, such as subset f1-score, precision, recall, accuracy, and hamming loss.

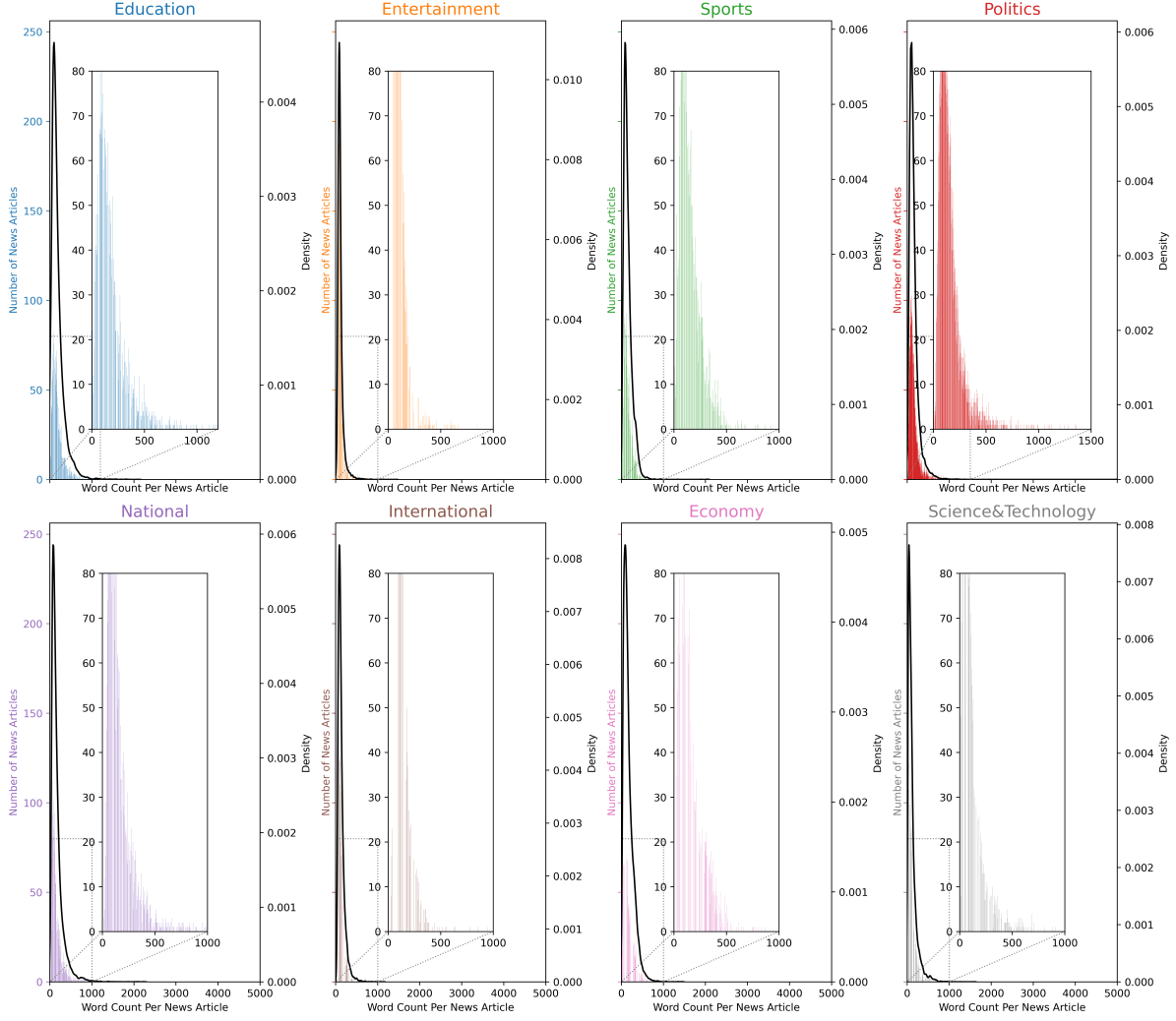


Figure 9: Distribution of Document Word Counts (Dominant Topic)

**Algorithm 7** News Article Classification with Automatic Labeling for Multi-Label**Input:** *autoLabelDF*, *cleanArticles*, [*th*]**Output:** *Evaluation of multi-label news article classification*

- 1:  $mulLabelDF \leftarrow getMultiLabelBinarizer(autoLabelDF.classProb, [th])$
- 2:  $manMulLblClass \leftarrow getMultiLabelBinarizer(autoLabelDF.class)$
- 3:  $evaluation \leftarrow manMulLblClass, mulLabelDF$
- 4:  $xtrain, ytrain, xtest, ytest \leftarrow mlAlgo\_d2v(autoLabelDF.ngram, mulLabelDF)$
- 5:  $evaluation \leftarrow mlAlgorithms(xtrain, ytrain, xtest, ytest)$

Figure 11 describes the number of news articles with multiple classes for different probability thresholds such as 0.1, 0.2, 0.3, 0.4, and 0.5. For a threshold of 0.1, 56,217 news articles have 1 class, 28,081 news articles have 2 classes, 27,939 news articles have 3 classes, 6,789 news articles have 4 classes, 688 news articles have 5 classes, 286 news articles have 6 classes and no article has no class. For a threshold of 0.5, most of the articles have 1 class and 22,420 articles have no class because the probability of each class is less than 0.5. To evaluate multi-label classification, we use the threshold of 0.3, where 88,394 articles have 1 class, 30,385 articles have 2 classes, 188 articles have 3 classes, and 1,033 articles have no class.

The number of news articles for each class is shown in Figure 12. The International and Politics classes have the most articles, with 29,245 and 31,741 respectively.

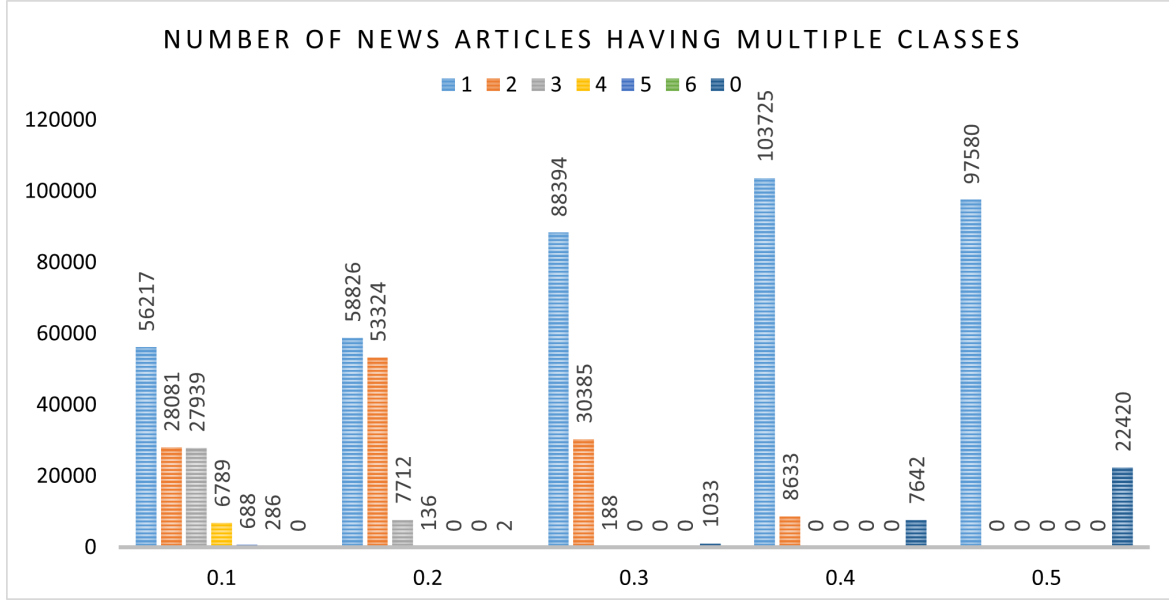


Figure 11: Number of News Articles Having Multiple Classes

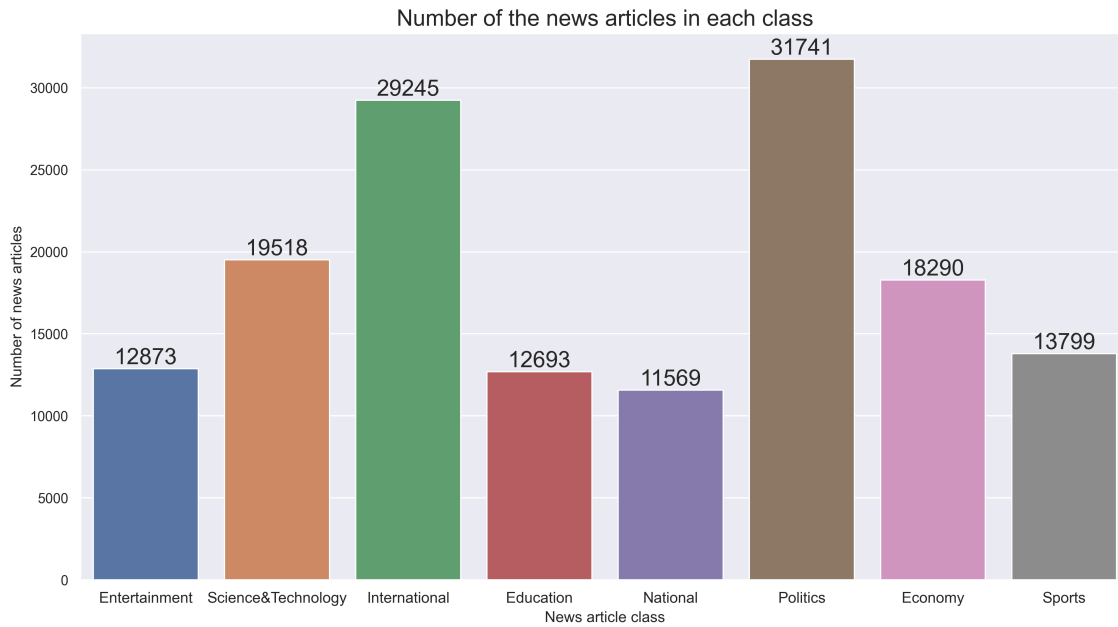


Figure 12: Number of News Articles for Each Class

#### 4 Manually labelled News Article Classification

This section analyses the performance of supervised machine and deep learning algorithms, which were discussed in Section 3.6 and 3.7, to establish a detailed comparison of the news article classification models.

Figure 13 shows the accuracy in percentage of the news article classification for machine learning algorithms using embedding techniques, i.e., BOW, TFIDF, and Doc2Vec. We used a maximum of 300 feature vectors for each word embedding technique. The highest accuracy was achieved 87.14% by logistic regression algorithm with the Doc2Vec technique. Moreover, the Doc2Vec technique obtained more than 80% accuracy for all the ML algorithms, whereas the BOW and TFIDF techniques obtained very low accuracy. The Doc2Vec is a word2vec extension that works as a memory for what is lacking in the present context. On the other hand, BOW and TFIDF do not capture the context of

the text. As a result, in large documents with extensive textual variation and complexity, the ML models are overfitted for BOW and TFIDF approaches. With the BOW and TFIDF techniques, we achieved the highest accuracy of 49.39% and 54.72% for SGD, respectively. In summary, Doc2Vec outperforms the BOW and TFIDF approaches in terms of accuracy.

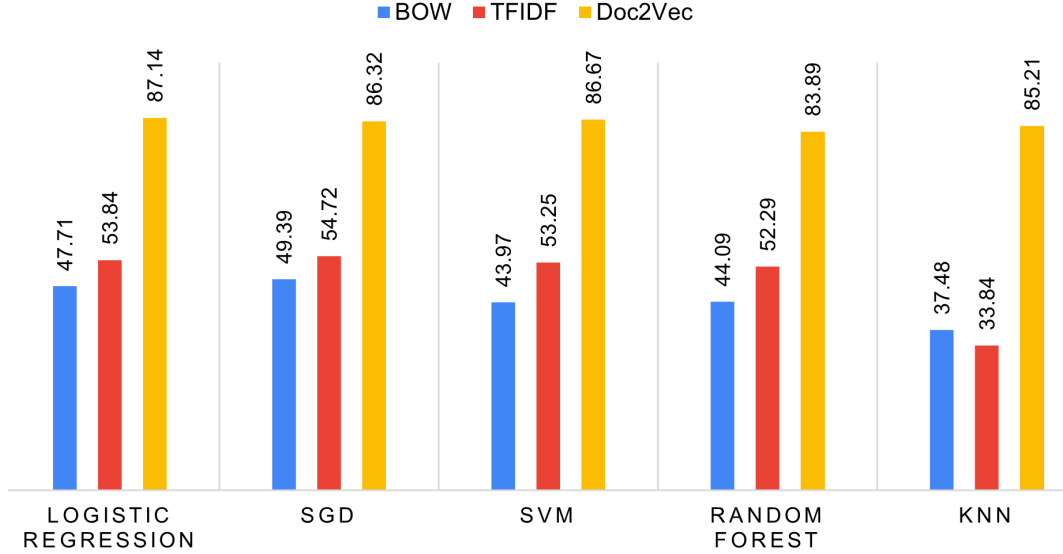


Figure 13: Manually labelled Article Classification Accuracy for ML Algorithms

Figure 14 depicts the execution time for ML algorithms where SVM and KNN have the highest and lowest execution times compared to other ML algorithms. For the SVM algorithm, BOW has the highest execution time of about 5411.86 seconds, whereas TFIDF gets less time. The execution time is high in SVM because of the kernel parameter which slows down the process.

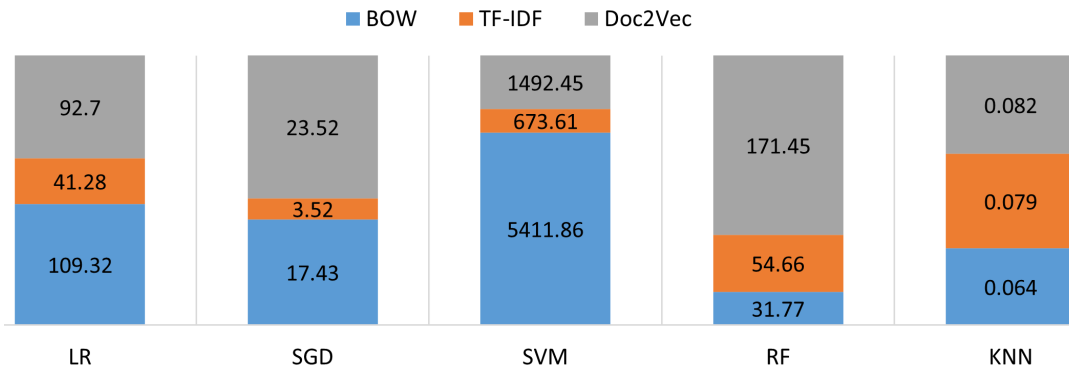


Figure 14: Execution Time (seconds) for ML Algorithms

We used three word embedding techniques (Word2Vec, Fasttext, and Glove) to evaluate the performance of the DL (CNN, LSTM, BiLSTM, and GRU) models. Additionally, we analysed the performance of article classification models in two ways: with stemmer and without stemmer (WS). We noticed that the stemmer has no significant impact on DL accuracy. Figure 15 demonstrates the accuracy of the DL algorithms for each word embedding technique. We obtained the highest accuracy of 91.83% for GRU with Fasttext and stemmer.

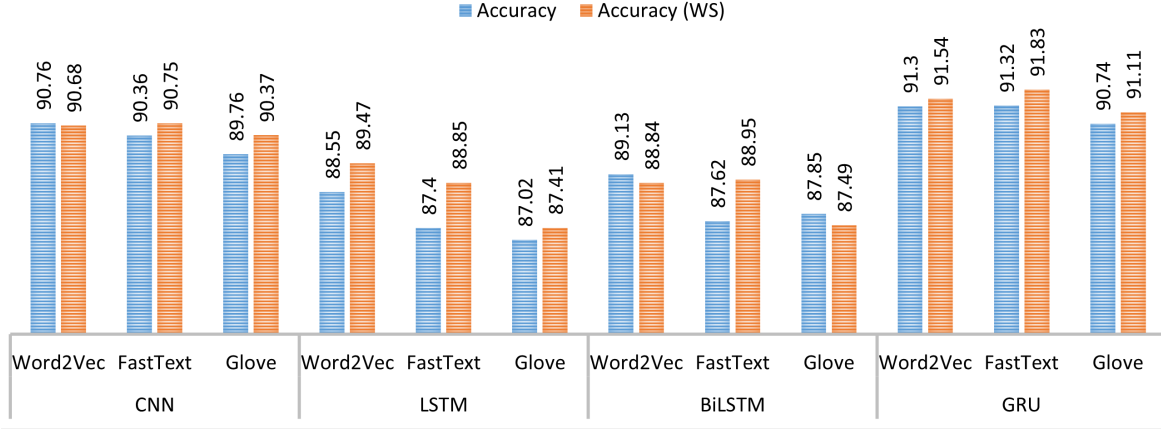


Figure 15: Manually labelled Article Classification Accuracy for DL Algorithms

Figure 16 shows that the LSTM and BiLSTM algorithms take the most execution time, whereas CNN and GRU have significantly less execution time. LSTM has less execution time than BiLSTM, as in BiLSTM, the input sequence flow is both forward and backward, which captures the context from both the past and present sequence. For example, BiLSTM execution times for word2vec and FastText techniques are 13.24 and 13.29 hour, whereas LSTM execution times for word2vec and FastText techniques are 9.25 and 9.55 hour. However, with the Glove technique, LSTM has a longer execution time than BiLSTM.

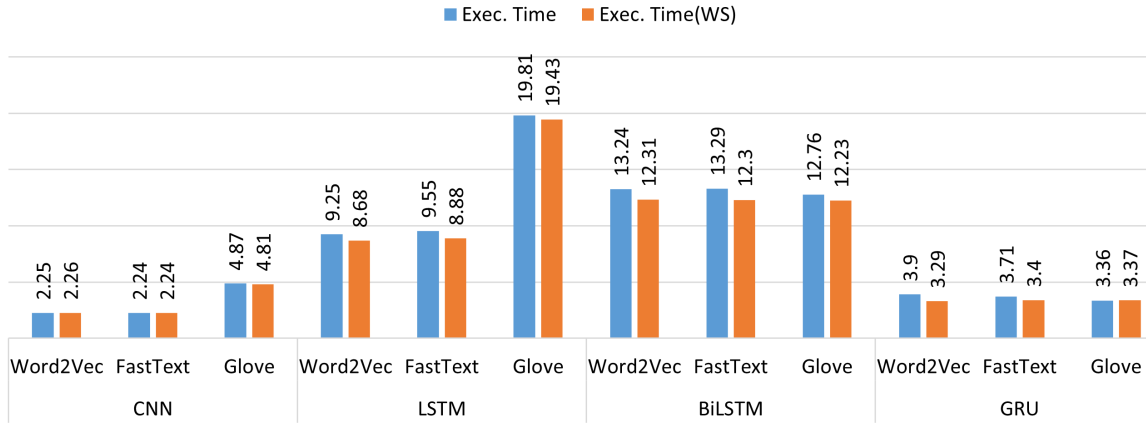


Figure 16: Execution Time (hour) for DL Algorithms

Figure 17, 18, 19, and 20 depict DL model performance for each class using f1-score, precision, and recall, as evaluation metrics. The harmonic average of recall and precision is the F1-score. The F1-score value will be 1 if both precision and recall values are 1. Precision is defined as the ratio of accurately predicted positive article classes to the total number of positively predicted classes. The recall, which is also known as the true positive rate or sensitivity, is defined as the ratio of the accurately predicted positive article classes to the total number of accurately predicted article classes. Ideally, a precision and recall value close to 1 indicates that the model has the best performance. The horizontal lines of the graph indicate the article class (see Table 1).

Figure 17 shows the CNN algorithm results (i.e., precision, recall, and F1-score) with word embedding techniques for each article class. The average precision, recall, and F1-score value for all word embedding techniques is about 0.9, demonstrating the high performance of the CNN article classification model. For the National class, the F1-score is below 0.8 for all word embedding techniques because the recall value is below 0.8, although the precision value is more than 0.8. The F1-score is about 0.9 for the remaining article classes.

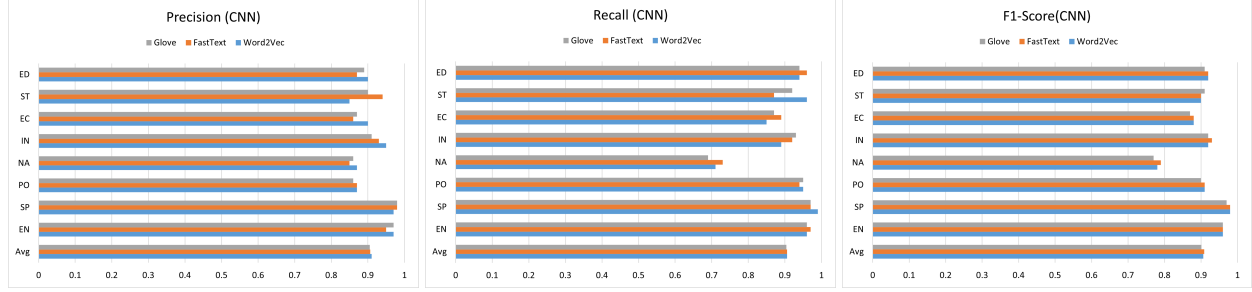


Figure 17: Evaluation Metrics for Each Article Class using CNN

Figure 18 shows the LSTM algorithm results (i.e., precision, recall, and F1-score) with word embedding techniques for each article class. The average precision, recall, and F1-score value for all word embedding techniques is approximately 0.9, presenting the high performance of the LSTM article classification model. For the National class, the F1-score is below 0.8 for all word embedding techniques because the recall and precision values are below 0.8. For the Economy class, the F1-score is below 0.9 for all word embedding techniques, as the recall and precision values are below 0.9. The F1-score, only for word2vec, is about 0.9 for the remaining article classes. The Glove and FastText do not perform well compared to Word2Vec for the LSTM article classification model.

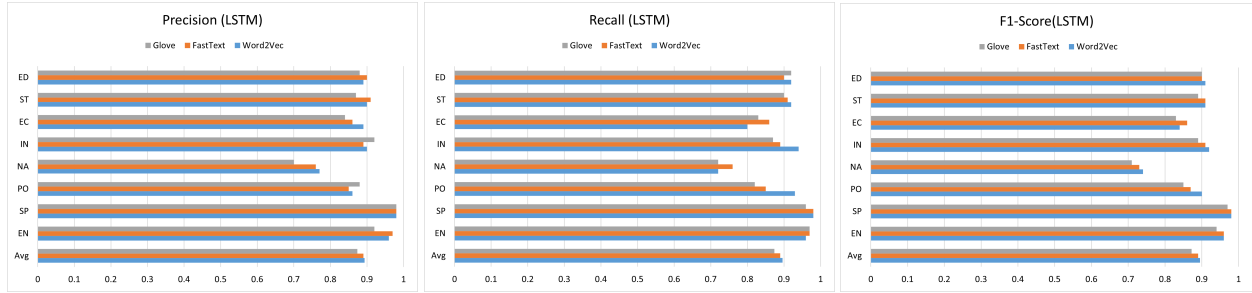


Figure 18: Evaluation Metrics for Each Article Class using LSTM

Figure 19 demonstrates the BiLSTM algorithm results (i.e., precision, recall and F1-score) with word embedding techniques for each article class. The average precision, recall and F1-score value for all word embedding techniques are approximately 0.9, showing the high performance of the BiLSTM article classification model. For the National class, the F1-score is below 0.8 for all word embedding techniques because the recall and precision values are below 0.8. For the Economy, and Politics class, the F1-score is below 0.9 for all word embedding techniques. For the remaining classes, the F1-score is about 0.9, which is demonstrating high performance.

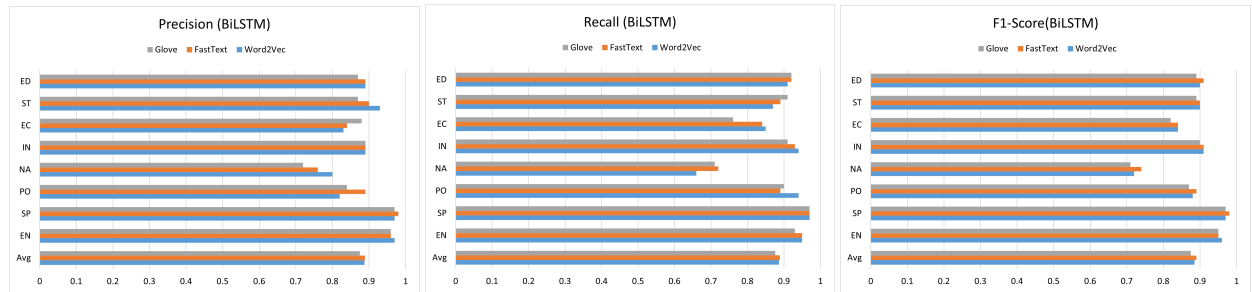


Figure 19: Evaluation Metrics for Each Article Class using BiLSTM

Figure 20 represents the GRU algorithm results (i.e., precision, recall, and F1-score) with word embedding techniques for each article class. The average precision, recall, and F1-score value for all word embedding techniques is about 0.9, presenting the heightened performance of the GRU article classification model. For the National class, the F1-score is below 0.8 for all word embedding techniques because the recall value is below 0.8, although the precision value is more than 0.8. For the Economy class, the F1-score is below 0.9 for all word embedding techniques, except the

precision of Word2Vec, as the recall and precision values are below 0.9. The F1-score is about 0.9 for the remaining article classes.

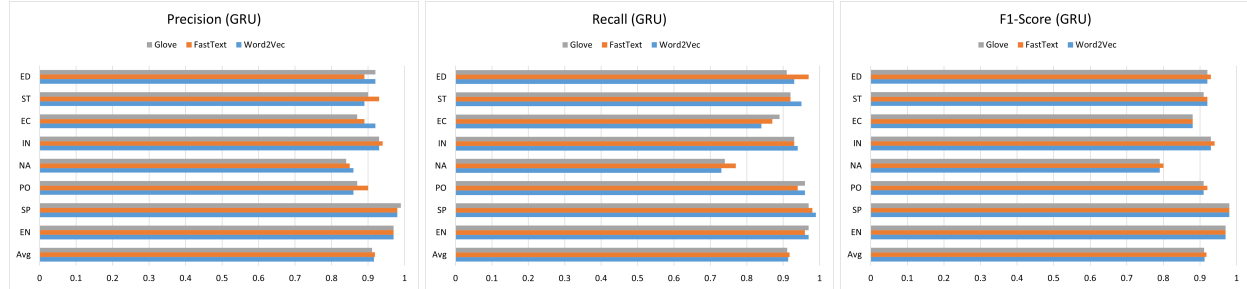


Figure 20: Evaluation Metrics for Each Article Class using GRU

By following the above consequences, we infer that GRU with Word2Vec outperforms other DL algorithms in terms of accuracy, and evaluation metrics (i.e., precision, recall, and F1-score).

## 5 Automatically labelled News Article Classification

This section goes through our findings related to the automatically labelled dataset in depth. We evaluate automatic single labelled dataset in two ways: first, we consider the most dominant automatic cluster label (see Section 3.8) as the predicted class and the original class as the tested class. After that, we evaluate the predicted and tested classes by using precision, recall, and the F1-score. Secondly, we divide the automatically labelled dataset into training and testing sets based on threshold values (i.e., 0.5, 0.6, 0.7, 0.8, and 0.9). Furthermore, we use ML algorithms to train the models and evaluate them.

Table 21 lists three examples to analyse the result of automatic labelling. The first column contains the original news article text in Bengali, and the second column translates the news article into English for Bengali non-speakers to understand the text. Column 3 contains the article class, which is divided into two sub-columns: one column contains the original class, and the other column contains the cluster class. We perceive that the first example describes the opinions of several political ministers about national mourning day. This news article is mostly political and also related to national issue though it was originally labelled as National by the news reporter. After automatic labeling, we observed that this example is 75% Politics and 23% National. Consequently, we labelled this example as Politics based on the highest probability. The second example is describing the opinion of the election committee about the use of the digital voting machine in the election, which cost about 300k. This example is originally labelled as National and the automatic labeling marked it as 53% National, 22% Politics and 23% Economic. Finally, we labelled that as National because of the highest probability. The third example is about International football news, which is originally labelled as Sports but after automatic labeling, we discover that it's 5% Entertainment, 18% International, 77% Sports and finally labelled it as Sports. After reviewing the above examples, we ascertain that single-labelled news articles can't properly describe the news as compared to multi-labelled news articles. Consequently, we extended our research to the multi-label classification.



News Article	English Translation	Class	
		Original	Cluster
বঙ্গবন্ধুর ৪৩তম শাহাদাতবার্ষিকী ও জাতীয় শোক দিবস ... প্রধানমন্ত্রীর সঙ্গে ছিলেন উপদেষ্টামণ্ডলীর সদস্য ..., আওয়ামী লীগ সাধারণ সম্পাদক, সড়ক পরিবহন ও সেতুমন্ত্রী ... প্রমুখ।	Bangabandhu's 43rd Martyrdom Anniversary and National Mourning Day ... The Prime Minister was accompanied by members of the Advisory Council ..., Awami League General Secretary, Road Transport and Bridges Minister ... etc.	National	PO: 75% NA: 23%
প্রধান নির্বাচন কমিশনার (সিইসি) কে এম নূরুল হুদা বলেছেন, রাজনৈতিক দলগুলোর সম্মতি পেলেই আসন্ন জাতীয় সংসদ নির্বাচনে ডিজিটাল ভোটিং মেশিন (ডিভিএম) ব্যবহার করা হবে। ... ভোটিংগ্রহণ করতে চাইলে ৩ লাখ মেশিনের প্রয়োজন হবে। ... একাদশ সংসদ নির্বাচনকে সামনে রেখে রোডম্যাপ চূড়ান্ত করেই রাজনৈতিক দলসহ অংশীজনের সঙ্গে সংলাপের দিনস্ফূর্ণ ঠিক করা হবে...।	Chief Election Commissioner (CEC) KM Nurul Huda has said that the digital voting machine (DVM) will be used in the upcoming parliamentary elections only if the political parties agree. ... 3 lakh machines will be needed to vote. ... With the Eleventh Parliamentary Election in mind, the roadmap will be finalized and the date of dialogue with the political parties and partners will be fixed ...	National	NA: 53% EC: 23% PO: 22%
এই নিয়ে টানা দ্বিতীয় মেজর টুর্নামেন্টের ফাইনাল খেলছে ফ্রান্স... এই প্রতিভা ব্যাপারটার অভাব নেই আরেক ফাইনালিস্ট ক্রোয়েশিয়ারও। ... যদি একটা গোল আসে, তাহলে আমি খুব খুশী হবে।	France is playing the final of the second major tournament in a row ... There is no shortage of talent in another finalist Croatia. ... if a goal comes, I'll be happy.	Sports	SP: 77% IN: 18% EN: 5%

Figure 21: News Article Example of Original and Cluster Class

Figure 22 demonstrates the classification report for each class where we consider the most dominant automatic cluster label as the predicted class and the original class as the tested class. The horizontal lines of the graph indicate the article label class (see Table 1). We achieved 66% average precision, recall, and F1-score for automatic labelling, but remarkably low precision, recall, and F1-score for the National class. As we see in the first example in Table 21, the news is more related to the Politics about 75% and 23% is National. The news reporter labelled this example as National although the news is more about politics. Generally, the National class can include other classes, for example, Politics, Sports, Education, and so on. As a result, our automatic labelled evaluation outcomes are low for the National class. The Education class achieved a 50% F1-score, whereas the remaining classes achieved more than a 60% F1-score. Apart from the National class, our automatically labelled model performs well for the remaining classes.

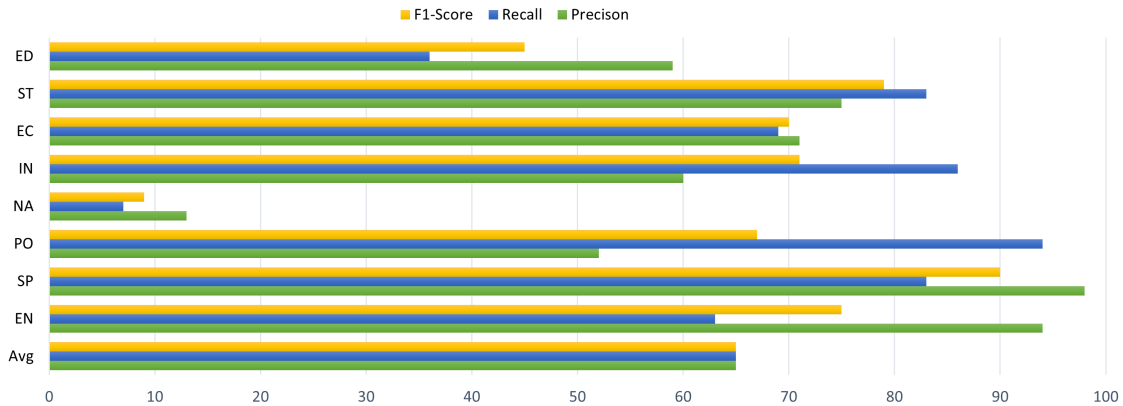


Figure 22: Compare Result between Original Class and Cluster Class

Figure 23 describes the distribution of training and testing sets for the different threshold values. More details of the training and testing set are given in Figure 24. For threshold 0.5, the dataset has 82% training and 18% testing set,

furthermore, the proportion of training and testing set continuously changes for different threshold values. The training set is about 30% and the testing set is 70% for a threshold of 0.9.

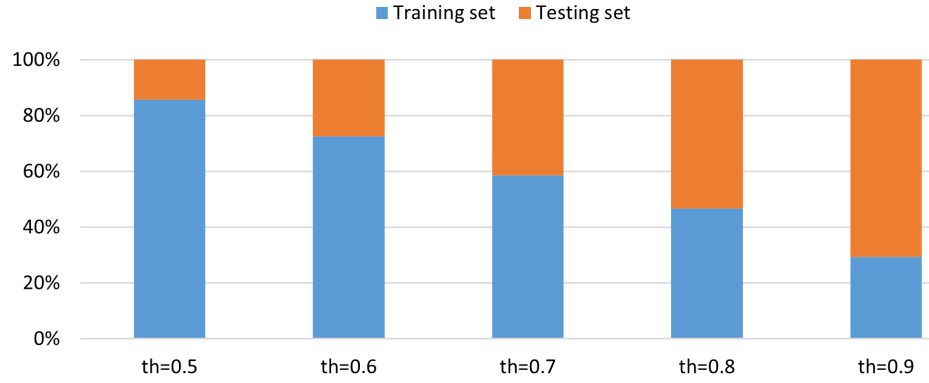


Figure 23: Training and Testing Set for ML Algorithms using Topic Modeling

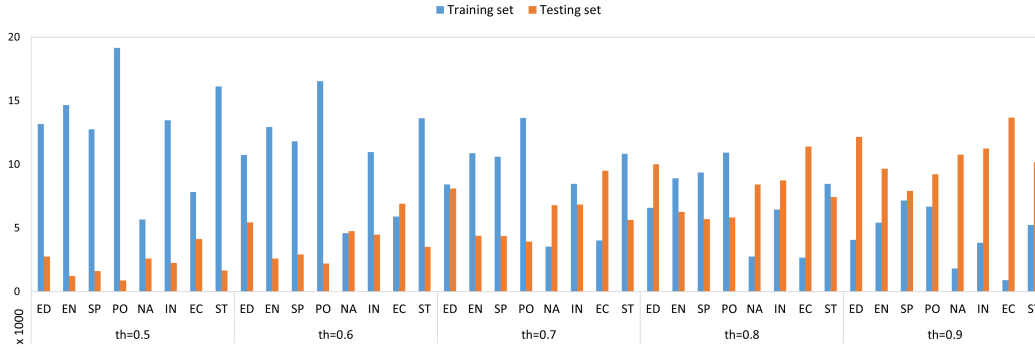


Figure 24: Details of Training and Testing Set for ML Algorithms using Topic Modeling

In Section 4, we achieved the highest accuracy for the Doc2Vec word embedding method using several machine learning algorithms for the originally labelled dataset. Here, we also use the Doc2Vec method for the same machine learning algorithms. We calculate the accuracy of the model using the auto-labelled test set news articles' predicted class and the manually (originally) labelled class of the following test set news articles.

Figure 25 shows the accuracy of several machine learning algorithms with different threshold values. The highest accuracy of 57.72% was achieved by the KNN algorithm for threshold 0.9.

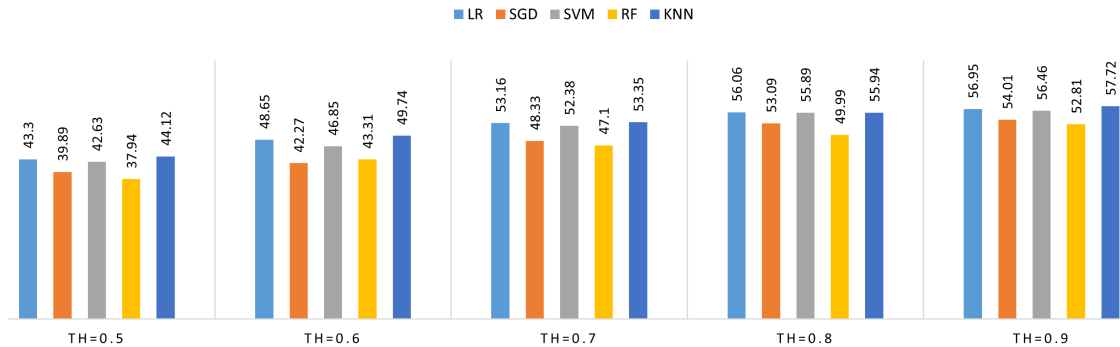


Figure 25: Result of ML Algorithms using Topic Modeling

The precision, recall, and f1-score for each class using KNN algorithms are shown in Figures 26, 27, and 28 at various threshold values. The automatically labelled dataset performed better for almost all classes excluding National and Education. As we see in Table 21, some news articles are labelled as National, although they are strongly or partially related to other topics.

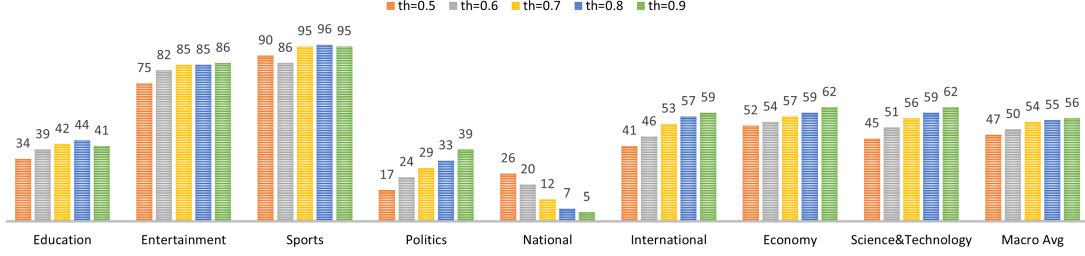


Figure 26: Precision of ML Algorithms using Topic Modeling

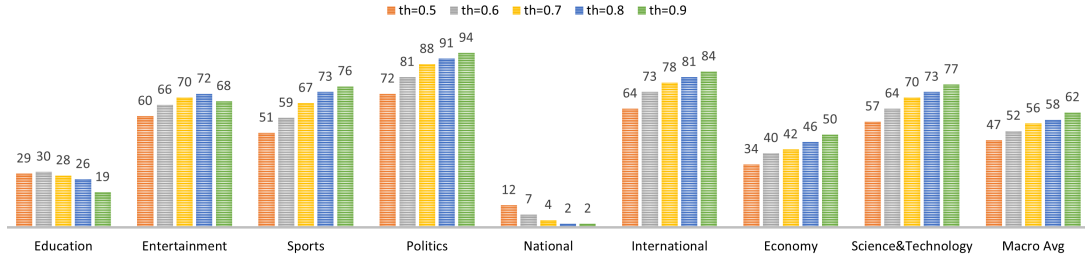


Figure 27: Recall of ML Algorithms using Topic Modeling

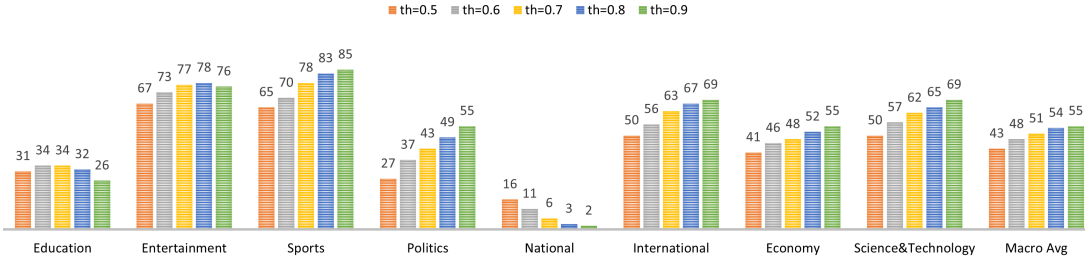


Figure 28: F1-Score of ML Algorithms using Topic Modeling

We evaluated automated single-labelled news articles and concluded that a single label is insufficient to accurately classify a news article. As a consequence, we look into the performance of multi-label news article classification in this section. The evaluation of multi-label news classification is more difficult than that of single-label text classification. The multi-label evaluation metrics are discussed further below.

Let a multi-label dataset  $T$  has a label set  $L$ ,  $|L| = k$ , and  $n$  multi-label instances  $(x_i, Y_i)$ ,  $1 \leq i \leq n$ ,  $(x_i \in X, Y_i \in y = \{0, 1\}^k)$ . The multi-label classifier  $h$  and  $Z_i = h(x_i) = \{0, 1\}^k$  be the set of predicted labels by  $h$  for the instance  $x_i$  [63].

Accuracy (A): The percentage of predicted accurate labels to the total number (predicted and actual) of labels for each occurrence is defined as accuracy. The overall accuracy is calculated as the average of all cases.

$$Accuracy, A = \frac{1}{n} \sum_{i=1}^n \frac{|Y_i \cap Z_i|}{|Y_i \cup Z_i|} \quad (4)$$

**Precision (P):** Precision is defined as the ratio of expected accurate labels to total number of actual labels, averaged over all cases.

$$Precision, P = \frac{1}{n} \sum_{i=1}^n \frac{|Y_i \cap Z_i|}{|Z_i|} \quad (5)$$

**Recall (R):** The dimension of correctly predicted labels to the cumulative number of predicted labels, averaged over all examples, is referred to as recall.

$$Recall, R = \frac{1}{n} \sum_{i=1}^n \frac{|Y_i \cap Z_i|}{|Y_i|} \quad (6)$$

**F1-Measure (F):** It is the harmonic mean of precision and recall.

$$F1 = \frac{1}{n} \sum_{i=1}^n \frac{2|Y_i \cap Z_i|}{|Y_i| + |Z_i|} \quad (7)$$

**Hamming Loss (HL):** The Hamming Loss statistic calculates the average number of times an example's relevance to a class label is predicted incorrectly. As a consequence, hamming loss takes into account both the prediction error (when an incorrect label is predicted) and the missing error (when a relevant label is not predicted), standardized across all classes and instances.

$$HammingLoss, HL = \frac{1}{kn} \sum_{i=1}^n \sum_{l=1}^k [I(l \in Z_i \wedge l \notin Y_i) + I(l \notin Z_i \wedge l \in Y_i)] \quad (8)$$

where the indicator function is 'I'. In terms of multi-label classification model performance, the less the hamming loss, the better the learning algorithm performs. The range of hamming loss is between 0 to 1.

The multi-label classification was performed using the same ML algorithms and the Doc2Vec word embedding method (see Section 3.9), which has been shown to be effective for single-label classification. We use a minimum threshold value of 0.3 for multi-labelling. We evaluate the multi-label news classification model in two testing set: using the original label, and predicted multi-label. The original label is a single label that cannot be evaluated with a multi-label evaluation process. As a result, we transform the single label into a multi-label by assigning 1 to the original label and 0 to the remainder of the label. For example, in multi-label, we have 8 labels for 8 classes and in the single label, we have only one label. Consequently, we assign 1 to the original label, and 0 to the remaining label, like this (1,0,0,0,0,0,0). Figure 29 shows the multi-label classification evaluation metrics. The OL refers to the original label and MUL refers to automatic multi-label. The highest accuracy, about 75%, is achieved by the KNN algorithm. It also performs the best with 90% precision, 88% recall, and 87% F1-score.

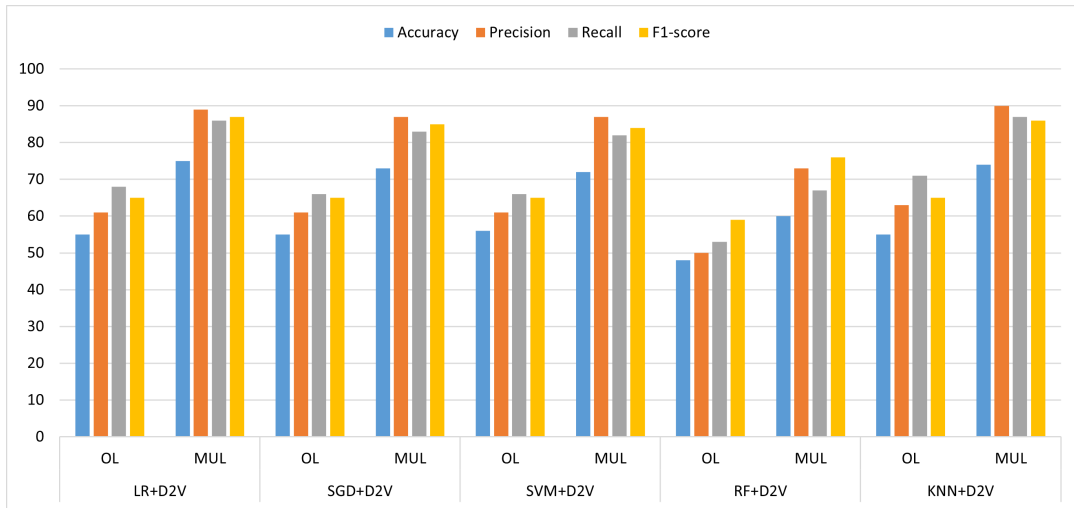


Figure 29: Comparison of Single Label and Multi-label Text Classification

Figure 30 shows that both testing sets (original label class and multi-label class) have a very low hamming loss of 0.09 and 0.03 for logistics regression, respectively.

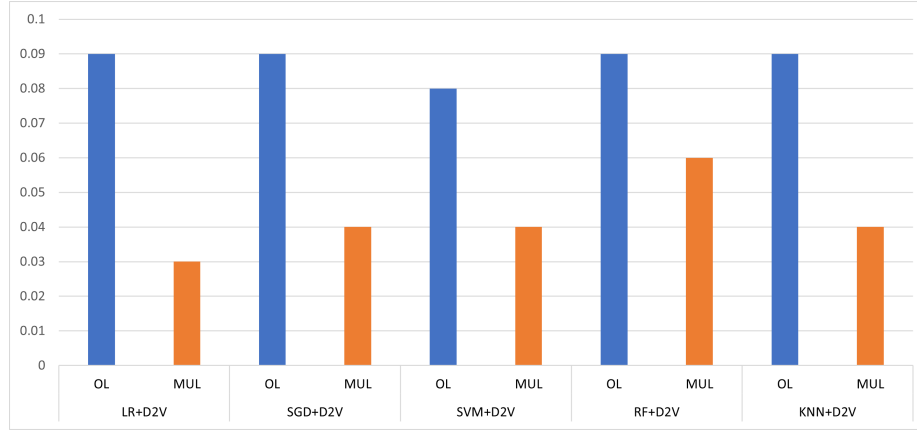


Figure 30: Comparison of Single Label and Multi-label Text Classification using Hamming Loss

## 6 Discussion

The broader aim of our work is to investigate how to name the parameter with the best output. In this paper, we investigate the following studies: we develop a comprehensive Bangla news article dataset using web scraping technique [6]; present a case study to investigate the performance of machine learning and deep learning algorithms with various word embedding techniques (see Section 3, and 4); propose a method for developing automatically labelled datasets, i.e., single-label and multi-label, from manually labelled datasets (see Section 3), and investigate automatic labelling (see Section 5). Automatic labelling is one of the issues which is addressed in Deep Journalism and Deep-Journal V1.0 [49].

For news article classification, we used a manually labelled Potrika dataset, which is labelled by the news reporter. We analysed the performance of word embedding techniques (i.e., BOW, TF-IDF, Doc2vec, word2Vec, fasttext, and glove), machine learning (i.e., logistic regression, SGD, SVM, random forest, and KNN), and deep learning algorithms (CNN, LSTM, BiLSTM, and GRU). We conclude that the deep learning algorithm's performance is much better compared to machine learning algorithms as deep learning algorithms hold the contextual meaning of the articles. On the other hand, machine learning algorithms are based on statistical learning theory and are unable to capture the context of the articles. Word embedding techniques, such as BOW and TF-IDF, achieved very low accuracy, whereas Doc2vec got good accuracy of 87.14% for logistic regression because the Doc2vec model is performing as a memory that recalls what is lacking from the present context. Deep learning algorithms achieved high accuracy for all word embedding techniques. Additionally, we investigate the performance of deep learning algorithms in two ways: using a stemmer and without a stemmer. We discovered that the use of stemmer has no significant impact in terms of the performance of the algorithms. GRU achieved the highest accuracy of about 92% using the fasttext word embedding technique.

## 7 Conclusion and Future Work

In this research, we have used the most comprehensive Bangla newspaper dataset called Potrika and performed extended experimentation to compare the performance of several machine learning and deep learning algorithms using several word embedding techniques. We further investigate the possibility of using topic modeling algorithms for automatic labeling of the classification dataset. Besides, we evaluate the multi-label news article classification with the state-of-the-art evaluation metrics. In the future, we will adopt hybrid deep learning algorithms with several word embedding techniques to improve the accuracy of Bangla news article classification and also improve the performance of automatic labeling techniques.

## Acknowledgments

The work reported in this paper is supported by the High Performance Computing Centre (HPC Center) at King Abdulaziz University, Saudi Arabia. The experiments reported in this paper were performed on the Aziz supercomputer at the HPC Center, King Abdulaziz University.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## References

- [1] Walaa Medhat, Ahmed Hassan, and Hoda Korashy. Sentiment analysis algorithms and applications: A survey. *Ain Shams engineering journal*, 5(4):1093–1113, 2014.
- [2] Sanjay K Dwivedi and Chandrakala Arya. Automatic text classification in information retrieval: A survey. In *Proceedings of the Second International Conference on Information and Communication Technology for Competitive Strategies*, pages 1–6, 2016.
- [3] Ziqiang Cao, Wenjie Li, Sujian Li, and Furu Wei. Improving multi-document summarization via text classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- [4] Asif Karim, Sami Azam, Bharanidharan Shanmugam, Krishnan Kannoorpatti, and Mamoun Alazab. A comprehensive survey for intelligent spam email detection. *IEEE Access*, 7:168261–168295, 2019.
- [5] <https://www.statista.com/statistics/266808/the-most-spoken-languages-worldwide/>. Accessed August 10, 2021.
- [6] Istiak Ahmad, Fahad AlQurashi, and Rashid Mehmood. Potrika: Raw and balanced newspaper datasets in the bangla language with eight topics and five attributes, 2022.
- [7] Shervin Minaee, Nal Kalchbrenner, Erik Cambria, Narjes Nikzad, Meysam Chenaghlu, and Jianfeng Gao. Deep learning based text classification: A comprehensive review. *arXiv preprint arXiv:2004.03705*, 2020.
- [8] Muhammad Mahmudun Nabi, Md Tanzir Altaf, and Sabir Ismail. Detecting sentiment from bangla text using machine learning technique and feature analysis. *International Journal of Computer Applications*, 153(11):28–34, 2016.
- [9] Shamsul Arafin Mahtab, Nazmul Islam, and Md Mahfuzur Rahaman. Sentiment analysis on bangladesh cricket with support vector machine. In *2018 International Conference on Bangla Speech and Language Processing (ICBSLP)*, pages 1–4. IEEE, 2018.
- [10] Nusrath Tabassum and Muhammad Ibrahim Khan. Design an empirical framework for sentiment analysis from bangla text using machine learning. In *2019 International Conference on Electrical, Computer and Communication Engineering (ECCE)*, pages 1–5. IEEE, 2019.
- [11] Md Al-Amin, Md Saiful Islam, and Shapan Das Uzzal. Sentiment analysis of bengali comments with word2vec and sentiment information of words. In *2017 international conference on electrical, computer and communication engineering (ECCE)*, pages 186–190. IEEE, 2017.
- [12] Md Asimuzzaman, Pinku Deb Nath, Farah Hossain, Asif Hossain, and Rashedur M Rahman. Sentiment analysis of bangla microblogs using adaptive neuro fuzzy system. In *2017 13th International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery (ICNC-FSKD)*, pages 1631–1638. IEEE, 2017.
- [13] KM Azharul Hasan, Mosiur Rahman, et al. Sentiment detection from bangla text using contextual valency analysis. In *2014 17th International Conference on Computer and Information Technology (ICCIT)*, pages 292–295. IEEE, 2014.
- [14] Rashedul Amin Tuhin, Bechitra Kumar Paul, Faria Nawrine, Mahbuba Akter, and Amit Kumar Das. An automated system of sentiment analysis from bangla text using supervised learning techniques. In *2019 IEEE 4th International Conference on Computer and Communication Systems (ICCCS)*, pages 360–364. IEEE, 2019.
- [15] Md Rahman, Md Seddiqui, et al. Comparison of classical machine learning approaches on bangla textual emotion analysis. *arXiv preprint arXiv:1907.07826*, 2019.

- [16] Animesh Kumar Paul and Pintu Chandra Shill. Sentiment mining from bangla data using mutual information. In *2016 2nd international conference on electrical, computer & telecommunication engineering (ICECTE)*, pages 1–4. IEEE, 2016.
- [17] SM Abu Taher, Kazi Afsana Akhter, and KM Azharul Hasan. N-gram based sentiment mining for bangla text using support vector machine. In *2018 International Conference on Bangla Speech and Language Processing (ICBSLP)*, pages 1–5. IEEE, 2018.
- [18] Al Amin, Imran Hossain, Aysha Akther, and Kazi Masudul Alam. Bengali vader: A sentiment analysis approach using modified vader. In *2019 International Conference on Electrical, Computer and Communication Engineering (ECCE)*, pages 1–6. IEEE, 2019.
- [19] Asif Hassan, Mohammad Rashedul Amin, N Mohammed, and AKA Azad. Sentiment analysis on bangla and romanized bangla text (brbt) using deep recurrent models. *arXiv preprint arXiv:1610.00369*, 2016.
- [20] Md Habibul Alam, Md-Mizanur Rahoman, and Md Abul Kalam Azad. Sentiment analysis for bangla sentences using convolutional neural network. In *2017 20th International Conference of Computer and Information Technology (ICCIT)*, pages 1–6. IEEE, 2017.
- [21] Md Tanvir Alam and Md Mofijul Islam. Bard: Bangla article classification using a new comprehensive dataset. In *2018 International Conference on Bangla Speech and Language Processing (ICBSLP)*, pages 1–5. IEEE, 2018.
- [22] Adnan Ahmad and Mohammad Ruhul Amin. Bengali word embeddings and it’s application in solving document classification problem. In *2016 19th International Conference on Computer and Information Technology (ICCIT)*, pages 425–430. IEEE, 2016.
- [23] Ankita Dhar, NiladriSekhar Dash, and Kaushik Roy. Classification of text documents through distance measurement: An experiment with multi-domain bangla text documents. In *2017 3rd International Conference on Advances in Computing, Communication & Automation (ICACCA)(Fall)*, pages 1–6. IEEE, 2017.
- [24] Ankita Dhar, Niladri Sekhar Dash, and Kaushik Roy. Application of tf-idf feature for categorizing documents of online bangla web text corpus. In *Intelligent Engineering Informatics*, pages 51–59. Springer, 2018.
- [25] Pritom Mojumder, Mahmudul Hasan, Md Faruque Hossain, and KM Azharul Hasan. A study of fasttext word embedding effects in document classification in bangla language. In *International Conference on Cyber Security and Computer Science*, pages 441–453. Springer, 2020.
- [26] Tanvirul Alam, Akib Khan, and Firoj Alam. Bangla text classification using transformers. *arXiv preprint arXiv:2011.04446*, 2020.
- [27] MA Helal and Malek Mouhoub. Topic modelling in bangla language: An lda approach to optimize topics and news classification. *Computer and Information Science*, 11(4):77–83, 2018.
- [28] Kazi Masudul Alam, Md Tanvir Hossain Hemel, SM Muhaiminul Islam, and Aysha Akther. Bangla news trend observation using lda based topic modeling. In *2020 23rd International Conference on Computer and Information Technology (ICCIT)*, pages 1–6. IEEE, 2020.
- [29] Kanish Shah, Henil Patel, Devanshi Sanghvi, and Manan Shah. A comparative analysis of logistic regression, random forest and knn models for the text classification. *Augmented Human Research*, 5(1):1–16, 2020.
- [30] Kumarshankar Raychaudhuri, Manoj Kumar, and Sanjana Bhanu. A comparative study and performance analysis of classification techniques: support vector machine, neural networks and decision trees. In *International Conference on Advances in Computing and Data Sciences*, pages 13–21. Springer, 2016.
- [31] Jason Wei and Kai Zou. Eda: Easy data augmentation techniques for boosting performance on text classification tasks. *arXiv preprint arXiv:1901.11196*, 2019.
- [32] Xinjie Zhou, Xiaojun Wan, and Jianguo Xiao. Attention-based lstm network for cross-lingual sentiment classification. In *Proceedings of the 2016 conference on empirical methods in natural language processing*, pages 247–256, 2016.
- [33] Jiachen Du, Lin Gui, Ruifeng Xu, and Yulan He. A convolutional attention model for text classification. In *National CCF conference on natural language processing and Chinese computing*, pages 183–195. Springer, 2017.
- [34] Gang Liu and Jiabao Guo. Bidirectional lstm with attention mechanism and convolutional layer for text classification. *Neurocomputing*, 337:325–338, 2019.
- [35] Weijiang Li, Fang Qi, Ming Tang, and Zhengtao Yu. Bidirectional lstm with self-attention mechanism and multi-channel features for sentiment classification. *Neurocomputing*, 387:63–77, 2020.
- [36] Liang Yao, Chengsheng Mao, and Yuan Luo. Graph convolutional networks for text classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 7370–7377, 2019.

- [37] Xien Liu, Xinxin You, Xiao Zhang, Ji Wu, and Ping Lv. Tensor graph convolutional networks for text classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8409–8416, 2020.
- [38] Min Yang, Wei Zhao, Lei Chen, Qiang Qu, Zhou Zhao, and Ying Shen. Investigating the transferring capability of capsule networks for text classification. *Neural Networks*, 118:247–261, 2019.
- [39] Jaeyoung Kim, Sion Jang, Eunjeong Park, and Sungchul Choi. Text classification using capsules. *Neurocomputing*, 376:214–221, 2020.
- [40] Deepak Kumar Jain, Rachna Jain, Yash Upadhyay, Abhishek Kathuria, and Xiangyuan Lan. Deep refinement: capsule network with attention mechanism-based system for text classification. *Neural Computing and Applications*, 32(7):1839–1856, 2020.
- [41] Ximing Li, Changchun Li, Jinjin Chi, Jihong Ouyang, and Chenliang Li. Dataless text classification: A topic modeling approach with document manifold. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pages 973–982, 2018.
- [42] Daochen Zha and Chenliang Li. Multi-label dataless text classification with topic modeling. *Knowledge and Information Systems*, 61(1):137–160, 2019.
- [43] Aditya Anantharaman, Arpit Jadiya, Chandana Tulasi Sai Siri, Bharath NVS Adikar, and Biju Mohan. Performance evaluation of topic modeling algorithms for text classification. In *2019 3rd International Conference on Trends in Electronics and Informatics (ICOEI)*, pages 704–708. IEEE, 2019.
- [44] Pinaki Prasad Guha Neogi, Amit Kumar Das, Saptarsi Goswami, and Joy Mustafi. Topic modeling for text classification. In *Emerging technology in modelling and graphics*, pages 395–407. Springer, 2020.
- [45] Miha Pavlinek and Vili Podgorelec. Text classification method based on self-training and lda topic models. *Expert Systems with Applications*, 80:83–93, 2017.
- [46] Istiak Ahmad, Ehab Abozinadah, Fahad Al Qurashi, and Rashid Mehmood. Potrika: Raw and balanced newspaper datasets in the bangla language with eight topics and five attributes. <https://doi.org/10.17632/v362rp78dc.2>, 2021.
- [47] Ebtesam Alomari, Iyad Katib, and Rashid Mehmood. Iktishaf: A big data road-traffic event detection tool using twitter and spark machine learning. *Mobile Networks and Applications*, pages 1–16, 2020.
- [48] Ebtesam Alomari, Iyad Katib, Aiiad Albeshri, Tan Yigitcanlar, and Rashid Mehmood. Iktishaf+: A big data tool with automatic labeling for road traffic social sensing and event detection using distributed machine learning. *Sensors*, 21(9):2993, 2021.
- [49] Istiak Ahmad, Fahad Alqurashi, Ehab Abozinadah, and Rashid Mehmood. Deep journalism and deepjournal v1.0: A data-driven deep learning approach to discover parameters for transportation. *Sustainability*, 14(9):5711, 2022.
- [50] Shoayee Alotaibi, Rashid Mehmood, and Iyad Katib. Sentiment Analysis of Arabic Tweets in Smart Cities: A Review of Saudi Dialect. In *2019 Fourth International Conference on Fog and Mobile Edge Computing (FMEC)*, pages 330–335. IEEE, 2019.
- [51] Nala Alahmari, Sarah Alswedani, Ahmed Alzahrani, Iyad Katib, Aiiad Albeshri, and Rashid Mehmood. Musawah: A data-driven ai approach and tool to co-create healthcare services with a case study on cancer disease in saudi arabia. *Sustainability*, 14(6), 2022.
- [52] Sarah Alswedani, Iyad Katib, Ehab Abozinadah, and Rashid Mehmood. Discovering urban governance parameters for online learning in saudi arabia during covid-19 using topic modeling of twitter data. *Frontiers in Sustainable Cities*, 4, 2022.
- [53] Sarah Alswedani, Rashid Mehmood, and Iyad Katib. Sustainable participatory governance: Data-driven discovery of parameters for planning online and in-class education in saudi arabia during covid-19. *Frontiers in Sustainable Cities*, 4, 2022.
- [54] Ebtesam Alomari, Iyad Katib, Aiiad Albeshri, and Rashid Mehmood. Covid-19: Detecting government pandemic measures and public concerns from twitter arabic data using distributed machine learning. *International Journal of Environmental Research and Public Health*, 18(1):282, 2021.
- [55] Sugimiyanto Suma, Rashid Mehmood, Nasser Albugami, Iyad Katib, and Aiiad Albeshri. Enabling next generation logistics and planning for smarter societies. *Procedia Computer Science*, 109:1122–1127, 2017.
- [56] Sugimiyanto Suma, Rashid Mehmood, and Aiiad Albeshri. Automatic detection and validation of smart city events using hpc and apache spark platforms. In *Smart Infrastructure and Applications*, pages 55–78. Springer, 2020.



- [57] Eman Alqahtani, Nourah Janbi, Sanaa Sharaf, and Rashid Mehmood. Smart homes and families to enable sustainable societies: A data-driven approach for multi-perspective parameter discovery using bert modelling. *Sustainability*, 14(20), 2022.
- [58] Quoc Le and Tomas Mikolov. Distributed representations of sentences and documents. In *International conference on machine learning*, pages 1188–1196. PMLR, 2014.
- [59] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. *arXiv preprint arXiv:1310.4546*, 2013.
- [60] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [61] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [62] Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, H erve J egou, and Tomas Mikolov. Fasttext. zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*, 2016.
- [63] Mohammad S Sorower. A literature survey on algorithms for multi-label learning. *Oregon State University, Corvallis*, 18:1–25, 2010.